





CONTENTS

1 WELCOME TO SCIKITLEARN 1			
11 INSTALLING SCIKITLEARN		1	
12 FREQUENTLY ASKED QUESTIONS			2
13 SUPPORT	8		
14 RELATED PROJECTS		9	
15 ABOUT US	12		
16 WHO IS USING SCIKITLEARN			17
17 RELEASE HISTORY		26	
18 VERSION 0213	27		
19 VERSION 0212	28		
110 VERSION 0211	29		
111 VERSION 0210	29		
112 VERSION 0204	41		
113 VERSION 0203	42		
114 VERSION 0202	43		
115 VERSION 0201	44		
116 VERSION 0200	48		
117 PREVIOUS RELEASES		63	
118 ROADMAP	145		
119 SCIKITLEARN GOVERNANCE AND DECISIONMAKING			148
2 SCIKITLEARN TUTORIALS 151			
21 AN INTRODUCTION TO MACHINE LEARNING WITH SCIKITLEARN			151
22 A TUTORIAL ON STATISTICALLEARNING FOR SCIENTIFIC DATA PROCESSING			157
23 WORKING WITH TEXT DATA		186	
24 CHOOSING THE RIGHT ESTIMATOR			193
25 EXTERNAL RESOURCES VIDEOS AND TALKS			194
3 USER GUIDE 197			
31 SUPERVISED LEARNING		197	
32 UNSUPERVISED LEARNING		336	
33 MODEL SELECTION AND EVALUATION			434
34 INSPECTION	573		
35 DATASET TRANSFORMATIONS		575	
36 DATASET LOADING UTILITIES		625	
37 COMPUTING WITH SCIKITLEARN		650	
4 GLOSSARY OF COMMON TERMS AND API ELEMENTS 663			
41 GENERAL CONCEPTS	663		
42 CLASS APIS AND ESTIMATOR TYPES			672

43 TARGET TYPES	674		
44 METHODS	676		
45 PARAMETERS	678		
46 ATTRIBUTES	681		
47 DATA AND SAMPLE PROPERTIES		682	
5 EXAMPLES	683		
51 MISCELLANEOUS EXAMPLES		683	
52 EXAMPLES BASED ON REAL WORLD DATASETS			716
53 BICLUSTERING	775		
54 CALIBRATION	787		
55 CLASSIFICATION	805		
56 CLUSTERING	820		
57 PIPELINES AND COMPOSITE ESTIMATORS		908	
58 COVARIANCE ESTIMATION	942		
59 CROSS DECOMPOSITION	957		
510 DATASET EXAMPLES	961		
511 DECOMPOSITION	970		
512 ENSEMBLE METHODS	1016		
513 TUTORIAL EXERCISES	1070		
514 FEATURE SELECTION	1079		
515 GAUSSIAN PROCESS FOR MACHINE LEARNING			1090
516 MISSING VALUE IMPUTATION		1119	
517 INSPECTION	1125		
518 GENERALIZED LINEAR MODELS		1130	
519 MANIFOLD LEARNING	1214		
520 GAUSSIAN MIXTURE MODELS		1244	
521 MODEL SELECTION	1261		
522 MULTIOUTPUT METHODS	1311		
523 NEAREST NEIGHBORS	1314		
524 NEURAL NETWORKS	1347		
525 PREPROCESSING	1360		
526 SEMI SUPERVISED CLASSIFICATION		1386	
527 SUPPORT VECTOR MACHINES		1399	
528 WORKING WITH TEXT DOCUMENTS		1431	
529 DECISION TREES	1447		
6 API REFERENCE	1459		
61SKLEARNBASE BASE CLASSES AND UTILITY FUNCTIONS			1459
62SKLEARNCALIBRATION PROBABILITY CALIBRATION			1467
63SKLEARNCLUSTER CLUSTERING	1471		
64SKLEARNCLUSTERBICLUSTER BICLUSTERING		1517	
65SKLEARNCOMPOSE COMPOSITE ESTIMATORS		1524	
66SKLEARNCOVARIANCE COVARIANCE ESTIMATORS		1532	
67SKLEARNCROSSDECOMPOSITION CROSS DECOMPOSITION			1563
68SKLEARNDATASETS DATASETS	1577		
69SKLEARNDECOMPOSITION MATRIX DECOMPOSITION		1622	
610SKLEARNDISCRIMINANTANALYSIS DISCRIMINANT ANALYSIS			1678
611SKLEARNDUMMY DUMMY ESTIMATORS		1685	
612SKLEARNENSEMBLE ENSEMBLE METHODS		1690	
613SKLEARNEXCEPTIONS EXCEPTIONS AND WARNINGS			1730
614SKLEARNEXPERIMENTAL EXPERIMENTAL	1735		
615SKLEARNFEATUREEXTRACTION FEATURE EXTRACTION			1736
616SKLEARNFEATURESELECTION FEATURE SELECTION			1764

617SKLEARN	GAUSSIANPROCESS	GAUSSIAN PROCESSES	1798
618SKLEARN	ISOTONIC	ISOTONIC REGRESSION	1838
619SKLEARN	IMPUTE	IMPUTE	1843
620SKLEARN	KERNELAPPROXIMATION	KERNEL APPROXIMATION	1851
621SKLEARN	KERNELRIDGE	KERNEL RIDGE REGRESSION	1861
622SKLEARN	LINEARMODEL	GENERALIZED LINEAR MODELS	1864
623SKLEARN	MANIFOLD	MANIFOLD LEARNING	1965
624SKLEARN	METRICS	METRICS	1983
625SKLEARN	MIXTURE	GAUSSIAN MIXTURE MODELS	2056
626SKLEARN	MODELSELECTION	MODEL SELECTION	2068
627SKLEARN	MULTICLASS	MULTICLASS AND MULTILABEL CLASSIFICATION	2122
628SKLEARN	MULTIOUTPUT	MULTIOUTPUT REGRESSION AND CLASSIFICATION	2130
629SKLEARN	NNAIVEBAYES	NAIVE BAYES	2140
630SKLEARN	NEIGHBORS	NEAREST NEIGHBORS	2153
631SKLEARN	NEURALNETWORK	NEURAL NETWORK MODELS	2205
632SKLEARN	PIPELINE	PIPELINE	2218
633SKLEARN	INSPECTION	INSPECTION	2227
634SKLEARN	PREPROCESSING	PREPROCESSING AND NORMALIZATION	2231
635SKLEARN	RANDOMPROJECTION	RANDOM PROJECTION	2286
636SKLEARN	SEMISUPERVISED	SEMISUPERVISED LEARNING	2293
637SKLEARN	SVM	SUPPORT VECTOR MACHINES	2299
638SKLEARN	TREE	DECISION TREES	2331
639SKLEARN	UTILS	UTILITIES	2358
640	RECENTLY DEPRECATED		2386
7	DEVELOPER'S GUIDE	2407	
71	CONTRIBUTING		
72	DEVELOPERS' TIPS AND TRICKS		2429
73	UTILITIES FOR DEVELOPERS		2433
74	HOW TO OPTIMIZE FOR SPEED		2436
75	ADVANCED INSTALLATION INSTRUCTIONS		2443
76	MAINTAINER	COREDEVELOPER INFORMATION	2447
BIBLIOGRAPHY 2451			
INDEX 2459			
III			



CHAPTER  
ONE  
WELCOME TO SCIKITLEARN  
11 INSTALLING SCIKITLEARN  
NOTE IF YOU WISH TO CONTRIBUTE TO THE PROJECT IT’S RECOMMENDED YOU INSTALL THE LATEST DEVELOPMENT VERSION  
111 INSTALLING THE LATEST RELEASE  
SCIKITLEARN REQUIRES  
• PYTHON 35  
• NUMPY 1110  
• SCIPY 0170  
• JOBLIB 011  
SCIKITLEARN PLOTTING CAPABILITIES IE FUNCTIONS START WITH “ PLOT ” REQUIRE MATPLOTLIB 151 SOME OF THE SCIKIT  
LEARN EXAMPLES MIGHT REQUIRE ONE OR MORE EXTRA DEPENDENCIES SCIKITIMAGE 0123 PANDAS 0180  
WARNING SCIKITLEARN 020 WAS THE LAST VERSION TO SUPPORT PYTHON 27 AND PYTHON 34 SCIKITLEARN NOW REQUIRES  
PYTHON 35 OR NEWER  
IF YOU ALREADY HAVE A WORKING INSTALLATION OF NUMPY AND SCIPY THE EASIEST WAY TO INSTALL SCIKITLEARN IS USING PIP  
PIP INSTALL U SCIKITLEARN  
ORCONDA  
CONDA INSTALL SCIKITLEARN  
IF YOU HAVE NOT INSTALLED NUMPY OR SCIPY YET YOU CAN ALSO INSTALL THESE USING CONDA OR PIP WHEN USING PIP PLEASE  
ENSURE THAT BINARY WHEELS ARE USED AND NUMPY AND SCIPY ARE NOT RECOMPILED FROM SOURCE WHICH CAN HAPPEN WHEN  
USING PARTICULAR CONFIGURATIONS OF OPERATING SYSTEM AND HARDWARE SUCH AS LINUX ON A RASPBERRY PI BUILDING NUMPY  
AND SCIPY FROM SOURCE CAN BE COMPLEX ESPECIALLY ON WINDOWS AND REQUIRES CAREFUL CONFIGURATION TO ENSURE THAT THEY  
LINK AGAINST AN OPTIMIZED IMPLEMENTATION OF LINEAR ALGEBRA ROUTINES INSTEAD USE A THIRDPARTY DISTRIBUTION AS DESCRIBED  
BELOW  
1

SCIKITLEARN USER GUIDE RELEASE 0213

IF YOU MUST INSTALL SCIKITLEARN AND ITS DEPENDENCIES WITH PIP YOU CAN INSTALL IT AS SCIKITLEARNALLDEPS THE MOST COMMON USE CASE FOR THIS IS IN A REQUIREMENTSTXT FILE USED AS PART OF AN AUTOMATED BUILD PROCESS FOR A PAAS APPLICATION OR A DOCKER IMAGE THIS OPTION IS NOT INTENDED FOR MANUAL INSTALLATION FROM THE COMMAND LINE NOTE FOR INSTALLING ON PYPY PYPY3V510 NUMPY 1140 AND SCIPY 110 ARE REQUIRED FOR INSTALLATION INSTRUCTIONS FOR MORE DISTRIBUTIONS SEE OTHER DISTRIBUTIONS FOR COMPILING THE DEVELOPMENT VERSION FROM SOURCE OR BUILDING THE PACKAGE IF NO DISTRIBUTION IS AVAILABLE FOR YOUR ARCHITECTURE SEE THE ADVANCED INSTALLATION INSTRUCTIONS

112 THIRDPARTY DISTRIBUTIONS

IF YOU DON'T ALREADY HAVE A PYTHON INSTALLATION WITH NUMPY AND SCIPY WE RECOMMEND TO INSTALL EITHER VIA YOUR PACKAGE MANAGER OR VIA A PYTHON BUNDLE THESE COME WITH NUMPY SCIPY SCIKITLEARN MATPLOTLIB AND MANY OTHER HELPFUL SCIENTIFIC AND DATA PROCESSING LIBRARIES

AVAILABLE OPTIONS ARE

CANOPY AND ANACONDA FOR ALL SUPPORTED PLATFORMS

CANOPY AND ANACONDA BOTH SHIP A RECENT VERSION OF SCIKITLEARN IN ADDITION TO A LARGE SET OF SCIENTIFIC PYTHON LIBRARY FOR WINDOWS MAC OSX AND LINUX

ANACONDA OFFERS SCIKITLEARN AS PART OF ITS FREE DISTRIBUTION

WARNING TO UPGRADE OR UNINSTALL SCIKITLEARN INSTALLED WITH ANACONDA OR CONDA YOUSHOULD NOT USE THE PIP COMMAND INSTEAD

TO UPGRADE SCIKITLEARN

CONDA UPDATE SCIKITLEARN

TO UNINSTALL SCIKITLEARN

CONDA REMOVE SCIKITLEARN

UPGRADING WITH PIP INSTALL U SCIKITLEARN OR UNINSTALLING PIP UNINSTALL SCIKITLEARN IS LIKELY FAIL TO PROPERLY REMOVE FILES INSTALLED BY THE CONDA COMMAND PIP UPGRADE AND UNINSTALL OPERATIONS ONLY WORK ON PACKAGES INSTALLED VIA PIP INSTALL

WINPYTHON FOR WINDOWS

THE WINPYTHON PROJECT DISTRIBUTES SCIKITLEARN AS AN ADDITIONAL PLUGIN

12 FREQUENTLY ASKED QUESTIONS

HERE WE TRY TO GIVE SOME ANSWERS TO QUESTIONS THAT REGULARLY POP UP ON THE MAILING LIST

2 CHAPTER 1 WELCOME TO SCIKITLEARN



SCIKITLEARN USER GUIDE RELEASE 0213

121 WHAT IS THE PROJECT NAME A LOT OF PEOPLE GET IT WRONG  
SCIKITLEARN BUT NOT SCIKIT OR SCIKIT NOR SCIKIT LEARN ALSO NOT SCIKITSLEARN OR SCIKITSLEARN WHICH WERE PREVIOUSLY USED

122 HOW DO YOU PRONOUNCE THE PROJECT NAME  
SYKIT LEARN SCI STANDS FOR SCIENCE

123 WHY SCIKIT  
THERE ARE MULTIPLE SCIKITS WHICH ARE SCIENTIFIC TOOLBOXES BUILT AROUND SCIPY YOU CAN FIND A LIST AT [HTTP://SCIKITS](http://scikit-learn.org)  
APPSPOTCOMSCIKITS APART FROM SCIKITLEARN ANOTHER POPULAR ONE IS SCIKITIMAGE

124 HOW CAN I CONTRIBUTE TO SCIKITLEARN  
SEECONTRIBUTING BEFORE WANTING TO ADD A NEW ALGORITHM WHICH IS USUALLY A MAJOR AND LENGTHY UNDERTAKING IT IS RECOMMENDED TO START WITH KNOWN ISSUES PLEASE DO NOT CONTACT THE CONTRIBUTORS OF SCIKITLEARN DIRECTLY REGARDING CONTRIBUTING TO SCIKITLEARN

125 WHAT’S THE BEST WAY TO GET HELP ON SCIKITLEARN USAGE  
FOR GENERAL MACHINE LEARNING QUESTIONS PLEASE USE CROSS VALIDATED WITH THE MACHINELEARNING TAG  
FOR SCIKITLEARN USAGE QUESTIONS PLEASE USE STACK OVERFLOW WITH THE SCIKITLEARN ANDPYTHON TAGS YOU CAN ALTERNATIVELY USE THE MAILING LIST

PLEASE MAKE SURE TO INCLUDE A MINIMAL REPRODUCTION CODE SNIPPET IDEALLY SHORTER THAN 10 LINES THAT HIGHLIGHTS YOUR PROBLEM ON A TOY DATASET FOR INSTANCE FROM SKLEARNDATASETS OR RANDOMLY GENERATED WITH FUNCTIONS OF NUMPY  
RANDOM WITH A FIXED RANDOM SEED PLEASE REMOVE ANY LINE OF CODE THAT IS NOT NECESSARY TO REPRODUCE YOUR PROBLEM  
THE PROBLEM SHOULD BE REPRODUCIBLE BY SIMPLY COPYPASTING YOUR CODE SNIPPET IN A PYTHON SHELL WITH SCIKITLEARN  
INSTALLED DO NOT FORGET TO INCLUDE THE IMPORT STATEMENTS  
MORE GUIDANCE TO WRITE GOOD REPRODUCTION CODE SNIPPETS CAN BE FOUND AT  
[HTTPS://STACKOVERFLOW.COM/HELP/MCVE](https://stackoverflow.com/help/mcve)

IF YOUR PROBLEM RAISES AN EXCEPTION THAT YOU DO NOT UNDERSTAND EVEN AFTER GOOGLING IT PLEASE MAKE SURE TO INCLUDE THE FULL TRACEBACK THAT YOU OBTAIN WHEN RUNNING THE REPRODUCTION SCRIPT  
FOR BUG REPORTS OR FEATURE REQUESTS PLEASE MAKE USE OF THE ISSUE TRACKER ON GITHUB  
THERE IS ALSO A SCIKITLEARN GITTER CHANNEL WHERE SOME USERS AND DEVELOPERS MIGHT BE FOUND  
PLEASE DO NOT EMAIL ANY AUTHORS DIRECTLY TO ASK FOR ASSISTANCE REPORT BUGS OR FOR ANY OTHER ISSUE RELATED TO SCIKITLEARN

126 HOW SHOULD I SAVE EXPORT OR DEPLOY ESTIMATORS FOR PRODUCTION  
SEEMODEL PERSISTENCE

12 FREQUENTLY ASKED QUESTIONS 3

SCIKITLEARN USER GUIDE RELEASE 0213

127 HOW CAN I CREATE A BUNCH OBJECT

DON'T MAKE A BUNCH OBJECT THEY ARE NOT PART OF THE SCIKITLEARN API BUNCH OBJECTS ARE JUST A WAY TO PACKAGE SOME NUMPY ARRAYS AS A SCIKITLEARN USER YOU ONLY EVER NEED NUMPY ARRAYS TO FEED YOUR MODEL WITH DATA FOR INSTANCE TO TRAIN A CLASSIFIER ALL YOU NEED IS A 2D ARRAY XFOR THE INPUT VARIABLES AND A 1D ARRAY YFOR THE TARGET VARIABLES THE ARRAY XHOLDS THE FEATURES AS COLUMNS AND SAMPLES AS ROWS THE ARRAY YCONTAINS INTEGER VALUES TO ENCODE THE CLASS MEMBERSHIP OF EACH SAMPLE IN X

128 HOW CAN I LOAD MY OWN DATASETS INTO A FORMAT USABLE BY SCIKITLEARN

GENERALLY SCIKITLEARN WORKS ON ANY NUMERIC DATA STORED AS NUMPY ARRAYS OR SCIPY SPARSE MATRICES OTHER TYPES THAT ARE CONVERTIBLE TO NUMERIC ARRAYS SUCH AS PANDAS DATAFRAME ARE ALSO ACCEPTABLE

FOR MORE INFORMATION ON LOADING YOUR DATA FILES INTO THESE USABLE DATA STRUCTURES PLEASE REFER TO LOADING EXTERNAL DATASETS

129 WHAT ARE THE INCLUSION CRITERIA FOR NEW ALGORITHMS

WE ONLY CONSIDER WELLESTABLISHED ALGORITHMS FOR INCLUSION A RULE OF THUMB IS AT LEAST 3 YEARS SINCE PUBLICATION 200 CITATIONS AND WIDE USE AND USEFULNESS A TECHNIQUE THAT PROVIDES A CLEARCUT IMPROVEMENT EG AN ENHANCED DATA STRUCTURE OR A MORE EFFICIENT APPROXIMATION TECHNIQUE ON A WIDELYUSED METHOD WILL ALSO BE CONSIDERED FOR INCLUSION FROM THE ALGORITHMS OR TECHNIQUES THAT MEET THE ABOVE CRITERIA ONLY THOSE WHICH FIT WELL WITHIN THE CURRENT API OF SCIKITLEARN THAT IS A FITPREDICTTTRANSFORM INTERFACE AND ORDINARILY HAVING INPUTOUTPUT THAT IS A NUMPY ARRAY OR SPARSE MATRIX ARE ACCEPTED

THE CONTRIBUTOR SHOULD SUPPORT THE IMPORTANCE OF THE PROPOSED ADDITION WITH RESEARCH PAPERS ANDOR IMPLEMENTATIONS IN OTHER SIMILAR PACKAGES DEMONSTRATE ITS USEFULNESS VIA COMMON USECASESAPPLICATIONS AND CORROBORATE PERFORMANCE IMPROVEMENTS IF ANY WITH BENCHMARKS ANDOR PLOTS IT IS EXPECTED THAT THE PROPOSED ALGORITHM SHOULD OUTPERFORM THE METHODS THAT ARE ALREADY IMPLEMENTED IN SCIKITLEARN AT LEAST IN SOME AREAS

INCLUSION OF A NEW ALGORITHM SPEEDING UP AN EXISTING MODEL IS EASIER IF

- IT DOES NOT INTRODUCE NEW HYPERPARAMETERS AS IT MAKES THE LIBRARY MORE FUTUREPROOF
- IT IS EASY TO DOCUMENT CLEARLY WHEN THE CONTRIBUTION IMPROVES THE SPEED AND WHEN IT DOES NOT FOR INSTANCE “WHEN NFEATURES NSAMPLES”
- BENCHMARKS CLEARLY SHOW A SPEED UP

ALSO NOTE THAT YOUR IMPLEMENTATION NEED NOT BE IN SCIKITLEARN TO BE USED TOGETHER WITH SCIKITLEARN TOOLS YOU CAN IMPLEMENT YOUR FAVORITE ALGORITHM IN A SCIKITLEARN COMPATIBLE WAY UPLOAD IT TO GITHUB AND LET US KNOW WE WILL BE HAPPY TO LIST IT UNDER RELATED PROJECTS IF YOU ALREADY HAVE A PACKAGE ON GITHUB FOLLOWING THE SCIKITLEARN API YOU MAY ALSO BE INTERESTED TO LOOK AT SCIKITLEARNCONTRIB

1210 WHY ARE YOU SO SELECTIVE ON WHAT ALGORITHMS YOU INCLUDE IN SCIKITLEARN

CODE IS MAINTENANCE COST AND WE NEED TO BALANCE THE AMOUNT OF CODE WE HAVE WITH THE SIZE OF THE TEAM AND ADD TO THIS THE FACT THAT COMPLEXITY SCALES NON LINEARLY WITH THE NUMBER OF FEATURES THE PACKAGE RELIES ON CORE DEVELOPERS USING THEIR FREE TIME TO FIX BUGS MAINTAIN CODE AND REVIEW CONTRIBUTIONS ANY ALGORITHM THAT IS ADDED NEEDS FUTURE ATTENTION BY THE DEVELOPERS AT WHICH POINT THE ORIGINAL AUTHOR MIGHT LONG HAVE LOST INTEREST SEE ALSO WHAT ARE THE INCLUSION CRITERIA FOR NEW ALGORITHMS FOR A GREAT READ ABOUT LONGTERM MAINTENANCE ISSUES IN OPENSOURCE SOFTWARE LOOK AT THE EXECUTIVE SUMMARY OF ROADS AND BRIDGES

4 CHAPTER 1 WELCOME TO SCIKITLEARN

SCIKITLEARN USER GUIDE RELEASE 0213

1211 WHY DID YOU REMOVE HMMS FROM SCIKITLEARN

SEEWILL YOU ADD GRAPHICAL MODELS OR SEQUENCE PREDICTION TO SCIKITLEARN

1212 WILL YOU ADD GRAPHICAL MODELS OR SEQUENCE PREDICTION TO SCIKITLEARN

NOT IN THE FORESEEABLE FUTURE SCIKITLEARN TRIES TO PROVIDE A UNIFIED API FOR THE BASIC TASKS IN MACHINE LEARNING WITH PIPELINES AND METAALGORITHMS LIKE GRID SEARCH TO TIE EVERYTHING TOGETHER THE REQUIRED CONCEPTS APIS ALGORITHMS AND EXPERTISE REQUIRED FOR STRUCTURED LEARNING ARE DIFFERENT FROM WHAT SCIKITLEARN HAS TO OFFER IF WE STARTED DOING ARBITRARY STRUCTURED LEARNING WE'D NEED TO REDESIGN THE WHOLE PACKAGE AND THE PROJECT WOULD LIKELY COLLAPSE UNDER ITS OWN WEIGHT

THERE ARE TWO PROJECT WITH API SIMILAR TO SCIKITLEARN THAT DO STRUCTURED PREDICTION

- PYSTRUCT HANDLES GENERAL STRUCTURED LEARNING FOCUSES ON SSVMS ON ARBITRARY GRAPH STRUCTURES WITH APPROXIMATE INFERENCE DEFINES THE NOTION OF SAMPLE AS AN INSTANCE OF THE GRAPH STRUCTURE
- SEQLEARN HANDLES SEQUENCES ONLY FOCUSES ON EXACT INFERENCE HAS HMMS BUT MOSTLY FOR THE SAKE OF COMPLETE NESS TREATS A FEATURE VECTOR AS A SAMPLE AND USES AN OFFSET ENCODING FOR THE DEPENDENCIES BETWEEN FEATURE VECTORS

1213 WILL YOU ADD GPU SUPPORT

NO OR AT LEAST NOT IN THE NEAR FUTURE THE MAIN REASON IS THAT GPU SUPPORT WILL INTRODUCE MANY SOFTWARE DEPENDENCIES AND INTRODUCE PLATFORM SPECIFIC ISSUES SCIKITLEARN IS DESIGNED TO BE EASY TO INSTALL ON A WIDE VARIETY OF PLATFORMS OUTSIDE OF NEURAL NETWORKS GPUS DON'T PLAY A LARGE ROLE IN MACHINE LEARNING TODAY AND MUCH LARGER GAINS IN SPEED CAN OFTEN BE ACHIEVED BY A CAREFUL CHOICE OF ALGORITHMS

1214 DO YOU SUPPORT PYPY

IN CASE YOU DIDN'T KNOW PYPY IS AN ALTERNATIVE PYTHON IMPLEMENTATION WITH A BUILTIN JUSTINTIME COMPILER EXPERIMENTAL SUPPORT FOR PYPY3V510 HAS BEEN ADDED WHICH REQUIRES NUMPY 1140 AND SCIPY 110

1215 HOW DO I DEAL WITH STRING DATA OR TREES GRAPHS

SCIKITLEARN ESTIMATORS ASSUME YOU'LL FEED THEM REALVALUED FEATURE VECTORS THIS ASSUMPTION IS HARDCODED IN PRETTY MUCH ALL OF THE LIBRARY HOWEVER YOU CAN FEED NONNUMERICAL INPUTS TO ESTIMATORS IN SEVERAL WAYS

IF YOU HAVE TEXT DOCUMENTS YOU CAN USE A TERM FREQUENCY FEATURES SEE TEXT FEATURE EXTRACTION FOR THE BUILTIN TEXT VECTORIZERS FOR MORE GENERAL FEATURE EXTRACTION FROM ANY KIND OF DATA SEE LOADING FEATURES FROM DICTS ANDFEATURE HASHING

ANOTHER COMMON CASE IS WHEN YOU HAVE NONNUMERICAL DATA AND A CUSTOM DISTANCE OR SIMILARITY METRIC ON THESE DATA EXAMPLES INCLUDE STRINGS WITH EDIT DISTANCE AKA LEVENSHTAIN DISTANCE EG DNA OR RNA SEQUENCES THESE CAN BE ENCODED AS NUMBERS BUT DOING SO IS PAINFUL AND ERRORPRONE WORKING WITH DISTANCE METRICS ON ARBITRARY DATA CAN BE DONE IN TWO WAYS

FIRSTLY MANY ESTIMATORS TAKE PRECOMPUTED DISTANCESIMILARITY MATRICES SO IF THE DATASET IS NOT TOO LARGE YOU CAN COMPUTE DISTANCES FOR ALL PAIRS OF INPUTS IF THE DATASET IS LARGE YOU CAN USE FEATURE VECTORS WITH ONLY ONE "FEATURE" WHICH IS AN INDEX INTO A SEPARATE DATA STRUCTURE AND SUPPLY A CUSTOM METRIC FUNCTION THAT LOOKS UP THE ACTUAL DATA IN THIS DATA STRUCTURE EG TO USE DBSCAN WITH LEVENSHTAIN DISTANCES

12 FREQUENTLY ASKED QUESTIONS 5

```
SCIKITLEARN USER GUIDE RELEASE 0213
FROM LEVEN IMPORT LEVENSHTAIN
IMPORT NUMPY AS NP
FROM SKLEARNCLUSTER IMPORT DBSCAN
DATA ACCTCCTAGAAG ACCTACTAGAAGTT GAATATTAGGCCGA
DEF LEVMETRICX Y
I J INTX0 INTY0 EXTRACT INDICES
RETURN LEVENSHTAINDATAI DATAJ
```

```
X NPARANGELENDATARESHAPE1 1
X
ARRAY0
1
2
WE NEED TO SPECIFY ALGORITUMBRUTE AS THE DEFAULT ASSUMES
A CONTINUOUS FEATURE SPACE
DBSCANX METRICLEVMEETRIC EPS5 MINSAMPLES2 ALGORITHMBRUTE
```

```
0 1 ARRAY 0 0 1
THIS USES THE THIRDPARTY EDIT DISTANCE PACKAGE LEVEN
SIMILAR TRICKS CAN BE USED WITH SOME CARE FOR TREE KERNELS GRAPH KERNELS ETC
1216 WHY DO I SOMETIME GET A CRASHFREEZE WITH NJOBS 1 UNDER OSX OR LINUX
SEVERAL SCIKITLEARN TOOLS SUCH AS GRIDSEARCHCV ANDCROSSVALSCORE RELY INTERNALLY ON PYTHON'S
MULTIPROCESSING MODULE TO PARALLELIZE EXECUTION ONTO SEVERAL PYTHON PROCESSES BY PASSING NJOBS 1 AS
ARGUMENT
```

THE PROBLEM IS THAT PYTHON MULTIPROCESSING DOES AFORK SYSTEM CALL WITHOUT FOLLOWING IT WITH AN EXEC SYSTEM CALL FOR PERFORMANCE REASONS MANY LIBRARIES LIKE SOME VERSIONS OF ACCELERATE VECLIB UNDER OSX SOME VERSIONS OF MKL THE OPENMP RUNTIME OF GCC NVIDIA'S CUDA AND PROBABLY MANY OTHERS MANAGE THEIR OWN INTERNAL THREAD POOL UPON A CALL TO FORK THE THREAD POOL STATE IN THE CHILD PROCESS IS CORRUPTED THE THREAD POOL BELIEVES IT HAS MANY THREADS WHILE ONLY THE MAIN THREAD STATE HAS BEEN FORKED IT IS POSSIBLE TO CHANGE THE LIBRARIES TO MAKE THEM DETECT WHEN A FORK HAPPENS AND REINITIALIZE THE THREAD POOL IN THAT CASE WE DID THAT FOR OPENBLAS MERGED UPSTREAM IN MASTER SINCE 0210 AND WE CONTRIBUTED A PATCH TO GCC'S OPENMP RUNTIME NOT YET REVIEWED BUT IN THE END THE REAL CULPRIT IS PYTHON'S MULTIPROCESSING THAT DOESFORK WITHOUTEXEC TO REDUCE THE OVERHEAD OF STARTING AND USING NEW PYTHON PROCESSES FOR PARALLEL COMPUTING UNFORTUNATELY THIS IS A VIOLATION OF THE POSIX STANDARD AND THEREFORE SOME SOFTWARE EDITORS LIKE APPLE REFUSE TO CONSIDER THE LACK OF FORKSAFETY IN ACCELERATE VECLIB AS A BUG

IN PYTHON 34 IT IS NOW POSSIBLE TO CONFIGURE MULTIPROCESSING TO USE THE 'FORKSERVER' OR 'SPAWN' START METHODS INSTEAD OF THE DEFAULT 'FORK' TO MANAGE THE PROCESS POOLS TO WORK AROUND THIS ISSUE WHEN USING SCIKITLEARN YOU CAN SET THE JOBLIBSTARTMETHOD ENVIRONMENT VARIABLE TO 'FORKSERVER' HOWEVER THE USER SHOULD BE AWARE THAT USING THE 'FORKSERVER' METHOD PREVENTS JOBLIBPARALLEL TO CALL FUNCTION INTERACTIVELY DEFINED IN A SHELL SESSION IF YOU HAVE CUSTOM CODE THAT USES MULTIPROCESSING DIRECTLY INSTEAD OF USING IT VIA JOBLIB YOU CAN ENABLE THE 'FORKSERVER' MODE GLOBALLY FOR YOUR PROGRAM INSERT THE FOLLOWING INSTRUCTIONS IN YOUR MAIN SCRIPT

```
IMPORT MULTIPROCESSING
OTHER IMPORTS CUSTOM CODE LOAD DATA DEFINE MODEL
IFNAME MAIN
MULTIPROCESSINGSETSTARTMETHODFORKSERVER
6 CHAPTER 1 WELCOME TO SCIKITLEARN
```

SCIKITLEARN USER GUIDE RELEASE 0213

CALL SCIKITLEARN UTILS WITH NJOBS 1 HERE

YOU CAN FIND MORE DEFAULT ON THE NEW START METHODS IN THE MULTIPROCESSING DOCUMENTATION

1217 WHY DOES MY JOB USE MORE CORES THAN SPECIFIED WITH NJOBS UNDER OSX OR LINUX

THIS HAPPENS WHEN VECTORIZED NUMPY OPERATIONS ARE HANDLED BY LIBRARIES SUCH AS MKL OR OPENBLAS WHILE SCIKITLEARN ADHERES TO THE LIMIT SET BY NJOBS NUMPY OPERATIONS VECTORIZED USING MKL OR OPENBLAS WILL MAKE USE OF MULTIPLE THREADS WITHIN EACH SCIKITLEARN JOB THREAD OR PROCESS THE NUMBER OF THREADS USED BY THE BLAS LIBRARY CAN BE SET VIA AN ENVIRONMENT VARIABLE FOR EXAMPLE TO SET THE MAXIMUM NUMBER OF THREADS TO SOME INTEGER VALUE N THE FOLLOWING ENVIRONMENT VARIABLES SHOULD BE SET

- FOR MKL EXPORT MKLNUMTHREADSN
- FOR OPENBLAS EXPORT OPENBLASNUMTHREADSN

1218 WHY IS THERE NO SUPPORT FOR DEEP OR REINFORCEMENT LEARNING WILL THERE BE SUPPORT FOR DEEP OR REINFORCEMENT LEARNING IN SCIKITLEARN

DEEP LEARNING AND REINFORCEMENT LEARNING BOTH REQUIRE A RICH VOCABULARY TO DEFINE AN ARCHITECTURE WITH DEEP LEARNING ADDITIONALLY REQUIRING GPUS FOR EFFICIENT COMPUTING HOWEVER NEITHER OF THESE FIT WITHIN THE DESIGN CONSTRAINTS OF SCIKITLEARN AS A RESULT DEEP LEARNING AND REINFORCEMENT LEARNING ARE CURRENTLY OUT OF SCOPE FOR WHAT SCIKITLEARN SEEKS TO ACHIEVE

YOU CAN FIND MORE INFORMATION ABOUT ADDITION OF GPU SUPPORT AT WILL YOU ADD GPU SUPPORT

1219 WHY IS MY PULL REQUEST NOT GETTING ANY ATTENTION

THE SCIKITLEARN REVIEW PROCESS TAKES A SIGNIFICANT AMOUNT OF TIME AND CONTRIBUTORS SHOULD NOT BE DISCOURAGED BY A LACK OF ACTIVITY OR REVIEW ON THEIR PULL REQUEST WE CARE A LOT ABOUT GETTING THINGS RIGHT THE FIRST TIME AS MAINTENANCE AND LATER CHANGE COMES AT A HIGH COST WE RARELY RELEASE ANY “EXPERIMENTAL” CODE SO ALL OF OUR CONTRIBUTIONS WILL BE SUBJECT TO HIGH USE IMMEDIATELY AND SHOULD BE OF THE HIGHEST QUALITY POSSIBLE INITIALLY BEYOND THAT SCIKITLEARN IS LIMITED IN ITS REVIEWING BANDWIDTH MANY OF THE REVIEWERS AND CORE DEVELOPERS ARE WORKING ON SCIKITLEARN ON THEIR OWN TIME IF A REVIEW OF YOUR PULL REQUEST COMES SLOWLY IT IS LIKELY BECAUSE THE REVIEWERS ARE BUSY WE ASK FOR YOUR UNDERSTANDING AND REQUEST THAT YOU NOT CLOSE YOUR PULL REQUEST OR DISCONTINUE YOUR WORK SOLELY BECAUSE OF THIS REASON

1220 HOW DO I SET A RANDOMSTATE FOR AN ENTIRE EXECUTION

FOR TESTING AND REPLICABILITY IT IS OFTEN IMPORTANT TO HAVE THE ENTIRE EXECUTION CONTROLLED BY A SINGLE SEED FOR THE PSEUDO RANDOM NUMBER GENERATOR USED IN ALGORITHMS THAT HAVE A RANDOMIZED COMPONENT SCIKITLEARN DOES NOT USE ITS OWN GLOBAL RANDOM STATE WHENEVER A RANDOMSTATE INSTANCE OR AN INTEGER RANDOM SEED IS NOT PROVIDED AS AN ARGUMENT IT RELIES ON THE NUMPY GLOBAL RANDOM STATE WHICH CAN BE SET USING NUMPYRANDOMSEED FOR EXAMPLE TO SET AN EXECUTION’S NUMPY GLOBAL RANDOM STATE TO 42 ONE COULD EXECUTE THE FOLLOWING IN HIS OR HER SCRIPT

```
import numpy as np
np.random.seed(42)
```

12 FREQUENTLY ASKED QUESTIONS 7

SCIKITLEARN USER GUIDE RELEASE 0213

HOWEVER A GLOBAL RANDOM STATE IS PRONE TO MODIFICATION BY OTHER CODE DURING EXECUTION THUS THE ONLY WAY TO ENSURE REPLICABILITY IS TO PASS RANDOMSTATE INSTANCES EVERYWHERE AND ENSURE THAT BOTH ESTIMATORS AND CROSSVALIDATION SPLITTERS HAVE THEIR RANDOMSTATE PARAMETER SET

1221 WHY DO CATEGORICAL VARIABLES NEED PREPROCESSING IN SCIKITLEARN COMPARED TO OTHER TOOLS

MOST OF SCIKITLEARN ASSUMES DATA IS IN NUMPY ARRAYS OR SCIPY SPARSE MATRICES OF A SINGLE NUMERIC DTYPE THESE DO NOT EXPLICITLY REPRESENT CATEGORICAL VARIABLES AT PRESENT THUS UNLIKE R'S DATAFRAMES OR PANDASDATAFRAME WE REQUIRE EXPLICIT CONVERSION OF CATEGORICAL FEATURES TO NUMERIC VALUES AS DISCUSSED IN ENCODING CATEGORICAL FEATURES SEE ALSO COLUMN TRANSFORMER WITH MIXED TYPES FOR AN EXAMPLE OF WORKING WITH HETEROGENEOUS EG CATEGORICAL AND NUMERIC DATA

1222 WHY DOES SCIKITLEARN NOT DIRECTLY WORK WITH FOR EXAMPLE PAN  
DASDATAFRAME

THE HOMOGENEOUS NUMPY AND SCIPY DATA OBJECTS CURRENTLY EXPECTED ARE MOST EFFICIENT TO PROCESS FOR MOST OPERATIONS EXTENSIVE WORK WOULD ALSO BE NEEDED TO SUPPORT PANDAS CATEGORICAL TYPES RESTRICTING INPUT TO HOMOGENEOUS TYPES THEREFORE REDUCES MAINTENANCE COST AND ENCOURAGES USAGE OF EFFICIENT DATA STRUCTURES

13 SUPPORT

THERE ARE SEVERAL WAYS TO GET IN TOUCH WITH THE DEVELOPERS

131 MAILING LIST

- THE MAIN MAILING LIST IS SCIKITLEARN
- THERE IS ALSO A COMMIT LIST SCIKITLEARNCOMMITTS WHERE UPDATES TO THE MAIN REPOSITORY AND TEST FAILURES GET NOTIFIED

132 USER QUESTIONS

- SOME SCIKITLEARN DEVELOPERS SUPPORT USERS ON STACKOVERFLOW USING THE SCIKITLEARN TAG
- FOR GENERAL THEORETICAL OR METHODOLOGICAL MACHINE LEARNING QUESTIONS STACK EXCHANGE IS PROBABLY A MORE SUITABLE VENUE

IN BOTH CASES PLEASE USE A DESCRIPTIVE QUESTION IN THE TITLE FIELD EG NO "PLEASE HELP WITH SCIKITLEARN" AS THIS IS NOT A QUESTION AND PUT DETAILS ON WHAT YOU TRIED TO ACHIEVE WHAT WERE THE EXPECTED RESULTS AND WHAT YOU OBSERVED INSTEAD IN THE DETAILS FIELD

CODE AND DATA SNIPPETS ARE WELCOME MINIMALISTIC UP TO 20 LINES LONG REPRODUCTION SCRIPT VERY HELPFUL PLEASE DESCRIBE THE NATURE OF YOUR DATA AND THE HOW YOU PREPROCESSED IT WHAT IS THE NUMBER OF SAMPLES WHAT IS THE NUMBER AND TYPE OF FEATURES ID CATEGORICAL OR NUMERICAL AND FOR SUPERVISED LEARNING TASKS WHAT TARGET ARE YOUR TRYING TO PREDICT BINARY MULTICLASS 1 OUT OF NCLASSES OR MULTILABEL KOUT OFNCLASSES CLASSIFICATION OR

CONTINUOUS VARIABLE REGRESSION

8 CHAPTER 1 WELCOME TO SCIKITLEARN

SCIKITLEARN USER GUIDE RELEASE 0213

133 BUG TRACKER

IF YOU THINK YOU'VE ENCOUNTERED A BUG PLEASE REPORT IT TO THE ISSUE TRACKER

[HTTPSGITHUBCOMSCIKITLEARNSCIKITLEARNISSUES](https://github.com/scikitlearn/scikitlearn/issues)

DON'T FORGET TO INCLUDE

- STEPS OR BETTER SCRIPT TO REPRODUCE
- EXPECTED OUTCOME
- OBSERVED OUTCOME OR PYTHON OR GDB TRACEBACKS

TO HELP DEVELOPERS FIX YOUR BUG FASTER PLEASE LINK TO A [HTTPSGISTGITHUBCOM](https://gist.github.com) HOLDING A STANDALONE MINIMALISTIC PYTHON SCRIPT THAT REPRODUCES YOUR BUG AND OPTIONALLY A MINIMALISTIC SUBSAMPLE OF YOUR DATASET FOR INSTANCE EXPORTED AS CSV FILES USING `NUMPYSAVETXT`

NOTE GISTS ARE GIT CLONEABLE REPOSITORIES AND THUS YOU CAN USE GIT TO PUSH DATAFILES TO THEM

134 IRC

SOME DEVELOPERS LIKE TO HANG OUT ON CHANNEL [#scikitlearn](https://irc.freenode.net) ON [IRC](https://freenode.net) FREENODENET

IF YOU DO NOT HAVE AN IRC CLIENT OR ARE BEHIND A FIREWALL THIS WEB CLIENT WORKS FINE [HTTPSWEBCHATFREENODENET](https://webchat.freenode.net)

135 DOCUMENTATION RESOURCES

THIS DOCUMENTATION IS RELATIVE TO 0213 DOCUMENTATION FOR OTHER VERSIONS CAN BE FOUND [HERE](#)

PRINTABLE PDF DOCUMENTATION FOR OLD VERSIONS CAN BE FOUND [HERE](#)

14 RELATED PROJECTS

PROJECTS IMPLEMENTING THE SCIKITLEARN ESTIMATOR API ARE ENCOURAGED TO USE THE SCIKITLEARNCONTRIB TEMPLATE WHICH FACILITATES BEST PRACTICES FOR TESTING AND DOCUMENTING ESTIMATORS THE SCIKITLEARNCONTRIB GITHUB ORGANISATION ALSO ACCEPTS HIGHQUALITY CONTRIBUTIONS OF REPOSITORIES CONFORMING TO THIS TEMPLATE

BELOW IS A LIST OF SISTERPROJECTS EXTENSIONS AND DOMAIN SPECIFIC PACKAGES

141 INTEROPERABILITY AND FRAMEWORK ENHANCEMENTS

THESE TOOLS ADAPT SCIKITLEARN FOR USE WITH OTHER TECHNOLOGIES OR OTHERWISE ENHANCE THE FUNCTIONALITY OF SCIKITLEARN'S ESTIMATORS

DATA FORMATS

- SKLEARNPANDAS BRIDGE FOR SCIKITLEARN PIPELINES AND PANDAS DATA FRAME WITH DEDICATED TRANSFORMERS
- SKLEARNXARRAY PROVIDES COMPATIBILITY OF SCIKITLEARN ESTIMATORS WITH XARRAY DATA STRUCTURES

AUTOML

- AUTOML AUTOMATED MACHINE LEARNING FOR PRODUCTION AND ANALYTICS BUILT ON SCIKITLEARN AND RELATED PROJECTS TRAINS A PIPELINE WITH ALL THE STANDARD MACHINE LEARNING STEPS TUNED FOR PREDICTION SPEED AND EASE OF TRANSFER TO PRODUCTION ENVIRONMENTS
- AUTOSKLEARN AN AUTOMATED MACHINE LEARNING TOOLKIT AND A DROPIN REPLACEMENT FOR A SCIKITLEARN ESTIMATOR

14 RELATED PROJECTS 9

SCIKITLEARN USER GUIDE RELEASE 0213

- TPOT AN AUTOMATED MACHINE LEARNING TOOLKIT THAT OPTIMIZES A SERIES OF SCIKITLEARN OPERATORS TO DESIGN A MACHINE LEARNING PIPELINE INCLUDING DATA AND FEATURE PREPROCESSORS AS WELL AS THE ESTIMATORS WORKS AS A DROPIN REPLACEMENT FOR A SCIKITLEARN ESTIMATOR
  - SCIKITOPTIMIZE A LIBRARY TO MINIMIZE VERY EXPENSIVE AND NOISY BLACKBOX FUNCTIONS IT IMPLEMENTS SEVERAL METHODS FOR SEQUENTIAL MODELBASED OPTIMIZATION AND INCLUDES A REPLACEMENT FOR GRIDSEARCHCV OR RANDOMIZEDSEARCHCV TO DO CROSSVALIDATED PARAMETER SEARCH USING ANY OF THESE STRATEGIES
  - EXPERIMENTATION FRAMEWORKS
  - REP ENVIRONMENT FOR CONDUCTING DATADRIVEN RESEARCH IN A CONSISTENT AND REPRODUCIBLE WAY
  - ML FRONTEND PROVIDES DATASET MANAGEMENT AND SVM FITTINGPREDICTION THROUGH WEBBASED AND PROGRAMMATIC INTERFACES
  - SCIKITLEARN LABORATORY A COMMANDLINE WRAPPER AROUND SCIKITLEARN THAT MAKES IT EASY TO RUN MACHINE LEARNING EXPERIMENTS WITH MULTIPLE LEARNERS AND LARGE FEATURE SETS
  - XCESSIV IS A NOTEBOOKLIKE APPLICATION FOR QUICK SCALABLE AND AUTOMATED HYPERPARAMETER TUNING AND STACKED ENSEMBLING PROVIDES A FRAMEWORK FOR KEEPING TRACK OF MODELHYPERPARAMETER COMBINATIONS
  - MODEL INSPECTION AND VISUALISATION
  - ELI5 A LIBRARY FOR DEBUGGINGINSPECTING MACHINE LEARNING MODELS AND EXPLAINING THEIR PREDICTIONS
  - MLXTEND INCLUDES MODEL VISUALIZATION UTILITIES
  - SCIKITPLOT A VISUALIZATION LIBRARY FOR QUICK AND EASY GENERATION OF COMMON PLOTS IN DATA ANALYSIS AND MACHINE LEARNING
  - YELLOWBRICK A SUITE OF CUSTOM MATPLOTLIB VISUALIZERS FOR SCIKITLEARN ESTIMATORS TO SUPPORT VISUAL FEATURE ANALYSIS
  - MODEL SELECTION EVALUATION AND DIAGNOSTICS
  - MODEL EXPORT FOR PRODUCTION
  - ONNXMLTOOLS SERIALIZES MANY SCIKITLEARN PIPELINES TO ONNX FOR INTERCHANGE AND PREDICTION
  - SKLEARN2PMML SERIALIZATION OF A WIDE VARIETY OF SCIKITLEARN ESTIMATORS AND TRANSFORMERS INTO PMML WITH THE HELP OF JPMMLSKLEARN LIBRARY
  - SKLEARNPORTER TRANSPILE TRAINED SCIKITLEARN MODELS TO C JAVA JAVASCRIPT AND OTHERS
  - SKLEARNCOMPILEDTREES GENERATE A C IMPLEMENTATION OF THE PREDICT FUNCTION FOR DECISION TREES AND ENSEMBLES TRAINED BY SKLEARN USEFUL FOR LATENCYSENSITIVE PRODUCTION ENVIRONMENTS
  - 142 OTHER ESTIMATORS AND TASKS
  - NOT EVERYTHING BELONGS OR IS MATURE ENOUGH FOR THE CENTRAL SCIKITLEARN PROJECT THE FOLLOWING ARE PROJECTS PROVIDING INTERFACES SIMILAR TO SCIKITLEARN FOR ADDITIONAL LEARNING ALGORITHMS INFRASTRUCTURES AND TASKS
  - STRUCTURED LEARNING
  - SEQLEARN SEQUENCE CLASSIFICATION USING HMMS OR STRUCTURED PERCEPTRON
  - HMMLEARN IMPLEMENTATION OF HIDDEN MARKOV MODELS THAT WAS PREVIOUSLY PART OF SCIKITLEARN
  - PYSTRUCT GENERAL CONDITIONAL RANDOM FIELDS AND STRUCTURED PREDICTION
  - POMEGRANATE PROBABILISTIC MODELLING FOR PYTHON WITH AN EMPHASIS ON HIDDEN MARKOV MODELS
  - SKLEARNCRFSUITE LINEARCHAIN CONDITIONAL RANDOM FIELDS CRFSUITE WRAPPER WITH SKLEARNLIKE API
  - DEEP NEURAL NETWORKS ETC
  - PYLEARN2 A DEEP LEARNING AND NEURAL NETWORK LIBRARY BUILD ON THEANO WITH SCIKITLEARN LIKE INTERFACE
- 10 CHAPTER 1 WELCOME TO SCIKITLEARN



SCIKITLEARN USER GUIDE RELEASE 0213

- SKLEARNTHEANO SCIKITLEARN COMPATIBLE ESTIMATORS TRANSFORMERS AND DATASETS WHICH USE THEANO INTERNALLY
- NOLEARN A NUMBER OF WRAPPERS AND ABSTRACTIONS AROUND EXISTING NEURAL NETWORK LIBRARIES
- KERAS DEEP LEARNING LIBRARY CAPABLE OF RUNNING ON TOP OF EITHER TENSORFLOW OR THEANO
- LASAGNE A LIGHTWEIGHT LIBRARY TO BUILD AND TRAIN NEURAL NETWORKS IN THEANO
- SKORCH A SCIKITLEARN COMPATIBLE NEURAL NETWORK LIBRARY THAT WRAPS PYTORCH

BROAD SCOPE

- MLXTEND INCLUDES A NUMBER OF ADDITIONAL ESTIMATORS AS WELL AS MODEL VISUALIZATION UTILITIES
- SPARKITLEARN SCIKITLEARN API AND FUNCTIONALITY FOR PYSARK’S DISTRIBUTED MODELLING

OTHER REGRESSION AND CLASSIFICATION

- XGBOOST OPTIMISED GRADIENT BOOSTED DECISION TREE LIBRARY
- MLENSEMBLE GENERALIZED ENSEMBLE LEARNING STACKING BLENDING SUBSEMBLE DEEP ENSEMBLES ETC
- LIGHTNING FAST STATEOFHEART LINEAR MODEL SOLVERS SDCA ADAGRAD SVRG SAG ETC
- PYEARTH MULTIVARIATE ADAPTIVE REGRESSION SPLINES
- KERNEL REGRESSION IMPLEMENTATION OF NADARAYAWATSON KERNEL REGRESSION WITH AUTOMATIC BANDWIDTH SELECTION
- GPLEARN GENETIC PROGRAMMING FOR SYMBOLIC REGRESSION TASKS
- MULTIISOTONIC ISOTONIC REGRESSION ON MULTIDIMENSIONAL FEATURES
- SCIKITMULTILEARN MULTILABEL CLASSIFICATION WITH FOCUS ON LABEL SPACE MANIPULATION
- SEGLEARN TIME SERIES AND SEQUENCE LEARNING USING SLIDING WINDOW SEGMENTATION

DECOMPOSITION AND CLUSTERING

- LDA FAST IMPLEMENTATION OF LATENT DIRICHLET ALLOCATION IN CYTHON WHICH USES GIBBS SAMPLING
- TO SAMPLE FROM THE TRUE POSTERIOR DISTRIBUTION SCIKITLEARN’S SKLEARNDECOMPOSITION
- LATENTDIRICHLETALLOCATION IMPLEMENTATION USES VARIATIONAL INFERENCE TO SAMPLE FROM A TRACTABLE
- APPROXIMATION OF A TOPIC MODEL’S POSTERIOR DISTRIBUTION
- SPARSE FILTERING UNSUPERVISED FEATURE LEARNING BASED ON SPARSEFILTERING
- KMODES KMODES CLUSTERING ALGORITHM FOR CATEGORICAL DATA AND SEVERAL OF ITS VARIATIONS
- HDBSCAN HDBSCAN AND ROBUST SINGLE LINKAGE CLUSTERING ALGORITHMS FOR ROBUST VARIABLE DENSITY CLUSTERING
- SPHERECLUSTER SPHERICAL KMEANS AND MIXTURE OF VON MISES FISHER CLUSTERING ROUTINES FOR DATA ON THE UNIT HYPER SPHERE

PREPROCESSING

- CATEGORICAENCODING A LIBRARY OF SKLEARN COMPATIBLE CATEGORICAL VARIABLE ENCODERS
- IMBALANCEDLEARN VARIOUS METHODS TO UNDER AND OVERSAMPLE DATASETS

143 STATISTICAL LEARNING WITH PYTHON

OTHER PACKAGES USEFUL FOR DATA ANALYSIS AND MACHINE LEARNING

- PANDAS TOOLS FOR WORKING WITH HETEROGENEOUS AND COLUMNAR DATA RELATIONAL QUERIES TIME SERIES AND BASIC STATISTICS
- THEANO A CPUGPU ARRAY PROCESSING FRAMEWORK GEARED TOWARDS DEEP LEARNING RESEARCH

14 RELATED PROJECTS 11

SCIKITLEARN USER GUIDE RELEASE 0213

• STATSMODELS ESTIMATING AND ANALYSING STATISTICAL MODELS MORE FOCUSED ON STATISTICAL TESTS AND LESS ON PREDICTION THAN SCIKITLEARN

- PYMC BAYESIAN STATISTICAL MODELS AND FITTING ALGORITHMS
- SACRED TOOL TO HELP YOU CONFIGURE ORGANIZE LOG AND REPRODUCE EXPERIMENTS
- SEABORN VISUALIZATION LIBRARY BASED ON MATPLOTLIB IT PROVIDES A HIGHLEVEL INTERFACE FOR DRAWING ATTRACTIVE STATISTICAL GRAPHICS

• DEEP LEARNING A CURATED LIST OF DEEP LEARNING SOFTWARE LIBRARIES

DOMAIN SPECIFIC PACKAGES

- SCIKITIMAGE IMAGE PROCESSING AND COMPUTER VISION IN PYTHON
- NATURAL LANGUAGE TOOLKIT NLTK NATURAL LANGUAGE PROCESSING AND SOME MACHINE LEARNING
- GENSIM A LIBRARY FOR TOPIC MODELLING DOCUMENT INDEXING AND SIMILARITY RETRIEVAL
- NILEARN MACHINE LEARNING FOR NEUROIMAGING
- ASTROML MACHINE LEARNING FOR ASTRONOMY
- MSMBUILDER MACHINE LEARNING FOR PROTEIN CONFORMATIONAL DYNAMICS TIME SERIES
- SCIKITSURPRISE A SCIKIT FOR BUILDING AND EVALUATING RECOMMENDER SYSTEMS

144 SNIPPETS AND TIDBITS

THE WIKI HAS MORE

15 ABOUT US

151 HISTORY

THIS PROJECT WAS STARTED IN 2007 AS A GOOGLE SUMMER OF CODE PROJECT BY DAVID COURNAPEAU LATER THAT YEAR MATTHIEU BRUCHER STARTED WORK ON THIS PROJECT AS PART OF HIS THESIS

IN 2010 FABIAN PEDREGOSA GAELE VAROQUAUX ALEXANDRE GRAMFORT AND VINCENT MICHEL OF INRIA TOOK LEADERSHIP OF THE PROJECT AND MADE THE FIRST PUBLIC RELEASE FEBRUARY THE 1ST 2010 SINCE THEN SEVERAL RELEASES HAVE APPEARED FOLLOWING A 3 MONTH CYCLE AND A THRIVING INTERNATIONAL COMMUNITY HAS BEEN LEADING THE DEVELOPMENT

152 GOVERNANCE

THE DECISION MAKING PROCESS AND GOVERNANCE STRUCTURE OF SCIKITLEARN IS LAID OUT IN THE GOVERNANCE DOCUMENT

153 AUTHORS

THE FOLLOWING PEOPLE ARE CURRENTLY CORE CONTRIBUTORS TO SCIKITLEARN'S DEVELOPMENT AND MAINTENANCE

PLEASE DO NOT EMAIL THE AUTHORS DIRECTLY TO ASK FOR ASSISTANCE OR REPORT ISSUES INSTEAD PLEASE SEE WHAT'S THE BEST WAY TO ASK QUESTIONS ABOUT SCIKITLEARN IN THE FAQ

SEE ALSO

12 CHAPTER 1 WELCOME TO SCIKITLEARN

SCIKITLEARN USER GUIDE RELEASE 0213  
HOW YOU CAN CONTRIBUTE TO THE PROJECT  
154 EMERITUS CORE DEVELOPERS  
THE FOLLOWING PEOPLE HAVE BEEN ACTIVE CONTRIBUTORS IN THE PAST BUT ARE NO LONGER ACTIVE IN THE PROJECT

- ALEXANDER FABISCH
- ALEXANDRE PASSOS
- ANGEL SOLER GOLLONET
- ARNAUD JOLY
- CHRIS GORGOLEWSKI
- DAVID COURNAPEAU
- DAVID WARDEFARLEY
- EDUARD DUCHESNAY
- FABIAN PEDRAGOSA
- GILLES LOUPPE
- JACOB SCHREIBER
- JAKE VANDERPLAS
- JAKUES GROBLER
- JARROD MILLMAN
- KYLE KASTNER
- LARS BUITINCK
- MANOJ KUMAR
- MATHIEU BLONDEL
- MATTHIEU BRUCHER
- NOEL DAWE
- PAOLO LOSI
- PETER PRETTENHOFER
- RAGHAV RAJAGOPALAN
- ROBERT LAYTON
- RON WEISS
- SATRAJIT GHOSH
- SHIQIAO DU
- THOUIS RAY JONES
- VINCENT DUBOURG
- VINCENT MICHEL
- VIRGILE FRITSCH
- WEI LI

15 ABOUT US 13

SCIKITLEARN USER GUIDE RELEASE 0213  
155 CITING SCIKITLEARN  
IF YOU USE SCIKITLEARN IN A SCIENTIFIC PUBLICATION WE WOULD APPRECIATE CITATIONS TO THE FOLLOWING PAPER  
SCIKITLEARN MACHINE LEARNING IN PYTHON PEDREGOSA ET AL JMLR 12 PP 28252830 2011  
BIBTEX ENTRY  
ARTICLE SCIKITLEARN  
TITLES  
SCIKITLEARN MACHINE LEARNING IN PYTHON  
AUTHOR PEDREGOSA F ANDVAROQUAUX G ANDGRAMFORT A ANDMICHEL V  
ANDTHIRION B ANDGRISEL O ANDBLONDEL M ANDPRETTENHOFER P  
ANDWEISS R ANDDUBOURG V ANDVANDERPLAS J ANDPASSOS A AND  
COURNAPEAU D ANDBRUCHER M ANDPERROT M ANDDUCHESNAY E  
JOURNAL JOURNAL OF MACHINE LEARNING RESEARCH  
VOLUME 12  
PAGES 28252830  
YEAR 2011

IF YOU WANT TO CITE SCIKITLEARN FOR ITS API OR DESIGN YOU MAY ALSO WANT TO CONSIDER THE FOLLOWING PAPER  
API DESIGN FOR MACHINE LEARNING SOFTWARE EXPERIENCES FROM THE SCIKITLEARN PROJECT BUITINCK ET AL 2013  
BIBTEX ENTRY  
IN PROCEEDINGS SKLEARN API  
AUTHOR LARS BUITINCK ANDGILLES LOUPPE ANDMATHIEU BLONDEL AND  
FABIAN PEDREGOSA ANDANDREAS MUELLER ANDOLIVIER GRISEL AND  
VLAD NICULAE ANDPETER PRETTENHOFER ANDALEXANDRE GRAMFORT  
ANDJAKES GROBLER ANDROBERT LAYTON ANDJAKE VANDERPLAS AND  
ARNAUD JOLY ANDBRIAN HOLT ANDGAEL VAROQUAUX  
TITLE API DESIGN FOR MACHINE LEARNING SOFTWARE EXPERIENCES FROM  
SCIKITLEARN  
PROJECT  
BOOK TITLE ECML PKDD WORKSHOP LANGUAGES FOR DATA MINING AND MACHINE  
LEARNING  
YEAR 2013  
PAGES 108122

156 ARTWORK  
HIGH QUALITY PNG AND SVG LOGOS ARE AVAILABLE IN THE DOCLOGOS SOURCE DIRECTORY  
14 CHAPTER 1 WELCOME TO SCIKITLEARN

SCIKITLEARN USER GUIDE RELEASE 0213

157 FUNDING

INRIA ACTIVELY SUPPORTS THIS PROJECT IT HAS PROVIDED FUNDING FOR FABIAN PEDREGOSA 20102012 JAQUES GROB  
LER 20122013 AND OLIVIER GRISEL 20132017 TO WORK ON THIS PROJECT FULLTIME IT ALSO HOSTS CODING SPRINTS AND  
OTHER EVENTS

PARISSACLAY CENTER FOR DATA SCIENCE FUNDED ONE YEAR  
FOR A DEVELOPER TO WORK ON THE PROJECT FULLTIME 20142015 AND 50 OF THE TIME OF GUILLAUME LEMAITRE 2016  
2017

NYU MOORESLOAN DATA SCIENCE ENVIRONMENT FUNDED AN  
DREAS MUELLER 20142016 TO WORK ON THIS PROJECT THE MOORESLOAN DATA SCIENCE ENVIRONMENT ALSO FUNDS SEV  
ERAL STUDENTS TO WORK ON THE PROJECT PARTTIME

TÉLÉCOM PARIS

TECH FUNDED MANOJ KUMAR 2014 TOM DUPRÉ LA TOUR 2015 RAGHAV RV 20152017 THIERRY GUILLEMOT 2016  
2017 AND ALBERT THOMAS 2017 TO WORK ON SCIKITLEARN

COLUMBIA UNIVERSITY FUNDS AN  
DREAS MÜLLER SINCE 2016

ANDREAS MÜLLER ALSO RECEIVED A GRANT TO IMPROVE SCIKITLEARN  
FROM THE ALFRED P SLOAN FOUNDATION IN 2017

THE UNIVERSITY OF

15 ABOUT US 15

SCIKITLEARN USER GUIDE RELEASE 0213  
SYDNEY FUNDS JOEL NOTHMAN SINCE JULY 2017  
THE LABEX DIGI  
COSME FUNDED NICOLAS GOIX 20152016 TOM DUPRÉ LA TOUR 20152016 AND 20172018 MATHURIN MASSIAS 2018  
2019 TO WORK PART TIME ON SCIKITLEARN DURING THEIR PHDS IT ALSO FUNDED A SCIKITLEARN CODING SPRINT IN 2015  
THE FOLLOWING STUDENTS WERE SPONSORED BY GOOGLE TO WORK ON  
SCIKITLEARN THROUGH THE GOOGLE SUMMER OF CODE PROGRAM

- 2007 DAVID COURNAPEAU
- 2011 VLAD NICULAE
- 2012 VLAD NICULAE IMMANUEL BAYER
- 2013 KEMAL EREN NICOLAS TRÉSEGNIE
- 2014 HAMZEH ALSALHI ISSAM LARADJI MAHESHAKYA WIJEWARDENA MANOJ KUMAR
- 2015 RAGHAV RV WEI XUE
- 2016 NELSON LIU YENCHEN LIN

IT ALSO PROVIDED FUNDING FOR SPRINTS AND EVENTS AROUND SCIKITLEARN IF YOU WOULD LIKE TO PARTICIPATE IN THE NEXT GOOGLE  
SUMMER OF CODE PROGRAM PLEASE SEE THIS PAGE  
THE NEURODEBIAN PROJECT PROVIDING DEBIAN PACKAGING AND CONTRIBUTIONS IS SUPPORTED BY DR JAMES V HAXBY DART  
MOUTH COLLEGE  
THE PSF HELPED FIND AND MANAGE FUNDING FOR OUR 2011 GRANADA SPRINT MORE INFORMATION CAN BE FOUND HERE  
TINYCLUES FUNDED THE 2011 INTERNATIONAL GRANADA SPRINT  
DONATING TO THE PROJECT  
IF YOU ARE INTERESTED IN DONATING TO THE PROJECT OR TO ONE OF OUR CODESPRINTS YOU CAN USE THE PAYPAL BUTTON BELOW OR THE  
NUMFOCUS DONATIONS PAGE IF YOU USE THE LATTER PLEASE INDICATE THAT YOU ARE DONATING FOR THE SCIKITLEARN PROJECT  
ALL DONATIONS WILL BE HANDLED BY NUMFOCUS A NONPROFITORGANIZATION WHICH IS MANAGED BY A BOARD OF SCIPY  
COMMUNITY MEMBERS NUMFOCUS'S MISSION IS TO FOSTER SCIENTIFIC COMPUTING SOFTWARE IN PARTICULAR IN PYTHON AS  
A FISCAL HOME OF SCIKITLEARN IT ENSURES THAT MONEY IS AVAILABLE WHEN NEEDED TO KEEP THE PROJECT FUNDED AND AVAILABLE  
WHILE IN COMPLIANCE WITH TAX REGULATIONS  
THE RECEIVED DONATIONS FOR THE SCIKITLEARN PROJECT MOSTLY WILL GO TOWARDS COVERING TRAVELEXPENSES FOR CODE SPRINTS AS  
WELL AS TOWARDS THE ORGANIZATION BUDGET OF THE PROJECT1  
1REGARDING THE ORGANIZATION BUDGET IN PARTICULAR WE MIGHT USE SOME OF THE DONATED FUNDS TO PAY FOR OTHER PROJECT EXPENSES  
HOSTING OR CONTINUOUS INTEGRATION SERVICES  
16 CHAPTER 1 WELCOME TO SCIKITLEARN

SCIKITLEARN USER GUIDE RELEASE 0213

NOTES

THE 2013 PARIS INTERNATIONAL SPRINT

FIG 11 IAP VII19 DYSO

FOR MORE INFORMATION ON THIS SPRINT SEE HERE

158 INFRASTRUCTURE SUPPORT

- WE WOULD LIKE TO THANK RACKSPACE FOR PROVIDING US WITH A FREE RACKSPACE CLOUD ACCOUNT TO AUTOMATICALLY BUILD THE DOCUMENTATION AND THE EXAMPLE GALLERY FROM FOR THE DEVELOPMENT VERSION OF SCIKITLEARN USING THIS TOOL
- WE WOULD ALSO LIKE TO THANK SHINING PANDA FOR FREE CPU TIME ON THEIR CONTINUOUS INTEGRATION SERVER

16 WHO IS USING SCIKITLEARN

161 JPMORGAN

SCIKITLEARN IS AN INDISPENSABLE PART OF THE PYTHON MACHINE LEARNING TOOLKIT AT JPMORGAN IT IS VERY WIDELY USED ACROSS ALL PARTS OF THE BANK FOR CLASSIFICATION PREDICTIVE ANALYTICS AND VERY MANY OTHER MACHINE LEARNING TASKS ITS STRAIGHTFORWARD API ITS BREADTH OF ALGORITHMS AND THE QUALITY OF ITS DOCUMENTATION COMBINE TO MAKE SCIKITLEARN SIMULTANEOUSLY VERY APPROACHABLE AND VERY POWERFUL  
STEPHEN SIMMONS VP ATHENA RESEARCH JPMORGAN

16 WHO IS USING SCIKITLEARN 17

SCIKITLEARN USER GUIDE RELEASE 0213

162 SPOTIFY

SCIKITLEARN PROVIDES A TOOLBOX WITH SOLID IMPLEMENTATIONS OF A BUNCH OF STATEOFTHEART MODELS AND MAKES IT EASY TO PLUG THEM INTO EXISTING APPLICATIONS WE'VE BEEN USING IT QUITE A LOT FOR MUSIC RECOMMENDATIONS AT SPOTIFY AND I THINK IT'S THE MOST WELLDESIGNED ML PACKAGE I'VE SEEN SO FAR

ERIK BERNHARDSSON ENGINEERING MANAGER MUSIC DISCOVERY MACHINE LEARNING SPOTIFY

163 INRIA

AT INRIA WE USE SCIKITLEARN TO SUPPORT LEADINGEDGE BASIC RESEARCH IN MANY TEAMS PARIETAL FOR NEUROIMAGING LEAR FOR COMPUTER VISION VISAGES FOR MEDICAL IMAGE ANALYSIS PRIVATICS FOR SECURITY THE PROJECT IS A FANTASTIC TOOL TO ADDRESS DIFFICULT APPLICATIONS OF MACHINE LEARNING IN AN ACADEMIC ENVIRONMENT AS IT IS PERFORMANT AND VERSATILE BUT ALL EASYTOUSE AND WELL DOCUMENTED WHICH MAKES IT WELL SUITED TO GRAD STUDENTS

GAËL VAROQUAUX RESEARCH AT PARIETAL

164 BETAWORKS

BETAWORKS IS A NYCBASED STARTUP STUDIO THAT BUILDS NEW PRODUCTS GROWS COMPANIES AND INVESTS IN OTHERS OVER THE PAST 8 YEARS WE'VE LAUNCHED A HANDFUL OF SOCIAL DATA ANALYTICSDRIVEN SERVICES SUCH AS BITLY CHARTBEAT DIGG AND SCALE MODEL CONSISTENTLY THE BETAWORKS DATA SCIENCE TEAM USES SCIKITLEARN FOR A VARIETY OF TASKS FROM EXPLORATORY ANALYSIS TO PRODUCT DEVELOPMENT IT IS AN ESSENTIAL PART OF OUR TOOLKIT RECENT USES ARE INCLUDED IN DIGG'S NEW VIDEO RECOMMENDER SYSTEM AND PONCHO'S DYNAMIC HEURISTIC SUBSPACE CLUSTERING

GILAD LOTAN CHIEF DATA SCIENTIST

18 CHAPTER 1 WELCOME TO SCIKITLEARN



SCIKITLEARN USER GUIDE RELEASE 0213

165 HUGGING FACE

AT HUGGING FACE WE'RE USING NLP AND PROBABILISTIC MODELS TO GENERATE CONVERSATIONAL ARTIFICIAL INTELLIGENCES THAT ARE FUN TO CHAT WITH DESPITE USING DEEP NEURAL NETS FOR A FEW OF OUR NLP TASKS SCIKITLEARN IS STILL THE BREADANDBUTTER OF OUR DAILY MACHINE LEARNING ROUTINE THE EASE OF USE AND PREDICTABILITY OF THE INTERFACE AS WELL AS THE STRAIGHTFORWARD MATHEMATICAL EXPLANATIONS THAT ARE HERE WHEN YOU NEED THEM IS THE KILLER FEATURE WE USE A VARIETY OF SCIKITLEARN MODELS IN PRODUCTION AND THEY ARE ALSO OPERATIONALLY VERY PLEASANT TO WORK WITH

JULIEN CHAUMOND CHIEF TECHNOLOGY OFFICER

166 EVERNOTE

BUILDING A CLASSIFIER IS TYPICALLY AN ITERATIVE PROCESS OF EXPLORING THE DATA SELECTING THE FEATURES THE ATTRIBUTES OF THE DATA BELIEVED TO BE PREDICTIVE IN SOME WAY TRAINING THE MODELS AND FINALLY EVALUATING THEM FOR MANY OF THESE TASKS WE RELIED ON THE EXCELLENT SCIKITLEARN PACKAGE FOR PYTHON

READ MORE

MARK AYZENSHTAT VP AUGMENTED INTELLIGENCE

167 TÉLÉCOM PARISTECH

AT TELECOM PARISTECH SCIKITLEARN IS USED FOR HANDSON SESSIONS AND HOME ASSIGNMENTS IN INTRODUCTORY AND ADVANCED MACHINE LEARNING COURSES THE CLASSES ARE FOR UNDERGRADS AND MASTERS STUDENTS THE GREAT BENEFIT OF SCIKITLEARN IS ITS FAST LEARNING CURVE THAT ALLOWS STUDENTS TO QUICKLY START WORKING ON INTERESTING AND MOTIVATING PROBLEMS

16 WHO IS USING SCIKITLEARN 19

SCIKITLEARN USER GUIDE RELEASE 0213  
ALEXANDRE GRAMFORT ASSISTANT PROFESSOR  
168 BOOKINGCOM

AT BOOKINGCOM WE USE MACHINE LEARNING ALGORITHMS FOR MANY DIFFERENT APPLICATIONS SUCH AS RECOMMENDING HOTELS AND DESTINATIONS TO OUR CUSTOMERS DETECTING FRAUDULENT RESERVATIONS OR SCHEDULING OUR CUSTOMER SERVICE AGENTS SCIKITLEARN IS ONE OF THE TOOLS WE USE WHEN IMPLEMENTING STANDARD ALGORITHMS FOR PREDICTION TASKS ITS API AND DOCUMENTATIONS ARE EXCELLENT AND MAKE IT EASY TO USE THE SCIKITLEARN DEVELOPERS DO A GREAT JOB OF INCORPORATING STATE OF THE ART IMPLEMENTATIONS AND NEW ALGORITHMS INTO THE PACKAGE THUS SCIKITLEARN PROVIDES CONVENIENT ACCESS TO A WIDE SPECTRUM OF ALGORITHMS AND ALLOWS US TO READILY FIND THE RIGHT TOOL FOR THE RIGHT JOB  
MELANIE MUELLER DATA SCIENTIST  
169 AWEBER

THE SCIKITLEARN TOOLKIT IS INDISPENSABLE FOR THE DATA ANALYSIS AND MANAGEMENT TEAM AT AWEBER IT ALLOWS US TO DO AWESOME STUFF WE WOULD NOT OTHERWISE HAVE THE TIME OR RESOURCES TO ACCOMPLISH THE DOCUMENTATION IS EXCELLENT ALLOWING NEW ENGINEERS TO QUICKLY EVALUATE AND APPLY MANY DIFFERENT ALGORITHMS TO OUR DATA THE TEXT FEATURE EXTRACTION UTILITIES ARE USEFUL WHEN WORKING WITH THE LARGE VOLUME OF EMAIL CONTENT WE HAVE AT AWEBER THE RANDOMIZEDPCA IMPLEMENTATION ALONG WITH PIPELINING AND FEATUREUNIONS ALLOWS US TO DEVELOP COMPLEX MACHINE LEARNING ALGORITHMS EFFICIENTLY AND RELIABLY  
ANYONE INTERESTED IN LEARNING MORE ABOUT HOW AWEBER DEPLOYS SCIKITLEARN IN A PRODUCTION ENVIRONMENT SHOULD CHECK OUT TALKS FROM PYDATA BOSTON BY AWEBER'S MICHAEL BECKER AVAILABLE AT [HTTPSGITHUBCOMMDBECKERPYDATA2013](https://github.com/mdbecker/pydata2013)  
MICHAEL BECKER SOFTWARE ENGINEER DATA ANALYSIS AND MANAGEMENT NINJAS  
1610 YHAT

THE COMBINATION OF CONSISTENT APIS THOROUGH DOCUMENTATION AND TOP NOTCH IMPLEMENTATION MAKE SCIKITLEARN OUR FAVORITE MACHINE LEARNING PACKAGE IN PYTHON SCIKITLEARN MAKES DOING ADVANCED ANALYSIS IN PYTHON ACCESSIBLE TO ANYONE AT YHAT WE MAKE IT EASY TO INTEGRATE THESE MODELS INTO YOUR PRODUCTION APPLICATIONS THUS ELIMINATING THE UNNECESSARY DEV TIME ENCOUNTERED PRODUCTIONIZING ANALYTICAL WORK  
GREG LAMP COFOUNDER YHAT  
20 CHAPTER 1 WELCOME TO SCIKITLEARN

SCIKITLEARN USER GUIDE RELEASE 0213

1611 RANGESPAN

THE PYTHON SCIKITLEARN TOOLKIT IS A CORE TOOL IN THE DATA SCIENCE GROUP AT RANGESPAN ITS LARGE COLLECTION OF WELL DOCUMENTED MODELS AND ALGORITHMS ALLOW OUR TEAM OF DATA SCIENTISTS TO PROTOTYPE FAST AND QUICKLY ITERATE TO FIND THE RIGHT SOLUTION TO OUR LEARNING PROBLEMS WE FIND THAT SCIKITLEARN IS NOT ONLY THE RIGHT TOOL FOR PROTOTYPING BUT ITS CAREFUL AND WELL TESTED IMPLEMENTATION GIVE US THE CONFIDENCE TO RUN SCIKITLEARN MODELS IN PRODUCTION

JURGEN VAN GAEL DATA SCIENCE DIRECTOR AT RANGESPAN LTD

1612 BIRCHBOX

AT BIRCHBOX WE FACE A RANGE OF MACHINE LEARNING PROBLEMS TYPICAL TO ECOMMERCE PRODUCT RECOMMENDATION USER CLUSTERING INVENTORY PREDICTION TRENDS DETECTION ETC SCIKITLEARN LETS US EXPERIMENT WITH MANY MODELS ESPECIALLY IN THE EXPLORATION PHASE OF A NEW PROJECT THE DATA CAN BE PASSED AROUND IN A CONSISTENT WAY MODELS ARE EASY TO SAVE AND REUSE UPDATES KEEP US INFORMED OF NEW DEVELOPMENTS FROM THE PATTERN DISCOVERY RESEARCH COMMUNITY SCIKITLEARN IS AN IMPORTANT TOOL FOR OUR TEAM BUILT THE RIGHT WAY IN THE RIGHT LANGUAGE

THIERRY BERTINMAHIEUX BIRCHBOX DATA SCIENTIST

1613 BESTOFMEDIA GROUP

SCIKITLEARN IS OUR 1 TOOLKIT FOR ALL THINGS MACHINE LEARNING AT BESTOFMEDIA WE USE IT FOR A VARIETY OF TASKS EG SPAM FIGHTING AD CLICK PREDICTION VARIOUS RANKING MODELS THANKS TO THE VARIED STATEOFHEART ALGORITHM IMPLEMENTATIONS PACKAGED INTO IT IN THE LAB IT ACCELERATES PROTOTYPING OF COMPLEX PIPELINES IN PRODUCTION I CAN SAY IT HAS PROVEN TO BE ROBUST AND EFFICIENT ENOUGH TO BE DEPLOYED FOR BUSINESS CRITICAL COMPONENTS

EUSTACHE DIEMERT LEAD SCIENTIST BESTOFMEDIA GROUP

16 WHO IS USING SCIKITLEARN 21

SCIKITLEARN USER GUIDE RELEASE 0213

1614 CHANGEORG

AT CHANGEORG WE AUTOMATE THE USE OF SCIKITLEARN’S RANDOMFORESTCLASSIFIER IN OUR PRODUCTION SYSTEMS TO DRIVE EMAIL TARGETING THAT REACHES MILLIONS OF USERS ACROSS THE WORLD EACH WEEK IN THE LAB SCIKITLEARN’S EASEOFUSE PERFORMANCE AND OVERALL VARIETY OF ALGORITHMS IMPLEMENTED HAS PROVED INVALUABLE IN GIVING US A SINGLE RELIABLE SOURCE TO TURN TO FOR OUR MACHINELEARNING NEEDS

VIJAY RAMESH SOFTWARE ENGINEER IN DATASCIENCE AT CHANGEORG

1615 PHIMECA ENGINEERING

AT PHIMECA ENGINEERING WE USE SCIKITLEARN ESTIMATORS AS SURROGATES FOR EXPENSIVETOEVALUATE NUMERICAL MODELS MOSTLY BUT NOT EXCLUSIVELY FINITEELEMENT MECHANICAL MODELS FOR SPEEDING UP THE INTENSIVE POSTPROCESSING OPERATIONS INVOLVED IN OUR SIMULATIONBASED DECISION MAKING FRAMEWORK SCIKITLEARN’S FITPREDICT API TOGETHER WITH ITS EFFICIENT CROSSVALIDATION TOOLS CONSIDERABLY EASES THE TASK OF SELECTING THE BESTFIT ESTIMATOR WE ARE ALSO USING SCIKITLEARN FOR ILLUSTRATING CONCEPTS IN OUR TRAINING SESSIONS TRAINEES ARE ALWAYS IMPRESSED BY THE EASEOFUSE OF SCIKITLEARN DESPITE THE APPARENT THEORETICAL COMPLEXITY OF MACHINE LEARNING

VINCENT DUBOURG PHIMECA ENGINEERING PHD ENGINEER

1616 HOWABOUTWE

AT HOWABOUTWE SCIKITLEARN LETS US IMPLEMENT A WIDE ARRAY OF MACHINE LEARNING TECHNIQUES IN ANALYSIS AND IN PRODUCTION DESPITE HAVING A SMALL TEAM WE USE SCIKITLEARN’S CLASSIFICATION ALGORITHMS TO PREDICT USER BEHAVIOR ENABLING US TO FOR EXAMPLE ESTIMATE THE VALUE OF LEADS FROM A GIVEN TRAFFIC SOURCE EARLY IN THE LEAD’S TENURE ON OUR SITE ALSO OUR 22 CHAPTER 1 WELCOME TO SCIKITLEARN

SCIKITLEARN USER GUIDE RELEASE 0213

USERS' PROFILES CONSIST OF PRIMARILY UNSTRUCTURED DATA ANSWERS TO OPENENDED QUESTIONS SO WE USE SCIKITLEARN'S FEATURE EXTRACTION AND DIMENSIONALITY REDUCTION TOOLS TO TRANSLATE THESE UNSTRUCTURED DATA INTO INPUTS FOR OUR MATCHMAKING SYSTEM

DANIEL WEITZENFELD SENIOR DATA SCIENTIST AT HOWABOUTWE

1617 PEERINDEX

AT PEERINDEX WE USE SCIENTIFIC METHODOLOGY TO BUILD THE INFLUENCE GRAPH A UNIQUE DATASET THAT ALLOWS US TO IDENTIFY WHO'S REALLY INFLUENTIAL AND IN WHICH CONTEXT TO DO THIS WE HAVE TO TACKLE A RANGE OF MACHINE LEARNING AND PREDICTIVE MODELING PROBLEMS SCIKITLEARN HAS EMERGED AS OUR PRIMARY TOOL FOR DEVELOPING PROTOTYPES AND MAKING QUICK PROGRESS FROM PREDICTING MISSING DATA AND CLASSIFYING TWEETS TO CLUSTERING COMMUNITIES OF SOCIAL MEDIA USERS SCIKITLEARN PROVED USEFUL IN A VARIETY OF APPLICATIONS ITS VERY INTUITIVE INTERFACE AND EXCELLENT COMPATIBILITY WITH OTHER PYTHON TOOLS MAKES IT AN INDISPENSABLE TOOL IN OUR DAILY RESEARCH EFFORTS

FERENC HUSZAR SENIOR DATA SCIENTIST AT PEERINDEX

1618 DATAROBOT

DATAROBOT IS BUILDING NEXT GENERATION PREDICTIVE ANALYTICS SOFTWARE TO MAKE DATA SCIENTISTS MORE PRODUCTIVE AND SCIKITLEARN IS AN INTEGRAL PART OF OUR SYSTEM THE VARIETY OF MACHINE LEARNING TECHNIQUES IN COMBINATION WITH THE SOLID IMPLEMENTATIONS THAT SCIKITLEARN OFFERS MAKES IT A ONESTOPSHOPPING LIBRARY FOR MACHINE LEARNING IN PYTHON MOREOVER ITS CONSISTENT API WELLTESTED CODE AND PERMISSIVE LICENSING ALLOW US TO USE IT IN A PRODUCTION ENVIRONMENT SCIKITLEARN HAS LITERALLY SAVED US YEARS OF WORK WE WOULD HAVE HAD TO DO OURSELVES TO BRING OUR PRODUCT TO MARKET JEREMY ACHIN CEO COFOUNDER DATAROBOT INC

1619 OKCUPID

WE'RE USING SCIKITLEARN AT OKCUPID TO EVALUATE AND IMPROVE OUR MATCHMAKING SYSTEM THE RANGE OF FEATURES IT HAS ESPECIALLY PREPROCESSING UTILITIES MEANS WE CAN USE IT FOR A WIDE VARIETY OF PROJECTS AND IT'S PERFORMANT ENOUGH TO HANDLE THE VOLUME OF DATA THAT WE NEED TO SORT THROUGH THE DOCUMENTATION IS REALLY THOROUGH AS WELL WHICH MAKES THE LIBRARY QUITE EASY TO USE

DAVID KOH SENIOR DATA SCIENTIST AT OKCUPID

16 WHO IS USING SCIKITLEARN 23

SCIKITLEARN USER GUIDE RELEASE 0213

1620 LOVELY

AT LOVELY WE STRIVE TO DELIVER THE BEST APARTMENT MARKETPLACE WITH RESPECT TO OUR USERS AND OUR LISTINGS FROM UNDERSTANDING USER BEHAVIOR IMPROVING DATA QUALITY AND DETECTING FRAUD SCIKITLEARN IS A REGULAR TOOL FOR GATHERING INSIGHTS PREDICTIVE MODELING AND IMPROVING OUR PRODUCT THE EASYTOREAD DOCUMENTATION AND INTUITIVE ARCHITECTURE OF THE API MAKES MACHINE LEARNING BOTH EXPLORABLE AND ACCESSIBLE TO A WIDE RANGE OF PYTHON DEVELOPERS I’M CONSTANTLY RECOMMENDING THAT MORE DEVELOPERS AND SCIENTISTS TRY SCIKITLEARN

SIMON FRID DATA SCIENTIST LEAD AT LOVELY

1621 DATA PUBLICA

DATA PUBLICA BUILDS A NEW PREDICTIVE SALES TOOL FOR COMMERCIAL AND MARKETING TEAMS CALLED CRADAR WE EXTENSIVELY USE SCIKITLEARN TO BUILD SEGMENTATIONS OF CUSTOMERS THROUGH CLUSTERING AND TO PREDICT FUTURE CUSTOMERS BASED ON PAST PARTNERSHIPS SUCCESS OR FAILURE WE ALSO CATEGORIZE COMPANIES USING THEIR WEBSITE COMMUNICATION THANKS TO SCIKITLEARN AND ITS MACHINE LEARNING ALGORITHM IMPLEMENTATIONS EVENTUALLY MACHINE LEARNING MAKES IT POSSIBLE TO DETECT WEAK SIGNALS THAT TRADITIONAL TOOLS CANNOT SEE ALL THESE COMPLEX TASKS ARE PERFORMED IN AN EASY AND STRAIGHTFORWARD WAY THANKS TO THE GREAT QUALITY OF THE SCIKITLEARN FRAMEWORK

GUILLAUME LEBOURGEOIS SAMUEL CHARRON DATA SCIENTISTS AT DATA PUBLICA

1622 MACHINALIS

SCIKITLEARN IS THE CORNERSTONE OF ALL THE MACHINE LEARNING PROJECTS CARRIED AT MACHINALIS IT HAS A CONSISTENT API A WIDE SELECTION OF ALGORITHMS AND LOTS OF AUXILIARY TOOLS TO DEAL WITH THE BOILERPLATE WE HAVE USED IT IN PRODUCTION ENVIRONMENTS ON A VARIETY OF PROJECTS INCLUDING CLICKTHROUGH RATE PREDICTION INFORMATION EXTRACTION AND EVEN COUNTING SHEEP

IN FACT WE USE IT SO MUCH THAT WE’VE STARTED TO FREEZE OUR COMMON USE CASES INTO PYTHON PACKAGES SOME OF THEM OPENSOURCED LIKE FEATUREFORGE SCIKITLEARN IN ONE WORD AWESOME

RAFAEL CARRASCOSA LEAD DEVELOPER

24 CHAPTER 1 WELCOME TO SCIKITLEARN

SCIKITLEARN USER GUIDE RELEASE 0213

1623 SOLIDO

SCIKITLEARN IS HELPING TO DRIVE MOORE’S LAW VIA SOLIDO SOLIDO CREATES COMPUTERAIDED DESIGN TOOLS USED BY THE MAJORITY OF TOP20 SEMICONDUCTOR COMPANIES AND FABs TO DESIGN THE BLEEDINGEDGE CHIPS INSIDE SMARTPHONES AUTO MOBILES AND MORE SCIKITLEARN HELPS TO POWER SOLIDO’S ALGORITHMS FOR RAREEVENT ESTIMATION WORSTCASE VERIFICATION OPTIMIZATION AND MORE AT SOLIDO WE ARE PARTICULARLY FOND OF SCIKITLEARN’S LIBRARIES FOR GAUSSIAN PROCESS MODELS LARGESCALE REGULARIZED LINEAR REGRESSION AND CLASSIFICATION SCIKITLEARN HAS INCREASED OUR PRODUCTIVITY BECAUSE FOR MANY ML PROBLEMS WE NO LONGER NEED TO “ROLL OUR OWN” CODE THIS PYDATA 2014 TALK HAS DETAILS TRENT MCCONAGHY FOUNDER SOLIDO DESIGN AUTOMATION INC

1624 INFONEA

WE EMPLOY SCIKITLEARN FOR RAPID PROTOTYPING AND CUSTOMMADE DATA SCIENCE SOLUTIONS WITHIN OUR INMEMORY BASED BUSINESS INTELLIGENCE SOFTWARE INFONEA® AS A WELLDOCUMENTED AND COMPREHENSIVE COLLECTION OF STATEOFTHEART ALGORITHMS AND PIPELINING METHODS SCIKITLEARN ENABLES US TO PROVIDE FLEXIBLE AND SCALABLE SCIENTIFIC ANALYSIS SOLUTIONS THUS SCIKITLEARN IS IMMENSELY VALUABLE IN REALIZING A POWERFUL INTEGRATION OF DATA SCIENCE TECHNOLOGY WITHIN SELF SERVICE BUSINESS ANALYTICS THORSTEN KRANZ DATA SCIENTIST COMA SOFT AG

1625 DATAKU

OUR SOFTWARE DATA SCIENCE STUDIO DSS ENABLES USERS TO CREATE DATA SERVICES THAT COMBINE ETL WITH MACHINE LEARNING OUR MACHINE LEARNING MODULE INTEGRATES MANY SCIKITLEARN ALGORITHMS THE SCIKITLEARN LIBRARY IS A PERFECT INTEGRATION WITH DSS BECAUSE IT OFFERS ALGORITHMS FOR VIRTUALLY ALL BUSINESS CASES OUR GOAL IS TO OFFER A TRANSPARENT AND FLEXIBLE TOOL THAT MAKES IT EASIER TO OPTIMIZE TIME CONSUMING ASPECTS OF BUILDING A DATA SERVICE PREPARING DATA AND TRAINING MACHINE LEARNING ALGORITHMS ON ALL TYPES OF DATA FLORIAN DOUETTEAU CEO DATAKU

1626 OTTO GROUP

HERE AT OTTO GROUP ONE OF GLOBAL BIG FIVE B2C ONLINE RETAILERS WE ARE USING SCIKITLEARN IN ALL ASPECTS OF OUR DAILY WORK FROM DATA EXPLORATION TO DEVELOPMENT OF MACHINE LEARNING APPLICATION TO THE PRODUCTIVE DEPLOYMENT OF THOSE SERVICES IT HELPS US TO TACKLE MACHINE LEARNING PROBLEMS RANGING FROM ECOMMERCE TO LOGISTICS IT CONSISTENT APIS ENABLED US TO BUILD THE PALLADIUM RESTAPI FRAMEWORK AROUND IT AND CONTINUOUSLY DELIVER SCIKITLEARN BASED SERVICES 16 WHO IS USING SCIKITLEARN 25

SCIKITLEARN USER GUIDE RELEASE 0213  
CHRISTIAN RAMMIG HEAD OF DATA SCIENCE OTTO GROUP  
1627 ZOPA

AT ZOPA THE FIRST EVER PEERTOPEER LENDING PLATFORM WE EXTENSIVELY USE SCIKITLEARN TO RUN THE BUSINESS AND OPTIMIZE OUR USERS' EXPERIENCE IT POWERS OUR MACHINE LEARNING MODELS INVOLVED IN CREDIT RISK FRAUD RISK MARKETING AND PRICING AND HAS BEEN USED FOR ORIGINATING AT LEAST 1 BILLION GBP WORTH OF ZOPA LOANS IT IS VERY WELL DOCUMENTED POWERFUL AND SIMPLE TO USE WE ARE GRATEFUL FOR THE CAPABILITIES IT HAS PROVIDED AND FOR ALLOWING US TO DELIVER ON OUR MISSION OF MAKING MONEY SIMPLE AND FAIR  
VLASIOS VASILEIOU HEAD OF DATA SCIENCE ZOPA

1628 MARS  
SCIKITLEARN IS INTEGRAL TO THE MACHINE LEARNING ECOSYSTEM AT MARS WHETHER WE'RE DESIGNING BETTER RECIPES FOR PETFOOD OR CLOSELY ANALYSING OUR COCOA SUPPLY CHAIN SCIKITLEARN IS USED AS A TOOL FOR RAPIDLY PROTOTYPING IDEAS AND TAKING THEM TO PRODUCTION THIS ALLOWS US TO BETTER UNDERSTAND AND MEET THE NEEDS OF OUR CONSUMERS WORLDWIDE SCIKITLEARN'S FEATURERICH TOOLSET IS EASY TO USE AND EQUIPS OUR ASSOCIATES WITH THE CAPABILITIES THEY NEED TO SOLVE THE BUSINESS CHALLENGES THEY FACE EVERY DAY  
MICHAEL FITZKE NEXT GENERATION TECHNOLOGIES SR LEADER MARS INC

17 RELEASE HISTORY  
RELEASE NOTES FOR CURRENT AND RECENT RELEASES ARE DETAILED ON THIS PAGE WITH PREVIOUS RELEASES LINKED BELOW  
TIP SUBSCRIBE TO SCIKITLEARN RELEASES ON LIBRARIESIO TO BE NOTIFIED WHEN NEW VERSIONS ARE RELEASED

171 LEGEND FOR CHANGELOGS

- MAJOR FEATURE SOMETHING BIG THAT YOU COULDN'T DO BEFORE
- FEATURE SOMETHING THAT YOU COULDN'T DO BEFORE
- EFFICIENCY AN EXISTING FEATURE NOW MAY NOT REQUIRE AS MUCH COMPUTATION OR MEMORY
- ENHANCEMENT A MISCELLANEOUS MINOR IMPROVEMENT
- FIX SOMETHING THAT PREVIOUSLY DIDN'T WORK AS DOCUMENTATED - OR ACCORDING TO REASONABLE EXPECTATIONS - SHOULD NOW WORK
- API C HANGE YOU WILL NEED TO CHANGE YOUR CODE TO HAVE THE SAME EFFECT IN THE FUTURE OR A FEATURE WILL BE REMOVED IN THE FUTURE

26 CHAPTER 1 WELCOME TO SCIKITLEARN



SCIKITLEARN USER GUIDE RELEASE 0213

18 VERSION 0213

JULY 30 2019

181 CHANGED MODELS

THE FOLLOWING ESTIMATORS AND FUNCTIONS WHEN FIT WITH THE SAME DATA AND PARAMETERS MAY PRODUCE DIFFERENT MODELS FROM THE PREVIOUS VERSION THIS OFTEN OCCURS DUE TO CHANGES IN THE MODELLING LOGIC BUG FIXES OR ENHANCEMENTS OR IN RANDOM SAMPLING PROCEDURES

- THE V0200 RELEASE NOTES FAILED TO MENTION A BACKWARDS INCOMPATIBILITY IN METRICSMAKESCORER WHENNEEDSPROBATRUE ANDYTRUE IS BINARY NOW THE SCORER FUNCTION IS SUPPOSED TO ACCEPT A 1D YPRED IE PROBABILITY OF THE POSITIVE CLASS SHAPE NSAMPLES INSTEAD OF A 2D YPRED IE SHAPE NSAMPLES 2

182 CHANGELOG

SKLEARNCLUSTER

- FIXFIXED A BUG IN CLUSTERKMEANS WHERE COMPUTATION WITH INITRANDOM WAS SINGLE THREADED FOR NJOBS 1 ORNJOBS 1 12955 BY PRABAKARAN KUMARESSHAN
- FIXFIXED A BUG IN CLUSTEROPTICS WHERE USERS WERE UNABLE TO PASS FLOAT MINSAMPLES AND MINCLUSTERSIZE 14496 BY FABIAN KLOPPER AND HANMIN QIN

SKLEARNCOMPOSE

- FIXFIXED AN ISSUE IN COMPOSECOLUMNTRANSFORMER WHERE USING DATAFRAMES WHOSE COLUMN ORDER DIFFERS BETWEEN FUNC FIT AND FUNCTransFORM COULD LEAD TO SILENTLY PASSING INCORRECT COLUMNS TO THE REMAINDER TRANSFORMER 14237 BY ANDREAS SCHUDERER

SKLEARNDATASETS

- FIXDATASETSFETCHCALIFORNIAHOUSING DATASETSFETCHCOVTYPE DATASETS FETCHKDDCUP99 DATASETSFETCHOLIVETTIFACES DATASETSFETCHRCV1 AND DATASETSFETCHSPECIESDISTRIBUTIONS TRY TO PERSIST THE PREVIOUSLY CACHE USING THE NEW JOBLIB IF THE CACHED DATA WAS PERSISTED USING THE DEPRECATED SKLEARNEXTERNALSJOLIB THIS BEHAVIOR IS SET TO BE DEPRECATED AND REMOVED IN V023 14197 BY ADRIN JALALI

SKLEARNENSEMBLE

- FIX FIX ZERO DIVISION ERROR IN HISTGRADIENTBOOSTINGCLASSIFIER AND HISTGRADIENTBOOSTINGREGRESSOR 14024 BY NICOLAS HUG

SKLEARNIMPUTE

- FIXFIXED A BUG IN IMPUTESIMPLEIMPUTER ANDIMPUTEITERATIVEIMPUTER SO THAT NO ERRORS ARE THROWN WHEN THERE ARE MISSING VALUES IN TRAINING DATA 13974 BY FRANK HOANG

18 VERSION 0213 27

SCIKITLEARN USER GUIDE RELEASE 0213

SKLEARNINSPECTION

- FIXFIXED A BUG IN INSPECTIONPLOTPARTIALDEPENDENCE WHERE TARGET PARAMETER WAS NOT BEING TAKEN INTO ACCOUNT FOR MULTICLASS PROBLEMS 14393 BY GUILLEM G SUBIES

SKLEARNLINEARMODEL

- FIXFIXED A BUG IN LINEARMODELLOGISTICREGRESSIONCV WHERE REFIT FALSE WOULD FAIL DEPENDING ON THE MULTICLASS AND PENALTY PARAMETERS REGRESSION INTRODUCED IN 021 14087 BY NICOLAS HUG
- FIXCOMPATIBILITY FIX FOR LINEARMODELARDREGRESSION AND SCIPY130 ADAPTS TO UPSTREAM CHANGES TO THE DEFAULT PINVH CUTOFF THRESHOLD WHICH OTHERWISE RESULTS IN POOR ACCURACY IN SOME CASES 14067 BY TIM STALEY

SKLEARNNEIGHBORS

- FIXFIXED A BUG IN NEIGHBORSNEIGHBORHOODCOMPONENTSANALYSIS WHERE THE VALIDATION OF INITIAL PARAMETERS NCOMPONENTS MAXITER AND TOL REQUIRED TOO STRICT TYPES 14092 BY JÉRÉMIE DU BOISBERANGER

SKLEARNTREE

- FIXFIXED BUG IN TREEEXPORTTEXT WHEN THE TREE HAS ONE FEATURE AND A SINGLE FEATURE NAME IS PASSED IN 14053 BY THOMAS FAN
- FIXFIXED AN ISSUE WITH PLOTTREE WHERE IT DISPLAYED ENTROPY CALCULATIONS EVEN FOR GINI CRITERION IN DECISIONTREECLASSIFIERS 13947 BY FRANK HOANG

19 VERSION 0212

24 MAY 2019

191 CHANGELOG

SKLEARNDECOMPOSITION

- FIXFIXED A BUG IN CROSSDECOMPOSITIONCCA IMPROVING NUMERICAL STABILITY WHEN Y IS CLOSE TO ZERO 13903 BY THOMAS FAN

SKLEARNMETRICS

- FIXFIXED A BUG IN METRICSPAIRWISEEUCLIDEANDISTANCES WHERE A PART OF THE DISTANCE MATRIX WAS LEFT UNINSTANTIATED FOR SUFFICIENTLY LARGE FLOAT32 DATASETS REGRESSION INTRODUCED IN 021 13910 BY JÉRÉMIE DU BOISBERANGER

28 CHAPTER 1 WELCOME TO SCIKITLEARN

SCIKITLEARN USER GUIDE RELEASE 0213  
SKLEARNPREPROCESSING  
•FIXFIXED A BUG IN PREPROCESSINGONEHOTENCODER WHERE THE NEW DROP PARAMETER WAS NOT REFLECTED  
INGETFEATURENAMES 13894 BY JAMES MYATT  
SKLEARNUTILSSPARSEFUNCS  
•FIXFIXED A BUG WHERE MINMAXAXIS WOULD FAIL ON 32BIT SYSTEMS FOR CERTAIN LARGE INPUTS THIS  
AFFECTSPREPROCESSINGMAXABSSCALER PREPROCESSINGNORMALIZE ANDPREPROCESSING  
LABELBINARIZER 13741 BY RODDY MACSWEEN  
110 VERSION 0211  
17 MAY 2019  
THIS IS A BUGFIX RELEASE TO PRIMARILY RESOLVE SOME PACKAGING ISSUES IN VERSION 0210 IT ALSO INCLUDES MINOR DOCU  
MENTATION IMPROVEMENTS AND SOME BUG FIXES  
1101 CHANGELOG  
SKLEARNINSPECTION  
•FIXFIXED A BUG IN INSPECTIONPARTIALDEPENDENCE TO ONLY CHECK CLASSIFIER AND NOT REGRESSOR FOR  
THE MULTICLASSMULTIOUTPUT CASE 14309 BY GUILLAUME LEMAITRE  
SKLEARNMETRICS  
•FIXFIXED A BUG IN METRICSPAIRWISEDISTANCES WHERE IT WOULD RAISE ATTRIBUTEERROR FOR  
BOOLEAN METRICS WHEN XHAD A BOOLEAN DTYPE AND Y NONE 13864 BY PARESH MATHUR  
•FIXFIXED TWO BUGS IN METRICSPAIRWISEDISTANCES WHENNJOBS 1 FIRST IT USED TO RETURN A  
DISTANCE MATRIX WITH SAME DTYPE AS INPUT EVEN FOR INTEGER DTYPE THEN THE DIAGONAL WAS NOT ZEROS FOR EUCLIDEAN  
METRIC WHEN YISX 13877 BY JÉRÉMIE DU BOISBERRANGER  
SKLEARNNEIGHBORS  
•FIXFIXED A BUG IN NEIGHBORSKERNELDENSITY WHICH COULD NOT BE RESTORED FROM A PICKLE IF  
SAMPLEWEIGHT HAD BEEN USED 13772 BY ADITYA VYAS  
111 VERSION 0210  
MAY 2019  
110 VERSION 0211 29

SCIKITLEARN USER GUIDE RELEASE 0213

1111 CHANGED MODELS

THE FOLLOWING ESTIMATORS AND FUNCTIONS WHEN FIT WITH THE SAME DATA AND PARAMETERS MAY PRODUCE DIFFERENT MODELS FROM THE PREVIOUS VERSION THIS OFTEN OCCURS DUE TO CHANGES IN THE MODELLING LOGIC BUG FIXES OR ENHANCEMENTS OR IN RANDOM SAMPLING PROCEDURES

- DISCRIMINANTANALYSISLINEARDISCRIMINANTANALYSIS FOR MULTICLASS CLASSIFICATION FIX
- DISCRIMINANTANALYSISLINEARDISCRIMINANTANALYSIS WITH ‘EIGEN’ SOLVER FIX
- LINEARMODELBAYESIANRIDGE FIX
- DECISION TREES AND DERIVED ENSEMBLES WHEN BOTH MAXDEPTH ANDMAXLEAFNODES ARE SET FIX
- LINEARMODELLOGISTICREGRESSION ANDLINEARMODELLOGISTICREGRESSIONCV WITH ‘SAGA’ SOLVER FIX
- ENSEMBLEGRADIENTBOOSTINGCLASSIFIER FIX
- SKLEARNFEATUREEXTRACTIONTEXTHASHINGVECTORIZER SKLEARN FEATUREEXTRACTIONTEXTTFIDFVECTORIZER AND SKLEARNFEATUREEXTRACTION TEXTCOUNTVECTORIZER FIX
- NEURALNETWORKMLPCLASSIFIER FIX
- SVMSCVDECISIONFUNCTION AND MULTICLASSONEVSONECLASSIFIER DECISIONFUNCTION FIX
- LINEARMODELSGDCLASSIFIER AND ANY DERIVED CLASSIFIERS FIX
- ANY MODEL USING THE LINEARMODELSAGSAGSOLVER FUNCTION WITH A 0SEED INCLUD INLINEARMODELLOGISTICREGRESSION LINEARMODELLOGISTICREGRESSIONCV LINEARMODELRIDGE ANDLINEARMODELRIDGECV WITH ‘SAG’ SOLVER FIX
- LINEARMODELRIDGECV WHEN USING GENERALIZED CROSSVALIDATION WITH SPARSE INPUTS FIX

DETAILS ARE LISTED IN THE CHANGELOG BELOW

WHILE WE ARE TRYING TO BETTER INFORM USERS BY PROVIDING THIS INFORMATION WE CANNOT ASSURE THAT THIS LIST IS COMPLETE

1112 KNOWN MAJOR BUGS

- THE DEFAULT MAXITER FORLINEARMODELLOGISTICREGRESSION IS TOO SMALL FOR MANY SOLVERS GIVEN THE DEFAULTTOL IN PARTICULAR WE ACCIDENTALLY CHANGED THE DEFAULT MAXITER FOR THE LIBLINEAR SOLVER FROM 1000 TO 100 ITERATIONS IN 3591 RELEASED IN VERSION 016 IN A FUTURE RELEASE WE HOPE TO CHOOSE BETTER DEFAULT MAXITER ANDTOL HEURISTICALLY DEPENDING ON THE SOLVER SEE 13317

1113 CHANGELOG

SUPPORT FOR PYTHON 34 AND BELOW HAS BEEN OFFICIALLY DROPPED

SKLEARNBASE

- API C HANGE THE R2 SCORE USED WHEN CALLING SCORE ON A REGRESSOR WILL USE MULTIOUTPUTUNIFORMAVERAGE FROM VERSION 023 TO KEEP CONSISTENT WITH METRICSR2SCORE THIS WILL INFLUENCE THE SCORE METHOD OF ALL THE MULTIOUTPUT REGRESSORS EXCEPT FOR MULTIOUTPUT MULTIOUTPUTREGRESSOR 13157 BY HANMIN QIN
- 30 CHAPTER 1 WELCOME TO SCIKITLEARN

SCIKITLEARN USER GUIDE RELEASE 0213

SKLEARNCALIBRATION

- ENHANCEMENT ADDED SUPPORT TO BIN THE DATA PASSED INTO CALIBRATIONCALIBRATIONCURVE BY QUAN TILES INSTEAD OF UNIFORMLY BETWEEN 0 AND 1 13086 BY SCOTT COLE
- ENHANCEMENT ALLOW NDIMENSIONAL ARRAYS AS INPUT FOR CALIBRATIONCALIBRATEDCLASSIFIERCV 13485 BY WILLIAM DE VAZELHES

SKLEARNCLUSTER

- MAJOR FEATURE A NEW CLUSTERING ALGORITHM CLUSTEROPTICS AN ALGORITM RELATED TO CLUSTER DBSCAN THAT HAS HYPERPARAMETERS EASIER TO SET AND THAT SCALES BETTER BY SHANE ADRIN JALALI ERICH SCHUBERT HANMIN QIN AND ASSIA BENBIHI
- FIXFIXED A BUG WHERE CLUSTERBIRCH COULD OCCASIONALLY RAISE AN ATTRIBUTEERROR 13651 BY JOEL NOTH MAN
- FIXFIXED A BUG IN CLUSTERKMEANS WHERE EMPTY CLUSTERS WEREN'T CORRECTLY RELOCATED WHEN USING SAMPLE WEIGHTS 13486 BY JÉRÉMIE DU BOISBERRANGER
- API C HANGE THENCOMPONENTS ATTRIBUTE IN CLUSTERAGGLOMERATIVECLUSTERING AND CLUSTERFEATUREAGGLOMERATION HAS BEEN RENAMED TO NCONNECTEDCOMPONENTS 13427 BY STEPHANE COUVREUR
- ENHANCEMENT CLUSTERAGGLOMERATIVECLUSTERING ANDCLUSTERFEATUREAGGLOMERATION NOW ACCEPT A DISTANCETHRESHOLD PARAMETER WHICH CAN BE USED TO FIND THE CLUSTERS INSTEAD OF NCLUSTERS 9069 BY VATHSALA ACHAR AND ADRIN JALALI

SKLEARNCOMPOSE

- API C HANGE COMPOSECOLUMNTRANSFORMER IS NO LONGER AN EXPERIMENTAL FEATURE 13835 BY HANMIN QIN

SKLEARNDATASETS

- FIXADDED SUPPORT FOR 64BIT GROUP IDS AND POINTERS IN SVMLIGHT FILES 10727 BY BRYAN K WOODS
- FIXDATASETSLOADSAMPLEIMAGES RETURNS IMAGES WITH A DETERMINISTIC ORDER 13250 BY THOMAS FAN

SKLEARNDECOMPOSITION

- ENHANCEMENT DECOMPOSITIONKERNELPCA NOW HAS DETERMINISTIC OUTPUT RESOLVED SIGN AMBIGUITY IN EIGENVALUE DECOMPOSITION OF THE KERNEL MATRIX 13241 BY AURÉLIEN BELLET
  - FIXFIXED A BUG IN DECOMPOSITIONKERNELPCA FITTRANSFORM NOW PRODUCES THE CORRECT OUTPUT THE SAME AS FITTRANSFORM IN CASE OF NONREMOVED ZERO EIGENVALUES REMOVEZEROEIGFALSE FITINVERSETRANSFORM WAS ALSO ACCELERATED BY USING THE SAME TRICK ASFITTRANSFORM TO COMPUTE THE TRANSFORM OF X 12143 BY SYLVAIN MARIÉ
  - FIXFIXED A BUG IN DECOMPOSITIONNMF WHEREINIT NNDSVD INIT NNDSVDA AND INIT NNDSVDAR ARE ALLOWED WHEN NCOMPONENTS NFEATURES INSTEAD OFNCOMPONENTS MINNSAMPLES NFEATURES 11650 BY HOSSEIN POURBOZORG AND ZIJIE ZJ POH
- 111 VERSION 0210 31

SCIKITLEARN USER GUIDE RELEASE 0213

- API C HANGE THE DEFAULT VALUE OF THE INIT ARGUMENT IN DECOMPOSITION  
NONNEGATIVEFACTORIZATION WILL CHANGE FROM RANDOM TONONE IN VERSION 023 TO MAKE IT  
CONSISTENT WITH DECOMPOSITIONNMF A FUTUREWARNING IS RAISED WHEN THE DEFAULT VALUE IS USED 12988  
BY ZIJIE ZJ POH
- SKLEARNDISCRIMINANTANALYSIS
- ENHANCEMENT DISCRIMINANTANALYSISLINEARDISCRIMINANTANALYSIS NOW PRESERVES  
FLOAT32 ANDFLOAT64 DTYPES 8769 AND 11000 BY THIBAUT SEJOURNE
- FIXACHANGEDBEHAVIOURWARNING IS NOW RAISED WHEN DISCRIMINANTANALYSIS  
LINEARDISCRIMINANTANALYSIS IS GIVEN AS PARAMETER NCOMPONENTS MINNFEATURES  
NCLASSES 1 ANDNCOMPONENTS IS CHANGED TO MINNFEATURES NCLASSES 1 IF SO  
PREVIOUSLY THE CHANGE WAS MADE BUT SILENTLY 11526 BY WILLIAM DE VAZELHES
- FIXFIXED A BUG IN DISCRIMINANTANALYSISLINEARDISCRIMINANTANALYSIS WHERE THE PRE  
DICTED PROBABILITIES WOULD BE INCORRECTLY COMPUTED IN THE MULTICLASS CASE 6848 BY AGAMEMNON KRASOULIS  
ANDGUILLAUME LEMAITRE
- FIXFIXED A BUG IN DISCRIMINANTANALYSISLINEARDISCRIMINANTANALYSIS WHERE THE PRE  
DICTED PROBABILITIES WOULD BE INCORRECTLY COMPUTED WITH EIGEN SOLVER 11727 BY AGAMEMNON KRASOULIS  
SKLEARNDUMMY
- FIXFIXED A BUG IN DUMMYDUMMYCLASSIFIER WHERE THEPREDICTPROBA METHOD WAS RETURNING INT32  
ARRAY INSTEAD OF FLOAT64 FOR THE STRATIFIED STRATEGY 13266 BY CHRISTOS ARIDAS
- FIXFIXED A BUG IN DUMMYDUMMYCLASSIFIER WHERE IT WAS THROWING A DIMENSION MISMATCH ERROR IN  
PREDICTION TIME IF A COLUMN VECTOR YWITHSHAPEN 1 WAS GIVEN AT FIT TIME 13545 BY NICK SORROS AND  
ADRIN JALALI
- SKLEARNENSEMBLE
- MAJOR FEATURE ADD TWO NEW IMPLEMENTATIONS OF GRADIENT BOOSTING TREES ENSEMBLE  
HISTGRADIENTBOOSTINGCLASSIFIER ANDENSEMBLEHISTGRADIENTBOOSTINGREGRESSOR  
THE IMPLEMENTATION OF THESE ESTIMATORS IS INSPIRED BY LIGHTGBM AND CAN BE ORDERS OF MAGNITUDE FASTER THAN  
ENSEMBLEGRADIENTBOOSTINGREGRESSOR ANDENSEMBLEGRADIENTBOOSTINGCLASSIFIER  
WHEN THE NUMBER OF SAMPLES IS LARGER THAN TENS OF THOUSANDS OF SAMPLES THE API OF THESE NEW ESTIMATORS  
IS SLIGHTLY DIFFERENT AND SOME OF THE FEATURES FROM ENSEMBLEGRADIENTBOOSTINGCLASSIFIER AND  
ENSEMBLEGRADIENTBOOSTINGREGRESSOR ARE NOT YET SUPPORTED  
THESE NEW ESTIMATORS ARE EXPERIMENTAL WHICH MEANS THAT THEIR RESULTS OR THEIR API MIGHT CHANGE WITHOUT ANY  
DEPRECATION CYCLE TO USE THEM YOU NEED TO EXPLICITLY IMPORT ENABLEHISTGRADIENTBOOSTING  
EXPLICITLY REQUIRE THIS EXPERIMENTAL FEATURE  
FROM SKLEARNEXPERIMENTAL IMPORT ENABLEHISTGRADIENTBOOSTING NOQA  
NOW YOU CAN IMPORT NORMALLY FROM SKLEARNENSEMBLE  
FROM SKLEARNENSEMBLE IMPORT HISTGRADIENTBOOSTINGCLASSIFIER  
12807 BY NICOLAS HUG
- FEATURE ADDENSEMBLEVOTINGREGRESSOR WHICH PROVIDES AN EQUIVALENT OF ENSEMBLE  
VOTINGCLASSIFIER FOR REGRESSION PROBLEMS 12513 BY RAMIL NUGMANOV AND MOHAMED ALI JAMAOU

SCIKITLEARN USER GUIDE RELEASE 0213

- EFFICIENCY MAKEENSEMBLEISOLATIONFOREST PREFER THREADS OVER PROCESSES WHEN RUNNING WITH NJOBS 1 AS THE UNDERLYING DECISION TREE FIT CALLS DO RELEASE THE GIL THIS CHANGES REDUCES MEMORY USAGE AND COMMUNICATION OVERHEAD 12543 BY ISAAC STORCH AND OLIVIER GRISEL
  - EFFICIENCY MAKEENSEMBLEISOLATIONFOREST MORE MEMORY EFFICIENT BY AVOIDING KEEPING IN MEMORY EACH TREE PREDICTION 13260 BY NICOLAS GOIX
  - EFFICIENCY ENSEMBLEISOLATIONFOREST NOW USES CHUNKS OF DATA AT PREDICTION STEP THUS CAPPING THE MEMORY USAGE 13283 BY NICOLAS GOIX
  - EFFICIENCY SKLEARNENSEMBLEGRADIENTBOOSTINGCLASSIFIER ANDSKLEARNENSEMBLE GRADIENTBOOSTINGREGRESSOR NOW KEEP THE INPUT YASFLOAT64 TO AVOID IT BEING COPIED INTERNALLY BY TREES 13524 BY ADRIN JALALI
  - ENHANCEMENT MINIMIZED THE VALIDATION OF X IN ENSEMBLEADABOOSTCLASSIFIER ANDENSEMBLE ADABOOSTREGRESSOR 13174 BY CHRISTOS ARIDAS
  - ENHANCEMENT ENSEMBLEISOLATIONFOREST NOW EXPOSES WARMSTART PARAMETER ALLOWING ITERATIVE ADDITION OF TREES TO AN ISOLATION FOREST 13496 BY PETER MARKO
  - FIXTHE VALUES OF FEATUREIMPORTANCES IN ALL RANDOM FOREST BASED MODELS IE ENSEMBLERANDOMFORESTCLASSIFIER ENSEMBLERANDOMFORESTREGRESSOR ENSEMBLEEXTRATREESCLASSIFIER ENSEMBLEEXTRATREESREGRESSOR ENSEMBLE RANDOMTREESEMBEDDING ENSEMBLEGRADIENTBOOSTINGCLASSIFIER ANDENSEMBLE GRADIENTBOOSTINGREGRESSOR NOW -SUM UP TO1
  - ALL THE SINGLE NODE TREES IN FEATURE IMPORTANCE CALCULATION ARE IGNORED
  - IN CASE ALL TREES HAVE ONLY ONE SINGLE NODE IE A ROOT NODE FEATURE IMPORTANCES WILL BE AN ARRAY OF ALL ZEROS 13636 AND 13620 BY ADRIN JALALI
  - FIXFIXED A BUG IN ENSEMBLEGRADIENTBOOSTINGCLASSIFIER ANDENSEMBLE GRADIENTBOOSTINGREGRESSOR WHICH DIDN'T SUPPORT SCIKITLEARN ESTIMATORS AS THE INITIAL ESTIMATOR ALSO ADDED SUPPORT OF INITIAL ESTIMATOR WHICH DOES NOT SUPPORT SAMPLE WEIGHTS 12436 BY JÉRÉMIE DU BOISBERRANGER AND 12983 BY NICOLAS HUG
  - FIXFIXED THE OUTPUT OF THE AVERAGE PATH LENGTH COMPUTED IN ENSEMBLEISOLATIONFOREST WHEN THE INPUT IS EITHER 0 1 OR 2 13251 BY ALBERT THOMAS AND JOSHUAKENNETHJONES
  - FIXFIXED A BUG IN ENSEMBLEGRADIENTBOOSTINGCLASSIFIER WHERE THE GRADIENTS WOULD BE INCORRECTLY COMPUTED IN MULTICLASS CLASSIFICATION PROBLEMS 12715 BY NICOLAS HUG
  - FIXFIXED A BUG IN ENSEMBLEGRADIENTBOOSTINGCLASSIFIER WHERE VALIDATION SETS FOR EARLY STOPPING WERE NOT SAMPLED WITH STRATIFICATION 13164 BY NICOLAS HUG
  - FIXFIXED A BUG IN ENSEMBLEGRADIENTBOOSTINGCLASSIFIER WHERE THE DEFAULT INITIAL PREDICTION OF A MULTICLASS CLASSIFIER WOULD PREDICT THE CLASSES PRIORS INSTEAD OF THE LOG OF THE PRIORS 12983 BY NICOLAS HUG
  - FIXFIXED A BUG IN ENSEMBLERANDOMFORESTCLASSIFIER WHERE THEPREDICT METHOD WOULD ERROR FOR MULTICLASS MULTIOUTPUT FORESTS MODELS IF ANY TARGETS WERE STRINGS 12834 BY ELIZABETH SANDER
  - FIXFIXED A BUG IN ENSEMBLEGRADIENTBOOSTINGLOSSFUNCTION ANDENSEMBLE GRADIENTBOOSTINGLEASTSQUARESERROR WHERE THE DEFAULT VALUE OF LEARNINGRATE IN UPDATETERMINALREGIONS IS NOT CONSISTENT WITH THE DOCUMENT AND THE CALLER FUNCTIONS NOTE HOWEVER THAT DIRECTLY USING THESE LOSS FUNCTIONS IS DEPRECATED 6463 BY MOVELIKERIVER
  - FIXENSEMBLEPARTIALDEPENDENCE AND CONSEQUENTLY THE NEW VERSION SKLEARNINSPECTION PARTIALDEPENDENCE NOW TAKES SAMPLE WEIGHTS INTO ACCOUNT FOR THE PARTIAL DEPENDENCE COMPUTATION WHEN THE GRADIENT BOOSTING MODEL HAS BEEN TRAINED WITH SAMPLE WEIGHTS 13193 BY SAMUEL O RONSIN
- 111 VERSION 0210 33

SCIKITLEARN USER GUIDE RELEASE 0213

- API C HANGE ENSEMBLEPARTIALDEPENDENCE ANDENSEMBLEPLOTPARTIALDEPENDENCE ARE NOW DEPRECATED IN FAVOR OF INSPECTIONPARTIALDEPENDENCE ANDINSPECTION PLOTPARTIALDEPENDENCE 12599 BY TREVOR STEPHENS AND NICOLAS HUG
- FIXENSEMBLEVOTINGCLASSIFIER ANDENSEMBLEVOTINGREGRESSOR WERE FAILING DURING FIT IN ONE OF THE ESTIMATORS WAS SET TO NONE ANDSAMPLEWEIGHT WAS NOTNONE 13779 BY GUILLAUME LEMAITRE
- API C HANGE ENSEMBLEVOTINGCLASSIFIER ANDENSEMBLEVOTINGREGRESSOR ACCEPTDROP TO DISABLE AN ESTIMATOR IN ADDITION TO NONE TO BE CONSISTENT WITH OTHER ESTIMATORS IE PIPELINE FEATUREUNION ANDCOMPOSECOLUMNTRANSFORMER 13780 BY GUILLAUME LEMAITRE
- SKLEARNEXTERNALS
- API C HANGE DEPRECATED EXTERNALSSIX SINCE WE HAVE DROPPED SUPPORT FOR PYTHON 27 12916 BY HAN MIN QIN
- SKLEARNFEATUREEXTRACTION
- FIXINPUTFILE ORINPUTFILENAME AND A CALLABLE IS GIVEN AS THE ANALYZER SKLEARN FEATUREEXTRACTIONTEXTHASHINGVECTORIZER SKLEARNFEATUREEXTRACTIONTEXT TFIDFVECTORIZER ANDSKLEARNFEATUREEXTRACTIONTEXTCOUNTVECTORIZER NOW READ THE DATA FROM THE FILES AND THEN PASS IT TO THE GIVEN ANALYZER INSTEAD OF PASSING THE FILE NAMES OR THE FILE OBJECTS TO THE ANALYZER 13641 BY ADRIN JALALI
- SKLEARNIMPUTE
- MAJOR FEATURE ADDEDIMPUTEITERATIVEIMPUTER WHICH IS A STRATEGY FOR IMPUTING MISSING VALUES BY MODELING EACH FEATURE WITH MISSING VALUES AS A FUNCTION OF OTHER FEATURES IN A ROUNDROBIN FASHION 8478 AND 12177 BY SERGEY FELDMAN AND BEN LAWSON
- THE API OF ITERATIVEIMPUTER IS EXPERIMENTAL AND SUBJECT TO CHANGE WITHOUT ANY DEPRECATION CYCLE TO USE THEM YOU NEED TO EXPLICITLY IMPORT ENABLEITERATIVEIMPUTER
- FROM SKLEARNEXPERIMENTAL IMPORT ENABLEITERATIVEIMPUTER NOQA
- NOW YOU CAN IMPORT NORMALLY FROM SKLEARNIMPUTE
- FROM SKLEARNIMPUTE IMPORT ITERATIVEIMPUTER
- FEATURE THEIMPUTESIMPLEIMPUTER ANDIMPUTEITERATIVEIMPUTER HAVE A NEW PARAMETER ADDINDICATOR WHICH SIMPLY STACKS A IMPUTEMISSINGINDICATOR TRANSFORM INTO THE OUTPUT OF THE IMPUTER'S TRANSFORM THAT ALLOWS A PREDICTIVE ESTIMATOR TO ACCOUNT FOR MISSINGNESS 12583 13601 BY DANYLO BAIBAK
- FIXINIMPUTEMISSINGINDICATOR AVOID IMPLICIT DENSIFICATION BY RAISING AN EXCEPTION IF INPUT IS SPARSE ADDMISSINGVALUES PROPERTY IS SET TO 0 13240 BY BARTOSZ TELENCZUK
- FIXFIXED TWO BUGS IN IMPUTEMISSINGINDICATOR FIRST WHEN XIS SPARSE ALL THE NONZERO NON MISSING VALUES USED TO BECOME EXPLICIT FALSE IN THE TRANSFORMED DATA THEN WHEN FEATURESMISSINGONLY ALL FEATURES USED TO BE KEPT IF THERE WERE NO MISSING VALUES AT ALL 13562 BY JÉRÉMIE DU BOISBERRANGER
- SKLEARNINSPECTION
- NEW SUBPACKAGE
- 34 CHAPTER 1 WELCOME TO SCIKITLEARN



SCIKITLEARN USER GUIDE RELEASE 0213

- FEATURE PARTIAL DEPENDENCE PLOTS INSPECTIONPLOTPARTIALDEPENDENCE ARE NOW SUPPORTED FOR ANY REGRESSOR OR CLASSIFIER PROVIDED THAT THEY HAVE A PREDICTPROBA METHOD 12599 BY TREVOR STEPHENS AND NICOLAS HUG
- SKLEARNISOTONIC
- FEATURE ALLOW DIFFERENT DTYPES SUCH AS FLOAT32 IN ISOTONICISOTONICREGRESSION 8769 BY VLAD NICULAE
- SKLEARNLINEARMODEL
- ENHANCEMENT LINEARMODELRIDGE NOW PRESERVES FLOAT32 ANDFLOAT64 DTYPES 8769 AND 11000 BY GUILLAUME LEMAITRE AND JOAN MASSICH
- FEATURE LINEARMODELLOGISTICREGRESSION AND LINEARMODEL LOGISTICREGRESSIONCV NOW SUPPORT ELASTICNET PENALTY WITH THE ‘SAGA’ SOLVER 11646 BY NICOLAS HUG
- FEATURE ADDEDLINEARMODELLARSPATHGRAM WHICH IS LINEARMODELLARSPATH IN THE SUFFICIENT STATS MODE ALLOWING USERS TO COMPUTE LINEARMODELLARSPATH WITHOUT PROVIDING XANDY 11699 BY KUAI YU
- EFFICIENCY LINEARMODELMAKEDATASET NOW PRESERVES FLOAT32 ANDFLOAT64 DTYPES REDUCING MEMORY CONSUMPTION IN STOCHASTIC GRADIENT SAG AND SAGA SOLVERS 8769 AND 11000 BY NELLE VAROQUAUX ARTHUR IMBERT GUILLAUME LEMAITRE AND JOAN MASSICH
- ENHANCEMENT LINEARMODELLOGISTICREGRESSION NOW SUPPORTS AN UNREGULARIZED OBJECTIVE WHEN PENALTYNONE IS PASSED THIS IS EQUIVALENT TO SETTING CNPINF WITH L2 REGULARIZATION NOT SUPPORTED BY THE LIBLINEAR SOLVER 12860 BY NICOLAS HUG
- ENHANCEMENT SPARSECG SOLVER INLINEARMODELRIDGE NOW SUPPORTS FITTING THE INTERCEPT IE FITINTERCEPTTRUE WHEN INPUTS ARE SPARSE 13336 BY BARTOSZ TELENCZUK
- ENHANCEMENT THE COORDINATE DESCENT SOLVER USED IN LASSO ELASTICNET ETC NOW ISSUES A CONVERGENCEWARNING WHEN IT COMPLETES WITHOUT MEETING THE DESIRED TOLERANBCE 11754 AND 13397 BY BRENT FAGAN AND ADRIN JALALI
- FIXFIXED A BUG IN LINEARMODELLOGISTICREGRESSION ANDLINEARMODEL LOGISTICREGRESSIONCV WITH ‘SAGA’ SOLVER WHERE THE WEIGHTS WOULD NOT BE CORRECTLY UPDATED IN SOME CASES 11646 BY TOM DUPRE LA TOUR
- FIXFIXED THE POSTERIOR MEAN POSTERIOR COVARIANCE AND RETURNED REGULARIZATION PARAMETERS IN LINEARMODELBAYESIANRIDGE THE POSTERIOR MEAN AND THE POSTERIOR COVARIANCE WERE NOT THE ONES COMPUTED WITH THE LAST UPDATE OF THE REGULARIZATION PARAMETERS AND THE RETURNED REGULARIZATION PARAMETERS WERE NOT THE FINAL ONES ALSO FIXED THE FORMULA OF THE LOG MARGINAL LIKELIHOOD USED TO COMPUTE THE SCORE WHEN COMPUTESCORETRUE 12174 BY ALBERT THOMAS
- FIXFIXED A BUG IN LINEARMODELLASSOLARSIC WHERE USER INPUT COPYXFALSE AT INSTANCE CREATION WOULD BE OVERRIDDEN BY DEFAULT PARAMETER VALUE COPYXTRUE INFIT 12972 BY LUCIO FERNANDEZ ARJONA
- FIXFIXED A BUG IN LINEARMODELLINEARREGRESSION THAT WAS NOT RETURNING THE SAME COEFFECIENTS AND INTERCEPTS WITH FITINTERCEPTTRUE IN SPARSE AND DENSE CASE 13279 BY ALEXANDRE GRAMFORT
- FIXFIXED A BUG IN LINEARMODELHUBERREGRESSOR THAT WAS BROKEN WHEN XWAS OF DTYPE BOOL 13328 BY ALEXANDRE GRAMFORT

SCIKITLEARN USER GUIDE RELEASE 0213

- FIXFIXED A PERFORMANCE ISSUE OF SAGA ANDSAG SOLVERS WHEN CALLED IN A JOBLIBPARALLEL SETTING WITH NJOBS 1 ANDBACKENDTHREADING CAUSING THEM TO PERFORM WORSE THAN IN THE SEQUENTIAL CASE 13389 BY PIERRE GLASER
  - FIXFIXED A BUG IN LINEARMODELSTOCHASTICGRADIENTBASESGDCLASSIFIER THAT WAS NOT DETERMINISTIC WHEN TRAINED IN A MULTICLASS SETTING ON SEVERAL THREADS 13422 BY CLÉMENT DOUMOIRO
  - FIX FIXED BUG IN LINEARMODELRIDGEREGRESSION LINEARMODELRIDGE ANDLINEARMODELRIDGECLASSIFIER THAT CAUSED UNHANDLED EXCEPTION FOR ARGUMENTS RETURNINTERCEPTTRUE ANDSOLVERAUTO DEFAULT OR ANY OTHER SOLVER DIFFERENT FROM SAG 13363 BY BARTOSZ TELENCZUK
  - FIXLINEARMODELRIDGEREGRESSION WILL NOW RAISE AN EXCEPTION IF RETURNINTERCEPTTRUE AND SOLVER IS DIFFERENT FROM SAG PREVIOUSLY ONLY WARNING WAS ISSUED 13363 BY BARTOSZ TELENCZUK
  - FIXLINEARMODELRIDGEREGRESSION WILL CHOOSE SPARSECG SOLVER FOR SPARSE INPUTS WHEN SOLVERAUTO ANDSAMPLEWEIGHT IS PROVIDED PREVIOUSLY CHOLESKY SOLVER WAS SELECTED 13363 BY BARTOSZ TELENCZUK
  - API C HANGE THE USE OFLINEARMODELLARSPATH WITHXNONE WHILE PASSING GRAM IS DEPRECATED IN VERSION 021 AND WILL BE REMOVED IN VERSION 023 USE LINEARMODELLARSPATHGRAM INSTEAD 11699 BY KUAI YU
  - API C HANGE LINEARMODELLOGISTICREGRESSIONPATH IS DEPRECATED IN VERSION 021 AND WILL BE REMOVED IN VERSION 023 12821 BY NICOLAS HUG
  - FIXLINEARMODELRIDGECV WITH GENERALIZED CROSSVALIDATION NOW CORRECTLY FITS AN INTERCEPT WHEN FITINTERCEPTTRUE AND THE DESIGN MATRIX IS SPARSE 13350 BY JÉRÔME DOCKÈS
- SKLEARNMANIFOLD
- EFFICIENCY MAKEMANIFOLDTSNETRUSTWORTHINESS USE AN INVERTED INDEX INSTEAD OF AN NPWHERE LOOKUP TO FIND THE RANK OF NEIGHBORS IN THE INPUT SPACE THIS IMPROVES EFFICIENCY IN PARTICULAR WHEN COMPUTED WITH LOTS OF NEIGHBORS ANDOR SMALL DATASETS 9907 BY WILLIAM DE VAZELHES
- SKLEARNMETRICS
- FEATURE ADDED THEMETRICSMAXERROR METRIC AND A CORRESPONDING MAXERROR SCORER FOR SINGLE OUTPUT REGRESSION 12232 BY KRISHNA SANGEETH
  - FEATURE ADDMETRICSMULTILABELCONFUSIONMATRIX WHICH CALCULATES A CONFUSION MATRIX WITH TRUE POSITIVE FALSE POSITIVE FALSE NEGATIVE AND TRUE NEGATIVE COUNTS FOR EACH CLASS THIS FACILITATES THE CALCULATION OF SETWISE METRICS SUCH AS RECALL SPECIFICITY FALL OUT AND MISS RATE 11179 BY SHANGWU YAO AND JOEL NOTHMAN
  - FEATURE METRICSJACCARDScore HAS BEEN ADDED TO CALCULATE THE JACCARD COEFFICIENT AS AN EVALUATION METRIC FOR BINARY MULTILABEL AND MULTICLASS TASKS WITH AN INTERFACE ANALOGOUS TO METRICSF1SCORE 13151 BY GAURAV DHINGRA AND JOEL NOTHMAN
  - FEATURE ADDEDMETRICSPAIRWISEHAVERSINEDISTANCES WHICH CAN BE ACCESSED WITH METRICPAIRWISE THROUGHMETRICSPAIRWISEDISTANCES AND ESTIMATORS HAVERSINE DISTANCE WAS PREVIOUSLY AVAILABLE FOR NEAREST NEIGHBORS CALCULATION 12568 BY WEI XUE EMMANUEL ARIAS AND JOEL NOTHMAN
  - EFFICIENCY FASTERMETRICSPAIRWISEDISTANCES WITH NJOBS 1 BY USING A THREADBASED BACKEND INSTEAD OF PROCESSBASED BACKENDS 8216 BY PIERRE GLASER AND ROMUALD MENUET
  - EFFICIENCY THE PAIRWISE MANHATTAN DISTANCES WITH SPARSE INPUT NOW USES THE BLAS SHIPPED WITH SCIPY INSTEAD OF THE BUNDLED BLAS 12732 BY JÉRÉMIE DU BOISBERRANGER
- 36 CHAPTER 1 WELCOME TO SCIKITLEARN

SCIKITLEARN USER GUIDE RELEASE 0213

- ENHANCEMENT USE LABEL ACCURACY INSTEAD OF MICROAVERAGE ONMETRICS CLASSIFICATIONREPORT TO AVOID CONFUSION MICROAVERAGE IS ONLY SHOWN FOR MULTILABEL OR MULTICLASS WITH A SUBSET OF CLASSES BECAUSE IT IS OTHERWISE IDENTICAL TO ACCURACY 12334 BY EMMANUEL ARIAS JOEL NOTHMAN AND ANDREAS MÜLLER
  - ENHANCEMENT ADDEDBETA PARAMETER TO METRICSHOMOGENEITYCOMPLETENESSVMEASURE AND METRICSVMEASURESCORE TO CONFIGURE THE TRADEOFF BETWEEN HOMOGENEITY AND COMPLETENESS 13607 BY STEPHANE COUVREUR AND AND IVAN SANCHEZ
  - FIXTHE METRIC METRICSR2SCORE IS DEGENERATE WITH A SINGLE SAMPLE AND NOW IT RETURNS NAN AND RAISES EXCEPTIONSUNDEFINEDMETRICWARNING 12855 BY PAWEL SENDYK
  - FIXFIXED A BUG WHERE METRICSBRIERSCORELOSS WILL SOMETIMES RETURN INCORRECT RESULT WHEN THERE'S ONLY ONE CLASS IN YTRUE 13628 BY HANMIN QIN
  - FIXFIXED A BUG IN METRICSLABELRANKINGAVERAGEPRECISIONSCORE WHERE SAMPLEWEIGHT WASN'T TAKEN INTO ACCOUNT FOR SAMPLES WITH DEGENERATE LABELS 13447 BY DAN ELLIS
  - API C HANGE THE PARAMETER LABELS INMETRICSHAMMINGLOSS IS DEPRECATED IN VERSION 021 AND WILL BE REMOVED IN VERSION 023 10580 BY RESHAMA SHAIKH AND SANDRA MITROVIC
  - FIXTHE FUNCTION METRICSPAIRWISEEUCLIDEANDISTANCES AND THEREFORE SEVERAL ESTIMATORS WITH METRICEUCLIDEAN SUFFERED FROM NUMERICAL PRECISION ISSUES WITH FLOAT32 FEATURES PRECISION HAS BEEN INCREASED AT THE COST OF A SMALL DROP OF PERFORMANCE 13554 BY CELELIBI AND JÉRÉMIE DU BOISBERRANGER
  - API C HANGE METRICSJACCARDSIMILARITYSCORE IS DEPRECATED IN FAVOUR OF THE MORE CONSISTENT METRICSJACCARDSCORE THE FORMER BEHAVIOR FOR BINARY AND MULTICLASS TARGETS IS BROKEN 13151 BY JOEL NOTHMAN
  - SKLEARNMIXTURE
  - FIXFIXED A BUG IN MIXTUREBASEMIXTURE AND THEREFORE ON ESTIMATORS BASED ON IT IE MIXTURE GAUSSIANMIXTURE ANDMIXTUREBAYESIANGAUSSIANMIXTURE WHEREFITPREDICT ANDFIT PREDICT WERE NOT EQUIVALENT 13142 BY JÉRÉMIE DU BOISBERRANGER
  - SKLEARNMODELSELECTION
  - FEATURE CLASSESGRIDSEARCHCV ANDRANDOMIZEDSEARCHCV NOW ALLOW FOR REFITCALLABLE TO ADD FLEX IBIITY IN IDENTIFYING THE BEST ESTIMATOR SEE BALANCE MODEL COMPLEXITY AND CROSSVALIDATED SCORE 11354 BY WENHAO ZHANG JOEL NOTHMAN AND ADRIN JALALI
  - ENHANCEMENT CLASSESGRIDSEARCHCV RANDOMIZEDSEARCHCV AND METHODS CROSSVALSCORE CROSSVALPREDICT CROSSVALIDATE NOW PRINT TRAIN SCORES WHEN RETURNTRAINSCORES IS TRUE AND VERBOSE 2 FORLEARNINGCURVE ANDVALIDATIONCURVE ONLY THE LATTER IS REQUIRED 12613 AND 12669 BY MARC TORRELLAS
  - ENHANCEMENT SOME CV SPLITTER CLASSES AND MODELSELECTIONTRAINTESTSPLIT NOW RAISE VALUEERROR WHEN THE RESULTING TRAINING SET IS EMPTY 12861 BY NICOLAS HUG
  - FIXFIXED A BUG WHERE MODELSELECTIONSTRATIFIEDKFOLD SHUFFLES EACH CLASS'S SAMPLES WITH THE SAMERANDOMSTATE MAKINGSHUFFLETRUE INEFFECTIVE 13124 BY HANMIN QIN
  - FIXADDED ABILITY FOR MODELSELECTIONCROSSVALPREDICT TO HANDLE MULTILABEL AND MULTIOUTPUTMULTICLASS TARGETS WITH PREDICTPROBA TYPE METHODS 8773 BY STEPHEN HOOVER
  - FIXFIXED AN ISSUE IN CROSSVALPREDICT WHEREMETHODPREDICTPROBA RETURNED ALWAYS 00 WHEN ONE OF THE CLASSES WAS EXCLUDED IN A CROSSVALIDATION FOLD 13366 BY GUILLAUME FOURNIER
- 111 VERSION 0210 37

SCIKITLEARN USER GUIDE RELEASE 0213

SKLEARNMULTICLASS

- FIXFIXED AN ISSUE IN MULTICLASSONEVSONECLASSIFIERDECISIONFUNCTION WHERE THE DECISIONFUNCTION VALUE OF A GIVEN SAMPLE WAS DIFFERENT DEPENDING ON WHETHER THE DECISIONFUNCTION WAS EVALUATED ON THE SAMPLE ALONE OR ON A BATCH CONTAINING THIS SAME SAMPLE DUE TO THE SCALING USED IN DECISIONFUNCTION 10440 BY JONATHAN OHAYON

SKLEARNMULTIOUTPUT

- FIXFIXED A BUG IN MULTIOUTPUTMULTIOUTPUTCLASSIFIER WHERE THE PREDICTPROBA METHOD INCORRECTLY CHECKED FOR PREDICTPROBA ATTRIBUTE IN THE ESTIMATOR OBJECT 12222 BY REBEKAH KIM

SKLEARNNEIGHBORS

- MAJOR FEATURE ADDEDNEIGHBORSNEIGHBORHOODCOMPONENTSANALYSIS FOR METRIC LEARNING WHICH IMPLEMENTS THE NEIGHBORHOOD COMPONENTS ANALYSIS ALGORITHM 10058 BY WILLIAM DE VAZELHES AND JOHN CHIOTELLIS

- API C HANGE METHODS IN NEIGHBORSNEARESTNEIGHBORS KNEIGHBORS RADIUSNEIGHBORS KNEIGHBORSGRAPH RADIUSNEIGHBORSGRAPH NOW RAISE NOTFITTEDERROR RATHER THAN ATTRIBUTEERROR WHEN CALLED BEFORE FIT 12279 BY KRISHNA SANGEETH

SKLEARNNEURALNETWORK

- FIXFIXED A BUG IN NEURALNETWORKMLPCLASSIFIER ANDNEURALNETWORKMLPREGRESSOR WHERE THE OPTION SHUFFLEFALSE WAS BEING IGNORED 12582 BY SAM WATERBURY

- FIXFIXED A BUG IN NEURALNETWORKMLPCLASSIFIER WHERE VALIDATION SETS FOR EARLY STOPPING WERE NOT SAMPLED WITH STRATIFICATION IN THE MULTILABEL CASE HOWEVER SPLITS ARE STILL NOT STRATIFIED 13164 BY NICOLAS HUG

SKLEARNPIPELINE

- FEATURE PIPELINEPIPELINE CAN NOW USE INDEXING NOTATION EG MYPIPELINE01 TO EXTRACT A SUBSEQUENCE OF STEPS AS ANOTHER PIPELINE INSTANCE A PIPELINE CAN ALSO BE INDEXED DIRECTLY TO EXTRACT A PARTICULAR STEP EGMYPipelineSVC RATHER THAN ACCESSING NAMEDSTEPS 2568 BY JOEL NOTHMAN

- FEATURE ADDED OPTIONAL PARAMETER VERBOSE INPIPELINEPIPELINE COMPOSE COLUMNTRANSFORMER ANDPIPELINEFEATUREUNION AND CORRESPONDING MAKE HELPERS FOR SHOWING PROGRESS AND TIMING OF EACH STEP 11364 BY BAZE PETRUSHEV KARAN DESAI JOEL NOTHMAN AND THOMAS FAN

- ENHANCEMENT PIPELINEPIPELINE NOW SUPPORTS USING PASSTHROUGH AS A TRANSFORMER WITH THE SAME EFFECT AS NONE 11144 BY THOMAS FAN

- ENHANCEMENT PIPELINEPIPELINE IMPLEMENTS LEN AND THEREFORE LENPIPELINE RETURNS THE NUMBER OF STEPS IN THE PIPELINE 13439 BY LAKSHYA KD

SKLEARNPREPROCESSING

- FEATURE PREPROCESSINGONEHOTENCODER NOW SUPPORTS DROPPING ONE FEATURE PER CATEGORY WITH A NEW DROP PARAMETER 12908 BY DREW JOHNSTON

- EFFICIENCY PREPROCESSINGONEHOTENCODER ANDPREPROCESSINGORDINALENCODER NOW HANDLE PANDAS DATAFRAMES MORE EFFICIENTLY 13253 BY MAIKIA

SCIKITLEARN USER GUIDE RELEASE 0213

- EFFICIENCY MAKEPREPROCESSINGMULTILABELBINARIZER CACHE CLASS MAPPINGS INSTEAD OF CALCULATING IT EVERY TIME ON THE FLY 12116 BY EKATERINA KRIVICH AND JOEL NOTHMAN
- EFFICIENCY PREPROCESSINGPOLYNOMIALFEATURES NOW SUPPORTS COMPRESSED SPARSE ROW CSR MATRICES AS INPUT FOR DEGREES 2 AND 3 THIS IS TYPICALLY MUCH FASTER THAN THE DENSE CASE AS IT SCALES WITH MATRIX DENSITY AND EXPANSION DEGREE ON THE ORDER OF DENSITYDEGREE AND IS MUCH MUCH FASTER THAN THE COMPRESSED SPARSE COLUMN CSC CASE 12197 BY ANDREW NYSTROM
- EFFICIENCY SPEED IMPROVEMENT IN PREPROCESSINGPOLYNOMIALFEATURES IN THE DENSE CASE ALSO ADDED A NEW PARAMETER ORDER WHICH CONTROLS OUTPUT ORDER FOR FURTHER SPEED PERFORMANCES 12251 BY TOM DUPRE LA TOUR
- FIXFIXED THE CALCULATION OVERFLOW WHEN USING A FLOAT16 DTYPE WITH PREPROCESSINGSTANDARDSCALER 13007 BY RAFFAELLO BALUYOT
- FIXFIXED A BUG IN PREPROCESSINGQUANTILETRANSFORMER ANDPREPROCESSINGQUANTILETRANSFORM TO FORCE NQUANTILES TO BE AT MOST EQUAL TO NSAMPLES VALUES OF NQUANTILES LARGER THAN NSAMPLES WERE EITHER USELESS OR RESULTING IN A WRONG APPROXIMATION OF THE CUMULATIVE DISTRIBUTION FUNCTION ESTIMATOR 13333 BY ALBERT THOMAS
- API C HANGE THE DEFAULT VALUE OF COPY INPREPROCESSINGQUANTILETRANSFORM WILL CHANGE FROM FALSE TO TRUE IN 023 IN ORDER TO MAKE IT MORE CONSISTENT WITH THE DEFAULT COPY VALUES OF OTHER FUNCTIONS IN PREPROCESSING AND PREVENT UNEXPECTED SIDE EFFECTS BY MODIFYING THE VALUE OF XINPLACE 13459 BY HUNTER MCGUSHION
- SKLEARN SVM
- FIXFIXED AN ISSUE IN SVM SVCDECISIONFUNCTION WHENDECISIONFUNCTIONSHAPEOVR THE DECISIONFUNCTION VALUE OF A GIVEN SAMPLE WAS DIFFERENT DEPENDING ON WHETHER THE DECISIONFUNCTION WAS EVALUATED ON THE SAMPLE ALONE OR ON A BATCH CONTAINING THIS SAME SAMPLE DUE TO THE SCALING USED IN DECISIONFUNCTION 10440 BY JONATHAN OHAYON
- SKLEARN TREE
- FEATURE DECISION TREES CAN NOW BE PLOTTED WITH MATPLOTLIB USING TREEPLOT TREE WITHOUT RELYING ON THE DOT LIBRARY REMOVING A HARDTOINSTALL DEPENDENCY 8508 BY ANDREAS MÜLLER
- FEATURE DECISION TREES CAN NOW BE EXPORTED IN A HUMAN READABLE TEXTUAL FORMAT USING TREEEXPORTTEXT 6261 BY GIUSEPPE VETTIGLI
- FEATURE GETNLEAVES ANDGETDEPTH HAVE BEEN ADDED TO TREEBASEDECISIONTREE AND CONSEQUENTLY ALL ESTIMATORS BASED ON IT INCLUDING TREEDECISIONTREECLASSIFIER TREEDECISIONTREEREgressor TREEEXTRATREECLASSIFIER ANDTREEEXTRATREEREgressor 12300 BY ADRIN JALALI
- FIXTREES AND FORESTS DID NOT PREVIOUSLY PREDICT MULTIOUTPUT CLASSIFICATION TARGETS WITH STRING LABELS DESPITE ACCEPTING THEM IN FIT 11458 BY MITAR MILUTINOVIC
- FIXFIXED AN ISSUE WITH TREEBASEDECISIONTREE AND CONSEQUENTLY ALL ESTIMATORS BASED ON IT INCLUDING TREEDECISIONTREECLASSIFIER TREEDECISIONTREEREgressor TREEEXTRATREECLASSIFIER ANDTREEEXTRATREEREgressor WHERE THEY USED TO EXCEED THE GIVEN MAXDEPTH BY 1 WHILE EXPANDING THE TREE IF MAXLEAFNODES ANDMAXDEPTH WERE BOTH SPECIFIED BY THE USER PLEASE NOTE THAT THIS ALSO AFFECTS ALL ENSEMBLE METHODS USING DECISION TREES 12344 BY ADRIN JALALI

SCIKITLEARN USER GUIDE RELEASE 0213

SKLEARNUTILS

- FEATURE UTILSRESAMPLE NOW ACCEPTS A STRATIFY PARAMETER FOR SAMPLING ACCORDING TO CLASS DISTRIBUTIONS 13549 BY NICOLAS HUG
- API CHANGE DEPRECATED WARNONDDTYPE PARAMETER FROM UTILSCHECKARRAY ANDUTILS CHECKXY ADDED EXPLICIT WARNING FOR DTYPE CONVERSION IN CHECKPAIRWISEARRAYS IF THEMETRIC BEING PASSED IS A PAIRWISE BOOLEAN METRIC 13382 BY PRATHMESH SAVALE
- MULTIPLE MODULES
- MAJOR FEATURE THEREPR METHOD OF ALL ESTIMATORS USED WHEN CALLING PRINTESTIMATOR HAS BEEN ENTIRELY REWRITTEN BUILDING ON PYTHON’S PRETTY PRINTING STANDARD LIBRARY ALL PARAMETERS ARE PRINTED BY DEFAULT BUT THIS CAN BE ALTERED WITH THE PRINTCHANGEDONLY OPTION INSKLEARNSETCONFIG 11705 BY NICOLAS HUG

- MAJOR FEATURE ADD ESTIMATORS TAGS THESE ARE ANNOTATIONS OF ESTIMATORS THAT ALLOW PROGRAMMATIC INSPECTION OF THEIR CAPABILITIES SUCH AS SPARSE MATRIX SUPPORT SUPPORTED OUTPUT TYPES AND SUPPORTED METHODS ESTIMATOR TAGS ALSO DETERMINE THE TESTS THAT ARE RUN ON AN ESTIMATOR WHEN CHECKESTIMATOR IS CALLED READ MORE IN THEUSER GUIDE 8022 BY ANDREAS MÜLLER
- EFFICIENCY MEMORY COPIES ARE AVOIDED WHEN CASTING ARRAYS TO A DIFFERENT DTYPE IN MULTIPLE ESTIMATORS 11973 BY ROMAN YURCHAK
- FIXFIXED A BUG IN THE IMPLEMENTATION OF THE OURRANDR HELPER FUNCTION THAT WAS NOT BEHAVING CONSISTENTLY ACROSS PLATFORMS 13422 BY MADHURA PARIKH AND CLÉMENT DOUMOIRO

MISCELLANEOUS

- ENHANCEMENT JOBLIB IS NO LONGER VENDORED IN SCIKITLEARN AND BECOMES A DEPENDENCY MINIMAL SUPPORTED VERSION IS JOBLIB 011 HOWEVER USING VERSION 013 IS STRONGLY RECOMMENDED 13531 BY ROMAN YURCHAK
- 1114 CHANGES TO ESTIMATOR CHECKS
- THESE CHANGES MOSTLY AFFECT LIBRARY DEVELOPERS
- ADDCHECKFITIDEMPOTENT TOCHECKESTIMATOR WHICH CHECKS THAT WHEN FITIS CALLED TWICE WITH THE SAME DATA THE OUPUT OF PREDICT PREDICTPROBA TRANSFORM AND DECISIONFUNCTION DOES NOT CHANGE 12328 BY NICOLAS HUG
- MANY CHECKS CAN NOW BE DISABLED OR CONFIGURED WITH ESTIMATOR TAGS 8022 BY ANDREAS MÜLLER

1115 CODE AND DOCUMENTATION CONTRIBUTORS

THANKS TO EVERYONE WHO HAS CONTRIBUTED TO THE MAINTENANCE AND IMPROVEMENT OF THE PROJECT SINCE VERSION 020 INCLUDING

ADANHAWTH ADITYA VYAS ADRIN JALALI AGAMEMNON KRASOULIS ALBERT THOMAS ALBERTO TORRES ALEXANDRE GRAMFORT AMOURAV ANDREA NAVARRETE ANDREAS MUELLER ANDREW NYSTROM ASSIABEN AURÉLIEN BELLET BARTOSZ MICHAŁOWSKI BARTOSZ TELENCZUK BAUKS BENJASTUDIO BERTRANDHAUT BHARAT RAGHUNATHAN BRENTFAGAN BRYAN WOODS CAT CHENAL CHEUK TING HO CHRIS CHOE CHRISTOS ARIDAS CLÉMENT DOUMOIRO COLE SMITH CONNOSSOR COREY LEVINSON DAN ELLIS DAN STINE DANYLO BAIBAK DATENKIEKER DENIS KATAEV DIDI BARZEV DILLON GARDNER DMITRY MOTTL DMITRY VUKOLOV DOUGAL J SUTHERLAND DOWON DREWMJOHNSTON DROR ATARIAH EDWARD J BROWN EKATERINA KRIVICH ELIZA BETH SANDER EMMANUEL ARIAS ERIC CHANG ERIC LARSON ERICH SCHUBERT ESVHD FALAK FEDA CURIC FEDERICO CASELLI

40 CHAPTER 1 WELCOME TO SCIKITLEARN

SCIKITLEARN USER GUIDE RELEASE 0213

FRANK HOANG FIBINSE XAVIER’ FINN O’SHEA GABRIEL MARZINOTTO GABRIEL VACALIUC GABRIELE CALVO GAEL VAROQUAUX GAURAV AHLAWAT GIUSEPPE VETTIGLI GREG GANDENBERGER GUILLAUME FOURNIER GUILLAUME LEMAITRE GUSTAVO DE MARI PEREIRA HANMIN QIN HAROLD FOX HHULUQI HUNTER MCGUSHION IAN SANDERS JACKLANGERMAN JACOPO NOTARSTEFANO JAKIRKHAM JAMES BOURBEAU JAN KOCH JAN S JANVANRIJN JARROD MILLMAN JDETHURENS JEREMIEDBB JF JOAAK JOAN MASSICH JOEL NOTHMAN JONATHAN OHAYON JORIS VAN DEN BOSSCHE JOSEPHSALMON JÉRÉMIE MÉHAULT KATRIN LEINWEBER KEN KMS15 KOEN KOSSORI ARUKU KRISHNA SANGEETH KUA I YU KULBEAR KUSHAL CHAUHAN KYLE JACKSON LAKSHYA KD LEANDRO HERMIDA LEE YI JIE JOEL LILY XIONG LISA SARAH THOMAS LOIC ESTEVE LOUIB LUKFA MAIKIA MAILLIAM MANIMARAN MANUEL LÓPEZIBÁÑEZ MARC TORRELLAS MARCO GAIDO MARCO GORELLI MARCOGORELLI MARINELM MARK HANNEL MARTIN GUBRI MASSTRAN MATHURINM MATTHEW ROESCHKE MAX COPELAND MELS Y T MFERRARI3 MICKAÉL SCHOENTGEN MING LI MITAR MOHAMMAD AFTAB MOHAMMED ABDELAAL MOHAMMED IBRAHEEM MUHAMMAD HASSAAN RAFIQUE MWESTT NAOYA IJIMA NICHOLAS SMITH NICOLAS GOIX NICOLAS HUG NIKOLAY SHEBANOV OLEKSANDR PAVLYK OLIVER RAUSCH OLIVIER GRISEL ORESTIS OSMAN OWEN FLANAGAN PAUL PACZUSKI PAVEL SORIANO PAVLOS KALLIS PAWEŁ SENDYK PEAY PETER PETER COCK PETER HAUSAMANN PETER MARKO PIERRE GLASER PIERRE TALLOTTE PIM DE HAAN PIOTR SZYMA ŃSKI PRABAKARAN KUMARESSHAN PRADEEP REDDY RAAMANA PRATHMESH SAVALE PULKIT MALOO QUENTIN BATISTA RADOSTIN STOYANOV RAF BALUYOT RAJDEEP DUA RAMIL NUGMANOV RAÚL GARCÍA CALVO REBEKAH KIM RESHAMA SHAIKH ROHAN LEKHWANI ROHAN SINGH ROHAN VARMA ROHIT KAPOOR ROMAN FELDBAUER ROMAN YURCHAK ROMUALD M ROOPAM SHARMA RYAN RÜDIGER BUSCHE SAM WATERBURY SAMUEL O RONSIN SANDROCASAGRANDE SCOTT COLE SCOTT LOWE SEBASTIAN RASCHKA SHANGWU YAO SHIVAM KOTWALIA SHIYU DUAN SMARIE SRIHARSHA HATWAR STEPHEN HOOVER STEPHEN TIERNEY STÉPHANE COUVREUR SURGAN12 SYLVAINLAN TAKINGITCASUAL TASHAY GREEN THIBSEJ THOMAS FAN THOMAS J FAN THOMAS MOREAU TOM DUPRÉ LA TOUR TOMMY TULIO CASAGRANDE UMAR FAROUK UMAR UTKARSH UPADHYAY VINAYAK MEHTA VISHAAL KAPOOR VIVEK KUMAR VLAD NICULAE VQEAN3 WENHAO ZHANG WILLIAM DE VAZELHES XHAN XING HAN LU XINYULIU12 YAROSLAV HALCHENKO ZACH GRIFFITH ZACH MILLER ZAYD HAMMOUDEH ZHUYI XUE ZIJIE ZJ POH

112 VERSION 0204

JULY 30 2019

THIS IS A BUGFIX RELEASE WITH SOME BUG FIXES APPLIED TO VERSION 0203

1121 CHANGELOG

THE BUNDLED VERSION OF JOBLIB WAS UPGRADED FROM 0130 TO 0132

SKLEARNCLUSTER

•FIXFIXED A BUG IN CLUSTERKMEANS WHERE KMEANS INITIALISATION COULD RARELY RESULT IN AN INDEXERROR

11756 BY JOEL NOTHMAN

SKLEARNCOMPOSE

•FIXFIXED AN ISSUE IN COMPOSECOLUMNSTRANSFORMER WHERE USING DATAFRAMES WHOSE COLUMN ORDER DIFFERS BETWEEN FUNC FIT AND FUNCTRANSFORM COULD LEAD TO SILENTLY PASSING INCORRECT COLUMNS TO THE REMAINDER TRANSFORMER 14237 BY ANDREAS SCHUDERER

SKLEARNMODELSELECTION

•FIXFIXED A BUG WHERE MODELSELECTIONSTRATIFIEDKFOLD SHUFFLES EACH CLASS’S SAMPLES WITH THE SAMERANDOMSTATE MAKINGSHUFFLETRUE INEFFECTIVE 13124 BY HANMIN QIN

112 VERSION 0204 41

SCIKITLEARN USER GUIDE RELEASE 0213

SKLEARNNEIGHBORS

- FIXFIXED A BUG IN NEIGHBORSKERNELDENSITY WHICH COULD NOT BE RESTORED FROM A PICKLE IF SAMPLEWEIGHT HAD BEEN USED 13772 BY ADITYA VYAS

113 VERSION 0203

MARCH 1 2019

THIS IS A BUGFIX RELEASE WITH SOME MINOR DOCUMENTATION IMPROVEMENTS AND ENHANCEMENTS TO FEATURES RELEASED IN 0200

1131 CHANGELOG

SKLEARNCLUSTER

- FIXFIXED A BUG IN CLUSTERKMEANS WHERE COMPUTATION WAS SINGLE THREADED WHEN NJOBS 1 OR NJOBS 1 12949 BY PRABAKARAN KUMARESSHAN

SKLEARNCOMPOSE

- FIXFIXED A BUG IN COMPOSECOLUMNTRANSFORMER TO HANDLE NEGATIVE INDEXES IN THE COLUMNS LIST OF THE TRANSFORMERS 12946 BY PIERRE TALLOTTE

SKLEARNCOVARIANCE

- FIXFIXED A REGRESSION IN COVARIANCEGRAPHICALASSO SO THAT THE CASE NFEATURES2 IS HANDLED CORRECTLY 13276 BY AURÉLIEN BELLET

SKLEARNDECOMPOSITION

- FIXFIXED A BUG IN DECOMPOSITIONSPARSEENCODE WHERE COMPUTATION WAS SINGLE THREADED WHEN NJOBS 1 ORNJOBS 1 13005 BY PRABAKARAN KUMARESSHAN

SKLEARNDATASETS

- EFFICIENCY SKLEARNDATASETSFETCHOPENML NOW LOADS DATA BY STREAMING AVOIDING HIGH MEMORY USAGE 13312 BY JORIS VAN DEN BOSSCHE

SKLEARNFEATUREEXTRACTION

- FIXFIXED A BUG IN FEATUREEXTRACTIONTEXTCOUNTVECTORIZER WHICH WOULD RESULT IN THE SPARSE FEATURE MATRIX HAVING CONFLICTING INDPTR ANDINDICES PRECISIONS UNDER VERY LARGE VOCABULARIES 11295 BY GABRIEL VACALIUC

42 CHAPTER 1 WELCOME TO SCIKITLEARN



SCIKITLEARN USER GUIDE RELEASE 0213

SKLEARNIMPUTE

- FIXADD SUPPORT FOR NONNUMERIC DATA IN SKLEARNIMPUTEMISSINGINDICATOR WHICH WAS NOT SUPPORTED WHILE SKLEARNIMPUTESIMPLEIMPUTER WAS SUPPORTING THIS FOR SOME IMPUTATION STRATEGIES 13046 BY GUILLAUME LEMAITRE

SKLEARNLINEARMODEL

- FIXFIXED A BUG IN LINEARMODELMULTITASKELASTICNET ANDLINEARMODEL MULTITASKLASSO WHICH WERE BREAKING WHEN WARMSTART TRUE 12360 BY AAKANKSHA JOSHI

SKLEARNPREPROCESSING

- FIXFIXED A BUG IN PREPROCESSINGKBINSDISCRETIZER WHERESTRATEGYKMEANS FAILS WITH AN ERROR DURING TRANSFORMATION DUE TO UNSORTED BIN EDGES 13134 BY SANDRO CASAGRANDE
- FIXFIXED A BUG IN PREPROCESSINGONEHOTENCODER WHERE THE DEPRECATION OF CATEGORICALFEATURES WAS HANDLED INCORRECTLY IN COMBINATION WITH HANDLEUNKNOWNIGNORE 12881 BY JORIS VAN DEN BOSSCHE
- FIXBINS WHOSE WIDTH ARE TOO SMALL IE 1E8 ARE REMOVED WITH A WARNING IN PREPROCESSING KBINSDISCRETIZER 13165 BY HANMIN QIN

SKLEARN SVM

- FIXFIXED A BUG IN SVM SVC SVMNU SVC SVM SVR SVMNU SVR AND SVM ONE CLASS SVM WHERE THE SCALE OPTION OF PARAMETER GAMMA IS ERRONEOUSLY DEFINED AS 1 / NFEATURES X STD IT'S NOW DEFINED AS 1 / NFEATURES X VAR 13221 BY HANMIN QIN

1132 CODE AND DOCUMENTATION CONTRIBUTORS

WITH THANKS TO

ADRIN JALALI AGAMEMNON KRASOULIS ALBERT THOMAS ANDREAS MUELLER AURÉLIEN BELLET BERTRANDHAUT BHARAT RAGHU NATHAN DOWON EMMANUEL ARIAS FIBINSE XAVIER FINN O'SHEA GABRIEL VACALIUC GAELE VAROQUAUX GUILLAUME LEMAITRE HANMIN QIN JOAQUIN JOEL NOTHMAN JORIS VAN DEN BOSSCHE JÉRÉMIE MÉHAULT KMS15 KOSSORI ARUKU LAK SHYA KD MAIKIA MANUEL LÓPEZIBÁÑEZ MARCO GORELLI MARCO GORELLI MFERRARI3 MICKAËL SCHOENTGEN NICOLAS HUG PAVLOS KALLIS PIERRE GLASER PIERRE TALLOTTE PRABAKARAN KUMARESSHAN RESHAMA SHAIKH ROHIT KAPOOR ROMAN YURCHAK SANDRO CASAGRANDE TASHAY GREEN THOMAS FAN VISHAAL KAPOOR ZHUYI XUE ZIJIE ZJ POH

114 VERSION 0202

DECEMBER 20 2018

THIS IS A BUGFIX RELEASE WITH SOME MINOR DOCUMENTATION IMPROVEMENTS AND ENHANCEMENTS TO FEATURES RELEASED IN 0200

114 VERSION 0202 43

SCIKITLEARN USER GUIDE RELEASE 0213

1141 CHANGED MODELS

THE FOLLOWING ESTIMATORS AND FUNCTIONS WHEN FIT WITH THE SAME DATA AND PARAMETERS MAY PRODUCE DIFFERENT MODELS FROM THE PREVIOUS VERSION THIS OFTEN OCCURS DUE TO CHANGES IN THE MODELLING LOGIC BUG FIXES OR ENHANCEMENTS OR IN RANDOM SAMPLING PROCEDURES

- SKLEARNNEIGHBORS WHENMETRICJACCARD BUG FIX
- USE OFSEUCLIDEAN ORMAHALANOBIS METRICS IN SOME CASES BUG FIX

1142 CHANGELOG

SKLEARNCOMPOSE

- FIXFIXED AN ISSUE IN COMPOSEMAKECOLUMNSTRANSFORMER WHICH RAISES UNEXPECTED ERROR WHEN COLUMNS IS PANDAS INDEX OR PANDAS SERIES 12704 BY HANMIN QIN

SKLEARNMETRICS

- FIX FIXED A BUG IN METRICSPAIRWISEDISTANCES ANDMETRICS

PAIRWISEDISTANCESCHUNKED WHERE PARAMETERS VOFSEUCLIDEAN ANDVIOFMAHALANOBIS METRICS WERE COMPUTED AFTER THE DATA WAS SPLIT INTO CHUNKS INSTEAD OF BEING PRECOMPUTED ON WHOLE DATA 12701 BY JEREMIE DU BOISBERRANGER

SKLEARNNEIGHBORS

- FIXFIXEDSKLEARNNEIGHBORSDISTANCEMETRIC JACCARD DISTANCE FUNCTION TO RETURN 0 WHEN TWO ALL ZERO VECTORS ARE COMPARED 12685 BY THOMAS FAN

SKLEARNUTILS

- FIXCALLINGUTILSCHECKARRAY ONPANDASSERIES WITH CATEGORICAL DATA WHICH RAISED AN ERROR IN 0200 NOW RETURNS THE EXPECTED OUTPUT AGAIN 12699 BY JORIS VAN DEN BOSSCHE

1143 CODE AND DOCUMENTATION CONTRIBUTORS

WITH THANKS TO

ADANHAWTH ADRIN JALALI ALBERT THOMAS ANDREAS MUELLER DAN STINE FEDA CURIC HANMIN QIN JAN S JEREMIEDBB JOEL NOTHMAN JORIS VAN DEN BOSSCHE JOSEPHSALMON KATRIN LEINWEBER LOIC ESTEVE MUHAMMAD HASSAAN RAFIQUE NICOLAS HUG OLIVIER GRISEL PAUL PACZUSKI RESHAMA SHAIKH SAM WATERBURY SHIVAM KOTWALIA THOMAS FAN

115 VERSION 0201

NOVEMBER 21 2018

THIS IS A BUGFIX RELEASE WITH SOME MINOR DOCUMENTATION IMPROVEMENTS AND ENHANCEMENTS TO FEATURES RELEASED IN 0200 NOTE THAT WE ALSO INCLUDE SOME API CHANGES IN THIS RELEASE SO YOU MIGHT GET SOME EXTRA WARNINGS AFTER UPDATING FROM 0200 TO 0201

44 CHAPTER 1 WELCOME TO SCIKITLEARN

SCIKITLEARN USER GUIDE RELEASE 0213

1151 CHANGED MODELS

THE FOLLOWING ESTIMATORS AND FUNCTIONS WHEN FIT WITH THE SAME DATA AND PARAMETERS MAY PRODUCE DIFFERENT MODELS FROM THE PREVIOUS VERSION THIS OFTEN OCCURS DUE TO CHANGES IN THE MODELLING LOGIC BUG FIXES OR ENHANCEMENTS OR IN RANDOM SAMPLING PROCEDURES

- DECOMPOSITIONINCREMENTALPCA BUG FIX

1152 CHANGELOG

SKLEARNCLUSTER

- EFFICIENCY MAKECLUSTERMEANSHIFT NO LONGER TRY TO DO NESTED PARALLELISM AS THE OVERHEAD WOULD HURT PERFORMANCE SIGNIFICANTLY WHEN NJOBS 1 12159 BY OLIVIER GRISEL

- FIXFIXED A BUG IN CLUSTERDBSCAN WITH PRECOMPUTED SPARSE NEIGHBORS GRAPH WHICH WOULD ADD EXPLICITLY ZEROS ON THE DIAGONAL EVEN WHEN ALREADY PRESENT 12105 BY TOM DUPRE LA TOUR

SKLEARNCOMPOSE

- FIXFIXED AN ISSUE IN COMPOSECOLUMNTRANSFORMER WHEN STACKING COLUMNS WITH TYPES NOT CONVERTIBLE TO A NUMERIC 11912 BY ADRIN JALALI

- API C HANGE COMPOSECOLUMNTRANSFORMER NOW APPLIES THE SPARSETHRESHOLD EVEN IF ALL TRANSFORMATION RESULTS ARE SPARSE 12304 BY ANDREAS MÜLLER

- API C HANGE COMPOSEMAKECOLUMNTRANSFORMER NOW EXPECTS TRANSFORMER COLUMNS INSTEAD OFCOLUMNS TRANSFORMER TO KEEP CONSISTENT WITH COMPOSECOLUMNTRANSFORMER 12339

BY ADRIN JALALI

SKLEARNDATASETS

- FIXDATASETSFETCHOPENML TO CORRECTLY USE THE LOCAL CACHE 12246 BY JAN N VAN RIJN

- FIXDATASETSFETCHOPENML TO CORRECTLY HANDLE IGNORE ATTRIBUTES AND ROW ID ATTRIBUTES 12330 BY JAN N VAN RIJN

- FIXFIXED INTEGER OVERFLOW IN DATASETSMAKECLASSIFICATION FOR VALUES OF NINFORMATIVE PARAMETER LARGER THAN 64 10811 BY ROMAN FELDBAUER

- FIXFIXED OLIVETTI FACES DATASET DESCR ATTRIBUTE TO POINT TO THE RIGHT LOCATION IN DATASETS FETCHOLIVETTIFACES 12441 BY JÉRÉMIE DU BOISBERRANGER

- FIXDATASETSFETCHOPENML TO RETRY DOWNLOADING WHEN READING FROM LOCAL CACHE FAILS 12517 BY THOMAS FAN

SKLEARNDECOMPOSITION

- FIXFIXED A REGRESSION IN DECOMPOSITIONINCREMENTALPCA WHERE 0200 RAISED AN ERROR IF THE NUMBER OF SAMPLES IN THE FINAL BATCH FOR FITTING INCREMENTALPCA WAS SMALLER THAN NCOMPONENTS 12234 BY MING LI

115 VERSION 0201 45

SCIKITLEARN USER GUIDE RELEASE 0213

SKLEARNENSEMBLE

- FIX FIXED A BUG MOSTLY AFFECTING ENSEMBLERANDOMFORESTCLASSIFIER WHERE CLASSWEIGHTBALANCEDSUBSAMPLE FAILED WITH MORE THAN 32 CLASSES 12165 BY JOEL NOTHMAN
- FIXFIXED A BUG AFFECTING ENSEMBLEBAGGINGCLASSIFIER ENSEMBLEBAGGINGREGRESSOR AND ENSEMBLEISOLATIONFOREST WHEREMAXFEATURES WAS SOMETIMES ROUNDED DOWN TO ZERO 12388 BY CONNOR TANN

SKLEARNFEATUREEXTRACTION

- FIXFIXED A REGRESSION IN V0200 WHERE FEATUREEXTRACTIONTEXTCOUNTVECTORIZER AND OTHER TEXT VECTORIZERS COULD ERROR DURING STOP WORDS VALIDATION WITH CUSTOM PREPROCESSORS OR TOKENIZERS 12393 BY ROMAN YURCHAK

SKLEARNLINEARMODEL

- FIXLINEARMODELSGDCLASSIFIER AND VARIANTS WITH EARLYSTOPPINGTRUE WOULD NOT USE A CONSISTENT VALIDATION SPLIT IN THE MULTICLASS CASE AND THIS WOULD CAUSE A CRASH WHEN USING THOSE ESTIMATORS AS PART OF PARALLEL PARAMETER SEARCH OR CROSSVALIDATION 12122 BY OLIVIER GRISEL
- FIXFIXED A BUG AFFECTING SGDCLASSIFIER IN THE MULTICLASS CASE EACH ONEVERSUSALL STEP IS RUN IN A JOBLIBPARALLEL CALL AND MUTATING A COMMON PARAMETER CAUSING A SEGMENTATION FAULT IF CALLED WITHIN A BACKEND USING PROCESSES AND NOT THREADS WE NOW USE REQUIRESHAREDMEM AT THEJOBLIBPARALLEL INSTANCE CREATION 12518 BY PIERRE GLASER AND OLIVIER GRISEL

SKLEARNMETRICS

- FIXFIXED A BUG IN METRICSPAIRWISEPAIRWISEDISTANCESARGMINMIN WHICH RETURNED THE SQUARE ROOT OF THE DISTANCE WHEN THE METRIC PARAMETER WAS SET TO “EUCLIDEAN” 12481 BY JÉRÉMIE DU BOISBER RANGER
- FIXFIXED A BUG IN METRICSPAIRWISEPAIRWISEDISTANCESCHUNKED WHICH DIDN’T ENSURE THE DIAGONAL IS ZERO FOR EUCLIDEAN DISTANCES 12612 BY ANDREAS MÜLLER
- API C HANGE THEMETRICSCALINSKI HARABAZSCORE HAS BEEN RENAMED TO METRICS CALINSKI HARABASZSCORE AND WILL BE REMOVED IN VERSION 023 12211 BY LISA THOMAS MARK HANNEL AND MELISSA FERRARI

SKLEARNMIXTURE

- FIXENSURE THAT THE FITPREDICT METHOD OF MIXTUREGAUSSIANMIXTURE ANDMIXTURE BAYESIANGAUSSIANMIXTURE ALWAYS YIELD ASSIGNMENTS CONSISTENT WITH FIT FOLLOWED BY PREDICT EVEN IF THE CONVERGENCE CRITERION IS TOO LOOSE OR NOT MET 12451 BY OLIVIER GRISEL

SKLEARNNEIGHBORS

- FIXFORCE THE PARALLELISM BACKEND TO THREADING FORNEIGHBORSKDTREE ANDNEIGHBORSBALLTREE IN PYTHON 27 TO AVOID PICKLING ERRORS CAUSED BY THE SERIALIZATION OF THEIR METHODS 12171 BY THOMAS MOREAU

SCIKITLEARN USER GUIDE RELEASE 0213

SKLEARNPREPROCESSING

- FIXFIXED BUG IN PREPROCESSINGORDINALENCODER WHEN PASSING MANUALLY SPECIFIED CATEGORIES 12365 BY JORIS VAN DEN BOSSCHE
- FIXFIXED BUG IN PREPROCESSINGKBINSDISCRETIZER WHERE THETRANSFORM METHOD MUTATES THE ENCODER ATTRIBUTE THE TRANSFORM METHOD IS NOW THREAD SAFE 12514 BY HANMIN QIN
- FIXFIXED A BUG IN PREPROCESSINGPOWERTRANSFORMER WHERE THE YEOJOHNSON TRANSFORM WAS INCORRECT FOR LAMBDA PARAMETERS OUTSIDE OF 0 2 12522 BY NICOLAS HUG
- FIXFIXED A BUG IN PREPROCESSINGONEHOTENCODER WHERE TRANSFORM FAILED WHEN SET TO IGNORE UNKNOWN NUMPY STRINGS OF DIFFERENT LENGTHS 12471 BY GABRIEL MARZINOTTO
- API C HANGE THE DEFAULT VALUE OF THE METHOD ARGUMENT IN PREPROCESSINGPOWERTRANSFORM WILL BE CHANGED FROM BOXCOX TOYEOJOHNSON TO MATCHPREPROCESSINGPOWERTRANSFORMER IN VER SION 023 A FUTUREWARNING IS RAISED WHEN THE DEFAULT VALUE IS USED 12317 BY ERIC CHANG

SKLEARNUTILS

- FIXUSE FLOAT64 FOR MEAN ACCUMULATOR TO AVOID FLOATING POINT PRECISION ISSUES IN PREPROCESSING STANDARDSCALER ANDDECOMPOSITIONINCREMENTALPCA WHEN USING FLOAT32 DATASETS 12338 BY BAUKS

- FIXCALLINGUTILSCHECKARRAY ONPANDASSERIES WHICH RAISED AN ERROR IN 0200 NOW RETURNS THE EXPECTED OUTPUT AGAIN 12625 BY ANDREAS MÜLLER

MISCELLANEOUS

- FIXWHEN USING SITE JOBLIB BY SETTING THE ENVIRONMENT VARIABLE SKLEARNSITEJOBLIB ADDED COMPATIBILITY WITH JOBLIB 011 IN ADDITION TO 012 12350 BY JOEL NOTHMAN AND ROMAN YURCHAK
- FIXMAKE SURE TO AVOID RAISING FUTUREWARNING WHEN CALLING NPVSTACK WITH NUMPY 116 AND LATER USE LIST COMPREHENSIONS INSTEAD OF GENERATOR EXPRESSIONS IN MANY LOCATIONS OF THE SCIKITLEARN CODE BASE 12467 BY OLIVIER GRISEL
- API C HANGE REMOVED ALL MENTIONS OF SKLEARNEXTERNALSJOBLIB AND DEPRECATED JOBLIB METHODS EXPOSED IN SKLEARNUTILS EXCEPT FOR UTILSPARALLELBACKEND ANDUTILS REGISTERPARALLELBACKEND WHICH ALLOW USERS TO CONFIGURE PARALLEL COMPUTATION IN SCIKITLEARN OTHER FUNCTIONALITIES ARE PART OF JOBLIB PACKAGE AND SHOULD BE USED DIRECTLY BY INSTALLING IT THE GOAL OF THIS CHANGE IS TO PREPARE FOR UNVENDORING JOBLIB IN FUTURE VERSION OF SCIKITLEARN 12345 BY THOMAS MOREAU

1153 CODE AND DOCUMENTATION CONTRIBUTORS

WITH THANKS TO

ADRIN JALALI ANDREA NAVARRETE ANDREAS MUELLER BAUKS BENJASTUDIO CHEUK TING HO CONNOSSOR COREY LEVIN SON DAN STINE DATENKIEKER DENIS KATAEV DILLON GARDNER DMITRY VUKOLOV DOUGAL J SUTHERLAND EDWARD J BROWN ERIC CHANG FEDERICO CASELLI GABRIEL MARZINOTTO GAELE VAROQUAUX GAURAVAH LAWAT GUSTAVO DE MARI PEREIRA HAN MIN QIN HAROLDFOX JACKLANGERMAN JACOPO NOTARSTEFANO JANVANRIJN JDETHURENS JEREMIEDBB JOEL NOTHMAN JORIS VAN DEN BOSSCHE KOEN KUSHAL CHAUHAN LEE YI JIE JOEL LILY XIONG MAILLIAM MARK HANNEL MELS YU TING MING LI NICHOLAS SMITH NICOLAS HUG NIKOLAY SHEBANOV OLEKSANDR PAVLYK OLIVIER GRISEL PETER HAUSAMANN PIERRE GLASER PULKIT MALOO QUENTIN BATISTA RADOSTIN STOYANOV RAMIL NUGMANOV REBEKAH KIM RESHAMA SHAIKH ROHAN SINGH ROMAN FELDBAUER ROMAN YURCHAK ROOPAM SHARMA SAM WATERBURY SCOTT LOWE SEBASTIAN RASCHKA STEPHEN TIER NEY SYLVAINLAN TAKINGITCASUAL THOMAS FAN THOMAS MOREAU TOM DUPRÉ LA TOUR TULIO CASAGRANDE UTKARSH UPADHYAY XING HAN LU YAROSLAV HALCHENKO ZACH MILLER

115 VERSION 0201 47

SCIKITLEARN USER GUIDE RELEASE 0213

116 VERSION 0200  
SEPTEMBER 25 2018

THIS RELEASE PACKS IN A MOUNTAIN OF BUG FIXES FEATURES AND ENHANCEMENTS FOR THE SCIKITLEARN LIBRARY AND IMPROVEMENTS TO THE DOCUMENTATION AND EXAMPLES THANKS TO OUR CONTRIBUTORS

THIS RELEASE IS DEDICATED TO THE MEMORY OF RAGHAV RAJAGOPALAN

WARNING VERSION 020 IS THE LAST VERSION OF SCIKITLEARN TO SUPPORT PYTHON 27 AND PYTHON 34 SCIKITLEARN 021 WILL REQUIRE PYTHON 35 OR HIGHER

1161 HIGHLIGHTS

WE HAVE TRIED TO IMPROVE OUR SUPPORT FOR COMMON DATASCIENCE USECASES INCLUDING MISSING VALUES CATEGORICAL VARIABLES HETEROGENEOUS DATA AND FEATURE TARGETS WITH UNUSUAL DISTRIBUTIONS MISSING VALUES IN FEATURES REPRESENTED BY NANS ARE NOW ACCEPTED IN COLUMNWISE PREPROCESSING SUCH AS SCALERS EACH FEATURE IS FITTED DISREGARDING NANS AND DATA CONTAINING NANS CAN BE TRANSFORMED THE NEW IMPUTE MODULE PROVIDES ESTIMATORS FOR LEARNING DESPITE MISSING DATA

COLUMN TRANSFORMER HANDLES THE CASE WHERE DIFFERENT FEATURES OR COLUMNS OF A PANDAS DATAFRAME NEED DIFFERENT PREPROCESSING STRING OR PANDAS CATEGORICAL COLUMNS CAN NOW BE ENCODED WITH ONEHOT ENCODER OR ORDINAL ENCODER

TRANSFORMED TARGET REGRESSOR HELPS WHEN THE REGRESSION TARGET NEEDS TO BE TRANSFORMED TO BE MODELED

POWER TRANSFORMER AND K BIN DISCRETIZER JOIN QUANTILE TRANSFORMER AS NONLINEAR TRANSFORMATIONS

BEYOND THIS WE HAVE ADDED SAMPLE WEIGHT SUPPORT TO SEVERAL ESTIMATORS INCLUDING K MEANS

BAYESIAN RIDGE AND KERNEL DENSITY AND IMPROVED STOPPING CRITERIA IN OTHERS INCLUDING MLP REGRESSOR

GRADIENT BOOSTING REGRESSOR AND SGD REGRESSOR

THIS RELEASE IS ALSO THE FIRST TO BE ACCOMPANIED BY A GLOSSARY OF COMMON TERMS AND API ELEMENTS DEVELOPED BY JOEL NOTHMAN THE GLOSSARY IS A REFERENCE RESOURCE TO HELP USERS AND CONTRIBUTORS BECOME FAMILIAR WITH THE TERMINOLOGY AND CONVENTIONS USED IN SCIKITLEARN

SORRY IF YOUR CONTRIBUTION DIDN'T MAKE IT INTO THE HIGHLIGHTS THERE'S A LOT HERE

1162 CHANGED MODELS

THE FOLLOWING ESTIMATORS AND FUNCTIONS WHEN FIT WITH THE SAME DATA AND PARAMETERS MAY PRODUCE DIFFERENT MODELS FROM THE PREVIOUS VERSION THIS OFTEN OCCURS DUE TO CHANGES IN THE MODELLING LOGIC BUG FIXES OR ENHANCEMENTS OR IN RANDOM SAMPLING PROCEDURES

- CLUSTER MEANS SHIFT BUG FIX
- DECOMPOSITION INCREMENTAL PCA IN PYTHON 2 BUG FIX
- DECOMPOSITION SPARSE PCA BUG FIX
- ENSEMBLE GRADIENT BOOSTING CLASSIFIER BUG FIX AFFECTING FEATURE IMPORTANCES
- ISOTONIC ISOTONIC REGRESSION BUG FIX
- LINEAR MODEL ARD REGRESSION BUG FIX
- LINEAR MODEL LOGISTIC REGRESSION CV BUG FIX
- LINEAR MODEL ORTHOGONAL MATCHING PURSUIT BUG FIX

48 CHAPTER 1 WELCOME TO SCIKITLEARN

SCIKITLEARN USER GUIDE RELEASE 0213

- LINEARMODELPASSIVEAGGRESSIVECLASSIFIER BUG FIX
- LINEARMODELPASSIVEAGGRESSIVEREGRESSOR BUG FIX
- LINEARMODELPERCEPTRON BUG FIX
- LINEARMODELSGDCLASSIFIER BUG FIX
- LINEARMODELSGDR regressor BUG FIX
- METRICSROCAUCSCORE BUG FIX
- METRICSROCCURVE BUG FIX
- NEURALNETWORKBASEMULTILAYERPERCEPTRON BUG FIX
- NEURALNETWORKMLPCLASSIFIER BUG FIX
- NEURALNETWORKMLPREGRESSOR BUG FIX

• THE V0190 RELEASE NOTES FAILED TO MENTION A BACKWARDS INCOMPATIBILITY WITH MODELSELECTION  
STRATIFIEDKFOLD WHENSHUFFLETRUE DUE TO 7823

DETAILS ARE LISTED IN THE CHANGELOG BELOW

WHILE WE ARE TRYING TO BETTER INFORM USERS BY PROVIDING THIS INFORMATION WE CANNOT ASSURE THAT THIS LIST IS COMPLETE  
1163 KNOWN MAJOR BUGS

- 11924 LINEARMODELLOGISTICREGRESSIONCV WITHSOLVERLBFGS AND  
MULTICLASSMULTINOMIAL MAY BE NONDETERMINISTIC OR OTHERWISE BROKEN ON MACOS THIS AP  
PEARS TO BE THE CASE ON TRAVIS CI SERVERS BUT HAS NOT BEEN CONFIRMED ON PERSONAL MACBOOKS THIS ISSUE HAS  
BEEN PRESENT IN PREVIOUS RELEASES
- 9354METRICSPAIRWISEEUCLIDEANDISTANCES WHICH IS USED SEVERAL TIMES THROUGHOUT THE LI  
BRARY GIVES RESULTS WITH POOR PRECISION WHICH PARTICULARLY AFFECTS ITS USE WITH 32BIT FLOAT INPUTS THIS BECAME  
MORE PROBLEMATIC IN VERSIONS 018 AND 019 WHEN SOME ALGORITHMS WERE CHANGED TO AVOID CASTING 32BIT DATA  
INTO 64BIT

1164 CHANGELOG

SUPPORT FOR PYTHON 33 HAS BEEN OFFICIALLY DROPPED

SKLEARNCLUSTER

- MAJOR FEATURE CLUSTERAGGLOMERATIVECLUSTERING NOW SUPPORTS SINGLE LINKAGE CLUSTERING VIA  
LINKAGESINGLE 9372 BY LELAND MCINNES AND STEVE ASTELS
- FEATURE CLUSTERKMEANS ANDCLUSTERMINIBATCHKMEANS NOW SUPPORT SAMPLE WEIGHTS VIA NEW  
PARAMETERSAMPLEWEIGHT INFIT FUNCTION 10933 BY JOHANNES HANSEN
- EFFICIENCY CLUSTERKMEANS CLUSTERMINIBATCHKMEANS ANDCLUSTERKMEANS PASSED WITH  
ALGORITHMFULL NOW ENFORCES ROWMAJOR ORDERING IMPROVING RUNTIME 10471 BY GAURAV DHINGRA
- EFFICIENCY CLUSTERDBSCAN NOW IS PARALLELIZED ACCORDING TO NJOBS REGARDLESS OF ALGORITHM 8003  
BY JOËL BILLAUD
- ENHANCEMENT CLUSTERKMEANS NOW GIVES A WARNING IF THE NUMBER OF DISTINCT CLUSTERS FOUND IS SMALLER  
THANNCLUSTERS THIS MAY OCCUR WHEN THE NUMBER OF DISTINCT POINTS IN THE DATA SET IS ACTUALLY SMALLER THAN  
THE NUMBER OF CLUSTER ONE IS LOOKING FOR 10059 BY CHRISTIAN BRAUNE

116 VERSION 0200 49

SCIKITLEARN USER GUIDE RELEASE 0213

- FIXFIXED A BUG WHERE THE FIT METHOD OFCLUSTERAFFINITYPROPAGATION STORED CLUSTER CENTERS AS 3D ARRAY INSTEAD OF 2D ARRAY IN CASE OF NONCONVERGENCE FOR THE SAME CLASS FIXED UNDEFINED AND ARBITRARY BEHAVIOR IN CASE OF TRAINING DATA WHERE ALL SAMPLES HAD EQUAL SIMILARITY 9612 BY JONATAN SAMOOCHA
  - FIXFIXED A BUG IN CLUSTERSPECTRALCLUSTERING WHERE THE NORMALIZATION OF THE SPECTRUM WAS USING A DIVISION INSTEAD OF A MULTIPLICATION 8129 BY JAN MARGETA GUILLAUME LEMAITRE AND DEVANSH D
  - FIXFIXED A BUG IN CLUSTERKMEANSELKAN WHERE THE RETURNED ITERATION WAS 1 LESS THAN THE CORRECT VALUE ALSO ADDED THE MISSING NITER ATTRIBUTE IN THE DOCSTRING OF CLUSTERKMEANS 11353 BY JEREMIE DU BOISBERRANGER
  - FIXFIXED A BUG IN CLUSTERMEANSHIFT WHERE THE ASSIGNED LABELS WERE NOT DETERMINISTIC IF THERE WERE MULTIPLE CLUSTERS WITH THE SAME INTENSITIES 11901 BY ADRIN JALALI
  - API C HANGE DEPRECATE POOLINGFUNC UNUSED PARAMETER IN CLUSTER AGGLOMERATIVECLUSTERING 9875 BY KUMAR ASHUTOSH
- SKLEARNCOMPOSE
- NEW MODULE
  - MAJOR FEATURE ADDEDCOMPOSECOLUMNSTRANSFORMER WHICH ALLOWS TO APPLY DIFFERENT TRANSFORMERS TO DIFFERENT COLUMNS OF ARRAYS OR PANDAS DATAFRAMES 9012 BY ANDREAS MÜLLER AND JORIS VAN DEN BOSSCHE AND 11315 BY THOMAS FAN
  - MAJOR FEATURE ADDED THECOMPOSETRANSFORMEDTARGETREGRESSOR WHICH TRANSFORMS THE TARGET Y BEFORE FITTING A REGRESSION MODEL THE PREDICTIONS ARE MAPPED BACK TO THE ORIGINAL SPACE VIA AN INVERSE TRANSFORM 9041 BY ANDREAS MÜLLER AND GUILLAUME LEMAITRE
- SKLEARNCOVARIANCE
- EFFICIENCY RUNTIME IMPROVEMENTS TO COVARIANCEGRAPHICALLASSO 9858 BY STEVEN BROWN
  - API C HANGE THECOVARIANCEGRAPHLASSO COVARIANCEGRAPHLASSO ANDCOVARIANCE GRAPHLASSOCV HAVE BEEN RENAMED TO COVARIANCEGRAPHICALLASSO COVARIANCE GRAPHICALLASSO ANDCOVARIANCEGRAPHICALLASSOCV RESPECTIVELY AND WILL BE REMOVED IN VERSION 022 9993 BY ARTIEM KRINITSYN

SKLEARNDATASETS

    - MAJOR FEATURE ADDEDDATASETSFETCHOPENML TO FETCH DATASETS FROM OPENML OPENML IS A FREE OPEN DATA SHARING PLATFORM AND WILL BE USED INSTEAD OF MLDATA AS IT PROVIDES BETTER SERVICE AVAILABILITY 9908 BY ANDREAS MÜLLER AND JAN N VAN RIJN
    - FEATURE INDATASETSMAKEBLOBS ONE CAN NOW PASS A LIST TO THE NSAMPLES PARAMETER TO INDICATE THE NUMBER OF SAMPLES TO GENERATE PER CLUSTER 8617 BY MASKANI FILALI MOHAMED AND KONSTANTINOS KATRIOPLAS
    - FEATURE ADDFILENAME ATTRIBUTE TO DATASETS THAT HAVE A CSV FILE 9101 BY ALEX33 AND MASKANI FILALI MOHAMED
    - FEATURE RETURNXY PARAMETER HAS BEEN ADDED TO SEVERAL DATASET LOADERS 10774 BY CHRIS CATALFO
    - FIXFIXED A BUG IN DATASETSLOADBOSTON WHICH HAD A WRONG DATA POINT 10795 BY TAKESHI YOSHIZAWA
    - FIXFIXED A BUG IN DATASETSLOADIRIS WHICH HAD TWO WRONG DATA POINTS 11082 BY SADHANA SRINI VASAN AND HANMIN QIN

50 CHAPTER 1 WELCOME TO SCIKITLEARN



SCIKITLEARN USER GUIDE RELEASE 0213

- FIXFIXED A BUG IN DATASETSFETCHKDDCUP99 WHERE DATA WERE NOT PROPERLY SHUFFLED 9731 BY NICO LAS GOIX
  - FIXFIXED A BUG IN DATASETSMKECIRCLES WHERE NO ODD NUMBER OF DATA POINTS COULD BE GENERATED 10045 BY CHRISTIAN BRAUNE
  - API CHANGE DEPRECATED SKLEARNDATASETSFETCHMLDATA TO BE REMOVED IN VERSION 022 ML DATAORG IS NO LONGER OPERATIONAL UNTIL REMOVAL IT WILL REMAIN POSSIBLE TO LOAD CACHED DATASETS 11466 BY JOEL NOTHMAN
  - SKLEARNDECOMPOSITION
    - FEATURE DECOMPOSITIONDICTLEARNING FUNCTIONS AND MODELS NOW SUPPORT POSITIVITY CONSTRAINTS THIS APPLIES TO THE DICTIONARY AND SPARSE CODE 6374 BY JOHN KIRKHAM
    - FEATURE FIXDECOMPOSITIONSPARSEPCA NOW EXPOSES NORMALIZECOMPONENTS WHEN SET TO TRUE THE TRAIN AND TEST DATA ARE CENTERED WITH THE TRAIN MEAN REPSECTIVELY DURING THE FIT PHASE AND THE TRANSFORM PHASE THIS FIXES THE BEHAVIOR OF SPARSEPCA WHEN SET TO FALSE WHICH IS THE DEFAULT THE PREVIOUS ABNORMAL BEHAVIOUR STILL HOLDS THE FALSE VALUE IS FOR BACKWARD COMPATIBILITY AND SHOULD NOT BE USED 11585 BY IVAN PANICO
    - EFFICIENCY EFFICIENCY IMPROVEMENTS IN DECOMPOSITIONDICTLEARNING 11420 AND OTHERS BY JOHN KIRKHAM
    - FIXFIX FOR UNINFORMATIVE ERROR IN DECOMPOSITIONINCREMENTALPCA NOW AN ERROR IS RAISED IF THE NUMBER OF COMPONENTS IS LARGER THAN THE CHOSEN BATCH SIZE THE NCOMPONENTSNONE CASE WAS ADAPTED ACCORDINGLY 6452 BY WALLY GAUZE
    - FIXFIXED A BUG WHERE THE PARTIALFIT METHOD OFDECOMPOSITIONINCREMENTALPCA USED INTEGER DIVISION INSTEAD OF FLOAT DIVISION ON PYTHON 2 9492 BY JAMES BOURBEAU
    - FIXINDECOMPOSITIONPCA SELECTING A NCOMPONENTS PARAMETER GREATER THAN THE NUMBER OF SAMPLES NOW RAISES AN ERROR SIMILARLY THE NCOMPONENTSNONE CASE NOW SELECTS THE MINIMUM OF NSAMPLES AND NFEATURES 8484 BY WALLY GAUZE
    - FIXFIXED A BUG IN DECOMPOSITIONPCA WHERE USERS WILL GET UNEXPECTED ERROR WITH LARGE DATASETS WHEN NCOMPONENTSMLE ON PYTHON 3 VERSIONS 9886 BY HANMIN QIN
    - FIXFIXED AN UNDERFLOW IN CALCULATING KLDIVERGENCE FOR DECOMPOSITIONNMF 10142 BY TOM DUPRE LA TOUR
    - FIXFIXED A BUG IN DECOMPOSITIONSPARSECODER WHEN RUNNING OMP SPARSE CODING IN PARALLEL USING READONLY MEMORY MAPPED DATASTRUCTURES 5956 BY VIGHNESH BIRODKAR AND OLIVIER GRISEL
  - SKLEARNDISCRIMINANTANALYSIS
    - EFFICIENCY MEMORY USAGE IMPROVEMENT FOR CLASSMEANS ANDCLASSCOV IN DISCRIMINANTANALYSIS 10898 BY NANXIN CHEN
  - SKLEARNDUMMY
    - FEATURE DUMMYDUMMYREGRESSOR NOW HAS ARETURNSTD OPTION IN ITS PREDICT METHOD THE RETURNED STANDARD DEVIATIONS WILL BE ZEROS
    - FEATURE DUMMYDUMMYCLASSIFIER ANDDUMMYDUMMYREGRESSOR NOW ONLY REQUIRE X TO BE AN OBJECT WITH FINITE LENGTH OR SHAPE 9832 BY VRISHANK BHARDWAJ
- 116 VERSION 0200 51

SCIKITLEARN USER GUIDE RELEASE 0213

- FEATURE DUMMYDUMMYCLASSIFIER ANDDDUMMYDUMMYREGRESSOR CAN NOW BE SCORED WITHOUT SUPPLYING TEST SAMPLES 11951 BY RÜDIGER BUSCHE
- SKLEARNENSEMBLE
- FEATURE ENSEMBLEBAGGINGREGRESSOR ANDENSEMBLEBAGGINGCLASSIFIER CAN NOW BE FIT WITH MISSINGNONFINITE VALUES IN X ANDOR MULTIOUTPUT Y TO SUPPORT WRAPPING PIPELINES THAT PERFORM THEIR OWN IMPUTATION 9707 BY JIMMY WAN
- FEATURE ENSEMBLEGRADIENTBOOSTINGCLASSIFIER AND ENSEMBLE GRADIENTBOOSTINGREGRESSOR NOW SUPPORT EARLY STOPPING VIA NITERNOCHANGE VALIDATIONFRACTION ANDTOL 7071 BY RAGHAV RV
- FEATURE ADDEDNAMEDESTIMATORS PARAMETER IN ENSEMBLEVOTINGCLASSIFIER TO ACCESS FITTED ESTIMATORS 9157 BY HERILALAINA RAKOTOARISON
- FIXFIXED A BUG WHEN FITTING ENSEMBLEGRADIENTBOOSTINGCLASSIFIER ORENSSEMBLE GRADIENTBOOSTINGREGRESSOR WITHWARMSTARTTRUE WHICH PREVIOUSLY RAISED A SEGMENTATION FAULT DUE TO A NONCONVERSION OF CSC MATRIX INTO CSR FORMAT EXPECTED BY DECISIONFUNCTION SIMILARLY FORTRANORDERED ARRAYS ARE CONVERTED TO CORDERED ARRAYS IN THE DENSE CASE 9991 BY GUILLAUME LEMAITRE
- FIXFIXED A BUG IN ENSEMBLEGRADIENTBOOSTINGREGRESSOR ANDENSEMBLE GRADIENTBOOSTINGCLASSIFIER TO HAVE FEATURE IMPORTANCES SUMMED AND THEN NORMALIZED RATHER THAN NORMALIZING ON A PERTREE BASIS THE PREVIOUS BEHAVIOR OVERWEIGHTED THE GINI IMPORTANCE OF FEATURES THAT APPEAR IN LATER STAGES THIS ISSUE ONLY AFFECTED FEATURE IMPORTANCES 11176 BY GIL FORSYTH
- API C HANGE THE DEFAULT VALUE OF THE NESTIMATORS PARAMETER OF ENSEMBLE RANDOMFORESTCLASSIFIER ENSEMBLERANDOMFORESTREGRESSOR ENSEMBLE EXTRATREESCLASSIFIER ENSEMBLEEXTRATREESREGRESSOR AND ENSEMBLE RANDOMTREESEMBEDDING WILL CHANGE FROM 10 IN VERSION 020 TO 100 IN 022 A FUTUREWARNING IS RAISED WHEN THE DEFAULT VALUE IS USED 11542 BY ANNA AYZENSHTAT
- API C HANGE CLASSES DERIVED FROM ENSEMBLEBASEBAGGING THE ATTRIBUTE ESTIMATORSSAMPLES WILL RETURN A LIST OF ARRAYS CONTAINING THE INDICES SELECTED FOR EACH BOOTSTRAP INSTEAD OF A LIST OF ARRAYS CONTAINING THE MASK OF THE SAMPLES SELECTED FOR EACH BOOTSTRAP INDICES ALLOWS TO REPEAT SAMPLES WHILE MASK DOES NOT ALLOW THIS FUNCTIONALITY 9524 BY GUILLAUME LEMAITRE
- FIXENSEMBLEBASEBAGGING WHERE ONE COULD NOT DETERMINISTICALLY REPRODUCE FIT RESULT USING THE OBJECT ATTRIBUTES WHEN RANDOMSTATE IS SET 9723 BY GUILLAUME LEMAITRE
- SKLEARNFEATUREEXTRACTION
- FEATURE ENABLE THE CALL TO GETFEATURENAMES IN UNFITTED FEATUREEXTRACTIONTEXT COUNTVECTORIZER INITIALIZED WITH A VOCABULARY 10908 BY MOHAMED MASKANI
- ENHANCEMENT IDF CAN NOW BE SET ON A FEATUREEXTRACTIONTEXTTFIDFTRANSFORMER 10899 BY SERGEY MELDERIS
- FIXFIXED A BUG IN FEATUREEXTRACTIONIMAGEEXTRACTPATCHES2D WHICH WOULD THROW AN EXCEPTION IF MAXPATCHES WAS GREATER THAN OR EQUAL TO THE NUMBER OF ALL POSSIBLE PATCHES RATHER THAN SIMPLY RETURNING THE NUMBER OF POSSIBLE PATCHES 10101 BY VARUN AGRAWAL
- FIXFIXED A BUG IN FEATUREEXTRACTIONTEXTCOUNTVECTORIZER FEATUREEXTRACTION TEXTTFIDFVECTORIZER FEATUREEXTRACTIONTEXTHASHINGVECTORIZER TO SUPPORT 64 BIT SPARSE ARRAY INDEXING NECESSARY TO PROCESS LARGE DATASETS WITH MORE THAN 2·109TOKENS WORDS OR NGRAMS 9147 BY CLAESFREDRIK MANNBY AND ROMAN YURCHAK
- 52 CHAPTER 1 WELCOME TO SCIKITLEARN

SCIKITLEARN USER GUIDE RELEASE 0213

- FIXFIXED BUG IN FEATUREEXTRACTIONTEXTTFIDFVECTORIZER WHICH WAS IGNORING THE PARAME  
TERDTYPE IN ADDITION FEATUREEXTRACTIONTEXTTFIDFTRANSFORMER WILL PRESERVE DTYPE FOR  
FLOATING AND RAISE A WARNING IF DTYPE REQUESTED IS INTEGER 10441 BY MAYUR KULKARNI AND GUILLAUME LEMAITRE

SKLEARNFEATURESELECTION

- FEATURE ADDED SELECT K BEST FEATURES FUNCTIONALITY TO FEATURESELECTIONSELECTFROMMODEL  
6689 BY NIHAR SHETH AND QUAZI RAHMAN
- FEATURE ADDEDMINFEATURESTOSELECT PARAMETER TO FEATURESELECTIONRFECV TO BOUND  
EVALUATED FEATURES COUNTS 11293 BY BRENT YI
- FEATURE FEATURESELECTIONRFECV 'S FIT METHOD NOW SUPPORTS GROUPS 9656 BY ADAM GREENHALL
- FIXFIXED COMPUTATION OF NFEATURESTOCOMPUTE FOR EDGE CASE WITH TIED CV SCORES IN  
FEATURESELECTIONRFECV 9222 BY NICK HOH

SKLEARNGAUSSIANPROCESS

- EFFICIENCY INGAUSSIANPROCESSGAUSSIANPROCESSREGRESSOR METHODPREDICT IS FASTER  
WHEN USING RETURNSTDTRUE IN PARTICULAR MORE WHEN CALLED SEVERAL TIMES IN A ROW 9234 BY ANDREWWW  
AND MINGHUI LIU

SKLEARNIMPUTE

- NEW MODULE ADOPTING PREPROCESSINGIMPUTER ASIMPUTESIMPLEIMPUTER WITH MINOR CHANGES  
SEE UNDER PREPROCESSING BELOW
- MAJOR FEATURE ADDEDIMPUTEMISSINGINDICATOR WHICH GENERATES A BINARY INDICATOR FOR MISSING  
VALUES 8075 BY MANITEJA NANDANA AND GUILLAUME LEMAITRE
- FEATURE THEIMPUTESIMPLEIMPUTER HAS A NEW STRATEGY CONSTANT TO COMPLETE MISSING VALUES  
WITH A FIXED ONE GIVEN BY THE FILLVALUE PARAMETER THIS STRATEGY SUPPORTS NUMERIC AND NONNUMERIC DATA  
AND SO DOES THE MOSTFREQUENT STRATEGY NOW 11211 BY JEREMIE DU BOISBERRANGER

SKLEARNISOTONIC

- FIXFIXED A BUG IN ISOTONICISOTONICREGRESSION WHICH INCORRECTLY COMBINED WEIGHTS WHEN FITTING  
A MODEL TO DATA INVOLVING POINTS WITH IDENTICAL X VALUES 9484 BY DALLAS CARD

SKLEARNLINEARMODEL

- FEATURE LINEARMODELSGDCCLASSIFIER LINEARMODELSGDREGRESSOR LINEARMODEL  
PASSIVEAGGRESSIVECLASSIFIER LINEARMODELPASSIVEAGGRESSIVEREGRESSOR AND  
LINEARMODELPERCEPTRON NOW EXPOSE EARLYSTOPPING VALIDATIONFRACTION AND  
NITERNOCHANGE PARAMETERS TO STOP OPTIMIZATION MONITORING THE SCORE ON A VALIDATION SET A NEW LEARN  
ING RATEADAPTIVE STRATEGY DIVIDES THE LEARNING RATE BY 5 EACH TIME NITERNOCHANGE CONSECUTIVE  
EPOCHS FAIL TO IMPROVE THE MODEL 9043 BY TOM DUPRE LA TOUR
- FEATURE ADD SAMPLEWEIGHT PARAMETER TO THE FIT METHOD OF LINEARMODELBAYESIANRIDGE FOR  
WEIGHTED LINEAR REGRESSION 10112 BY PETER ST JOHN

116 VERSION 0200 53

SCIKITLEARN USER GUIDE RELEASE 0213

- FIXFIXED A BUG IN LOGISTICLOGISTICREGRESSIONPATH TO ENSURE THAT THE RETURNED COEFFICIENTS ARE CORRECT WHEN MULTICLASSMULTINOMIAL PREVIOUSLY SOME OF THE COEFFICIENTS WOULD OVERRIDE EACH OTHER LEADING TO INCORRECT RESULTS IN LINEARMODELLOGISTICREGRESSIONCV 11724 BY NICOLAS HUG
  - FIXFIXED A BUG IN LINEARMODELLOGISTICREGRESSION WHERE WHEN USING THE PARAMETER MULTICLASSMULTINOMIAL THEPREDICTPROBA METHOD WAS RETURNING INCORRECT PROBABILITIES IN THE CASE OF BINARY OUTCOMES 9939 BY ROGER WESTOVER
  - FIXFIXED A BUG IN LINEARMODELLOGISTICREGRESSIONCV WHERE THESCORE METHOD ALWAYS COMPUTES ACCURACY NOT THE METRIC GIVEN BY THE SCORING PARAMETER 10998 BY THOMAS FAN
  - FIXFIXED A BUG IN LINEARMODELLOGISTICREGRESSIONCV WHERE THE 'OVR' STRATEGY WAS ALWAYS USED TO COMPUTE CROSSVALIDATION SCORES IN THE MULTICLASS SETTING EVEN IF MULTINOMIAL WAS SET 8720 BY WILLIAM DE VAZELHES
  - FIXFIXED A BUG IN LINEARMODELORTHOGONALMATCHINGPURSUIT THAT WAS BROKEN WHEN SETTING NORMALIZEFALSE 10071 BY ALEXANDRE GRAMFORT
  - FIXFIXED A BUG IN LINEARMODELARDREGRESSION WHICH CAUSED INCORRECTLY UPDATED ESTIMATES FOR THE STANDARD DEVIATION AND THE COEFFICIENTS 10153 BY JÖRG DÖPFERT
  - FIXFIXED A BUG IN LINEARMODELARDREGRESSION ANDLINEARMODELBAYESIANRIDGE WHICH CAUSED NAN PREDICTIONS WHEN FITTED WITH A CONSTANT TARGET 10095 BY JÖRG DÖPFERT
  - FIXFIXED A BUG IN LINEARMODELRIDGECLASSIFIERCV WHERE THE PARAMETER STORECVVALUES WAS NOT IMPLEMENTED THOUGH IT WAS DOCUMENTED IN CVVALUES AS A WAY TO SET UP THE STORAGE OF CROSS VALIDATION VALUES FOR DIFFERENT ALPHAS 10297 BY MABEL VILLALBAJIMÉNEZ
  - FIXFIXED A BUG IN LINEARMODELELASTICNET WHICH CAUSED THE INPUT TO BE OVERRIDDEN WHEN USING PARAMETERCOPYXTRUE ANDCHECKINPUTFALSE 10581 BY YACINE MAZARI
  - FIXFIXED A BUG IN SKLEARNLINEARMODELLASSO WHERE THE COEFFICIENT HAD WRONG SHAPE WHEN FITINTERCEPTFALSE 10687 BY MARTIN HAHN
  - FIXFIXED A BUG IN SKLEARNLINEARMODELLOGISTICREGRESSION WHERE THE MULTICLASSMULTINOMIAL WITH BINARY OUTPUT WITH WARMSTARTTRUE 10836 BY AISH WARYA SRINIVASAN
  - FIXFIXED A BUG IN LINEARMODELRIDGECV WHERE USING INTEGER ALPHAS RAISED AN ERROR 10397 BY MABEL VILLALBAJIMÉNEZ
  - FIXFIXED CONDITION TRIGGERING GAP COMPUTATION IN LINEARMODELLASSO ANDLINEARMODEL ELASTICNET WHEN WORKING WITH SPARSE MATRICES 10992 BY ALEXANDRE GRAMFORT
  - FIX FIXED A BUG IN LINEARMODELSGDCLASSIFIER LINEARMODEL SGDREGRESSOR LINEARMODELPASSIVEAGGRESSIVECLASSIFIER LINEARMODEL PASSIVEAGGRESSIVEREGRESSOR ANDLINEARMODELPERCEPTRON WHERE THE STOPPING CRITERION WAS STOPPING THE ALGORITHM BEFORE CONVERGENCE A PARAMETER NITERNOCHANGE WAS ADDED AND SET BY DEFAULT TO 5 PREVIOUS BEHAVIOR IS EQUIVALENT TO SETTING THE PARAMETER TO 1 9043 BY TOM DUPRE LA TOUR
  - FIXFIXED A BUG WHERE LIBLINEAR AND LIBSVMBASED ESTIMATORS WOULD SEGFAULT IF PASSED A SCIPYSPARSE MATRIX WITH 64BIT INDICES THEY NOW RAISE A VALUEERROR 11327 BY KARAN DHINGRA AND JOEL NOTHMAN
  - API C HANGE THE DEFAULT VALUES OF THE SOLVER ANDMULTICLASS PARAMETERS OF LINEARMODEL LOGISTICREGRESSION WILL CHANGE RESPECTIVELY FROM LIBLINEAR ANDOVR IN VERSION 020 TO LBFGS ANDAUTO IN VERSION 022 A FUTUREWARNING IS RAISED WHEN THE DEFAULT VALUES ARE USED 11905 BY TOM DUPRE LA TOUR AND JOEL NOTHMAN
  - API C HANGE DEPRECATEPOSITIVETRUE OPTION INLINEARMODELLARS AS THE UNDERLYING IMPLEMEN TATION IS BROKEN USE LINEARMODELLASSO INSTEAD 9837 BY ALEXANDRE GRAMFORT
- 54 CHAPTER 1 WELCOME TO SCIKITLEARN

SCIKITLEARN USER GUIDE RELEASE 0213

- API C HANGE NITER MAY VARY FROM PREVIOUS RELEASES IN LINEARMODELLOGISTICREGRESSION WITHSOLVERLBFGS ANDLINEARMODELHUBERREGRESSOR FOR SCIPY 100 THE OPTIMIZER COULD PERFORM MORE THAN THE REQUESTED MAXIMUM NUMBER OF ITERATIONS NOW BOTH ESTIMATORS WILL REPORT AT MOST MAXITER ITERATIONS EVEN IF MORE WERE PERFORMED 10723 BY JOEL NOTHMAN
- SKLEARNMANIFOLD
- EFFICIENCY SPEED IMPROVEMENTS FOR BOTH ‘EXACT’ AND ‘BARNESHUT’ METHODS IN MANIFOLDTSNE 10593 AND 10610 BY TOM DUPRE LA TOUR
- FEATURE SUPPORT SPARSE INPUT IN MANIFOLDISOMAPFIT 8554 BY LELAND MCINNES
- FEATURE MANIFOLDTSNETRUSTWORTHINESS ACCEPTS METRICS OTHER THAN EUCLIDEAN 9775 BY WILLIAM DE VAZELHES
- FIXFIXED A BUG IN MANIFOLDSPECTRALEMBEDDING WHERE THE NORMALIZATION OF THE SPECTRUM WAS USING A DIVISION INSTEAD OF A MULTIPLICATION 8129 BY JAN MARGETA GUILLAUME LEMAITRE AND DEVANSH D
- API C HANGE F EATURE DEPRECATE PRECOMPUTED PARAMETER IN FUNCTION MANIFOLDTSNE TRUSTWORTHINESS INSTEAD THE NEW PARAMETER METRIC SHOULD BE USED WITH ANY COMPATIBLE METRIC IN CLUDING ‘PRECOMPUTED’ IN WHICH CASE THE INPUT MATRIX XSHOULD BE A MATRIX OF PAIRWISE DISTANCES OR SQUARED DISTANCES 9775 BY WILLIAM DE VAZELHES
- API C HANGE DEPRECATEPRECOMPUTED PARAMETER IN FUNCTION MANIFOLDTSNETRUSTWORTHINESS INSTEAD THE NEW PARAMETER METRIC SHOULD BE USED WITH ANY COMPATIBLE METRIC INCLUDING ‘PRECOMPUTED’ IN WHICH CASE THE INPUT MATRIX XSHOULD BE A MATRIX OF PAIRWISE DISTANCES OR SQUARED DISTANCES 9775 BY WILLIAM DE VAZELHES
- SKLEARNMETRICS
- MAJOR FEATURE ADDED THEMETRICSDAVIESBOULDINSORE METRIC FOR EVALUATION OF CLUSTERING MOD ELS WITHOUT A GROUND TRUTH 10827 BY LUIS OSA
- MAJOR FEATURE ADDED THE METRICSBALANCEDACCURACYSORE METRIC AND A CORRESPONDING BALANCEDACCURACY SCORER FOR BINARY AND MULTICLASS CLASSIFICATION 8066 BY XYGUO AND AMAN DALMIA AND 10587 BY JOEL NOTHMAN
- FEATURE PARTIAL AUC IS AVAILABLE VIA MAXFPR PARAMETER IN METRICSROCAUCSCORE 3840 BY ALEXANDER NIEDERBÜHL
- FEATURE A SCORER BASED ON METRICSBRIERSCORELOSS IS ALSO AVAILABLE 9521 BY HANMIN QIN
- FEATURE ADDED CONTROL OVER THE NORMALIZATION IN METRICSNORMALIZEDMUTUALINFOSORE AND METRICSAJUSTEDMUTUALINFOSORE VIA THEAVERAGEMETHOD PARAMETER IN VERSION 022 THE DEFAULT NORMALIZER FOR EACH WILL BECOME THE ARITHMETIC MEAN OF THE ENTROPIES OF EACH CLUSTERING 11124 BY ARYA MCCARTHY
- FEATURE ADDEDOUTPUTDICT PARAMETER IN METRICSCCLASSIFICATIONREPORT TO RETURN CLASSIFI CATION STATISTICS AS DICTIONARY 11160 BY DAN BARKHORN
- FEATURE METRICSCCLASSIFICATIONREPORT NOW REPORTS ALL APPLICABLE AVERAGES ON THE GIVEN DATA IN CLUDING MICRO MACRO AND WEIGHTED AVERAGE AS WELL AS SAMPLES AVERAGE FOR MULTILABEL DATA 11679 BY ALEXANDER PACHA
- FEATURE METRICSAVERAGEPRECISIONSCORE NOW SUPPORTS BINARY YTRUE OTHER THAN0 1 OR 1 1 THROUGHPOSLABEL PARAMETER 9980 BY HANMIN QIN
- FEATURE METRICSLABELRANKINGAVERAGEPRECISIONSCORE NOW SUPPORTS SAMPLEWEIGHT 10845 BY JOSE PEREZPARRAS TOLEDANO
- 116 VERSION 0200 55

SCIKITLEARN USER GUIDE RELEASE 0213

- FEATURE ADDDENSEOUTPUT PARAMETER TO METRICSPAIRWISELINEARKERNEL WHEN FALSE AND BOTH INPUTS ARE SPARSE WILL RETURN A SPARSE MATRIX 10999 BY TAYLOR G SMITH
  - EFFICIENCY METRICSSILHOUETTESCORE ANDMETRICSSILHOUETTESAMPLES ARE MORE MEMORY EFFICIENT AND RUN FASTER THIS AVOIDS SOME REPORTED FREEZES AND MEMORYERRORS 11135 BY JOEL NOTHMAN
  - FIXFIXED A BUG IN METRICSPRECISIONRECALLFScoresSUPPORT WHEN TRUNCATED RANGENLABELS IS PASSED AS VALUE FOR LABELS 10377 BY GAURAV DHINGRA
  - FIXFIXED A BUG DUE TO FLOATING POINT ERROR IN METRICSRUCAUCSCORE WITH NONINTEGER SAMPLE WEIGHTS 9786 BY HANMIN QIN
  - FIXFIXED A BUG WHERE METRICSRUCCURVE SOMETIMES STARTS ON YAXIS INSTEAD OF 0 0 WHICH IS INCONSISTENT WITH THE DOCUMENT AND OTHER IMPLEMENTATIONS NOTE THAT THIS WILL NOT INFLUENCE THE RESULT FROM METRICSRUCAUCSCORE 10093 BY ALEXRYNDIN AND HANMIN QIN
  - FIXFIXED A BUG TO AVOID INTEGER OVERFLOW CASTED PRODUCT TO 64 BITS INTEGER IN METRICSMUTUALINFOSCORE 9772 BY KUMAR ASHUTOSH
  - FIXFIXED A BUG WHERE METRICSAVERAGEPRECISIONSCORE WILL SOMETIMES RETURN NAN WHEN SAMPLEWEIGHT CONTAINS 0 9980 BY HANMIN QIN
  - FIXFIXED A BUG IN METRICSFOWLKESMALLOWSSCORE TO AVOID INTEGER OVERFLOW CASTED RETURN VALUE OFCONTINGENCY MATRIX TOINT64 AND COMPUTED PRODUCT OF SQUARE ROOTS RATHER THAN SQUARE ROOT OF PRODUCT 9515 BY ALAN LIDDELL AND MANH DAO
  - API C HANGE DEPRECATEREORDER PARAMETER IN METRICSAUC AS IT'S NO LONGER REQUIRED FOR METRICSRUCAUCSCORE MOREOVER USING REORDERTRUE CAN HIDE BUGS DUE TO FLOATING POINT ERROR IN THE INPUT 9851 BY HANMIN QIN
  - API C HANGE INMETRICSNORMALIZEDMUTUALINFOSCORE ANDMETRICSAJUSTEDMUTUALINFOSCORE WARN THAT AVERAGEMETHOD WILL HAVE A NEW DEFAULT VALUE IN VERSION 022 THE DEFAULT NORMALIZER FOR EACH WILL BECOME THE ARITHMETIC MEAN OF THE ENTROPIES OF EACH CLUSTERING CURRENTLY METRICSNORMALIZEDMUTUALINFOSCORE USES THE DEFAULT OF AVERAGEMETHODGEOMETRIC ANDMETRICSAJUSTEDMUTUALINFOSCORE USES THE DEFAULT OFAVERAGEMETHODMAX TO MATCH THEIR BEHAVIORS IN VERSION 019 11124 BY ARYA MCCARTHY
  - API C HANGE THEBATCHSIZE PARAMETER TO METRICSPAIRWISEDISTANCESARGMINMIN AND METRICSPAIRWISEDISTANCESARGMIN IS DEPRECATED TO BE REMOVED IN V022 IT NO LONGER HAS ANY EFFECT AS BATCH SIZE IS DETERMINED BY GLOBAL WORKINGMEMORY CONFIG SEE LIMITING WORKING MEMORY 10280 BY JOEL NOTHMAN AND AMAN DALMIA
- SKLEARNMIXTURE
- FEATURE ADDED FUNCTION FITPREDICT TOMIXTUREGAUSSIANMIXTURE ANDMIXTUREGAUSSIANMIXTURE WHICH IS ESSENTIALLY EQUIVALENT TO CALLING FITAND PREDICT 10336 BY SHU HAORAN AND ANDREW PENG
  - FIXFIXED A BUG IN MIXTUREBASEMIXTURE WHERE THE REPORTED NITER WAS MISSING AN ITERATION IT AFFECTEDMIXTUREGAUSSIANMIXTURE ANDMIXTUREBAYESIANGAUSSIANMIXTURE 10740 BY ERICH SCHUBERT AND GUILLAUME LEMAITRE
  - FIXFIXED A BUG IN MIXTUREBASEMIXTURE AND ITS SUBCLASSES MIXTUREGAUSSIANMIXTURE AND MIXTUREBAYESIANGAUSSIANMIXTURE WHERE THELOWERBOUND WAS NOT THE MAX LOWER BOUND ACROSS ALL INITIALIZATIONS WHEN NINIT 1 BUT JUST THE LOWER BOUND OF THE LAST INITIALIZATION 10869 BY AURÉLIE GÉRON
- 56 CHAPTER 1 WELCOME TO SCIKITLEARN

SCIKITLEARN USER GUIDE RELEASE 0213

SKLEARNMODELSELECTION

- FEATURE ADDRETURNESTIMATOR PARAMETER IN MODELSELECTIONCROSSVALIDATE TO RETURN ESTIMATORS FITTED ON EACH SPLIT 9686 BY AURÉLIEN BELLET
- FEATURE NEWREFITTIME ATTRIBUTE WILL BE STORED IN MODELSELECTIONGRIDSEARCHCV AND MODELSELECTIONRANDOMIZEDSEARCHCV IFREFIT IS SET TOTRUE THIS WILL ALLOW MEASURING THE COMPLETE TIME IT TAKES TO PERFORM HYPERPARAMETER OPTIMIZATION AND REFITTING THE BEST MODEL ON THE WHOLE DATASET 11310 BY MATTHIAS FEURER
- FEATURE EXPOSEERRORSCORE PARAMETER IN MODELSELECTIONCROSSVALIDATE MODELSELECTIONCROSSVALSCORE MODELSELECTIONLEARNINGCURVE AND MODELSELECTIONVALIDATIONCURVE TO CONTROL THE BEHAVIOR TRIGGERED WHEN AN ERROR OCCURS IN MODELSELECTIONFITANDSCORE 11576 BY SAMUEL O RONSIN
- FEATURE BASESEARCHCV NOW HAS AN EXPERIMENTAL PRIVATE INTERFACE TO SUPPORT CUSTOMIZED PARAMETER SEARCH STRATEGIES THROUGH ITS RUNSEARCH METHOD SEE THE IMPLEMENTATIONS IN MODELSELECTION GRIDSEARCHCV ANDMODELSELECTIONRANDOMIZEDSEARCHCV AND PLEASE PROVIDE FEEDBACK IF YOU USE THIS NOTE THAT WE DO NOT ASSURE THE STABILITY OF THIS API BEYOND VERSION 020 9599 BY JOEL NOTHMAN
- ENHANCEMENT ADD IMPROVED ERROR MESSAGE IN MODELSELECTIONCROSSVALSCORE WHEN MULTIPLE METRICS ARE PASSED IN SCORING KEYWORD 11006 BY MING LI
- API CHANGE THE DEFAULT NUMBER OF CROSSVALIDATION FOLDS CVAND THE DEFAULT NUMBER OF SPLITS NSPLITS IN THEMODELSELECTIONKFOLD LIKE SPLITTERS WILL CHANGE FROM 3 TO 5 IN 022 AS 3FOLD HAS A LOT OF VARIANCE 11557 BY ALEXANDRE BOUCAUD
- API CHANGE THE DEFAULT OF IID PARAMETER OF MODELSELECTIONGRIDSEARCHCV AND MODELSELECTIONRANDOMIZEDSEARCHCV WILL CHANGE FROM TRUE TOFALSE IN VERSION 022 TO CORRESPOND TO THE STANDARD DEFINITION OF CROSSVALIDATION AND THE PARAMETER WILL BE REMOVED IN VERSION 024 ALTOGETHER THIS PARAMETER IS OF GREATEST PRACTICAL SIGNIFICANCE WHERE THE SIZES OF DIFFERENT TEST SETS IN CROSSVALIDATION WERE VERY UNEQUAL IE IN GROUPTBASED CV STRATEGIES 9085 BY LAURENT DIRER AND ANDREAS MÜLLER
- API CHANGE THE DEFAULT VALUE OF THE ERRORSCORE PARAMETER IN MODELSELECTIONGRIDSEARCHCV ANDMODELSELECTIONRANDOMIZEDSEARCHCV WILL CHANGE TO NPNAN IN VERSION 022 10677 BY KIRILL ZHDANOVICH
- API CHANGE CHANGED VALUEERROR EXCEPTION RAISED IN MODELSELECTIONPARAMETERSAMPLER TO A USERWARNING FOR CASE WHERE THE CLASS IS INSTANTIATED WITH A GREATER VALUE OF NITER THAN THE TOTAL SPACE OF PARAMETERS IN THE PARAMETER GRID NITER NOW ACTS AS AN UPPER BOUND ON ITERATIONS 10982 BY JULIET LAWTON
- API CHANGE INVALID INPUT FOR MODELSELECTIONPARAMETERGRID NOW RAISES TYPEERROR 10928 BY SOLUTUS IMMENSUS

SKLEARNMULTIOUTPUT

- MAJOR FEATURE ADDEDMULTIOUTPUTREGRESSORCHAIN FOR MULTITARGET REGRESSION 9257 BY KUMAR ASHUTOSH

SKLEARNNAIVEBAYES

- MAJOR FEATURE ADDEDNAIVEBAYESCOMPLEMENTNB WHICH IMPLEMENTS THE COMPLEMENT NAIVE BAYES CLASSIFIER DESCRIBED IN RENNIE ET AL 2003 8190 BY MICHAEL A ALCORN
- FEATURE ADDVARSMOOTHING PARAMETER IN NAIVEBAYESGAUSSIANNB TO GIVE A PRECISE CONTROL OVER VARIANCES CALCULATION 9681 BY DMITRY MOTTL

116 VERSION 0200 57

SCIKITLEARN USER GUIDE RELEASE 0213

- FIXFIXED A BUG IN NAIVEBAYESGAUSSIANNB WHICH INCORRECTLY RAISED ERROR FOR PRIOR LIST WHICH SUMMED TO 1 10005 BY GAURAV DHINGRA
  - FIXFIXED A BUG IN NAIVEBAYESMULTINOMIALNB WHICH DID NOT ACCEPT VECTOR VALUED PSEUDOCOUNTS ALPHA 10346 BY TOBIAS MADSEN
  - SKLEARNNEIGHBORS
    - EFFICIENCY NEIGHBORSRADIUSNEIGHBORSREGRESSOR AND NEIGHBORS RADIUSNEIGHBORSCLASSIFIER ARE NOW PARALLELIZED ACCORDING TO NJOBS REGARDLESS OF ALGORITHM 10887 BY JOËL BILLAUD
    - EFFICIENCY NEAREST NEIGHBORS QUERY METHODS ARE NOW MORE MEMORY EFFICIENT WHEN ALGORITHMBRUTE 11136 BY JOEL NOTHMAN AND AMAN DALMIA
    - FEATURE ADDSAMPLEWEIGHT PARAMETER TO THE FIT METHOD OF NEIGHBORSKERNELDENSITY TO ENABLE WEIGHTING IN KERNEL DENSITY ESTIMATION 4394 BY SAMUEL O RONSIN
    - FEATURE NOVELTY DETECTION WITH NEIGHBORSLOCALOUTLIERFACTOR ADD A NOVELTY PARAMETER TONEIGHBORSLOCALOUTLIERFACTOR WHEN NOVELTY IS SET TO TRUE NEIGHBORS LOCALOUTLIERFACTOR CAN THEN BE USED FOR NOVELTY DETECTION IE PREDICT ON NEW UNSEEN DATA AVAILABLE PREDICTION METHODS ARE PREDICT DECISIONFUNCTION ANDSCORESAMPLES BY DEFAULT NOVELTY IS SET TOFALSE AND ONLY THE FITPREDICT METHOD IS AVAIABLE BY ALBERT THOMAS
    - FIXFIXED A BUG IN NEIGHBORSNEARESTNEIGHBORS WHERE FITTING A NEARESTNEIGHBORS MODEL FAILS WHEN A THE DISTANCE METRIC USED IS A CALLABLE AND B THE INPUT TO THE NEARESTNEIGHBORS MODEL IS SPARSE 9579 BY THOMAS KOBER
    - FIXFIXED A BUG SO PREDICT INNEIGHBORSRADIUSNEIGHBORSREGRESSOR CAN HANDLE EMPTY NEIGHBOR SET WHEN USING NON UNIFORM WEIGHTS ALSO RAISES A NEW WARNING WHEN NO NEIGHBORS ARE FOUND FOR SAMPLES 9655 BY ANDREAS BJERRENIELSEN
    - FIX EFFICIENCY FIXED A BUG IN KDTREE CONSTRUCTION THAT RESULTS IN FASTER CONSTRUCTION AND QUERYING TIMES 11556 BY JAKE VANDERPLAS
    - FIXFIXED A BUG IN NEIGHBORSKDTEE ANDNEIGHBORSBALLTREE WHERE PICKLED TREE OBJECTS WOULD CHANGE THEIR TYPE TO THE SUPER CLASS BINARYTREE 11774 BY NICOLAS HUG
  - SKLEARNNEURALNETWORK
    - FEATURE ADD NITERNOCHANGE PARAMETER IN NEURALNETWORKBASEMULTILAYERPERCEPTRON NEURALNETWORKMLPREGRESSOR ANDNEURALNETWORKMLPCLASSIFIER TO GIVE CONTROL OVER MAXIMUM NUMBER OF EPOCHS TO NOT MEET TOL IMPROVEMENT 9456 BY NICHOLAS NADEAU
    - FIXFIXED A BUG IN NEURALNETWORKBASEMULTILAYERPERCEPTRON NEURALNETWORK MLPREGRESSOR ANDNEURALNETWORKMLPCLASSIFIER WITH NEWNITERNOCHANGE PARAMETER NOW AT 10 FROM PREVIOUSLY HARDCODED 2 9456 BY NICHOLAS NADEAU
    - FIXFIXED A BUG IN NEURALNETWORKMLPREGRESSOR WHERE FITTING QUIT UNEXPECTEDLY EARLY DUE TO LOCAL MINIMA OR FLUCTUATIONS 9456 BY NICHOLAS NADEAU
  - SKLEARNPIPELINE
    - FEATURE THEPREDICT METHOD OF PIPELINEPIPELINE NOW PASSES KEYWORD ARGUMENTS ON TO THE PIPELINE'S LAST ESTIMATOR ENABLING THE USE OF PARAMETERS SUCH AS RETURNSTD IN A PIPELINE WITH CAUTION 9304 BY BRENO FREITAS
- 58 CHAPTER 1 WELCOME TO SCIKITLEARN



SCIKITLEARN USER GUIDE RELEASE 0213

- API C HANGE PIPELINEFEATUREUNION NOW SUPPORTS DROP AS A TRANSFORMER TO DROP FEATURES 11144 BY THOMAS FAN
- SKLEARNPREPROCESSING
  - MAJOR FEATURE EXPANDEDPREPROCESSINGONEHOTENCODER TO ALLOW TO ENCODE CATEGORICAL STRING FEATURES AS A NUMERIC ARRAY USING A ONEHOT OR DUMMY ENCODING SCHEME AND ADDED PREPROCESSING ORDINALENCODER TO CONVERT TO ORDINAL INTEGERS THOSE TWO CLASSES NOW HANDLE ENCODING OF ALL FEATURE TYPES ALSO HANDLES STRINGVALUED FEATURES AND DERIVES THE CATEGORIES BASED ON THE UNIQUE VALUES IN THE FEATURES INSTEAD OF THE MAXIMUM VALUE IN THE FEATURES 9151 AND 10521 BY VIGHNESH BIRODKAR AND JORIS VAN DEN BOSSCHE
  - MAJOR FEATURE ADDEDPREPROCESSINGKBINSDISCRETIZER FOR TURNING CONTINUOUS FEATURES INTO CATEGORICAL OR ONEHOT ENCODED FEATURES 7668 9647 10195 10192 11272 11467 AND 11505 BY HENRY LIN HANMIN QIN TOM DUPRE LA TOUR AND GIOVANNI GIUSEPPE COSTA
  - MAJOR FEATURE ADDEDPREPROCESSINGPOWERTRANSFORMER WHICH IMPLEMENTS THE YEOJOHNSON AND BOXCOX POWER TRANSFORMATIONS POWER TRANSFORMATIONS TRY TO FIND A SET OF FEATUREWISE PARAMETRIC TRANSFORMATIONS TO APPROXIMATELY MAP DATA TO A GAUSSIAN DISTRIBUTION CENTERED AT ZERO AND WITH UNIT VARIANCE THIS IS USEFUL AS A VARIANCESTABILIZING TRANSFORMATION IN SITUATIONS WHERE NORMALITY AND HOMOSCEDASTICITY ARE DESIRABLE 10210 BY ERIC CHANG AND MANITEJA NANDANA AND 11520 BY NICOLAS HUG
  - MAJOR FEATURE NAN VALUES ARE IGNORED AND HANDLED IN THE FOLLOWING PREPROCESSING METHODS PREPROCESSINGMAXABSSCALER PREPROCESSINGMINMAXSCALER PREPROCESSING ROBUSTSCALER PREPROCESSINGSTANDARDSCALER PREPROCESSINGPOWERTRANSFORMER PREPROCESSINGQUANTILETRANSFORMER CLASSES AND PREPROCESSINGMAXABSSCALE PREPROCESSINGMINMAXSCALE PREPROCESSINGROBUSTSCALE PREPROCESSINGSCALE PREPROCESSINGPOWERTRANSFORM PREPROCESSINGQUANTILETRANSFORM FUNCTIONS RESPECTIVELY ADDRESSED IN ISSUES 11011 11005 11308 11206 11306 AND 10437 BY LUCIJA GREGOV AND GUILLAUME LEMAITRE
  - FEATURE PREPROCESSINGPOLYNOMIALFEATURES NOW SUPPORTS SPARSE INPUT 10452 BY AMAN DALMIA AND JOEL NOTHMAN
  - FEATURE PREPROCESSINGROBUSTSCALER ANDPREPROCESSINGROBUSTSCALE CAN BE FITTED USING SPARSE MATRICES 11308 BY GUILLAUME LEMAITRE
  - FEATURE PREPROCESSINGONEHOTENCODER NOW SUPPORTS THE GETFEATURENAMES METHOD TO OBTAIN THE TRANSFORMED FEATURE NAMES 10181 BY NIRVAN ANJIRBAG AND JORIS VAN DEN BOSSCHE
  - FEATURE A PARAMETER CHECKINVERSE WAS ADDED TO PREPROCESSINGFUNCTIONTRANSFORMER TO ENSURE THAT FUNC ANDINVERSEFUNC ARE THE INVERSE OF EACH OTHER 9399 BY GUILLAUME LEMAITRE
  - FEATURE THETRANSFORM METHOD OFSKLEARNPREPROCESSINGMULTILABELBINARIZER NOW IGNORES ANY UNKNOWN CLASSES A WARNING IS RAISED STATING THE UNKNOWN CLASSES CLASSES FOUND WHICH ARE IGNORED 10913 BY RODRIGO AGUNDEZ
  - FIXFIXED BUGS IN PREPROCESSINGLABELENCODER WHICH WOULD SOMETIMES THROW ERRORS WHEN TRANSFORM ORINVERSETRANSFORM WAS CALLED WITH EMPTY ARRAYS 10458 BY MAYUR KULKARNI
  - FIXFIX VALUEERROR IN PREPROCESSINGLABELENCODER WHEN USING INVERSETRANSFORM ON UNSEEN LABELS 9816 BY CHARLIE NEWWEY
  - FIXFIX BUG INPREPROCESSINGONEHOTENCODER WHICH DISCARDED THE DTYPE WHEN RETURNING A SPARSE MATRIX OUTPUT 11042 BY DANIEL MORALES
  - FIXFIXFIT ANDPARTIALFIT INPREPROCESSINGSTANDARDSCALER IN THE RARE CASE WHEN WITHMEANFALSE ANDWITHSTDFALSE WHICH WAS CRASHING BY CALLING FIT MORE THAN ONCE AND GIVING INCONSISTENT RESULTS FOR MEAN WHETHER THE INPUT WAS A SPARSE OR A DENSE MATRIX MEAN WILL BE SET TO NONE 116 VERSION 0200 59

SCIKITLEARN USER GUIDE RELEASE 0213

WITH BOTH SPARSE AND DENSE INPUTS NSAMPLESSEEN WILL BE ALSO REPORTED FOR BOTH INPUT TYPES 11235 BY GUILLAUME LEMAITRE

- API C HANGE DEPRECATE NVALUES ANDCATEGORICALFEATURES PARAMETERS AND ACTIVEFEATURES FEATUREINDICES ANDNVALUES ATTRIBUTES OF PREPROCESSING ONEHOTENCODER THENVALUES PARAMETER CAN BE REPLACED WITH THE NEW CATEGORIES PARAMETER AND THE ATTRIBUTES WITH THE NEW CATEGORIES ATTRIBUTE SELECTING THE CATEGORICAL FEATURES WITH THE CATEGORICALFEATURES PARAMETER IS NOW BETTER SUPPORTED USING THE COMPOSE COLUMNTRANSFORMER 10521 BY JORIS VAN DEN BOSSCHE

- API C HANGE DEPRECATEPREPROCESSINGIMPUTER AND MOVE THE CORRESPONDING MODULE TO IMPUTE SIMPLEIMPUTER 9726 BY KUMAR ASHUTOSH

- API C HANGE THEAXIS PARAMETER THAT WAS IN PREPROCESSINGIMPUTER IS NO LONGER PRESENT IN IMPUTESIMPLEIMPUTER THE BEHAVIOR IS EQUIVALENT TO AXIS0 IMPUTE ALONG COLUMNS ROW WISE IMPUTATION CAN BE PERFORMED WITH FUNCTIONTRANSFORMER EG FUNCTIONTRANSFORMERLAMBDA X SIMPLEIMPUTERFITTRANSFORMXTT 10829 BY GUILLAUME LEMAITRE AND GILBERTO OLIMPIO

- API C HANGE THE NAN MARKER FOR THE MISSING VALUES HAS BEEN CHANGED BETWEEN THE PREPROCESSING IMPUTER AND THE IMPUTESIMPLEIMPUTER MISSINGVALUESNAN SHOULD NOW BE

MISSINGVALUESNPNAN 11211 BY JEREMIE DU BOISBERRANGER

- API C HANGE INPREPROCESSINGFUNCTIONTRANSFORMER THE DEFAULT OF VALIDATE WILL BE FROM TRUE TOFALSE IN 022 10655 BY GUILLAUME LEMAITRE

SKLEARN SVM

- FIXFIXED A BUG IN SVM SVC WHERE WHEN THE ARGUMENT KERNEL IS UNICODE IN PYTHON2 THE PREDICTPROBA METHOD WAS RAISING AN UNEXPECTED TYPEERROR GIVEN DENSE INPUTS 10412 BY JIONGYAN ZHANG

- API C HANGE DEPRECATERANDOMSTATE PARAMETER IN SVMONECLASS SVM AS THE UNDERLYING IMPLEMENTATION IS NOT RANDOM 9497 BY ALBERT THOMAS

- API C HANGE THE DEFAULT VALUE OF GAMMA PARAMETER OF SVM SVC NUSVC SVRNUSVR ONECLASS SVM WILL CHANGE FROM AUTO TO SCALE IN VERSION 022 TO ACCOUNT BETTER FOR UNSCALED FEATURES 8361 BY GAURAV DHINGRA AND TING NEO

SKLEARN TREE

- ENHANCEMENT ALTHOUGH PRIVATE AND HENCE NOT ASSURED API STABILITY TREECRITERION CLASSIFICATIONCRITERION ANDTREECRITERIONREGRESSIONCRITERION MAY NOW BE CIM PORTED AND EXTENDED 10325 BY CAMIL STAPS

- FIXFIXED A BUG IN TREEBASEDECISIONTREE WITHSPLITTERBEST WHERE SPLIT THRESHOLD COULD BECOME INFINITE WHEN VALUES IN X WERE NEAR INFINITE 10536 BY JONATHAN OHAYON

- FIXFIXED A BUG IN TREEMAE TO ENSURE SAMPLE WEIGHTS ARE BEING USED DURING THE CALCULATION OF TREE MAE IMPURITY PREVIOUS BEHAVIOUR COULD CAUSE SUBOPTIMAL SPLITS TO BE CHOSEN SINCE THE IMPURITY CALCULATION CONSIDERED ALL SAMPLES TO BE OF EQUAL WEIGHT IMPORTANCE 11464 BY JOHN STOTT

SKLEARN UTILS

- FEATURE UTILSCHECKARRAY ANDUTILSCHECKXY NOW HAVEACCEPTLARGESPARSE TO CONTROL WHETHER SCIPYSPARSE MATRICES WITH 64BIT INDICES SHOULD BE REJECTED 11327 BY KARAN DHINGRA AND JOEL NOTHMAN

60 CHAPTER 1 WELCOME TO SCIKITLEARN

SCIKITLEARN USER GUIDE RELEASE 0213

- EFFICIENCY FIX AVOID COPYING THE DATA IN UTILSCHECKARRAY WHEN THE INPUT DATA IS A MEMMAP AND COPYFALSE 10663 BY ARTHUR MENSCH AND LOÏC ESTÈVE
  - API CHANGE UTILSCHECKARRAY YIELD A FUTURE WARNING INDICATING THAT ARRAYS OF BYTES STRINGS WILL BE INTERPRETED AS DECIMAL NUMBERS BEGINNING IN VERSION 022 10229 BY RYAN LEE
- MULTIPLE MODULES
- FEATURE API CHANGE MORE CONSISTENT OUTLIER DETECTION API ADD A SCORESAMPLES METHOD  
INSVMONECLASS SVM ENSEMBLE ISOLATION FOREST NEIGHBORS LOCAL OUTLIER FACTOR  
COVARIANCE ELLIPTIC ENVELOPE IT ALLOWS TO ACCESS RAW SCORE FUNCTIONS FROM ORIGINAL PERS  
A NEW OFFSET PARAMETER ALLOWS TO LINK SCORESAMPLES AND DECISION FUNCTION  
METHODS THE CONTAMINATION PARAMETER OF ENSEMBLE ISOLATION FOREST AND NEIGHBORS  
LOCAL OUTLIER FACTOR DECISION FUNCTION METHODS IS USED TO DEFINE THIS OFFSET SUCH THAT OUTLIERS  
RESP INLIERS HAVE NEGATIVE RESP POSITIVE DECISION FUNCTION VALUES BY DEFAULT CONTAMINATION  
IS KEPT UNCHANGED TO 0.1 FOR A DEPRECATION PERIOD IN 022 IT WILL BE SET TO "AUTO" THUS USING METHOD SPECIFIC  
SCORE OFFSETS IN COVARIANCE ELLIPTIC ENVELOPE DECISION FUNCTION METHOD THE RAW VALUES  
PARAMETER IS DEPRECATED AS THE SHIFTED MAHALANOBIS DISTANCE WILL BE ALWAYS RETURNED IN 022 9015 BY NICOLAS  
GOIX
  - FEATURE API CHANGE A BEHAVIOUR PARAMETER HAS BEEN INTRODUCED IN ENSEMBLE  
ISOLATION FOREST TO ENSURE BACKWARD COMPATIBILITY IN THE OLD BEHAVIOUR THE DECISION FUNCTION IS  
INDEPENDENT OF THE CONTAMINATION PARAMETER A THRESHOLD ATTRIBUTE DEPENDING ON THE CONTAMINATION  
PARAMETER IS THUS USED IN THE NEW BEHAVIOUR THE DECISION FUNCTION IS DEPENDENT ON THE  
CONTAMINATION PARAMETER IN SUCH A WAY THAT 0 BECOMES ITS NATURAL THRESHOLD TO DETECT OUTLIERS SET  
TING BEHAVIOUR TO "OLD" IS DEPRECATED AND WILL NOT BE POSSIBLE IN VERSION 022 BESIDE THE BEHAVIOUR PARAMETER  
WILL BE REMOVED IN 024 11553 BY NICOLAS GOIX
  - API CHANGE ADDED CONVERGENCE WARNING TO SVMLINEAR SVC AND LINEAR MODEL  
LOGISTIC REGRESSION WHEN VERBOSE IS SET TO 0 10881 BY ALEXANDRE SEVIN
  - API CHANGE CHANGED WARNING TYPE FROM USER WARNING TO EXCEPTIONS CONVERGENCE WARNING  
FOR FAILING CONVERGENCE IN LINEAR MODEL LOGISTIC REGRESSION PATH LINEAR MODEL  
RANSAC REGRESSOR LINEAR MODEL RIDGE REGRESSION GAUSSIAN PROCESS  
GAUSSIAN PROCESS REGRESSOR GAUSSIAN PROCESS GAUSSIAN PROCESS CLASSIFIER  
DECOMPOSITION FASTICA CROSS DECOMPOSITION PLSCANONICAL CLUSTER  
AFFINITY PROPAGATION AND CLUSTER BIRCH 10306 BY JONATHAN SIEBERT
- MISCELLANEOUS
- MAJOR FEATURE A NEW CONFIGURATION PARAMETER WORKING MEMORY WAS ADDED TO CONTROL MEMORY CON  
SUMPTION LIMITS IN CHUNKED OPERATIONS SUCH AS THE NEW METRIC PAIRWISE DISTANCES CHUNKED SEE  
LIMITING WORKING MEMORY 10280 BY JOEL NOTHMAN AND AMAN DALMIA
  - FEATURE THE VERSION OF JOBLIB BUNDLED WITH SCIKITLEARN IS NOW 0.12 THIS USES A NEW DEFAULT MULTIPROCESS  
ING IMPLEMENTATION NAMED LOKY WHILE THIS MAY INCUR SOME MEMORY AND COMMUNICATION OVERHEAD IT SHOULD  
PROVIDE GREATER CROSS PLATFORM STABILITY THAN RELYING ON PYTHON STANDARD LIBRARY MULTIPROCESSING 11741 BY THE  
JOBLIB DEVELOPERS ESPECIALLY THOMAS MOREAU AND OLIVIER GRISEL
  - FEATURE AN ENVIRONMENT VARIABLE TO USE THE SITE JOBLIB INSTEAD OF THE VENDOR ONE WAS ADDED ENVIRONMENT  
VARIABLES THE MAIN API OF JOBLIB IS NOW EXPOSED IN SKLEARN UTILS 11166 BY GAELE VAROQUAUX
  - FEATURE ADD ALMOST COMPLETE PYPY 3 SUPPORT KNOWN UNSUPPORTED FUNCTIONALITIES ARE DATASETS  
LOADS VMLIGHT FILE FEATURE EXTRACTION FEATURE HASHER AND FEATURE EXTRACTION  
TEXT HASHING VECTORIZER FOR RUNNING ON PYPY PYPY3V510 NUMPY 1.14.0 AND SCIPY 1.10.0 ARE  
REQUIRED 11010 BY ROMAN LAMY AND ROMAN YURCHAK
- 116 VERSION 0200 61

SCIKITLEARN USER GUIDE RELEASE 0213

- FEATURE A UTILITY METHOD SKLEARNSHOWVERSIONS WAS ADDED TO PRINT OUT INFORMATION RELEVANT FOR DEBUGGING IT INCLUDES THE USER SYSTEM THE PYTHON EXECUTABLE THE VERSION OF THE MAIN LIBRARIES AND BLAS BINDING INFORMATION 11596 BY ALEXANDRE BOUCAUD
  - FIXFIXED A BUG WHEN SETTING PARAMETERS ON METAESTIMATOR INVOLVING BOTH A WRAPPED ESTIMATOR AND ITS PARAMETER 9999 BY MARCUS V OSS AND JOEL NOTHMAN
  - FIXFIXED A BUG WHERE CALLING SKLEARNBASECLONE WAS NOT THREAD SAFE AND COULD RESULT IN A “POP FROM EMPTY LIST” ERROR 9569 BY ANDREAS MÜLLER
  - API CHANGE THE DEFAULT VALUE OF NJOBS IS CHANGED FROM 1TONONE IN ALL RELATED FUNCTIONS AND CLASSES NJOBSNONE MEANSUNSET IT WILL GENERALLY BE INTERPRETED AS NJOBS1 UNLESS THE CURRENT JOBLIB PARALLEL BACKEND CONTEXT SPECIFIES OTHERWISE SEE GLOSSARY FOR ADDITIONAL INFORMATION NOTE THAT THIS CHANGE HAPPENS IMMEDIATELY IE WITHOUT A DEPRECATION CYCLE 11741 BY OLIVIER GRISEL
  - FIXFIXED A BUG IN VALIDATION HELPERS WHERE PASSING A DASK DATAFRAME RESULTS IN AN ERROR 12462 BY ZACHARIAH MILLER
- 1165 CHANGES TO ESTIMATOR CHECKS
- THESE CHANGES MOSTLY AFFECT LIBRARY DEVELOPERS
- CHECKS FOR TRANSFORMERS NOW APPLY IF THE ESTIMATOR IMPLEMENTS TRANSFORM REGARDLESS OF WHETHER IT INHERITS FROM SKLEARNBASETRANSFORMERMIXIN 10474 BY JOEL NOTHMAN
  - CLASSIFIERS ARE NOW CHECKED FOR CONSISTENCY BETWEEN DECISIONFUNCTION AND CATEGORICAL PREDICTIONS 10500 BY NARINE KOKHLIKYAN
  - ALLOW TESTS IN UTILSESTIMATORCHECKSCHECKESTIMATOR TO TEST FUNCTIONS THAT ACCEPT PAIRWISE DATA 9701 BY KYLE JOHNSON
  - ALLOWUTILSESTIMATORCHECKSCHECKESTIMATOR TO CHECK THAT THERE IS NO PRIVATE SETTINGS APART FROM PARAMETERS DURING ESTIMATOR INITIALIZATION 9378 BY HERILALAINA RAKOTOARISON
  - THE SET OF CHECKS IN UTILSESTIMATORCHECKSCHECKESTIMATOR NOW INCLUDES A CHECKSETPARAMS TEST WHICH CHECKS THAT SETPARAMS IS EQUIVALENT TO PASSING PARAMETERS IN INIT AND WARNS IF IT ENCOUNTERS PARAMETER VALIDATION 7738 BY ALVIN CHIANG
  - ADD INVARIANCE TESTS FOR CLUSTERING METRICS 8102 BY ANKITA SINHA AND GUILLAUME LEMAITRE
  - ADDCHECKMETHODSSUBSETINVARIANCE TOCHECKESTIMATOR WHICH CHECKS THAT ESTIMATOR METHODS ARE INVARIANT IF APPLIED TO A DATA SUBSET 10428 BY JONATHAN OHAYON
  - ADD TESTS IN UTILSESTIMATORCHECKSCHECKESTIMATOR TO CHECK THAT AN ESTIMATOR CAN HANDLE READONLY MEMMAP INPUT DATA 10663 BY ARTHUR MENSCH AND LOÏC ESTÈVE
  - CHECKSAMPLEWEIGHTSPANDASSERIES NOW USES 8 RATHER THAN 6 SAMPLES TO ACCOMMODATE FOR THE DEFAULT NUMBER OF CLUSTERS IN CLUSTERKMEANS 10933 BY JOHANNES HANSEN
  - ESTIMATORS ARE NOW CHECKED FOR WHETHER SAMPLEWEIGHTNONE EQUATES TO SAMPLEWEIGHTNP ONES 11558 BY SERGUL AYDORE
- 1166 CODE AND DOCUMENTATION CONTRIBUTORS
- THANKS TO EVERYONE WHO HAS CONTRIBUTED TO THE MAINTENANCE AND IMPROVEMENT OF THE PROJECT SINCE VERSION 019 INCLUDING
- 211217613 AARSHAY JAIN ABSOLUTELYNOWARRANTY ADAM GREENHALL ADAM KLECZEWSKI ADAM RICHIEHALFORD ADEL R ADITYADAFLAPURKAR ADRIN JALALI AIDAN FITZGERALD AISHGRT1 AKASH SHIVRAM ALAN LIDDELL ALAN YEE ALBERT THOMAS
- 62 CHAPTER 1 WELCOME TO SCIKITLEARN

SCIKITLEARN USER GUIDE RELEASE 0213

ALEXANDER LENAIL ALEXANDERN ALEXANDRE BOUCAUD ALEXANDRE GRAMFORT ALEXANDRE SEVIN ALEX EGG ALVARO PEREZ DIAZ AMANDA AMAN DALMIA ANDREAS BJERRENIELSEN ANDREAS MUELLER ANDREW PENG ANGUS WILLIAMS ANIRUDDHA DAVE ANNAAYZENSHTAT ANTHONY GITTER ANTONIO QUINONEZ ANUBHAV MARWAHA ARIK PAMNANI ARTHUR OZGA ARTIEM K ARUNAVA ARYA MCCARTHY ATTRACTADORE AURÉLIEN BELLET AURÉLIEN GERON AYUSH GUPTA BALAKUMARAN MANOHA RAN BANGDA SUN BARRY HART BASTIAN VENTHUR BEN LAWSON BENN ROTH BRENO FREITAS BRENT YI BRETT KOONCE CAIO OLIVEIRA CAMIL STAPS CCLAUSS CHADY KAMAR CHARLIE BRUMMITT CHARLIE NEWWEY CHRIS CHRIS CHRIS CATALFO CHRIS FOSTER CHRIS HOLDGRAF CHRISTIAN BRAUNE CHRISTIAN HIRSCH CHRISTIAN HOGAN CHRISTOPHER JENNESS CLEMENT JOUDET CNX CWITTE DALLAS CARD DAN BARKHORN DANIEL DANIEL FERREIRA DANIEL GOMEZ DANIEL KLEVEBRING DANIELLE SHWED DANIEL MOHNS DANIL BAIBAK DARIUS MORAWIEC DAVID BEACH DAVID BURNS DAVID KIRKBY DAVID NICHOL SON DAVID PICKUP DEREK DIDI BARZEV DIEGODLH DILLON GARDNER DILLON NIEDERHUT DILUTEDSAUCE DLOVELL DMITRY MOTTTL DMITRY PETROV DOR COHEN DOUGLAS DUHAIME EKATERINA TUZOVA ERIC CHANG ERIC DEAN SANCHEZ ERICH SCHU BERT EUNJI FANGCHIEH CHOU FARAHSAAED FELIX FÉLIX RAIMUNDO FENX FILIPJ8 FRANKHUI FRANZ WOMPNER FREIJA DESCAMPS FRSI GABRIELE CALVO GAEL VAROQUAUX GAURAV DHINGRA GEORGI PEEV GIL FORSYTH GIOVANNI GIUSEPPE COSTA GKEVINYEN5418 GONCALORODRIGUES GRYLLOS PROKOPIS GUILLAUME LEMAITRE GUILLAUME “VERMEILLE” SANCHEZ GUSTAVO DE MARI PEREIRA HAKAA1 HANMIN QIN HENRY LIN HONG HONGHE HOSSEIN POURBOZORG HRISTO HUNAN ROS TOMYAN IAMPAT IVAN PANICO JAEWON CHUNG JAKE VANDERPLAS JAKIRKHAM JAMES BOURBEAU JAMES MALCOLM JAMIE COX JAN KOCH JAN MARGETA JAN SCHLÜTER JANVANRIJN JASON WOLOSONOVICH JC LIU JEB BEARER JEREMIEDBB JIMMY WAN JINKUN WANG JIONGYAN ZHANG JJABL JKLEINT JOAN MASSICH JOËL BILLAUD JOEL NOTHMAN JOHANNES HANSEN JOHNSTOTT JONATAN SAMOOCHA JONATHAN OHAYON JÖRG DÖPFERT JORIS VAN DEN BOSSCHE JOSE PEREZPARRAS TOLEDANO JOSEPHSALMON JOTASI JSCHENDEL JULIAN KUHLMANN JULIEN CHAUMOND JULIETCL JUSTIN SHENK KARL F KASPER PRIMDAL LAURITZEN KATRIN LEINWEBER KIRILL KSEMB KUAI YU KUMAR ASHUTOSH KYEONGPIL KANG KYE TAYLOR KYLEDROGO LELAND MCINNES LÉO DS LIAM GERON LIUTONG ZHOU LIZAO LI LKJCALC LOIC ESTEVE LOUIB LUCIANO VIOLA LUCIJA GREGOV LUIS OSA LUIS PEDRO COELHO LUKE M CRAIG LUKE PERSOLA MABEL MABEL VILLALBA MANITEJA NANDANA MARKI WANCHYSHYN MARK ROTH MARKUS MÜLLER MARSGUY MARTIN GUBRI MARTINHAHN MARTINKOKOS MATHURINM MATTHIAS FEURER MAX COPELAND MAYUR KULKARNI MEGHANN AGARWAL MELANIE GOETZ MICHAEL A ALCORN MINGHUI LIU MING LI MINH LE MOHAMED ALI JAMAOUI MOHAMED MASKANI MOHAMMAD SHAHEBAZ MUAYYAD ALSADI NABARUN PAL NA GARJUNA KUMAR NAOYA KANAI NARENDRAN SANTHANAM NARINEK NATHANIEL SAUL NATHAN SUH NICHOLAS NADEAU PENG A VS NICK HOH NICOLAS GOIX NICOLAS HUG NICOLAU WERNECK NIELSENMARKUS11 NIHAR SHETH NIKITA TITOV NILESH KEVLANI NIRVAN ANJIRBAG NOTMATTHANCOCK NZW OLEKSANDR PAVLYK OLIBLUM90 OLIVER RAUSCH OLIVIER GRISEL OREN MILMAN OSAID REHMAN NASIR PASBI PATRICK FERNANDES PATRICK OLDEN PAUL PACZUSKI PEDRO MORALES PETER PETER ST JOHN PIERREABLIN PIETRUH PINAKI NATH CHOWDHURY PIOTR SZYMA ´NSKI PRADEEP REDDY RAAMANA PRAVAR D MAHAJAN PRAVARMAHAJAN QINGYING CHEN RAGHAV RV RAJENDRA ARORA RAKOTOARISON HERILALAINA RAMESHWAR BHASKARAN RANKYLAU RASUL KERIMOV REIICHIRO NAKANO ROB ROMAN KOSOBRODOV ROMAN YURCHAK RONAN LAMY RRAGUNDEZ RÜDIGER BUSCHE RYAN SACHIN KELKAR SAGNIK BHATTACHARYA SAILESH CHOYAL SAM RADHAKRISHNAN SAM STEINGOLD SAMUEL BELL SAMUEL O RONSIN SAQIB NIZAM SHAMSI SATISH J SAURABH GUPTA SCOTT GIGANTE SEBASTIAN FLEN NERHAG SEBASTIAN RASCHKA SEBASTIEN DUBOIS SÉBASTIEN LERIQUE SEBASTIN SANTY SERGEY FELDMAN SERGEY MELDERIS SERGUL AYDORE SHAHEBAZ SHALIL AWALEY SHANGWU YAO SHARAD VIJALAPURAM SHARAN YALBURGI SHENHANC78 SHIVAM RASTOGI SHU HAORAN SIFTIKHA SINCLERT PÉREZ SOLUTUSIMMENSUS SOMYA ANAND SRAJAN PALIWAL SRIHARSHA HATWAR SRI KRISHNA STEFAN VAN DER WALT STEPHEN MCDOWELL STEVEN BROWN SYONEKURA TAEHOON LEE TAKANORI HAYASHI TARCUSX TAYLOR G SMITH THERILEY106 THOMAS THOMAS FAN THOMAS HEAVEY TOBIAS MADSEN TOBYCHEESE TOM AUGSPURGER TOM DUPRÉ LA TOUR TOMMY TREVOR STEPHENS TRISHNENDU GHORAI TULIO CASAGRANDE TWOSIGMAJAB UMAR FAROUK UMAR URVANG PATEL UTKARSH UPADHYAY VADIM MARKOVITSEV VARUN AGRAWAL VATHSALA ACHAR VILHELM VON EHREN HEIM VINAYAK MEHTA VINIT VINOD KUMAR L VIRAJ MAVANI VIRAJ NAVKAL VIVEK KUMAR VLAD NICULAE VQEAN3 VRIS HANK BHARDWAJ VUFG WALLYGAUZE WARUT VIJITBENJARONK WDEVAZELHES WENHAO ZHANG WES BARNETT WILL WILLIAM DE VAZELHES WILL ROSENFELD XIN XIONG YIMING PAUL LI YMAZARI YUFENG ZACH GRIFFITH ZÉ VINÍCIUS ZHENQING HU ZHIQING XIAO ZIJIE ZJ POH

117 PREVIOUS RELEASES  
1171 VERSION 0192  
JULY 2018  
117 PREVIOUS RELEASES 63

SCIKITLEARN USER GUIDE RELEASE 0213

THIS RELEASE IS EXCLUSIVELY IN ORDER TO SUPPORT PYTHON 3.7  
RELATED CHANGES

- NITER MAY VARY FROM PREVIOUS RELEASES IN LINEARMODELLOGISTICREGRESSION WITH SOLVERLBFGS AND LINEARMODELHUBERREGRESSOR FOR SCIPY 1.0.0 THE OPTIMIZER COULD PERFORM MORE THAN THE REQUESTED MAXIMUM NUMBER OF ITERATIONS NOW BOTH ESTIMATORS WILL REPORT AT MOST MAXITER ITERATIONS EVEN IF MORE WERE PERFORMED 10723 BY JOEL NOTHMAN

1.17.2 VERSION 0.19.1  
OCTOBER 23 2017

THIS IS A BUGFIX RELEASE WITH SOME MINOR DOCUMENTATION IMPROVEMENTS AND ENHANCEMENTS TO FEATURES RELEASED IN 0.19.0

NOTE THERE MAY BE MINOR DIFFERENCES IN TSNE OUTPUT IN THIS RELEASE DUE TO 9623 IN THE CASE WHERE MULTIPLE SAMPLES HAVE EQUAL DISTANCE TO SOME SAMPLE

CHANGELOG

API CHANGES

- REVERTED THE ADDITION OF METRICSNDGSCORE AND METRICSDCGSCORE WHICH HAD BEEN MERGED INTO VERSION 0.19.0 BY ERROR THE IMPLEMENTATIONS WERE BROKEN AND UNDOCUMENTED
- RETURN\_TRAINSCORE WHICH WAS ADDED TO MODELSELECTIONGRIDSEARCHCV  
MODELSELECTIONRANDOMIZEDSEARCHCV AND MODELSELECTIONCROSSVALIDATE IN VERSION 0.19.0 WILL BE CHANGING ITS DEFAULT VALUE FROM TRUE TO FALSE IN VERSION 0.21 WE FOUND THAT CALCULATING TRAINING SCORE COULD HAVE A GREAT EFFECT ON CROSS VALIDATION RUNTIME IN SOME CASES USERS SHOULD EXPLICITLY SET RETURN\_TRAINSCORE TO FALSE IF PREDICTION OR SCORING FUNCTIONS ARE SLOW RESULTING IN A DELETERIOUS EFFECT ON CV RUNTIME OR TO TRUE IF THEY WISH TO USE THE CALCULATED SCORES 9677 BY KUMAR ASHUTOSH AND JOEL NOTHMAN
- CORRELATIONMODELS AND REGRESSIONMODELS FROM THE LEGACY GAUSSIAN PROCESSES IMPLEMENTATION HAVE BEEN BELATEDLY DEPRECATED 9717 BY KUMAR ASHUTOSH

BUG FIXES

- AVOID INTEGER OVERFLOWS IN METRICSMATTHEWSCORRCOEFF 9693 BY SAM STEINGOLD
- FIXED A BUG IN THE OBJECTIVE FUNCTION FOR MANIFOLDTSNE BOTH EXACT AND WITH THE BARNESHUT APPROXIMATION WHEN N\_COMPONENTS > 3 9711 BY GONCALO RODRIGUES
- FIX REGRESSION IN MODELSELECTIONCROSSVALPREDICT WHERE IT RAISED AN ERROR WITH METHODPREDICTPROBA FOR SOME PROBABILISTIC CLASSIFIERS 9641 BY JAMES BOURBEAU
- FIXED A BUG WHERE DATASETSMAKECLASSIFICATION MODIFIED ITS INPUT WEIGHTS 9865 BY SACHIN KELKAR
- MODELSELECTIONSTRATIFIEDSHUFFLESPLIT NOW WORKS WITH MULTIOUTPUT MULTICLASS OR MULTILABEL DATA WITH MORE THAN 1000 COLUMNS 9922 BY CHARLIE BRUMMITT
- FIXED A BUG WITH NESTED AND CONDITIONAL PARAMETER SETTING EG SETTING A PIPELINE STEP AND ITS PARAMETER AT THE SAME TIME 9945 BY ANDREAS MÜLLER AND JOEL NOTHMAN

64 CHAPTER 1 WELCOME TO SCIKITLEARN

SCIKITLEARN USER GUIDE RELEASE 0213

REGRESSIONS IN 0190 FIXED IN 0191

- FIXED A BUG WHERE PARALLELISED PREDICTION IN RANDOM FORESTS WAS NOT THREADSAFE AND COULD RARELY RESULT IN ARBITRARY ERRORS 9830 BY JOEL NOTHMAN
  - FIX REGRESSION IN MODELSELECTIONCROSSVALPREDICT WHERE IT NO LONGER ACCEPTED XAS A LIST 9600 BY RASUL KERIMOV
  - FIXED HANDLING OF CROSSVALPREDICT FOR BINARY CLASSIFICATION WITH METHODDECISIONFUNCTION 9593 BY REIICHIRO NAKANO AND CORE DEVS
  - FIX REGRESSION IN PIPELINEPIPELINE WHERE IT NO LONGER ACCEPTED STEPS AS A TUPLE 9604 BY JORIS VAN DEN BOSSCHE
  - FIX BUG WHERE NITER WAS NOT PROPERLY DEPRECATED LEAVING NITER UNAVAILABLE FOR INTERIM USE
- INLINEARMODELSGDCLASSIFIER LINEARMODELSGDSREGRESSOR LINEARMODEL  
PASSIVEAGGRESSIVECLASSIFIER LINEARMODELPASSIVEAGGRESSIVEREGRESSOR AND  
LINEARMODELPERCEPTRON 9558 BY ANDREAS MÜLLER
- DATASET FETCHERS MAKE SURE TEMPORARY FILES ARE CLOSED BEFORE REMOVING THEM WHICH CAUSED ERRORS ON WINDOWS 9847 BY JOAN MASSICH
  - FIXED A REGRESSION IN MANIFOLDTSNE WHERE IT NO LONGER SUPPORTED METRICS OTHER THAN ‘EUCLIDEAN’ AND ‘PRE COMPUTED’ 9623 BY OLI BLUM
- ENHANCEMENTS
- OUR TEST SUITE AND UTILSESTIMATORCHECKSCHECKESTIMATORS CAN NOW BE RUN WITHOUT NOSE IN STALLED 9697 BY JOAN MASSICH
  - TO IMPROVE USABILITY OF VERSION 019’S PIPELINEPIPELINE CACHINGMEMORY NOW ALLOWS JOBLIB MEMORY INSTANCES THIS MAKE USE OF THE NEW UTILSVALIDATIONCHECKMEMORY HELPER ISSUE 9584 BY KUMAR ASHUTOSH
  - SOME FIXES TO EXAMPLES 9750 9788 9815
  - MADE A FUTUREWARNING IN SGDBASED ESTIMATORS LESS VERBOSE 9802 BY VRISHANK BHARDWAJ
- CODE AND DOCUMENTATION CONTRIBUTORS
- WITH THANKS TO

JOEL NOTHMAN LOIC ESTEVE ANDREAS MUELLER KUMAR ASHUTOSH VRISHANK BHARDWAJ HANMIN QIN RASUL KERIMOV  
JAMES BOURBEAU NAGARJUNA KUMAR NATHANIEL SAUL OLIVIER GRISEL ROMAN YURCHAK REIICHIRO NAKANO SACHIN KELKAR  
SAM STEINGOLD YAROSLAV HALCHENKO DIEGODLH FELIX GONCALORODRIGUES JKLEINT OLIBLUM90 PASBI ANTHONY GITTER BEN  
LAWSON CHARLIE BRUMMITT DIDI BARZEV GAELE VAROQUAUX JOAN MASSICH JORIS VAN DEN BOSSCHE NIELSENMARKUS11  
1173 VERSION 019

AUGUST 12 2017

HIGHLIGHTS

WE ARE EXCITED TO RELEASE A NUMBER OF GREAT NEW FEATURES INCLUDING NEIGHBORSLOCALOUTLIERFACTOR  
FOR ANOMALY DETECTION PREPROCESSINGQUANTILETRANSFORMER FOR ROBUST FEATURE TRANSFORMATION AND  
THEMULTIOUTPUTCLASSIFIERCHAIN METAESTIMATOR TO SIMPLY ACCOUNT FOR DEPENDENCIES BETWEEN CLASSES  
117 PREVIOUS RELEASES 65

SCIKITLEARN USER GUIDE RELEASE 0213

IN MULTILABEL PROBLEMS WE HAVE SOME NEW ALGORITHMS IN EXISTING ESTIMATORS SUCH AS MULTIPLICATIVE UP  
DATE INDECOMPOSITIONNMF AND MULTINOMIAL LINEARMODELLOGISTICREGRESSION WITH L1 LOSS USE

SOLVERSAGA

CROSS VALIDATION IS NOW ABLE TO RETURN THE RESULTS FROM MULTIPLE METRIC EVALUATIONS THE NEW MODELSELECTION  
CROSSVALIDATE CAN RETURN MANY SCORES ON THE TEST DATA AS WELL AS TRAINING SET PERFORMANCE AND TIMINGS AND WE  
HAVE EXTENDED THE SCORING ANDREFIT PARAMETERS FOR GRIDRANDOMIZED SEARCH TO HANDLE MULTIPLE METRICS  
YOU CAN ALSO LEARN FASTER FOR INSTANCE THE NEW OPTION TO CACHE TRANSFORMATIONS INPIPELINEPIPELINE MAKES  
GRID SEARCH OVER PIPELINES INCLUDING SLOW TRANSFORMATIONS MUCH MORE EFFICIENT AND YOU CAN PREDICT FASTER IF YOU'RE  
SURE YOU KNOW WHAT YOU'RE DOING YOU CAN TURN OFF VALIDATING THAT THE INPUT IS FINITE USING CONFIGCONTEXT  
WE'VE MADE SOME IMPORTANT FIXES TOO WE'VE FIXED A LONGSTANDING IMPLEMENTATION ERROR IN METRICS  
AVERAGEPRECISIONSCORE SO PLEASE BE CAUTIOUS WITH PRIOR RESULTS REPORTED FROM THAT FUNCTION A NUMBER  
OF ERRORS IN THE MANIFOLDTSNE IMPLEMENTATION HAVE BEEN FIXED PARTICULARLY IN THE DEFAULT BARNESHUT APPROX  
IMATIONSEMISUPERVISEDLABELSPREADING ANDSEMISUPERVISEDLABELPROPAGATION HAVE HAD  
SUBSTANTIAL FIXES LABELPROPAGATION WAS PREVIOUSLY BROKEN LABELSPREADING SHOULD NOW CORRECTLY RESPECT ITS ALPHA  
PARAMETER

CHANGED MODELS

THE FOLLOWING ESTIMATORS AND FUNCTIONS WHEN FIT WITH THE SAME DATA AND PARAMETERS MAY PRODUCE DIFFERENT MODELS  
FROM THE PREVIOUS VERSION THIS OFTEN OCCURS DUE TO CHANGES IN THE MODELLING LOGIC BUG FIXES OR ENHANCEMENTS OR IN  
RANDOM SAMPLING PROCEDURES

- CLUSTERKMEANS WITH SPARSE X AND INITIAL CENTROIDS GIVEN BUG FIX
- CROSSDECOMPOSITIONPLSREGRESSION WITHSCALETRUE BUG FIX
- ENSEMBLEGRADIENTBOOSTINGCLASSIFIER ANDENSEMBLEGRADIENTBOOSTINGREGRESSOR  
WHERE MINIMPURITYSPLIT IS USED BUG FIX
- GRADIENT BOOSTING LOSSQUANTILE BUG FIX
- ENSEMBLEISOLATIONFOREST BUG FIX
- FEATURESELECTIONSELECTFDR BUG FIX
- LINEARMODEL RANSACREGRESSOR BUG FIX
- LINEARMODEL LASSOLARS BUG FIX
- LINEARMODEL LASSOLARSIC BUG FIX
- MANIFOLDTSNE BUG FIX
- NEIGHBORSNEARESTCENTROID BUG FIX
- SEMISUPERVISEDLABELSPREADING BUG FIX
- SEMISUPERVISEDLABELPROPAGATION BUG FIX
- TREE BASED MODELS WHERE MINWEIGHTFRACTIONLEAF IS USED ENHANCEMENT
- MODELSELECTIONSTRATIFIEDKFOLD WITHSHUFFLETRUE THIS CHANGE DUE TO 7823 WAS NOT MEN

TIONED IN THE RELEASE NOTES AT THE TIME  
DETAILS ARE LISTED IN THE CHANGELOG BELOW

WHILE WE ARE TRYING TO BETTER INFORM USERS BY PROVIDING THIS INFORMATION WE CANNOT ASSURE THAT THIS LIST IS COMPLETE  
66 CHAPTER 1 WELCOME TO SCIKITLEARN



SCIKITLEARN USER GUIDE RELEASE 0213

CHANGELOG

NEW FEATURES

CLASSIFIERS AND REGRESSORS

- ADDEDMULTIOUTPUTCLASSIFIERCHAIN FOR MULTILABEL CLASSIFICATION BY ADAM KLECZEWSKI
- ADDED SOLVER SAGA THAT IMPLEMENTS THE IMPROVED VERSION OF STOCHASTIC AVERAGE GRADIENT IN LINEARMODELLOGISTICREGRESSION ANDLINEARMODELRIDGE IT ALLOWS THE USE OF L1 PENALTY WITH MULTINOMIAL LOGISTIC LOSS AND BEHAVES marginally better than ‘SAG’ DURING THE FIRST EPOCHS OF RIDGE AND LOGISTIC REGRESSION 8446 BY ARTHUR MENSCH

OTHER ESTIMATORS

- ADDED THE NEIGHBORSLOCALOUTLIERFACTOR CLASS FOR ANOMALY DETECTION BASED ON NEAREST NEIGHBORS 5279 BY NICOLAS GOIX AND ALEXANDRE GRAMFORT

• ADDED PREPROCESSINGQUANTILETRANSFORMER CLASS AND PREPROCESSING

QUANTILETRANSFORM FUNCTION FOR FEATURES NORMALIZATION BASED ON QUANTILES 8363 BY DENIS ENGEMANN GUILLAUME LEMAITRE OLIVIER GRISEL RAGHAV RV THIERRY GUILLEMOT AND GAEL VAROQUAUX

- THE NEW SOLVER MU IMPLEMENTS A MULTIPLICATE UPDATE IN DECOMPOSITIONNMF ALLOWING THE OPTIMIZATION OF ALL BETADIVERGENCES INCLUDING THE FROBENIUS NORM THE GENERALIZED KULLBACKLEIBLER DIVERGENCE AND THE ITAKURASAITO DIVERGENCE 5295 BY TOM DUPRE LA TOUR

MODEL SELECTION AND EVALUATION

- MODELSELECTIONGRIDSEARCHCV ANDMODELSELECTIONRANDOMIZEDSEARCHCV NOW SUPPORT SIMULTANEOUS EVALUATION OF MULTIPLE METRICS REFER TO THE SPECIFYING MULTIPLE METRICS FOR EVALUATION SECTION OF THE USER GUIDE FOR MORE INFORMATION 7388 BY RAGHAV RV
- ADDED THE MODELSELECTIONCROSSVALIDATE WHICH ALLOWS EVALUATION OF MULTIPLE METRICS THIS FUNC TION RETURNS A DICT WITH MORE USEFUL INFORMATION FROM CROSSVALIDATION SUCH AS THE TRAIN SCORES FIT TIMES AND SCORE TIMES REFER TO THE CROSSVALIDATE FUNCTION AND MULTIPLE METRIC EVALUATION SECTION OF THE USERGUIDE FOR MORE INFORMATION 7388 BY RAGHAV RV

- ADDEDMETRICSMEANSQUAREDLOGERROR WHICH COMPUTES THE MEAN SQUARE ERROR OF THE LOGARITHMIC TRANSFORMATION OF TARGETS PARTICULARLY USEFUL FOR TARGETS WITH AN EXPONENTIAL TREND 7655 BY KARAN DESAI

- ADDEDMETRICSDCGSCORE ANDMETRICSNDCGSCORE WHICH COMPUTE DISCOUNTED CUMULATIVE GAIN

DCG AND NORMALIZED DISCOUNTED CUMULATIVE GAIN NDCG 7739 BY DAVID GASQUEZ

- ADDED THE MODELSELECTIONREPEATEDKFOLD ANDMODELSELECTION REPEATEDSTRATIFIEDKFOLD 8120 BY NEERAJ GANGWAR

MISCELLANEOUS

- VALIDATION THAT INPUT DATA CONTAINS NO NAN OR INF CAN NOW BE SUPPRESSED USING CONFIGCONTEXT AT YOUR OWN RISK THIS WILL SAVE ON RUNTIME AND MAY BE PARTICULARLY USEFUL FOR PREDICTION TIME 7548 BY JOEL NOTHMAN
- ADDED A TEST TO ENSURE PARAMETER LISTING IN DOCSTRINGS MATCH THE FUNCTIONCLASS SIGNATURE 9206 BY ALEXANDRE GRAMFORT AND RAGHAV RV

ENHANCEMENTS

TREES AND ENSEMBLES

- THEMINWEIGHTFRACTIONLEAF CONSTRAINT IN TREE CONSTRUCTION IS NOW MORE EFFICIENT TAKING A FAST PATH TO DECLARE A NODE A LEAF IF ITS WEIGHT IS LESS THAN 2 THE MINIMUM NOTE THAT THE CONSTRUCTED TREE WILL BE DIFFERENT FROM PREVIOUS VERSIONS WHERE MINWEIGHTFRACTIONLEAF IS USED 7441 BY NELSON LIU

117 PREVIOUS RELEASES 67

SCIKITLEARN USER GUIDE RELEASE 0213

- ENSEMBLEGRADIENTBOOSTINGCLASSIFIER ANDENSEMBLEGRADIENTBOOSTINGREGRESSOR NOW SUPPORT SPARSE INPUT FOR PREDICTION 6101 BY IBRAIM GANIEV
  - ENSEMBLEVOTINGCLASSIFIER NOW ALLOWS CHANGING ESTIMATORS BY USING ENSEMBLE VOTINGCLASSIFIERSETPARAMS AN ESTIMATOR CAN ALSO BE REMOVED BY SETTING IT TO NONE 7674 BY YICHUAN LIU
  - TREEEXPORTGRAPHVIZ NOW SHOWS CONFIGURABLE NUMBER OF DECIMAL PLACES 8698 BY GUILLAUME LEMAITRE
  - ADDEDFLATTENTTRANSFORM PARAMETER TO ENSEMBLEVOTINGCLASSIFIER TO CHANGE OUTPUT SHAPE OF TRANSFORM METHOD TO 2 DIMENSIONAL 7794 BY IBRAIM GANIEV AND HERILALAINA RAKOTOARISON
  - LINEAR KERNELIZED AND RELATED MODELS
  - LINEARMODELSGDCLASSIFIER LINEARMODELSGDREGRESSOR LINEARMODEL PASSIVEAGGRESSIVECLASSIFIER LINEARMODELPASSIVEAGGRESSIVEREGRESSOR AND LINEARMODELPERCEPTRON NOW EXPOSE MAXITER ANDTOL PARAMETERS TO HANDLE CONVERGENCE MORE PRECISELYNITER PARAMETER IS DEPRECATED AND THE FITTED ESTIMATOR EXPOSES A NITER ATTRIBUTE WITH ACTUAL NUMBER OF ITERATIONS BEFORE CONVERGENCE 5036 BY TOM DUPRE LA TOUR
  - ADDED AVERAGE PARAMETER TO PERFORM WEIGHT AVERAGING IN LINEARMODEL PASSIVEAGGRESSIVECLASSIFIER 4939 BY ANDREA ESULI
  - LINEARMODELRANSACREGRESSOR NO LONGER THROWS AN ERROR WHEN CALLING FIT IF NO INLIERS ARE FOUND IN ITS FIRST ITERATION FURTHERMORE CAUSES OF SKIPPED ITERATIONS ARE TRACKED IN NEWLY ADDED ATTRIBUTES NSKIPS 7914 BY MICHAEL HORRELL
  - INGAUSSIANPROCESSGAUSSIANPROCESSREGRESSOR METHOD PREDICT IS A LOT FASTER WITH RETURNSTDTRUE 8591 BY HADRIEN BERTRAND
  - ADDEDRETURNSTD TOPREDICT METHOD OFLINEARMODELARDREGRESSION ANDLINEARMODEL BAYESIANRIDGE 7838 BY SERGEY FELDMAN
  - MEMORY USAGE ENHANCEMENTS PREVENT CAST FROM FLOAT32 TO FLOAT64 IN LINEARMODEL MULTITASKELASTICNET LINEARMODELLOGISTICREGRESSION WHEN USING NEWTONCG SOLVER AND LINEARMODELRIDGE WHEN USING SVD SPARSECG CHOLESKY OR LSQR SOLVERS 8835 8061 BY JOAN MASSICH AND NICOLAS CORDIER AND THIERRY GUILLEMOT
  - OTHER PREDICTORS
  - CUSTOM METRICS FOR THE NEIGHBORS BINARY TREES NOW HAVE FEWER CONSTRAINTS THEY MUST TAKE TWO 1DARRAYS AND RETURN A FLOAT 6288 BY JAKE VANDERPLAS
  - ALGORITHM AUTO INNEIGHBORS ESTIMATORS NOW CHOOSES THE MOST APPROPRIATE ALGORITHM FOR ALL INPUT TYPES AND METRICS 9145 BY HERILALAINA RAKOTOARISON AND REDDY CHINTHALA
  - DECOMPOSITION MANIFOLD LEARNING AND CLUSTERING
  - CLUSTERMINIBATCHKMEANS ANDCLUSTERKMEANS NOW USE SIGNIFICANTLY LESS MEMORY WHEN ASSIGNING DATA POINTS TO THEIR NEAREST CLUSTER CENTER 7721 BY JON CRALL
  - DECOMPOSITIONPCA DECOMPOSITIONINCREMENTALPCA ANDDECOMPOSITION TRUNCATEDSVD NOW EXPOSE THE SINGULAR VALUES FROM THE UNDERLYING SVD THEY ARE STORED IN THE ATTRIBUTESINGULARVALUES LIKE INDECOMPOSITIONINCREMENTALPCA 7685 BY TOMMY LÖFSTEDT
  - DECOMPOSITIONNMF NOW FASTER WHEN BETALOSSO 9277 BY HONGKAHJUN
  - MEMORY IMPROVEMENTS FOR METHOD BARNESHUT INMANIFOLDTSNE 7089 BY THOMAS MOREAU AND OLIVIER GRISEL
  - OPTIMIZATION SCHEDULE IMPROVEMENTS FOR BARNESHUT MANIFOLDTSNE SO THE RESULTS ARE CLOSER TO THE ONE FROM THE REFERENCE IMPLEMENTATION LVDMAATENBHTSNE BY THOMAS MOREAU AND OLIVIER GRISEL
- 68 CHAPTER 1 WELCOME TO SCIKITLEARN

SCIKITLEARN USER GUIDE RELEASE 0213

- MEMORY USAGE ENHANCEMENTS PREVENT CAST FROM FLOAT32 TO FLOAT64 IN DECOMPOSITIONPCA AND DECOMPOSITIONRANDOMIZEDSVDLOWRANK 9067 BY RAGHAV RV
  - PREPROCESSING AND FEATURE SELECTION
  - ADDEDNORMORDER PARAMETER TO FEATURESELECTIONSELECTFROMMODEL TO ENABLE SELECTION OF THE NORM ORDER WHEN COEF IS MORE THAN 1D 6181 BY ANTOINE WENDLINGER
  - ADDED ABILITY TO USE SPARSE MATRICES IN FEATURESELECTIONFREGRESSION WITHCENTERTRUE 8065 BY DANIEL LEJEUNE
  - SMALL PERFORMANCE IMPROVEMENT TO NGRAM CREATION IN FEATUREEXTRACTIONTEXT BY BINDING METHODS FOR LOOPS AND SPECIALCASING UNIGRAMS 7567 BY JAYE DOEPKÉ
  - RELAX ASSUMPTION ON THE DATA FOR THE KERNELAPPROXIMATIONSKEWEDCHI2SAMPLER SINCE THE SKEWEDCHI2 KERNEL IS DEFINED ON THE OPEN INTERVAL  $-\infty$   $\infty$  THE TRANSFORM FUNCTION SHOULD NOT CHECK WHETHER X 0 BUT WHETHER X SELF-SKEWEDNESS 7573 BY ROMAIN BRAULT
  - MADE DEFAULT KERNEL PARAMETERS KERNELDEPENDENT IN KERNELAPPROXIMATIONNNYSTROEM 5229 BY SAURABH BANSOD AND ANDREAS MÜLLER
  - MODEL EVALUATION AND METAESTIMATORS
  - PIPELINEPIPELINE IS NOW ABLE TO CACHE TRANSFORMERS WITHIN A PIPELINE BY USING THE MEMORY CONSTRUCTOR PARAMETER 7990 BY GUILLAUME LEMAITRE
  - PIPELINEPIPELINE STEPS CAN NOW BE ACCESSED AS ATTRIBUTES OF ITS NAMEDSTEPS ATTRIBUTE 8586 BY HERILALAINA RAKOTOARISON
  - ADDEDSAMPLEWEIGHT PARAMETER TO PIPELINEPIPELINESCORE 7723 BY MIKHAIL KOROBOV
  - ADDED ABILITY TO SET NJOBS PARAMETER TO PIPELINEMAKEUNION ATYPEERROR WILL BE RAISED FOR ANY OTHER KWARGS 8028 BY ALEXANDER BOOTH
  - MODELSELECTIONGRIDSEARCHCV MODELSELECTIONRANDOMIZEDSEARCHCV AND MODELSELECTIONCROSSVALSCORE NOW ALLOW ESTIMATORS WITH CALLABLE KERNELS WHICH WERE PREVIOUSLY PROHIBITED 8005 BY ANDREAS MÜLLER
  - MODELSELECTIONCROSSVALPREDICT NOW RETURNS OUTPUT OF THE CORRECT SHAPE FOR ALL VALUES OF THE ARGUMENTMETHOD 7863 BY AMAN DALMIA
  - ADDEDSHUFFLE ANDRANDOMSTATE PARAMETERS TO SHUFFLE TRAINING DATA BEFORE TAKING PREFIXES OF IT BASED ON TRAINING SIZES IN MODELSELECTIONLEARNINGCURVE 7506 BY NARINE KOKHLIKYAN
  - MODELSELECTIONSTRATIFIEDSHUFFLESPLIT NOW WORKS WITH MULTIOUTPUT MULTICLASS OR MULTILABEL DATA 9044 BY VLAD NICULAE
  - SPEED IMPROVEMENTS TO MODELSELECTIONSTRATIFIEDSHUFFLESPLIT 5991 BY ARTHUR MENSCH AND JOEL NOTHMAN
  - ADDSHUFFLE PARAMETER TO MODELSELECTIONTRAINTESTSPLIT 8845 BY THEM RMAX
  - MULTIOUTPUTMULTIOUTPUTREGRESSOR ANDMULTIOUTPUTMULTIOUTPUTCLASSIFIER NOW SUPPORT ONLINE LEARNING USING PARTIALFIT ISSUE8053 BY PENG YU
  - ADDMAXTRAIN SIZE PARAMETER TO MODELSELECTIONTIMESERIESSPLIT 8282 BY AMAN DALMIA
  - MORE CLUSTERING METRICS ARE NOW AVAILABLE THROUGH METRICSGETSCORER ANDSCORING PARAMETERS 8117 BY RAGHAV RV
  - A SCORER BASED ON METRICSEXPLAINEDVARIANCESCORE IS ALSO AVAILABLE 9259 BY HANMIN QIN
  - METRICS
  - METRICSMATTHEWSCORRCOEFF NOW SUPPORT MULTICLASS CLASSIFICATION 8094 BY JON CRALL
- 117 PREVIOUS RELEASES 69

SCIKITLEARN USER GUIDE RELEASE 0213

- ADDSAMPLEWEIGHT PARAMETER TO METRICSCOHENKAPPASCORE 8335 BY VICTOR PUGHON
- MISCELLANEOUS
- UTILSCHECKESTIMATOR NOW ATTEMPTS TO ENSURE THAT METHODS TRANSFORM PREDICT ETC DO NOT SET ATTRIBUTES ON THE ESTIMATOR 7533 BY EKATERINA KRIVICH
  - ADDED TYPE CHECKING TO THE ACCEPTSPARSE PARAMETER IN UTILSVALIDATION METHODS THIS PARAMETER NOW ACCEPTS ONLY BOOLEAN STRING OR LISTTUPLE OF STRINGS ACCEPTSPARSENONE IS DEPRECATED AND SHOULD BE REPLACED BY ACCEPTSPARSEFALSE 7880 BY JOSH KARNOFSKY
  - MAKE IT POSSIBLE TO LOAD A CHUNK OF AN SVMLIGHT FORMATTED FILE BY PASSING A RANGE OF BYTES TO DATASETS LOADSVMLIGHTFILE 935 BY OLIVIER GRISEL
  - DUMMYDUMMYCLASSIFIER ANDDUMMYDUMMYREGRESSOR NOW ACCEPT NONFINITE FEATURES 8931 BY ATTRACTADORE
- BUG FIXES
- TREES AND ENSEMBLES
- FIXED A MEMORY LEAK IN TREES WHEN USING TREES WITH CRITERIONMAE 8002 BY RAGHAV RV
  - FIXED A BUG WHERE ENSEMBLEISOLATIONFOREST USES AN AN INCORRECT FORMULA FOR THE AVERAGE PATH LENGTH 8549 BY PETER WANG
  - FIXED A BUG WHERE ENSEMBLEADABOOSTCLASSIFIER THROWSZERODIVISIONERROR WHILE FITTING DATA WITH SINGLE CLASS LABELS 7501 BY DOMINIK KRZEMINSKI
  - FIXED A BUG IN ENSEMBLEGRADIENTBOOSTINGCLASSIFIER ANDENSEMBLE GRADIENTBOOSTINGREGRESSOR WHERE A FLOAT BEING COMPARED TO 00 USINGCAUSED A DIVIDE BY ZERO ERROR 7970 BY HE CHEN
  - FIX A BUG WHERE ENSEMBLEGRADIENTBOOSTINGCLASSIFIER ANDENSEMBLE GRADIENTBOOSTINGREGRESSOR IGNORED THE MINIMPURITYSPLIT PARAMETER 8006 BY SEBASTIAN PÖLSTERL
  - FIXEDOOBSCORE INENSEMBLEBAGGINGCLASSIFIER 8936 BY MICHAEL LEWIS
  - FIXED EXCESSIVE MEMORY USAGE IN PREDICTION FOR RANDOM FORESTS ESTIMATORS 8672 BY MIKE BENFIELD
  - FIXED A BUG WHERE SAMPLEWEIGHT AS A LIST BROKE RANDOM FORESTS IN PYTHON 2 8068 BY XOR
  - FIXED A BUG WHERE ENSEMBLEISOLATIONFOREST FAILS WHEN MAXFEATURES IS LESS THAN 1 5732 BY ISHANK GULATI
  - FIX A BUG WHERE GRADIENT BOOSTING WITH LOSSQUANTILE COMPUTED NEGATIVE ERRORS FOR NEGATIVE VALUES OF YTRUE YPRED LEADING TO WRONG VALUES WHEN CALLING CALL 8087 BY ALEXIS MIGNON
  - FIX A BUG WHERE ENSEMBLEVOTINGCLASSIFIER RAISES AN ERROR WHEN A NUMPY ARRAY IS PASSED IN FOR WEIGHTS 7983 BY VINCENT PHAM
  - FIXED A BUG WHERE TREEEXPORTGRAPHVIZ RAISED AN ERROR WHEN THE LENGTH OF FEATURES NAMES DOES NOT MATCH NFEATURES IN THE DECISION TREE 8512 BY LI LI
- LINEAR KERNELIZED AND RELATED MODELS
- FIXED A BUG WHERE LINEARMODELTRANSACREGRESSORFIT MAY RUN UNTIL MAXITER IF IT FINDS A LARGE INLIER GROUP EARLY 8251 BY AIVISION2020
  - FIXED A BUG WHERE NAIVEBAYESMULTINOMIALNB ANDNAIVEBAYESBERNOULLINB FAILED WHEN ALPHA0 5814 BY YICHUAN LIU AND HERILALAINA RAKOTOARISON
- 70 CHAPTER 1 WELCOME TO SCIKITLEARN

SCIKITLEARN USER GUIDE RELEASE 0213

- FIXED A BUG WHERE LINEARMODELLASSOLARS DOES NOT GIVE THE SAME RESULT AS THE LASSOLARS IMPLEMENTATION AVAILABLE IN R LARS LIBRARY 7849 BY JAIR MONTOYA MARTINEZ
  - FIXED A BUG IN LINEARMODELRANDOMIZEDLASSO LINEARMODELLARS LINEARMODEL LASSOLARS LINEARMODELLARSCV ANDLINEARMODELLASSOLARSCV WHERE THE PARAMETER PRECOMPUTE WAS NOT USED CONSISTENTLY ACROSS CLASSES AND SOME VALUES PROPOSED IN THE DOCSTRING COULD RAISE ERRORS 5359 BY TOM DUPRE LA TOUR
  - FIX INCONSISTENT RESULTS BETWEEN LINEARMODELRIDGECV ANDLINEARMODEL RIDGE WHEN USING NORMALIZETRUE 9302 BY ALEXANDRE GRAMFORT
  - FIX A BUG WHERE LINEARMODELLASSOLARSFIT SOMETIMES LEFT COEF AS A LIST RATHER THAN AN NDARRAY 8160 BY CJ CAREY
  - FIXLINEARMODEL BAYESIANRIDGECV TO RETURN RIDGE PARAMETER ALPHA ANDLAMBDA CONSISTENT WITH CALCULATED COEFFICIENTS COEF ANDINTERCEPT 8224 BY PETER GEDECK
  - FIXED A BUG IN SVMONECLASS SVM WHERE IT RETURNED FLOATS INSTEAD OF INTEGER CLASSES 8676 BY VATHSALA ACHAR
  - FIX AICBIC CRITERION COMPUTATION IN LINEARMODELLASSOLARSIC 9022 BY ALEXANDRE GRAMFORT AND MEHMET BASBUG
  - FIXED A MEMORY LEAK IN OUR LIBLINEAR IMPLEMENTATION 9024 BY SERGEI LEBEDEV
  - FIX BUG WHERE STRATIFIED CV SPLITTERS DID NOT WORK WITH LINEARMODELLASSOCV 8973 BY PAULO HADDAD
  - FIXED A BUG IN GAUSSIANPROCESSGAUSSIANPROCESSREGRESSOR WHEN THE STANDARD DEVIATION AND COVARIANCE PREDICTED WITHOUT FIT WOULD FAIL WITH A UNMEANINGFUL ERROR BY DEFAULT 6573 BY QUAZI MARUFUR RAHMAN AND MANOJ KUMAR
- OTHER PREDICTORS
- FIXSEMISUPERVISED BASE LABELPROPAGATION TO CORRECTLY IMPLEMENT LABELPROPAGATION AND LABELSPREADING AS DONE IN THE REFERENCED PAPERS 9239 BY ANDRE AMBROSIO BOECHAT UTKARSH UPADHYAY AND JOEL NOTHMAN
- DECOMPOSITION MANIFOLD LEARNING AND CLUSTERING
- FIXED THE IMPLEMENTATION OF MANIFOLDTSNE
  - EARLYEXAGERATION PARAMETER HAD NO EFFECT AND IS NOW USED FOR THE FIRST 250 OPTIMIZATION ITERATIONS
  - FIXED THE ASSERTIONERROR TREE CONSISTENCY FAILED EXCEPTION REPORTED IN 8992
  - IMPROVE THE LEARNING SCHEDULE TO MATCH THE ONE FROM THE REFERENCE IMPLEMENTATION LVDMAATENBHTSNE BY THOMAS MOREAU AND OLIVIER GRISEL
  - FIX A BUG IN DECOMPOSITIONLATENTDIRICHLETALLOCATION WHERE THEPERPLEXITY METHOD WAS RETURNING INCORRECT RESULTS BECAUSE THE TRANSFORM METHOD RETURNS NORMALIZED DOCUMENT TOPIC DISTRIBUTIONS AS OF VERSION 018 7954 BY GARY FOREMAN
  - FIX OUTPUT SHAPE AND BUGS WITH NJOBS 1 IN DECOMPOSITIONSPARSECODER TRANSFORM AND DECOMPOSITIONSPARSEENCODE FOR ONEDIMENSIONAL DATA AND ONE COMPONENT THIS ALSO IMPACTS THE OUTPUT SHAPE OF DECOMPOSITIONDICTIONARYLEARNING 8086 BY ANDREAS MÜLLER
  - FIXED THE IMPLEMENTATION OF EXPLAINEDVARIANCE INDECOMPOSITIONPCA DECOMPOSITION RANDOMIZEDPCA ANDDECOMPOSITIONINCREMENTALPCA 9105 BY HANMIN QIN
  - FIXED THE IMPLEMENTATION OF NOISEVARIANCE INDECOMPOSITIONPCA 9108 BY HANMIN QIN
  - FIXED A BUG WHERE CLUSTERDBSCAN GIVES INCORRECT RESULT WHEN INPUT IS A PRECOMPUTED SPARSE MATRIX WITH INITIAL ROWS ALL ZERO 8306 BY AKSHAY GUPTA
- 117 PREVIOUS RELEASES 71

SCIKITLEARN USER GUIDE RELEASE 0213

- FIX A BUG REGARDING FITTING CLUSTERKMEANS WITH A SPARSE ARRAY X AND INITIAL CENTROIDS WHERE X’S MEANS WERE UNNECESSARILY BEING SUBTRACTED FROM THE CENTROIDS 7872 BY JOSH KARNOFSKY
  - FIXES TO THE INPUT VALIDATION IN COVARIANCEELLIPTICENVELOPE 8086 BY ANDREAS MÜLLER
  - FIXED A BUG IN COVARIANCEMINCOVDET WHERE INPUTTING DATA THAT PRODUCED A SINGULAR COVARIANCE MATRIX WOULD CAUSE THE HELPER METHOD CSTEP TO THROW AN EXCEPTION 3367 BY JEREMY STEWARD
  - FIXED A BUG IN MANIFOLDTSNE AFFECTING CONVERGENCE OF THE GRADIENT DESCENT 8768 BY DAVID DETOMASO
  - FIXED A BUG IN MANIFOLDTSNE WHERE IT STORED THE INCORRECT KLDIVERGENCE 6507 BY SEBASTIAN SAEGER
  - FIXED IMPROPER SCALING IN CROSSDECOMPOSITIONPLSREGRESSION WITHSCALETRUE 7819 BY JAYZED82
  - CLUSTERBICLUSTERSPECTRALCOCLUSTERING AND CLUSTERBICLUSTER SPECTRALBICLUSTERING FIT METHOD CONFORMS WITH API BY ACCEPTING YAND RETURNING THE OBJECT 6126 7814 BY LAURENT DIRER AND MANITEJA NANDANA
  - FIX BUG WHERE MIXTURESAMPLE METHODS DID NOT RETURN AS MANY SAMPLES AS REQUESTED 7702 BY LEVI JOHN WOLF
  - FIXED THE SHRINKAGE IMPLEMENTATION IN NEIGHBORSNEARESTCENTROID 9219 BY HANMIN QIN
  - PREPROCESSING AND FEATURE SELECTION
  - FOR SPARSE MATRICES PREPROCESSINGNORMALIZE WITHRETURNNORMTRUE WILL NOW RAISE A NOTIMPLEMENTEDERROR WITH ‘L1’ OR ‘L2’ NORM AND WITH NORM ‘MAX’ THE NORMS RETURNED WILL BE THE SAME AS FOR DENSE MATRICES 7771 BY ANG LU
  - FIX A BUG WHERE FEATURESELECTIONSELECTFDR DID NOT EXACTLY IMPLEMENT BENJAMINIHOCHBERG PROCEDURE IT FORMERLY MAY HAVE SELECTED FEWER FEATURES THAN IT SHOULD 7490 BY PENG MENG
  - FIXED A BUG WHERE LINEARMODELRANDOMIZEDLASSO ANDLINEARMODELRANDOMIZEDLOGISTICREGRESSION BREAKS FOR SPARSE INPUT 8259 BY AMAN DALMIA
  - FIX A BUG WHERE FEATUREEXTRACTIONFEATUREHASHER MANDATORILY APPLIED A SPARSE RANDOM PROJECTION TO THE HASHED FEATURES PREVENTING THE USE OF FEATUREEXTRACTIONTEXTHASHINGVECTORIZER IN A PIPELINE WITH FEATUREEXTRACTIONTEXTTFIDFTRANSFORMER 7565 BY ROMAN YURCHAK
  - FIX A BUG WHERE FEATURESELECTIONMUTUALINFOREGRESSION DID NOT CORRECTLY USE NNEIGHBORS 8181 BY GUILLAUME LEMAITRE
  - MODEL EVALUATION AND METAESTIMATORS
  - FIXED A BUG WHERE MODELSELECTIONBASESEARCHCVINVERSETRANSFORM RETURNSSSELFBESTESTIMATORTRANSFORM INSTEAD OF SELFBESTESTIMATOR INVERSETRANSFORM 8344 BY AKSHAY GUPTA AND RASMUS ERIKSSON
  - ADDED CLASSES ATTRIBUTE TO MODELSELECTIONGRIDSEARCHCV MODELSELECTIONRANDOMIZEDSEARCHCV GRIDSEARCHGRIDSEARCHCV AND GRIDSEARCHRANDOMIZEDSEARCHCV THAT MATCHES THE CLASSES ATTRIBUTE OF BESTESTIMATOR 7661 AND 8295 BY ALYSSA BATULA DYLAN WERNERMEIER AND STEPHEN HOOVER
  - FIXED A BUG WHERE MODELSELECTIONVALIDATIONCURVE REUSED THE SAME ESTIMATOR FOR EACH PARAMETER VALUE 7365 BY ALEKSANDR SANDROVSKII
  - MODELSELECTIONPERMUTATIONTESTSCORE NOW WORKS WITH PANDAS TYPES 5697 BY STIJN TONK
  - SEVERAL FIXES TO INPUT VALIDATION IN MULTICLASSOUTPUTCODECLASSIFIER 8086 BY ANDREAS MÜLLER
  - MULTICLASSONEVSONECLASSIFIER ‘SPARTIALFIT’ NOW ENSURES ALL CLASSES ARE PROVIDED UPFRONT 6250 BY ASISH PANDA
- 72 CHAPTER 1 WELCOME TO SCIKITLEARN

SCIKITLEARN USER GUIDE RELEASE 0213

- FIXMULTIOUTPUTMULTIOUTPUTCLASSIFIERPREDICTPROBA TO RETURN A LIST OF 2D ARRAYS RATHER THAN A 3D ARRAY IN THE CASE WHERE DIFFERENT TARGET COLUMNS HAD DIFFERENT NUMBERS OF CLASSES A VALUEERROR WOULD BE RAISED ON TRYING TO STACK MATRICES WITH DIFFERENT DIMENSIONS 8093 BY PETER BULL
  - CROSS VALIDATION NOW WORKS WITH PANDAS DATATYPES THAT THAT HAVE A READONLY INDEX 9507 BY LOIC ESTEVE METRICS
  - METRICSAVERAGEPRECISIONSCORE NO LONGER LINEARLY INTERPOLATES BETWEEN OPERATING POINTS AND IN STEAD WEIGHS PRECISIONS BY THE CHANGE IN RECALL SINCE THE LAST OPERATING POINT AS PER THE WIKIPEDIA ENTRY 7356 BY NICK DINGWALL AND GAEL VAROQUAUX
  - FIX A BUG IN METRICSClassificationCHECKTARGETS WHICH WOULD RETURN BINARY IFYTRUE ANDYPRED WERE BOTHBINARY BUT THE UNION OF YTRUE ANDYPRED WASMULTICLASS 8377 BY LOIC ESTEVE
  - FIXED AN INTEGER OVERFLOW BUG IN METRICSCONFUSIONMATRIX AND HENCE METRICS COHENKAPPA SCORE 8354 7929 BY JOEL NOTHMAN AND JON CRALL
  - FIXED PASSING OF GAMMA PARAMETER TO THE CHI2 KERNEL INMETRICSPAIRWISEPAIRWISEKERNELS 5211 BY NICK RHINEHART SAURABH BANSOD AND ANDREAS MÜLLER
- MISCELLANEOUS
- FIXED A BUG WHEN DATASETSMAKECLASSIFICATION FAILS WHEN GENERATING MORE THAN 30 FEATURES 8159 BY HERILALAINA RAKOTOARISON
  - FIXED A BUG WHERE DATASETSMAKEMOONS GIVES AN INCORRECT RESULT WHEN NSAMPLES IS ODD 8198 BY JOSH LEVY
  - SOMEFETCH FUNCTIONS IN DATASETS WERE IGNORING THE DOWNLOADIFMISSING KEYWORD 7944 BY RALF GOMMERS
  - FIX ESTIMATORS TO ACCEPT A SAMPLEWEIGHT PARAMETER OF TYPE PANDASSERIES IN THEIRFIT FUNCTION 7825 BY KATHLEEN CHEN
  - FIX A BUG IN CASES WHERE NUMPYCUMSUM MAY BE NUMERICALLY UNSTABLE RAISING AN EXCEPTION IF INSTABILITY IS IDENTIFIED 7376 AND 7331 BY JOEL NOTHMAN AND YANGARBITER
  - FIX A BUG WHERE BASEBASEESTIMATORGETSTATE OBSTRUCTED PICKLING CUSTOMIZATIONS OF CHILD CLASSES WHEN USED IN A MULTIPLE INHERITANCE CONTEXT 8316 BY HOLGER PETERS
  - UPDATE SPHINXGALLERY FROM 014 TO 017 FOR RESOLVING LINKS IN DOCUMENTATION BUILD WITH SPHINX15 8010 7986 BY OSCAR NAJERA
  - ADDDATAHOME PARAMETER TO SKLEARNDATASETSFETCHKDDCUP99 9289 BY LOIC ESTEVE
  - FIX DATASET LOADERS USING PYTHON 3 VERSION OF MAKEDIRS TO ALSO WORK IN PYTHON 2 9284 BY SEBASTIN SANTY
  - SEVERAL MINOR ISSUES WERE FIXED WITH THANKS TO THE ALERTS OF LGTMCOMHTTPSLGTMCOM 9278 BY JEAN HELIE AMONG OTHERS

API CHANGES SUMMARY

TREES AND ENSEMBLES

- GRADIENT BOOSTING BASE MODELS ARE NO LONGER ESTIMATORS BY ANDREAS MÜLLER
  - ALL TREE BASED ESTIMATORS NOW ACCEPT A MINIMPURITYDECREASE PARAMETER IN LIEU OF THE MINIMPURITYSPLIT WHICH IS NOW DEPRECATED THE MINIMPURITYDECREASE HELPS STOP SPLITTING THE NODES IN WHICH THE WEIGHTED IMPURITY DECREASE FROM SPLITTING IS NO LONGER AT LEAST MINIMPURITYDECREASE 8449 BY RAGHAV RV
- 117 PREVIOUS RELEASES 73

SCIKITLEARN USER GUIDE RELEASE 0213

LINEAR KERNELIZED AND RELATED MODELS

- NITER PARAMETER IS DEPRECATED IN LINEARMODELSGDCLASSIFIER LINEARMODEL SGDREGRESSOR LINEARMODELPASSIVEAGGRESSIVECLASSIFIER LINEARMODEL PASSIVEAGGRESSIVEREGRESSOR ANDLINEARMODELPERCEPTRON BY TOM DUPRE LA TOUR

OTHER PREDICTORS

- NEIGHBORSLSHFOREST HAS BEEN DEPRECATED AND WILL BE REMOVED IN 021 DUE TO POOR PERFORMANCE 9078 BY LAURENT DIRER
- NEIGHBORSNEARESTCENTROID NO LONGER PURPORTS TO SUPPORT METRICPRECOMPUTED WHICH NOW RAISES AN ERROR 8515 BY SERGUL AYDORE
- THEALPHA PARAMETER OF SEMISUPERVISEDLABELPROPAGATION NOW HAS NO EFFECT AND IS DEPRECATED TO BE REMOVED IN 021 9239 BY ANDRE AMBROSIO BOECHAT UTKARSH UPADHYAY AND JOEL NOTHMAN

DECOMPOSITION MANIFOLD LEARNING AND CLUSTERING

- DEPRECATE THE DOCTOPICDISTR ARGUMENT OF THE PERPLEXITY METHOD IN DECOMPOSITION LATENTDIRICHLETALLOCATION BECAUSE THE USER NO LONGER HAS ACCESS TO THE UNNORMALIZED DOCUMENT TOPIC DISTRIBUTION NEEDED FOR THE PERPLEXITY CALCULATION 7954 BY GARY FOREMAN
- THENTOPICS PARAMETER OF DECOMPOSITIONLATENTDIRICHLETALLOCATION HAS BEEN RENAMED TO NCOMPONENTS AND WILL BE REMOVED IN VERSION 021 8922 BY ATTRACTADORE
- DECOMPOSITIONSPARSEPCATTRANSFORM 'SRIDGEALPHA PARAMETER IS DEPRECATED IN PREFERENCE FOR CLASS PARAMETER 8137 BY NAOYA KANAI
- CLUSTERDBSCAN NOW HAS AMETRICPARAMS PARAMETER 8139 BY NAOYA KANAI

PREPROCESSING AND FEATURE SELECTION

- FEATURESELECTIONSELECTFROMMODEL NOW HAS APARTIALFIT METHOD ONLY IF THE UNDERLYING ESTIMATOR DOES BY ANDREAS MÜLLER
- FEATURESELECTIONSELECTFROMMODEL NOW VALIDATES THE THRESHOLD PARAMETER AND SETS THE THRESHOLD ATTRIBUTE DURING THE CALL TO FIT AND NO LONGER DURING THE CALL TO TRANSFORM BY ANDREAS MÜLLER
- THENONNEGATIVE PARAMETER IN FEATUREEXTRACTIONFEATUREHASHER HAS BEEN DEPRECATED AND REPLACED WITH A MORE PRINCIPLED ALTERNATIVE ALTERNATESIGN 7565 BY ROMAN YURCHAK
- LINEARMODELRANDOMIZEDLOGISTICREGRESSION ANDLINEARMODELRANDOMIZEDLASSO HAVE BEEN DEPRECATED AND WILL BE REMOVED IN VERSION 021 8995 BY RAMANAS

MODEL EVALUATION AND METAESTIMATORS

- DEPRECATE THE FITPARAMS CONSTRUCTOR INPUT TO THE MODELSELECTIONGRIDSEARCHCV AND MODELSELECTIONRANDOMIZEDSEARCHCV IN FAVOR OF PASSING KEYWORD PARAMETERS TO THE FIT METHODS

OF THOSE CLASSES DATADEPENDENT PARAMETERS NEEDED FOR MODEL TRAINING SHOULD BE PASSED AS KEYWORD ARGUMENTS TOFIT AND CONFORMING TO THIS CONVENTION WILL ALLOW THE HYPERPARAMETER SELECTION CLASSES TO BE USED WITH TOOLS SUCH ASMODELSELECTIONCROSSVALPREDICT 2879 BY STEPHEN HOOVER

- IN VERSION 021 THE DEFAULT BEHAVIOR OF SPLITTERS THAT USE THE TESTSIZE ANDTRAINSIZESIZE PARAMETER WILL CHANGE SUCH THAT SPECIFYING TRAINSIZESIZE ALONE WILL CAUSE TESTSIZE TO BE THE REMAINDER 7459 BY NELSON LIU
- MULTICLASSONEVSRESTCLASSIFIER NOW HAS PARTIALFIT DECISIONFUNCTION AND PREDICTPROBA METHODS ONLY WHEN THE UNDERLYING ESTIMATOR DOES 7812 BY ANDREAS MÜLLER AND MIKHAIL KOROBOV
- MULTICLASSONEVSRESTCLASSIFIER NOW HAS APARTIALFIT METHOD ONLY IF THE UNDERLYING ESTIMATOR DOES BY ANDREAS MÜLLER

74 CHAPTER 1 WELCOME TO SCIKITLEARN



SCIKITLEARN USER GUIDE RELEASE 0213

- THEDECISIONFUNCTION OUTPUT SHAPE FOR BINARY CLASSIFICATION IN MULTICLASS ONEVSRESTCLASSIFIER ANDMULTICLASSONEVSONECLASSIFIER IS NOWNSAMPLES TO CONFORM TO SCIKITLEARN CONVENTIONS 9100 BY ANDREAS MÜLLER
  - THEMULTIOUTPUTMULTIOUTPUTCLASSIFIERPREDICTPROBA FUNCTION USED TO RETURN A 3D ARRAY NSAMPLES NCLASSES NOUTPUTS IN THE CASE WHERE DIFFERENT TARGET COLUMNS HAD DIFFERENT NUMBERS OF CLASSES A VALUEERROR WOULD BE RAISED ON TRYING TO STACK MATRICES WITH DIFFERENT DIMENSIONS THIS FUNCTION NOW RETURNS A LIST OF ARRAYS WHERE THE LENGTH OF THE LIST IS NOUTPUTS AND EACH ARRAY IS NSAMPLES NCLASSES FOR THAT PARTICULAR OUTPUT 8093 BY PETER BULL
  - REPLACE ATTRIBUTE NAMEDSTEPS DICT TOUTILSBUNCH INPIPELINEPIPELINE TO ENABLE TAB COMPLETION IN INTERACTIVE ENVIRONMENT IN THE CASE CONFLICT VALUE ON NAMEDSTEPS ANDDICT ATTRIBUTEDICT BEHAVIOR WILL BE PRIORITIZED 8481 BY HERILALAINA RAKOTOARISON
- MISCELLANEOUS
- DEPRECATE THE YPARAMETER IN TRANSFORM ANDINVERSETRANSFORM THE METHOD SHOULD NOT ACCEPT Y PARAMETER AS IT'S USED AT THE PREDICTION TIME 8174 BY TAHAR ZANOUDA ALEXANDRE GRAMFORT AND RAGHAV RV
  - SCIPY 0133 AND NUMPY 182 ARE NOW THE MINIMUM SUPPORTED VERSIONS FOR SCIKITLEARN THE FOLLOWING BACKPORTED FUNCTIONS IN UTILS HAVE BEEN REMOVED OR DEPRECATED ACCORDINGLY 8854 AND 8874 BY NAOYA KANAI
  - THESTORECOVARIANCES ANDCOVARIANCES PARAMETERS OF DISCRIMINANTANALYSIS QUADRATICDISCRIMINANTANALYSIS HAS BEEN RENAMED TO STORECOVARIANCE ANDCOVARIANCE TO BE CONSISTENT WITH THE CORRESPONDING PARAMETER NAMES OF THE DISCRIMINANTANALYSIS LINEARDISCRIMINANTANALYSIS THEY WILL BE REMOVED IN VERSION 021 7998 BY JIACHENG

REMOVED IN 019

- UTILSFIXESARGPARTITION
  - UTILSFIXESARRAYEQUAL
  - UTILSFIXESASTYPE
  - UTILSFIXESBINCOUNT
  - UTILSFIXESEXPT
  - UTILSFIXESFROMBUFFEREMPTY
  - UTILSFIXESIN1D
  - UTILSFIXESNORM
  - UTILSFIXESRANKDATA
  - UTILSFIXESSAFECOPY
- DEPRECATED IN 019 TO BE REMOVED IN 021
- UTILSARPACKEIGS
  - UTILSARPACKEIGSH
  - UTILSARPACKSVDS
  - UTILSEXTMATHFASTDOT
  - UTILSEXTMATHLOGSUMEXP
  - UTILSEXTMATHNORM
  - UTILSEXTMATHPINVH
  - UTILSGRAPHGRAPHLAPLACIAN
- 117 PREVIOUS RELEASES 75

SCIKITLEARN USER GUIDE RELEASE 0213

-UTILSRANDOMCHOICE

-UTILSSPARSETOOLSCONNECTEDCOMPONENTS

-UTILSSTATSRANKDATA

- ESTIMATORS WITH BOTH METHODS DECISIONFUNCTION ANDPREDICTPROBA ARE NOW REQUIRED TO HAVE A MONOTONIC RELATION BETWEEN THEM THE METHOD CHECKDECISIONPROBACONSISTENCY HAS BEEN ADDED
- INUTILSESTIMATORCHECKS TO CHECK THEIR CONSISTENCY 7578 BY SHUBHAM BHARDWAJ
- ALL CHECKS IN UTILSESTIMATORCHECKS IN PARTICULAR UTILSESTIMATORCHECKS

CHECKESTIMATOR NOW ACCEPT ESTIMATOR INSTANCES MOST OTHER CHECKS DO NOT ACCEPT ESTIMATOR CLASSES ANY MORE 9019 BY ANDREAS MÜLLER

- ENSURE THAT ESTIMATORS' ATTRIBUTES ENDING WITH ARE NOT SET IN THE CONSTRUCTOR BUT ONLY IN THE FIT METHOD

MOST NOTABLY ENSEMBLE ESTIMATORS DERIVING FROM ENSEMBLEBASEENSEMBLE NOW ONLY HAVE SELF ESTIMATORS AVAILABLE AFTER FIT 7464 BY LARS BUITINCK AND LOIC ESTEVE

CODE AND DOCUMENTATION CONTRIBUTORS

THANKS TO EVERYONE WHO HAS CONTRIBUTED TO THE MAINTENANCE AND IMPROVEMENT OF THE PROJECT SINCE VERSION 018 IN CLUDING

JOEL NOTHMAN LOIC ESTEVE ANDREAS MUELLER GUILLAUME LEMAITRE OLIVIER GRISEL HANMIN QIN RAGHAV RV ALEXANDRE GRAMFORT THEMRRMAX AMAN DALMIA GAEL VAROQUAUX NAOYA KANAI TOM DUPRÉ LA TOUR RISHIKESH NELSON LIU TAE HOON LEE NELLE VAROQUAUX AASHIL MIKHAIL KOROBV SEBASTIN SANTY JOAN MASSICH ROMAN YURCHAK RAKOTOARI SON HERILALAINA THIERRY GUILLEMOT ALEXANDRE ABADIE CAROL WILLING BALAKUMARAN MANOHARAN JOSH KARNOFSKY VLAD NICULAE UTKARSH UPADHYAY DMITRY PETROV MINGHUI LIU SRIVATSAN VINCENT PHAM ALBERT THOMAS JAKE VAN DERPLAS ATTRACTADORE JC LIU ALEXANDERCBOOTH CHKOAR ÓSCAR NÁJERA AARSHAY JAIN KYLE GILLIAM RAMANA SUBRA MANYAM CJ CAREY CLEMENT JOUDET DAVID ROBLES HE CHEN JORIS VAN DEN BOSSCHE KARAN DESAI KATIE LUANGKOTE LELAND MCINNES MANITEJA NANDANA MICHELE LACCHIA SERGEI LEBEDEV SHUBHAM BHARDWAJ AKSHAY0724 OMTCYFZ RICKIEPARK WATERPONEY VATHSALA ACHAR JBDELAFOSSÉ RALF GOMMERS EKATERINA KRIVICH VIVEK KUMAR ISHANK GULATI DAVE ELLIOTT LDIRER REIICHIRO NAKANO LEVI JOHN WOLF MATHIEU BLONDEL SID KAPUR DOUGAL J SUTHERLAND MIDINAS MIKEBENFIELD SOURAV SINGH ASEEM BANSAL IBRAIM GANIEV STEPHEN HOOVER AISHWARYARK STEVEN C HOWELL GARY FOREMAN NEERAJ GANGWAR TAHAR JON CRALL DOKATO KATHY CHEN FERRIA THOMAS MOREAU CHARLIE BRUMMITT NICOLAS GOIX ADAM KLECZEWSKI SAM SHLEIFER NIKITA SINGH BASIL BEIROUTI GIORGIO PATRINI MANOJ KUMAR RAFAEL POSSAS JAMES BOURBEAU JAMES A BEDNAR JANINE HARPER JAYE JEAN HELIE JEREMY STEWARD ARTSIOM JOHN WEI JONATHAN LIGO JONATHAN RAHN SEANWILLIAMS ARTHUR MENSCH JOSH LEVY JULIAN KUHLMANN JULIEN AUBERT JÖRN HEES KAI SHIVAMGARGSYA KAT HEMPSTALK KAUSHIK LAKSHMIKANTH KENNEDY KENNETH LYONS KENNETH MYERS KEVIN YAP KIR ILL BOBYREV KONSTANTIN PODSHUMOK ARTHUR IMBERT LEE MURRAY TOASTEDCORNFALAKES LERA LI LI ARTHUR DOUILLARD MAINAK JAS TOBYCHEESE MANRAJ SINGH MANVENDRA SINGH MARC MEKETON MARCOFALKE MATTHEW BRETT MATTHIAS GILCH MEHUL AHUJA MELANIE GOETZ MENG PENG MICHAEL DEZUBE MICHAL BAUMGARTNER VIBRANTABHI19 ARTEM GOLU BIN MILEN PASKOV ANTONIN CARETTE MORIKKO MRMJAUH NALEPA EMMANUEL NAMIYA ANTOINE WENDLINGER NARINE KOKHLYKAN NARINEK NATE GUERIN ANGUS WILLIAMS ANG LU NICOLE VAVROVA NITISH PANDEY OKHLOPKOV DANIIL OLEGOVICH ANDY CRAZE OM PRAKASH PARMINDER SINGH PATRICK CARLSON PATRICK PEI PAUL GANSSE PAULO HADDAD PAWEŁ LOREK PENG YU PETE BACHANT PETER BULL PETER CSIZSEK PETER WANG PIETER ARTHUR DE JONG PINGYAO CHANG PRESTON PARRY PUNEET MATHUR QUENTIN HIBON ANDREW SMITH ANDREW JACKSON 1KASTNER RAMESHWAR BHASKARAN RE BECCA BILBRO REMI RAMPIN ANDREA ESULI ROB HALL ROBERT BRADSHAW ROMAIN BRAULT AMAN PRATIK RUIFENG ZHENG RUSSELL SMITH SACHIN AGARWAL SAILESH CHOYAL SAMSON TAN SAMUËL WEBER SARAH BROWN SEBASTIAN PÖLSTERL SE BASTIAN RASCHKA SEBASTIAN SAAGER ALYSSA BATULA ABHYUDAY PRATAP SINGH SERGEY FELDMAN SERGUL AYDORE SHARAN YALBURGI WILLDUAN SIDDHARTH GUPTA SRI KRISHNA ALMER STIJN TONK ALLEN RIDDELL THEOFILOS PAPAPANAGIOTOU ALISON ALEXIS MIGNON TOMMY BOUCHER TOMMY LÖFSTEDT TOSHIHIRO KAMISHIMA TYLER FOLKMAN TYLER LANIGAN ALEXANDER JUNGE VARUN SHENOY VICTOR POUGHON VILHELM VON EHRENHEIM ALEKSANDR SANDROVSKII ALAN YEE VLASIOS VASILEIOU WARUT VIJITBENJARONK YANG ZHANG YAROSLAV HALCHENKO YICHUAN LIU YUICHI FUJIKAWA AFFANV14 AIVISION2020 XOR ANDREH7 BRADY SALZ CAMPUSTRAMPUS AGAMEMNON KRASOULIS DITENBERG ELENASHAROVA FILIPJ8 FUKATANI GEDECK GUIN IOL GUOCI HAKAA1 HONGKAHJUN IAMXHY JAKIRKHAM JAROSLAWWEBER JAYZED82 JEROKO JMONTOYAM JONATHANSTRIEBEL

76 CHAPTER 1 WELCOME TO SCIKITLEARN

SCIKITLEARN USER GUIDE RELEASE 0213  
JOSEPHSALMON JSCHENDEL LEEREEVES MARTINHAHN MATHURINM MEHAKSACHDEVA MLEWIS1729 MLLIOU112 MTHORRELL  
NDINGWALL NUFFE YANGARBITER PLAGREE PLDTC325 BRENO FREITAS BRETT OLSEN BRIAN A ALFANO BRIAN BURNS POLMAURI  
BRANDON CARTER CHARLTON AUSTIN CHAYANT T15H CHINMAYA PANCHOLI CHRISTIAN DANIELSEN CHUNG YEN CHYIKWEI  
YAU PRAVARMHAJAN DOHMATOB ELVIS DANIEL LEJEUNE DANIEL HNYK DARIUS MORAWIEC DAVID DETOMASO DAVID  
GASQUEZ DAVID HABERTHÜR DAVID HERYANTO DAVID KIRKBY DAVID NICHOLSON RASHCHEDRIN DEBORAH GERTRUDE DIGGES  
DENIS ENGEMANN DEVANSH D DICKSON BOB BAXLEY DON86 E LYNCHKLARUP ED ROGERS ELIZABETH FERRISS ELLEN  
CO2 FABIAN EGLI FANGCHIEH CHOU BING TIAN DAI GREG STUPP GRZEGORZ SZPAK BERTRAND THIRION HADRIEN BERTRAND  
HARIZO RAJAONA ZXCVBNIUS HENRY LIN HOLGER PETERS ICYBLADE DAI IGOR ANDRIUSHCHENKO ILYA ISAAC LAUGHLIN IVÁN  
VALLÉS AURÉLIEN BELLET JPFRANCOIA JACOB SCHREIBER ASISH MAHAPATRA  
1174 VERSION 0182  
JUNE 20 2017  
LAST RELEASE WITH PYTHON 26 SUPPORT  
SCIKITLEARN 018 IS THE LAST MAJOR RELEASE OF SCIKITLEARN TO SUPPORT PYTHON 26 LATER VERSIONS OF SCIKITLEARN WILL  
REQUIRE PYTHON 27 OR ABOVE  
CHANGELOG  
• FIXES FOR COMPATIBILITY WITH NUMPY 1130 7946 8355 BY LOIC ESTEVE  
• MINOR COMPATIBILITY CHANGES IN THE EXAMPLES 9010 8040 9149  
CODE CONTRIBUTORS  
AMAN DALMIA LOIC ESTEVE NATE GUERIN SERGEI LEBEDEV  
1175 VERSION 0181  
NOVEMBER 11 2016  
CHANGELOG  
ENHANCEMENTS  
• IMPROVED SAMPLEWITHOUTREPLACEMENT SPEED BY UTILIZING NUMPYRANDOMPERMUTATION FOR MOST CASES  
AS A RESULT SAMPLES MAY DIFFER IN THIS RELEASE FOR A FIXED RANDOM STATE AFFECTED ESTIMATORS  
-ENSEMBLEBAGGINGCLASSIFIER  
-ENSEMBLEBAGGINGREGRESSOR  
-LINEARMODELRANSACREGRESSOR  
-MODELSELECTIONRANDOMIZEDSEARCHCV  
-RANDOMPROJECTIONSPPARSERANDOMPROJECTION  
THIS ALSO AFFECTS THE DATASETSMAKECLASSIFICATION METHOD  
117 PREVIOUS RELEASES 77

SCIKITLEARN USER GUIDE RELEASE 0213

BUG FIXES

- FIX ISSUE WHERE MINGRADNORM ANDNITERWITHOUTPROGRESS PARAMETERS WERE NOT BEING UTILISED BYMANIFOLDTSNE 6497 BY SEBASTIAN SÄGER
  - FIX BUG FOR SVM’S DECISION VALUES WHEN DECISIONFUNCTIONSHAPE ISOVR INSVMSVC SVM SVC ’S DECISIONFUNCTION WAS INCORRECT FROM VERSIONS 0170 THROUGH 0180 7724 BY BING TIAN DAI
  - ATTRIBUTE EXPLAINEDVARIANCERATIO OF DISCRIMINANTANALYSIS LINEARDISCRIMINANTANALYSIS CALCULATED WITH SVD AND EIGEN SOLVER ARE NOW OF THE SAME LENGTH 7632 BY JPFRANCOIA
  - FIXES ISSUE IN UNIVARIATE FEATURE SELECTION WHERE SCORE FUNCTIONS WERE NOT ACCEPTING MULTILABEL TARGETS 7676 BY MOHAMMED AFFAN
  - FIXED SETTING PARAMETERS WHEN CALLING FIT MULTIPLE TIMES ON FEATURESELECTIONSELECTFROMMODEL 7756 BY ANDREAS MÜLLER
  - FIXES ISSUE IN PARTIALFIT METHOD OFMULTICLASSONEVSRESTCLASSIFIER WHEN NUMBER OF CLASSES USED INPARTIALFIT WAS LESS THAN THE TOTAL NUMBER OF CLASSES IN THE DATA 7786 BY SRIVATSAN RAMESH
  - FIXES ISSUE IN CALIBRATIONCALIBRATEDCLASSIFIERCV WHERE THE SUM OF PROBABILITIES OF EACH CLASS FOR A DATA WAS NOT 1 AND CALIBRATEDCLASSIFIERCV NOW HANDLES THE CASE WHERE THE TRAINING SET HAS LESS NUMBER OF CLASSES THAN THE TOTAL DATA 7799 BY SRIVATSAN RAMESH
  - FIX A BUG WHERE SKLEARNFEATURESELECTIONSELECTFDR DID NOT EXACTLY IMPLEMENT BENJAMINI HOCHBERG PROCEDURE IT FORMERLY MAY HAVE SELECTED FEWER FEATURES THAN IT SHOULD 7490 BY PENG MENG
  - SKLEARNMANIFOLDLOCALLYLINEAREMBEDDING NOW CORRECTLY HANDLES INTEGER INPUTS 6282 BY JAKE VANDERPLAS
  - THEM INWEIGHTFRACTIONLEAF PARAMETER OF TREEBASED CLASSIFIERS AND REGRESSORS NOW ASSUMES UNIFORM SAMPLE WEIGHTS BY DEFAULT IF THE SAMPLEWEIGHT ARGUMENT IS NOT PASSED TO THE FIT FUNCTION PREVIOUSLY THE PARAMETER WAS SILENTLY IGNORED 7301 BY NELSON LIU
  - NUMERICAL ISSUE WITH LINEARMODELRIDGE CV ON CENTERED DATA WHEN NFEATURES NSAMPLES 6178 BY BERTRAND THIRION
  - TREE SPLITTING CRITERION CLASSES’ CLONINGPICKLING IS NOW MEMORY SAFE 7680 BY IBRAIM GANIEV
  - FIXED A BUG WHERE DECOMPOSITIONNMF SETS ITSNITERS ATTRIBUTE IN TRANSFORM 7553 BY EKATE RINA KRIVICH
  - SKLEARNLINEARMODELLOGISTICREGRESSIONCV NOW CORRECTLY HANDLES STRING LABELS 5874 BY RAGHAV RV
  - FIXED A BUG WHERE SKLEARNMODELSELECTIONTRAINTESTSPLIT RAISED AN ERROR WHEN STRATIFY IS A LIST OF STRING LABELS 7593 BY RAGHAV RV
  - FIXED A BUG WHERE SKLEARNMODELSELECTIONGRIDSEARCHCV ANDSKLEARN MODELSELECTIONRANDOMIZEDSEARCHCV WERE NOT PICKLEABLE BECAUSE OF A PICKLING BUG IN NP MAMASKEDARRAY 7594 BY RAGHAV RV
  - ALL CROSSVALIDATION UTILITIES IN SKLEARNMODELSELECTION NOW PERMIT ONE TIME CROSSVALIDATION SPLITTERS FOR THECVPARAMETER ALSO NONDETERMINISTIC CROSSVALIDATION SPLITTERS WHERE MULTIPLE CALLS TO SPLIT PRODUCE DISSIMILAR SPLITS CAN BE USED AS CVPARAMETER THE SKLEARNMODELSELECTIONGRIDSEARCHCV WILL CROSSVALIDATE EACH PARAMETER SETTING ON THE SPLIT PRODUCED BY THE FIRST SPLIT CALL TO THE CROSSVALIDATION SPLITTER 7660 BY RAGHAV RV
  - FIX BUG WHERE PREPROCESSINGMULTILABELBINARIZERFITTRANSFORM RETURNED AN INVALID CSR MATRIX 7750 BY CJ CAREY
- 78 CHAPTER 1 WELCOME TO SCIKITLEARN

SCIKITLEARN USER GUIDE RELEASE 0213

- FIXED A BUG WHERE METRICSPAIRWISECOSINEDISTANCES COULD RETURN A SMALL NEGATIVE DISTANCE 7732 BY ARTSION

API CHANGES SUMMARY

TREES AND FORESTS

- THEMINWEIGHTFRACTIONLEAF PARAMETER OF TREEBASED CLASSIFIERS AND REGRESSORS NOW ASSUMES UNIFORM SAMPLE WEIGHTS BY DEFAULT IF THE SAMPLEWEIGHT ARGUMENT IS NOT PASSED TO THE FIT FUNCTION PREVIOUSLY THE PARAMETER WAS SILENTLY IGNORED 7301 BY NELSON LIU

- TREE SPLITTING CRITERION CLASSES' CLONINGPICKLING IS NOW MEMORY SAFE 7680 BY IBRAIM GANIEV

LINEAR KERNELIZED AND RELATED MODELS

- LENGTH OF EXPLAINEDVARIANCERATIO OFDISCRIMINANTANALYSIS LINEARDISCRIMINANTANALYSIS CHANGED FOR BOTH EIGEN AND SVD SOLVERS THE ATTRIBUTE HAS NOW A LENGTH OF MINNCOMPONENTS NCLASSES 1 7632 BY JPFRANCOIA
- NUMERICAL ISSUE WITH LINEARMODELRIDGECV ON CENTERED DATA WHEN NFEATURES NSAMPLES 6178 BY BERTRAND THIRION

1176 VERSION 018

SEPTEMBER 28 2016

LAST RELEASE WITH PYTHON 26 SUPPORT

SCIKITLEARN 018 WILL BE THE LAST VERSION OF SCIKITLEARN TO SUPPORT PYTHON 26 LATER VERSIONS OF SCIKITLEARN WILL REQUIRE PYTHON 27 OR ABOVE

MODEL SELECTION ENHANCEMENTS AND API CHANGES

- THE MODELSELECTION MODULE

THE NEW MODULE SKLEARNMODELSELECTION WHICH GROUPS TOGETHER THE FUNCTIONALITIES OF FORMERLY SKLEARNCROSSVALIDATION SKLEARNGRIDSEARCH ANDSKLEARNLEARNINGCURVE INTRO DUCES NEW POSSIBILITIES SUCH AS NESTED CROSSVALIDATION AND BETTER MANIPULATION OF PARAMETER SEARCHES WITH PAN DAS

MANY THINGS WILL STAY THE SAME BUT THERE ARE SOME KEY DIFFERENCES READ BELOW TO KNOW MORE ABOUT THE CHANGES

- DATAINDEPENDENT CV SPLITTERS ENABLING NESTED CROSSVALIDATION

THE NEW CROSSVALIDATION SPLITTERS DEFINED IN THE SKLEARNMODELSELECTION ARE NO LONGER INITIALIZED WITH ANY DATADEPENDENT PARAMETERS SUCH AS Y INSTEAD THEY EXPOSE A SPLIT METHOD THAT TAKES IN THE DATA AND YIELDS A GENERATOR FOR THE DIFFERENT SPLITS

THIS CHANGE MAKES IT POSSIBLE TO USE THE CROSSVALIDATION SPLITTERS TO PERFORM NESTED CROSSVALIDATION FACILITATED BYMODELSELECTIONGRIDSEARCHCV ANDMODELSELECTIONRANDOMIZEDSEARCHCV UTILITIES

- THE ENHANCED CVRESULTS ATTRIBUTE

THE NEWCVRESULTS ATTRIBUTE OF MODELSELECTIONGRIDSEARCHCV ANDMODELSELECTION RANDOMIZEDSEARCHCV INTRODUCED IN LIEU OF THE GRIDSCORES ATTRIBUTE IS A DICT OF 1D ARRAYS WITH ELEMENTS IN EACH ARRAY CORRESPONDING TO THE PARAMETER SETTINGS IE SEARCH CANDIDATES

117 PREVIOUS RELEASES 79

SCIKITLEARN USER GUIDE RELEASE 0213

THECVRESULTS DICT CAN BE EASILY IMPORTED INTO PANDAS AS A DATAFRAME FOR EXPLORING THE SEARCH RESULTS

THECVRESULTS ARRAYS INCLUDE SCORES FOR EACH CROSSVALIDATION SPLIT WITH KEYS SUCH AS

SPLIT0TESTSCORE AS WELL AS THEIR MEAN MEANTESTSCORE AND STANDARD DEVIATION

STDTESTSCORE

THE RANKS FOR THE SEARCH CANDIDATES BASED ON THEIR MEAN CROSSVALIDATION SCORE IS AVAILABLE AT

CVRESULTSRANKTESTSCORE

THE PARAMETER VALUES FOR EACH PARAMETER IS STORED SEPARATELY AS NUMPY MASKED OBJECT ARRAYS THE VALUE FOR

THAT SEARCH CANDIDATE IS MASKED IF THE CORRESPONDING PARAMETER IS NOT APPLICABLE ADDITIONALLY A LIST OF ALL THE

PARAMETER DICTS ARE STORED AT CVRESULTSPARAMS

- PARAMETERS NFOLDS AND NITER RENAMED TO NSPLITS

SOME PARAMETER NAMES HAVE CHANGED THE NFOLDS PARAMETER IN NEW MODELSELECTIONKFOLD

MODELSELECTIONGROUPKFOLD SEE BELOW FOR THE NAME CHANGE AND MODELSELECTION

STRATIFIEDKFOLD IS NOW RENAMED TO NSPLITS THENITER PARAMETER IN MODELSELECTION

SHUFFLESPLIT THE NEW CLASS MODELSELECTIONGROUPSHUFFLESPLIT ANDMODELSELECTION

STRATIFIEDSHUFFLESPLIT IS NOW RENAMED TO NSPLITS

- RENAME OF SPLITTER CLASSES WHICH ACCEPTS GROUP LABELS ALONG WITH DATA

THE CROSSVALIDATION SPLITTERS LABELKFOLD LABELSHUFFLESPLIT LEAVEONELABELOUT AND

LEAVEPLABELOUT HAVE BEEN RENAMED TO MODELSELECTIONGROUPKFOLD MODELSELECTION

GROUPSHUFFLESPLIT MODELSELECTIONLEAVEONEGROUPOUT ANDMODELSELECTION

LEAVEPGROUPSOUT RESPECTIVELY

NOTE THE CHANGE FROM SINGULAR TO PLURAL FORM IN MODELSELECTIONLEAVEPGROUPSOUT

- FIT PARAMETER LABELS RENAMED TO GROUPS

THELABELS PARAMETER IN THE SPLIT METHOD OF THE NEWLY RENAMED SPLITTERS MODELSELECTION

GROUPKFOLD MODELSELECTIONLEAVEONEGROUPOUT MODELSELECTION

LEAVEPGROUPSOUT MODELSELECTIONGROUPSHUFFLESPLIT IS RENAMED TO GROUPS FOLLOWING THE

NEW NOMENCLATURE OF THEIR CLASS NAMES

- PARAMETER NLABELS RENAMED TO NGROUPS

THE PARAMETER NLABELS IN THE NEWLY RENAMED MODELSELECTIONLEAVEPGROUPSOUT IS CHANGED TO

NGROUPS

- TRAINING SCORES AND TIMING INFORMATION

CVRESULTS ALSO INCLUDES THE TRAINING SCORES FOR EACH CROSSVALIDATION SPLIT WITH KEYS SUCH

ASSPLIT0TRAINSCORE AS WELL AS THEIR MEAN MEANTRAINSCORE AND STAN

DARD DEVIATION STDTRAINSCORE TO AVOID THE COST OF EVALUATING TRAINING SCORE SET

RETURNTRAINSCOREFALSE

ADDITIONALLY THE MEAN AND STANDARD DEVIATION OF THE TIMES TAKEN TO SPLIT TRAIN AND SCORE THE MODEL ACROSS ALL THE

CROSSVALIDATION SPLITS IS AVAILABLE AT THE KEY MEANTIME ANDSTDTIME RESPECTIVELY

CHANGELOG

NEW FEATURES

CLASSIFIERS AND REGRESSORS

- THE GAUSSIAN PROCESS MODULE HAS BEEN REIMPLEMENTED AND NOW OFFERS CLASSIFICATION AND REGRESSION ESTI

MATORS THROUGH GAUSSIANPROCESSGAUSSIANPROCESSCLASSIFIER ANDGAUSSIANPROCESS

GAUSSIANPROCESSREGRESSOR AMONG OTHER THINGS THE NEW IMPLEMENTATION SUPPORTS KERNEL ENGINEERING

80 CHAPTER 1 WELCOME TO SCIKITLEARN

SCIKITLEARN USER GUIDE RELEASE 0213

GRADIENTBASED HYPERPARAMETER OPTIMIZATION OR SAMPLING OF FUNCTIONS FROM GP PRIOR AND GP POSTERIOR EXTENSIVE DOCUMENTATION AND EXAMPLES ARE PROVIDED BY JAN HENDRIK METZEN

- ADDED NEW SUPERVISED LEARNING ALGORITHM MULTILAYER PERCEPTRON 3204 BY ISSAM H LARADJI
- ADDEDLINEARMODELHUBERREGRESSOR A LINEAR MODEL ROBUST TO OUTLIERS 5291 BY MANOJ KUMAR
- ADDED THE MULTIOUTPUTMULTIOUTPUTREGRESSOR METAESTIMATOR IT CONVERTS SINGLE OUTPUT REGRESSORS TO MULTIOUTPUT REGRESSORS BY FITTING ONE REGRESSOR PER OUTPUT BY TIM HEAD

OTHER ESTIMATORS

- NEWMIXTUREGAUSSIANNMIXTURE ANDMIXTUREBAYESIANGAUSSIANNMIXTURE REPLACE FORMER MIXTURE MODELS EMPLOYING FASTER INFERENCE FOR SOUNDER RESULTS 7295 BY WEI XUE AND THIERRY GUILLEMOT
- CLASSDECOMPOSITIONRANDOMIZEDPCA IS NOW FACTORED INTO DECOMPOSITIONPCA AND IT IS AVAIL

ABLE CALLING WITH PARAMETER SVDSOLVERRANDOMIZED THE DEFAULT NUMBER OF NITER FOR RANDOMIZED HAS CHANGED TO 4 THE OLD BEHAVIOR OF PCA IS RECOVERED BY SVDSOLVERFULL AN ADDITIONAL SOLVER CALLS ARPACK AND PERFORMS TRUNCATED NONRANDOMIZED SVD BY DEFAULT THE BEST SOLVER IS SELECTED DEPENDING ON THE SIZE OF THE INPUT AND THE NUMBER OF COMPONENTS REQUESTED 5299 BY GIORGIO PATRINI

- ADDED TWO FUNCTIONS FOR MUTUAL INFORMATION ESTIMATION FEATURESELECTION

MUTUALINFOCLASSIF ANDFEATURESELECTIONMUTUALINFOREGRESSION THESE FUNCTIONS CAN BE USED IN FEATURESELECTIONSELECTKBEST ANDFEATURESELECTION SELECTPERCENTILE AS SCORE FUNCTIONS BY ANDREA BRAVI AND NIKOLAY MAYOROV

- ADDED THE ENSEMBLEISOLATIONFOREST CLASS FOR ANOMALY DETECTION BASED ON RANDOM FORESTS BY NICOLAS GOIX
- ADDEDALGORITHMELKAN TOCLUSTERKMEANS IMPLEMENTING ELKAN’S FAST KMEANS ALGORITHM BY ANDREAS MÜLLER

MODEL SELECTION AND EVALUATION

- ADDEDMETRICSCUSTERFOWLKESMALLOWSSCORE THE FOWLKES MALLOWS INDEX WHICH MEASURES THE SIMILARITY OF TWO CLUSTERINGS OF A SET OF POINTS BY ARNAUD FOUCHE AND THIERRY GUILLEMOT
- ADDEDMETRICSCALINSKI HARABAZSCORE WHICH COMPUTES THE CALINSKI AND HARABAZ SCORE TO EVALUATE THE RESULTING CLUSTERING OF A SET OF POINTS BY ARNAUD FOUCHE AND THIERRY GUILLEMOT
- ADDED NEW CROSSVALIDATION SPLITTER MODELSELECTIONTIMESERIESSPLIT TO HANDLE TIME SERIES DATA 6586 BY YENCHEN LIN
- THE CROSSVALIDATION ITERATORS ARE REPLACED BY CROSSVALIDATION SPLITTERS AVAILABLE FROM SKLEARN MODELSELECTION ALLOWING FOR NESTED CROSSVALIDATION SEE MODEL SELECTION ENHANCEMENTS AND API CHANGES FOR MORE INFORMATION 4294 BY RAGHAV RV

ENHANCEMENTS

TREES AND ENSEMBLES

- ADDED A NEW SPLITTING CRITERION FOR TREEDECISIONTREEREGRESSOR THE MEAN ABSOLUTE ERROR THIS CRITERION CAN ALSO BE USED IN ENSEMBLEEXTRATREESREGRESSOR ENSEMBLE RANDOMFORESTREGRESSOR AND THE GRADIENT BOOSTING ESTIMATORS 6667 BY NELSON LIU
- ADDED WEIGHTED IMPURITYBASED EARLY STOPPING CRITERION FOR DECISION TREE GROWTH 6954 BY NELSON LIU
- THE RANDOM FOREST EXTRA TREE AND DECISION TREE ESTIMATORS NOW HAS A METHOD DECISIONPATH WHICH RETURNS THE DECISION PATH OF SAMPLES IN THE TREE BY ARNAUD JOLY
- A NEW EXAMPLE HAS BEEN ADDED UNVEILING THE DECISION TREE STRUCTURE BY ARNAUD JOLY

117 PREVIOUS RELEASES 81

SCIKITLEARN USER GUIDE RELEASE 0213

- RANDOM FOREST EXTRA TREES DECISION TREES AND GRADIENT BOOSTING ESTIMATOR ACCEPT THE PARAMETER MINSAMPLESSPLIT ANDMINSAMPLESLEAF PROVIDED AS A PERCENTAGE OF THE TRAINING SAMPLES BY YELITE AND ARNAUD JOLY
  - GRADIENT BOOSTING ESTIMATORS ACCEPT THE PARAMETER CRITERION TO SPECIFY TO SPLITTING CRITERION USED IN BUILT DECISION TREES 6667 BY NELSON LIU
  - THE MEMORY FOOTPRINT IS REDUCED SOMETIMES GREATLY FOR ENSEMBLEBAGGINGBASEBAGGING AND CLASSES THAT INHERIT FROM IT IE ENSEMBLEBAGGINGCLASSIFIER ENSEMBLEBAGGINGREGRESSOR AND ENSEMBLEISOLATIONFOREST BY DYNAMICALLY GENERATING ATTRIBUTE ESTIMATORSSAMPLES ONLY WHEN IT IS NEEDED BY DAVID STAUB
  - ADDEDNJOBS ANDSAMPLEWEIGHT PARAMETERS FOR ENSEMBLEVOTINGCLASSIFIER TO FIT UNDERLYING ESTIMATORS IN PARALLEL 5805 BY IBRAIM GANIEV
  - LINEAR KERNELIZED AND RELATED MODELS
  - INLINEARMODELLOGISTICREGRESSION THE SAG SOLVER IS NOW AVAILABLE IN THE MULTINOMIAL CASE 5251 BY TOM DUPRE LA TOUR
  - LINEARMODELRANSACREGRESSOR SVMLINEARSVC ANDSVMLINEARSVR NOW SUPPORT SAMPLEWEIGHT BY IMACULATE
  - ADD PARAMETER LOSS TOLINEARMODELRANSACREGRESSOR TO MEASURE THE ERROR ON THE SAMPLES FOR EVERY TRIAL BY MANOJ KUMAR
  - PREDICTION OF OUTFOSAMPLE EVENTS WITH ISOTONIC REGRESSION ISOTONICISOTONICREGRESSION IS NOW MUCH FASTER OVER 1000X IN TESTS WITH SYNTHETIC DATA BY JONATHAN ARFA
  - ISOTONIC REGRESSION ISOTONICISOTONICREGRESSION NOW USES A BETTER ALGORITHM TO AVOID ON2 BEHAVIOR IN PATHOLOGICAL CASES AND IS ALSO GENERALLY FASTER 6691 BY ANTONY LEE
  - NAIVEBAYESGAUSSIANNB NOW ACCEPTS DATAINDEPENDENT CLASSPRIORS THROUGH THE PARAMETER PRIORS BY GUILLAUME LEMAITRE
  - LINEARMODELELASTICNET ANDLINEARMODELLASSO NOW WORKS WITH NPFLOAT32 INPUT DATA WITHOUT CONVERTING IT INTO NPFLOAT64 THIS ALLOWS TO REDUCE THE MEMORY CONSUMPTION 6913 BY YENCHEN LIN
  - SEMISUPERVISEDLABELPROPAGATION ANDSEMISUPERVISEDLABELSPREADING NOW ACCEPT ARBITRARY KERNEL FUNCTIONS IN ADDITION TO STRINGS KNN ANDRBF 5762 BY UTKARSH UPADHYAY
  - DECOMPOSITION MANIFOLD LEARNING AND CLUSTERING
  - ADDEDINVERSETRANSFORM FUNCTION TO DECOMPOSITIONNMF TO COMPUTE DATA MATRIX OF ORIGINAL SHAPE BY ANISH SHAH
  - CLUSTERKMEANS ANDCLUSTERMINIBATCHKMEANS NOW WORKS WITH NPFLOAT32 ANDNP FLOAT64 INPUT DATA WITHOUT CONVERTING IT THIS ALLOWS TO REDUCE THE MEMORY CONSUMPTION BY USING NP FLOAT32 6846 BY SEBASTIAN SÄGER AND YENCHEN LIN
  - PREPROCESSING AND FEATURE SELECTION
  - PREPROCESSINGROBUSTSCALER NOW ACCEPTS QUANTILERANGE PARAMETER 5929 BY KONSTANTIN POD SHUMOK
  - FEATUREEXTRACTIONFEATUREHASHER NOW ACCEPTS STRING VALUES 6173 BY RYAD ZENINE AND DEVASHISH DESHPANDE
  - KEYWORD ARGUMENTS CAN NOW BE SUPPLIED TO FUNC INPREPROCESSINGFUNCTIONTRANSFORMER BY MEANS OF THE KWARGS PARAMETER BY BRIAN MCFEE
  - FEATURESELECTIONSELECTKBEST ANDFEATURESELECTIONSELECTPERCENTILE NOW ACCEPT SCORE FUNCTIONS THAT TAKE X Y AS INPUT AND RETURN ONLY THE SCORES BY NIKOLAY MAYOROV
- 82 CHAPTER 1 WELCOME TO SCIKITLEARN



SCIKITLEARN USER GUIDE RELEASE 0213

MODEL EVALUATION AND METAESTIMATORS

- MULTICLASSONEVSONECLASSIFIER ANDMULTICLASSONEVSRESTCLASSIFIER NOW SUPPORT PARTIALFIT BY ASISH PANDA AND PHILIPP DOWLING
- ADDED SUPPORT FOR SUBSTITUTING OR DISABLING PIPELINEPIPELINE ANDPIPELINEFEATUREUNION COMPONENTS USING THE SETPARAMS INTERFACE THAT POWERS SKLEARNGRIDSEARCH SEE SELECTING DIMENSIONALITY REDUCTION WITH PIPELINE AND GRIDSEARCHCV BY JOEL NOTHMAN AND ROBERT MCGIBBON
- THE NEW CVRESULTS ATTRIBUTE OF MODELSELECTIONGRIDSEARCHCV ANDMODELSELECTIONRANDOMIZEDSEARCHCV CAN BE EASILY IMPORTED INTO PANDAS AS A DATAFRAME REF MODEL SELECTION ENHANCEMENTS AND API CHANGES FOR MORE INFORMATION 6697 BY RAGHAV RV
- GENERALIZATION OF MODELSELECTIONCROSSVALPREDICT ONE CAN PASS METHOD NAMES SUCH AS PREDICTION TO BE USED IN THE CROSS VALIDATION FRAMEWORK INSTEAD OF THE DEFAULT PREDICT BY ORI ZIV AND SEARS MERRITT
- THE TRAINING SCORES AND TIME TAKEN FOR TRAINING FOLLOWED BY SCORING FOR EACH SEARCH CANDIDATE ARE NOW AVAILABLE AT THECVRESULTS DICT SEE MODEL SELECTION ENHANCEMENTS AND API CHANGES FOR MORE INFORMATION 7325 BY EUGENE CHEN AND RAGHAV RV

METRICS

- ADDEDLABELS FLAG TOMETRICSLoss TO EXPLICITLY PROVIDE THE LABELS WHEN THE NUMBER OF CLASSES IN YTRUE ANDYPRED DIFFER 7239 BY HONG GUANGGUO WITH HELP FROM MADSENSEN AND NELSON LIU
- SUPPORT SPARSE CONTINGENCY MATRICES IN CLUSTER EVALUATION METRICSClustersUPERVISED TO SCALE TO A LARGE NUMBER OF CLUSTERS 7419 BY GREGORY STUPP AND JOEL NOTHMAN
- ADDSAMPLEWEIGHT PARAMETER TO METRICSMATTHEWSCORRcoef BY JATIN SHAH AND RAGHAV RV
- SPEED UP METRICSSILHOUETTEScore BY USING VECTORIZED OPERATIONS BY MANOJ KUMAR
- ADDSAMPLEWEIGHT PARAMETER TO METRICSCONFUSIONMATRIX BY BERNARDO STEIN

MISCELLANEOUS

- ADDEDNJOBS PARAMETER TO FEATURESELECTIONRFECV TO COMPUTE THE SCORE ON THE TEST FOLDS IN PARALLEL BY MANOJ KUMAR
- CODEBASE DOES NOT CONTAIN CC CYTHON GENERATED FILES THEY ARE GENERATED DURING BUILD DISTRIBUTION PACKAGES WILL STILL CONTAIN GENERATED CC FILES BY ARTHUR MENSCH
- REDUCE THE MEMORY USAGE FOR 32BIT FLOAT INPUT ARRAYS OF UTILSSPARSEFUNCMEANVARIANCEaxis ANDUTILSSPARSEFUNCINCRMEANVARIANCEaxis BY SUPPORTING CYTHON FUSED TYPES BY YENCHEN LIN
- THEIGNOREWARNINGS NOW ACCEPT A CATEGORY ARGUMENT TO IGNORE ONLY THE WARNINGS OF A SPECIFIED TYPE BY THIERRY GUILLEMOT
- ADDED PARAMETER RETURNXY AND RETURN TYPE DATA TARGET TUPLE OPTION TOLOADIRIS DATASET 7049 LOADBREASTCANCER DATASET 7152 LOADDIGITS DATASETLLOADDIABETES DATASETLLOADLINNERUD DATASETLLOADBOSTON DATASET 7154 BY MANVENDRA SINGH
- SIMPLIFICATION OF THE CLONE FUNCTION DEPRECATE SUPPORT FOR ESTIMATORS THAT MODIFY PARAMETERS IN INIT 5540 BY ANDREAS MÜLLER
- WHEN UNPICKLING A SCIKITLEARN ESTIMATOR IN A DIFFERENT VERSION THAN THE ONE THE ESTIMATOR WAS TRAINED WITH A USERWARNING IS RAISED SEE THE DOCUMENTATION ON MODEL PERSISTENCE FOR MORE DETAILS 7248 BY ANDREAS MÜLLER

117 PREVIOUS RELEASES 83

SCIKITLEARN USER GUIDE RELEASE 0213

BUG FIXES

TREES AND ENSEMBLES

- RANDOM FOREST EXTRA TREES DECISION TREES AND GRADIENT BOOSTING WON'T ACCEPT ANYMORE MINSAMPLESPLIT1 AS AT LEAST 2 SAMPLES ARE REQUIRED TO SPLIT A DECISION TREE NODE BY ARNAUD JOLY

- ENSEMBLEVOTINGCLASSIFIER NOW RAISES NOTFITTEDERROR IFPREDICT TRANSFORM OR PREDICTPROBA ARE CALLED ON THE NONFITTED ESTIMATOR BY SEBASTIAN RASCHKA

- FIX BUG WHERE ENSEMBLEADABOOSTCLASSIFIER ANDENSEMBLEADABOOSTREGRESSOR WOULD PERFORM POORLY IF THE RANDOMSTATE WAS FIXED 7411 BY JOEL NOTHMAN

- FIX BUG IN ENSEMBLES WITH RANDOMIZATION WHERE THE ENSEMBLE WOULD NOT SET RANDOMSTATE ON BASE ESTIMATORS IN A PIPELINE OR SIMILAR NESTING 7411 NOTE RESULTS FOR ENSEMBLE BAGGINGCLASSIFIER ENSEMBLEBAGGINGREGRESSOR ENSEMBLEADABOOSTCLASSIFIER AND ENSEMBLEADABOOSTREGRESSOR WILL NOW DIFFER FROM PREVIOUS VERSIONS BY JOEL NOTHMAN

LINEAR KERNELIZED AND RELATED MODELS

- FIXED INCORRECT GRADIENT COMPUTATION FOR LOSS SQUAREDEPSILONINSENSITIVE IN LINEARMODELSGDCLASSIFIER ANDLINEARMODELSGDREGRESSOR 6764 BY WENHUA YANG

- FIX BUG IN LINEARMODELLOGISTICREGRESSIONCV WHERE SOLVERLIBLINEAR DID NOT ACCEPT CLASSWEIGHTSBALANCED 6817 BY TOM DUPRE LA TOUR

- FIX BUG IN NEIGHBORSRADIUSNEIGHBORSCLASSIFIER WHERE AN ERROR OCCURRED WHEN THERE WERE OUTLIERS BEING LABELLED AND A WEIGHT FUNCTION SPECIFIED 6902 BY LEONIEBORNE

- FIXLINEARMODELELASTICNET SPARSE DECISION FUNCTION TO MATCH OUTPUT WITH DENSE IN THE MULTIOUTPUT CASE

DECOMPOSITION MANIFOLD LEARNING AND CLUSTERING

- DECOMPOSITIONRANDOMIZEDPCA DEFAULT NUMBER OF ITERATEDPOWER IS 4 INSTEAD OF 3 5141 BY GIORGIO PATRINI

- UTILSEXTMATHRANDOMIZEDSVD PERFORMS 4 POWER ITERATIONS BY DEFAULT INSTEAD OF 0 IN PRACTICE THIS IS ENOUGH FOR OBTAINING A GOOD APPROXIMATION OF THE TRUE EIGENVALUESVECTORS IN THE PRESENCE OF NOISE WHEN NCOMPONENTS IS SMALL 1MINXSHAPE NITER IS SET TO 7 UNLESS THE USER SPECIFIES A HIGHER NUMBER THIS IMPROVES PRECISION WITH FEW COMPONENTS 5299 BY GIORGIO PATRINI

- WHITENNONWHITEN INCONSISTENCY BETWEEN COMPONENTS OF DECOMPOSITIONPCA ANDDECOMPOSITIONRANDOMIZEDPCA NOW FACTORED INTO PCA SEE THE NEW FEATURES IS FIXED COMPONENTS ARE STORED WITH NO WHITENING 5299 BY GIORGIO PATRINI

- FIXED BUG IN MANIFOLDSPECTRALEMBEDDING WHERE DIAGONAL OF UNNORMALIZED LAPLACIAN MATRIX WAS INCORRECTLY SET TO 1 4995 BY PETER FISCHER

- FIXED INCORRECT INITIALIZATION OF UTILSARPACKEIGSH ON ALL OCCURRENCES AFFECTS CLUSTER BICLUSTERSPECTRALBICLUSTERING DECOMPOSITIONKERNELPCA MANIFOLD LOCALLYLINEAREMBEDDING ANDMANIFOLDSPECTRALEMBEDDING 5012 BY PETER FISCHER

- ATTRIBUTE EXPLAINEDVARIANCERATIO CALCULATED WITH THE SVD SOLVER OF DISCRIMINANTANALYSISLINEARDISCRIMINANTANALYSIS NOW RETURNS CORRECT RESULTS BY JPFRANCOIA

PREPROCESSING AND FEATURE SELECTION

- PREPROCESSINGDATATRANSFORMSELECTED NOW ALWAYS PASSES A COPY OF XTO TRANSFORM FUNCTION WHENCOPYTRUE 7194 BY CAIO OLIVEIRA

SCIKITLEARN USER GUIDE RELEASE 0213

MODEL EVALUATION AND METAESTIMATORS

- MODELSELECTIONSTRATIFIEDKFOLD NOW RAISES ERROR IF ALL NLABELS FOR INDIVIDUAL CLASSES IS LESS THAN NFOLDS 6182 BY DEVASHISH DESHPANDE

- FIXED BUG IN MODELSELECTIONSTRATIFIEDSHUFFLESPLIT WHERE TRAIN AND TEST SAMPLE COULD OVERLAP IN SOME EDGE CASES SEE 6121 FOR MORE DETAILS BY LOIC ESTEVE

- FIX IN SKLEARNMODELSELECTIONSTRATIFIEDSHUFFLESPLIT TO RETURN SPLITS OF SIZE TRAINSIZE ANDTESTSIZE IN ALL CASES 6472 BY ANDREAS MÜLLER

- CROSSVALIDATION OF ONEVSONECLASSIFIER ANDONEVSRESTCLASSIFIER NOW WORKS WITH PRECOMPUTED KERNELS 7350 BY RUSSELL SMITH

- FIX INCOMPLETE PREDICTPROBA METHOD DELEGATION FROM MODELSELECTIONGRIDSEARCHCV TO LINEARMODELSGDCCLASSIFIER 7159 BY YICHUAN LIU

METRICS

- FIX BUG IN METRICSSILHOUETTESCORE IN WHICH CLUSTERS OF SIZE 1 WERE INCORRECTLY SCORED THEY SHOULD GET A SCORE OF 0 BY JOEL NOTHMAN

- FIX BUG IN METRICSSILHOUETTESAMPLES SO THAT IT NOW WORKS WITH ARBITRARY LABELS NOT JUST THOSE RANGING FROM 0 TO NCLUSTERS 1

- FIX BUG WHERE EXPECTED AND ADJUSTED MUTUAL INFORMATION WERE INCORRECT IF CLUSTER CONTINGENCY CELLS EXCEEDED 216 BY JOEL NOTHMAN

- METRICSPAIRWISEPAIRWISEDISTANCES NOW CONVERTS ARRAYS TO BOOLEAN ARRAYS WHEN REQUIRED IN SCIPYSPATIALDISTANCE 5460 BY TOM DUPRE LA TOUR

- FIX SPARSE INPUT SUPPORT IN METRICSSILHOUETTESCORE AS WELL AS EXAMPLE EXAM PLETEXTDOCUMENTCLUSTERINGPY BY YENCHEN LIN

- METRICSPROCCURVE ANDMETRICSPRECISIONRECALLCURVE NO LONGER ROUND YSCORE VALUES WHEN CREATING ROC CURVES THIS WAS CAUSING PROBLEMS FOR USERS WITH VERY SMALL DIFFERENCES IN SCORES 7353

MISCELLANEOUS

- MODELSELECTIONTESTSSEARCHCHECKPARAMGRID NOW WORKS CORRECTLY WITH ALL TYPES THAT EXTENDSIMPLEMENTS SEQUENCE EXCEPT STRING INCLUDING RANGE PYTHON 3X AND XRANGE PYTHON 2X 7323 BY VIACHESLAV KOVALEVSKYI

- UTILSEXTMATHRANDOMIZEDDRANGEFINDER IS MORE NUMERICALLY STABLE WHEN MANY POWER ITERATIONS ARE REQUESTED SINCE IT APPLIES LU NORMALIZATION BY DEFAULT IF NITER2 NUMERICAL ISSUES ARE UNLIKELY THUS NO NORMALIZATION IS APPLIED OTHER NORMALIZATION OPTIONS ARE AVAILABLE NONE LU ANDQR 5141 BY GIORGIO PATRINI

- FIX A BUG WHERE SOME FORMATS OF SCIPYSPARSE MATRIX AND ESTIMATORS WITH THEM AS PARAMETERS COULD NOT BE PASSED TO BASECLONE BY LOIC ESTEVE

- DATASETSLOADSVMLIGHTFILE NOW IS ABLE TO READ LONG INT QID VALUES 7101 BY IBRAIM GANIEV

API CHANGES SUMMARY

LINEAR KERNELIZED AND RELATED MODELS

- RESIDUALMETRIC HAS BEEN DEPRECATED IN LINEARMODELRANSACREGRESSOR USELOSS INSTEAD BY MANOJ KUMAR

- ACCESS TO PUBLIC ATTRIBUTES X ANDY HAS BEEN DEPRECATED IN ISOTONICISOTONICREGRESSION BY JONATHAN ARFA

117 PREVIOUS RELEASES 85

SCIKITLEARN USER GUIDE RELEASE 0213

DECOMPOSITION MANIFOLD LEARNING AND CLUSTERING

- THE OLDMIXTUREDPGMM IS DEPRECATED IN FAVOR OF THE NEW MIXTUREBAYESIANGAUSSIANMIXTURE WITH THE PARAMETER WEIGHTCONCENTRATIONPRIORTYPEDIRICHLETPROCESS THE NEW CLASS SOLVES THE COMPUTATIONAL PROBLEMS OF THE OLD CLASS AND COMPUTES THE GAUSSIAN MIXTURE WITH A DIRICHLET PROCESS PRIOR FASTER THAN BEFORE 7295 BY WEI XUE AND THIERRY GUILLEMOT
- THE OLDMIXTUREVBGMM IS DEPRECATED IN FAVOR OF THE NEW MIXTUREBAYESIANGAUSSIANMIXTURE WITH THE PARAMETER WEIGHTCONCENTRATIONPRIORTYPEDIRICHLETDISTRIBUTION THE NEW CLASS SOLVES THE COMPUTATIONAL PROBLEMS OF THE OLD CLASS AND COMPUTES THE VARIATIONAL BAYESIAN GAUSSIAN MIXTURE FASTER THAN BEFORE 6651 BY WEI XUE AND THIERRY GUILLEMOT
- THE OLDMIXTUREGMM IS DEPRECATED IN FAVOR OF THE NEW MIXTUREGAUSSIANMIXTURE THE NEW CLASS COMPUTES THE GAUSSIAN MIXTURE FASTER THAN BEFORE AND SOME OF COMPUTATIONAL PROBLEMS HAVE BEEN SOLVED 6666 BY WEI XUE AND THIERRY GUILLEMOT

MODEL EVALUATION AND METAESTIMATORS

- THESKLEARNCROSSVALIDATION SKLEARNGRIDSEARCH ANDSKLEARNLEARNINGCURVE HAVE BEEN DEPRECATED AND THE CLASSES AND FUNCTIONS HAVE BEEN REORGANIZED INTO THE SKLEARN MODELSELECTION MODULE REF MODEL SELECTION ENHANCEMENTS AND API CHANGES FOR MORE INFORMATION 4294 BY RAGHAV RV
- THEGRIDSCORES ATTRIBUTE OF MODELSELECTIONGRIDSEARCHCV ANDMODELSELECTION RANDOMIZEDSEARCHCV IS DEPRECATED IN FAVOR OF THE ATTRIBUTE CVRESULTS REF MODEL SELECTION ENHANCEMENTS AND API CHANGES FOR MORE INFORMATION 6697 BY RAGHAV RV
- THE PARAMETERS NITER ORNFOLDS IN OLD CV SPLITTERS ARE REPLACED BY THE NEW PARAMETER NSPLITS SINCE IT CAN PROVIDE A CONSISTENT AND UNAMBIGUOUS INTERFACE TO REPRESENT THE NUMBER OF TRAINTEST SPLITS 7187 BY YENCHEN LIN
- CLASSES PARAMETER WAS RENAMED TO LABELS INMETRICSHAMMINGLOSS 7260 BY SEBASTIÁN VANRELL
- THE SPLITTER CLASSES LABELKFOLD LABELSHUFFLESPLIT LEAVEONELABELOUT AND LEAVEPLABELSOUT ARE RENAMED TO MODELSELECTIONGROUPKFOLD MODELSELECTION GROUPSHUFFLESPLIT MODELSELECTIONLEAVEONEGROUPOUT ANDMODELSELECTION LEAVEPGROUPSOUT RESPECTIVELY ALSO THE PARAMETER LABELS IN THESPLIT METHOD OF THE NEWLY RENAMED SPLITTERSMODELSELECTIONLEAVEONEGROUPOUT ANDMODELSELECTIONLEAVEPGROUPSOUT IS RENAMED TO GROUPS ADDITIONALLY IN MODELSELECTIONLEAVEPGROUPSOUT THE PARAMETER NLABELS IS RENAMED TO NGROUPS 6660 BY RAGHAV RV
- ERROR AND LOSS NAMES FOR SCORING PARAMETERS ARE NOW PREFIXED BY NEG SUCH AS NEGMEANSQUAREDERROR THE UNPREFIXED VERSIONS ARE DEPRECATED AND WILL BE REMOVED IN VERSION 020 7261 BY TIM HEAD

CODE CONTRIBUTORS

ADITYA JOSHI ALEJANDRO ALEXANDER FABISCH ALEXANDER LOGINOV ALEXANDER MINYUSHKIN ALEXANDER RUDY ALEXAN DRE ABADIE ALEXANDRE ABRAHAM ALEXANDRE GRAMFORT ALEXANDRE SAINT ALEXFIELDS ALVARO ULLOA ALYSSAQ AMLAN KAR ANDREAS MUELLER ANDREW GIESSEL ANDREW JACKSON ANDREW MCCULLOH ANDREW MURRAY ANISH SHAH ARAFAT ARCHIT SHARMA ARIEL ROKEM ARNAUD JOLY ARNAUD RACHEZ ARTHUR MENSCH ASH HOOVER ASNT BONOI BEHZAD TABIB IAN BERNARDO BERNHARD KRATZWALD BHARGAV MANGIPUDI BLAKEFLEI BOYUAN DENG BRANDON CARTER BRETT NAUL BRIAN MCFEE CAIO OLIVEIRA CAMILO LAMUS CAROL WILLING CASS CESHINE LEE CHARLES TRUONG CHYIKWEI YAU CJ CAREY CODEVIG COLIN NI DAN SHIEBLER DANIEL DANIEL HNYK DAVID ELLIS DAVID NICHOLSON DAVID STAUB DAVID THALER DAVID WARSHAW DAVIDE LASAGNA DEBORAH DEFINITELYUNCERTAIN DIDI BARZEV DJIPEY DSQUAREINDIA EDWINENSAE ELIAS KUTHE ELVIS DOHMATOB ETHAN WHITE FABIAN PEDREGOSA FABIO TICCONI FISACHE FLORIAN WILHELM FRANCIS FRANCIS O'DONOVAN GAELE VAROQUAUX GANIEV IBRAIM GHG GILLES LOUPPE GIORGIO PATRINI GIOVANNI CHERUBIN GIO VANNI LANZANI GLENN QIAN GORDON MOHR GOVINVATSAN GRAHAM CLENAGHAN GREG REDA GREG STUPP GUILLAUME 86 CHAPTER 1 WELCOME TO SCIKITLEARN

SCIKITLEARN USER GUIDE RELEASE 0213

LEMAITRE GUSTAV MÖRTBERG HALWAI HARIZO RAJAONA HARRY MAVROFORAKIS HASHCODE55 HDMETOR HENRY LIN HOB  
SON LANE HUGO BOWNEANDERSON IGOR ANDRIUSHCHENKO IMACULATE INKI HWANG ISAAC SIJARANAMUAL ISHANK GULATI  
ISSAM LARADJI IVER JORDAL JACKMARTIN JACOB SCHREIBER JAKE VANDERPLAS JAMES FIEDLER JAMES ROUTLEY JAN ZIKES  
JANNA BRETTINGEN JARFA JASON LASKA JBLACKBURNE JEFF LEVESQUE JEFFREY BLACKBURNE JEFFREY04 JEREMY HINTZ JERE  
MYNIXON JEROEN JESSICA YUNG JILLJÉNN VIE JIMMY JIA JIYUAN QIAN JOEL NOTHMAN JOHANNAH JOHN JOHN BOERSMA  
JOHN KIRKHAM JOHN MOELLER JONATHANSTRIEBEL JONCRALL JORDI JOSEPH MUNOZ JOSHUA COOK JPFrancoia JRFIEDLER  
JULIANKAHNERT JULIATHEBRAVE KAICHOGAMI KAMALAKERDADI KENNETH LYONS KEVIN WANG KINGJR KJELL KONSTANTIN  
PODSHUMOK KORNEL KIELCZEWSKI KRISHNA KALYAN KRISHNAKALYAN3 KVLE PUTNAM KYLE JACKSON LARS BUITINCK LDAVID  
LEIG LEIGHTONZHANG LELAND MCINNES LIANGCHI Hsieh LILIAN BESSON LIZSZ LOIC ESTEVE LOUIS TIAO LÉONIE BORNE  
MADS JENSEN MANITEJA NANDANA MANOJ KUMAR MANVENDRA SINGH MARCO MARIO KRELL MARK BAO MARK SZEPIENIEC  
MARTIN MADSEN MARTINBPR MARYANMOREL MASSIL MATHEUS MATHIEU BLONDEL MATHIEU DUBOIS MATTEO MATTHIAS EK  
MAN MAX MOROZ MICHAEL SCHERER MICHIAKI ARIGA MIKHAIL KOROBV MOUSSA TAIFI MRANDREWANDRADE MRIDUL SETH  
NADYAP NAOYA KANAI NATE GEORGE NELLE VAROQUAUX NELSON LIU NICK JAMES NICKLEDAVE NICO NICOLAS GOIX  
NIKOLAY MAYOROV NINGCHI NLATHIA OKBALEFTHANDED OKHLOPKOV OLIVIER GRISEL PANOS LOURIDAS PAUL STRICKLAND PER  
RINE LETELLIER PESTRICKLAND PETER FISCHER PIETER PINGYAO CHANG PRACTICALSWIFT PRESTON PARRY QIMU ZHENG RACHIT  
KANSAL RAGHAV RV RALF GOMMERS RAMANAS RAMMIG RANDY OLSON ROB ALEXANDER ROBERT LUTZ ROBIN SCHUCKER  
ROHAN JAIN RUIFENG ZHENG RYAN YU RÉMY LÉONE SAIHTTAM SAIWING YEUNG SAM SHLEIFER SAMUEL STJEAN SAR  
TAJ SINGH SASANK CHILAMKURTHY SAURABHBANSOD SCOTT ANDREWS SCOTT LOWE SEALES SEBASTIAN RASCHKA SEBASTIAN  
SAEGER SEBASTIÁN VANRELL SERGEI LEBEDEV SHAGUN SODHANI SHANMUGA CV SHASHANK SHEKHAR SHAWPAN SHENGXID  
UAN SHOTA SHUCKLE16 SKIPPER SEABOLD SKLEARNCI SMEDBERGM SRVANRELL SÉBASTIEN LERIQUE TARANJEET THEMRRMAX  
THIERRY THIERRY GUILLEMOT THOMAS THOMAS HALLOCK THOMAS MOREAU TIM HEAD TKAMMY TOASTEDCORNFLAKES TOM  
TOMDLT TOSHIHIRO KAMISHIMA TRACEROTONG TRENT HAUCK TREVORSTEPHENS TUE V O VARUN VARUN JEWALIKAR VIACH  
ESLAV VIGHNESH BIRODKAR VIKRAM VILLU RUUSMANN VINAYAK MEHTA WALTER WATERPONEY WENHUA YANG WENJIAN  
HUANG WILL WELCH WYSEGUY7 XYGUO YANLEND YAROSLAV HALCHENKO YELITE YEN YENCHENLIN YICHUAN LIU YOAV  
RAM YOSHIKI ZHENG RUIFENG ZIVORI ÓSCAR NÁJERA

1177 VERSION 0171

FEBRUARY 18 2016

CHANGELOG

BUG FIXES

- UPGRADE VENDORED JOBLIB TO VERSION 094 THAT FIXES AN IMPORTANT BUG IN JOBLIBPARALLEL THAT CAN SILENTLY YIELD TO WRONG RESULTS WHEN WORKING ON DATASETS LARGER THAN 1MB [HTTPS://GITHUB.COM/JOBLIB/JOBLIBBLOB094](https://github.com/joblib/joblib/blob/094/CHANGESRST)
  - FIXED READING OF BUNCH PICKLES GENERATED WITH SCIKITLEARN VERSION 016 THIS CAN AFFECT USERS WHO HAVE ALREADY DOWNLOADED A DATASET WITH SCIKITLEARN 016 AND ARE LOADING IT WITH SCIKITLEARN 017 SEE 6196 FOR HOW THIS AFFECTED DATASETSFETCH20NEWSGROUPS BY LOIC ESTEVE
  - FIXED A BUG THAT PREVENTED USING ROC AUC SCORE TO PERFORM GRID SEARCH ON SEVERAL CPU CORES ON LARGE ARRAYS SEE 6147 BY OLIVIER GRISEL
  - FIXED A BUG THAT PREVENTED TO PROPERLY SET THE PRESORT PARAMETER IN ENSEMBLE GRADIENTBOOSTINGREGRESSOR SEE 5857 BY ANDREW MCCULLOH
  - FIXED A JOBLIB ERROR WHEN EVALUATING THE PERPLEXITY OF A DECOMPOSITION LATENTDIRICHLETALLOCATION MODEL SEE 6258 BY CHYIKWEI YAU
- 117 PREVIOUS RELEASES 87

SCIKITLEARN USER GUIDE RELEASE 0213

1178 VERSION 017

NOVEMBER 5 2015

CHANGELOG

NEW FEATURES

- ALL THE SCALER CLASSES BUT PREPROCESSINGROBUSTSCALER CAN BE FITTED ONLINE BY CALLING PARTIALFIT BY GIORGIO PATRINI
  - THE NEW CLASS ENSEMBLEVOTINGCLASSIFIER IMPLEMENTS A “MAJORITY RULE” “SOFT VOTING” ENSEMBLE CLASSIFIER TO COMBINE ESTIMATORS FOR CLASSIFICATION BY SEBASTIAN RASCHKA
  - THE NEW CLASS PREPROCESSINGROBUSTSCALER PROVIDES AN ALTERNATIVE TO PREPROCESSING STANDARDSCALER FOR FEATUREWISE CENTERING AND RANGE NORMALIZATION THAT IS ROBUST TO OUTLIERS BY THOMAS UNTERTHINER
  - THE NEW CLASS PREPROCESSINGMAXABSSCALER PROVIDES AN ALTERNATIVE TO PREPROCESSING MINMAXSCALER FOR FEATUREWISE RANGE NORMALIZATION WHEN THE DATA IS ALREADY CENTERED OR SPARSE BY THOMAS UNTERTHINER
  - THE NEW CLASS PREPROCESSINGFUNCTIONTRANSFORMER TURNS A PYTHON FUNCTION INTO A PIPELINE COMPATIBLE TRANSFORMER OBJECT BY JOE JEVNIK
  - THE NEW CLASSES CROSSVALIDATIONLABELKFOLD ANDCROSSVALIDATION LABELSHUFFLESPLIT GENERATE TRAINTEST FOLDS RESPECTIVELY SIMILAR TO CROSSVALIDATIONKFOLD AND CROSSVALIDATIONSHUFFLESPLIT EXCEPT THAT THE FOLDS ARE CONDITIONED ON A LABEL ARRAY BY BRIAN MCFEE JEAN KOSSAIFI AND GILLES LOUPPE
  - DECOMPOSITIONLATENTDIRICHLETALLOCATION IMPLEMENTS THE LATENT DIRICHLET ALLOCATION TOPIC MODEL WITH ONLINE VARIATIONAL INFERENCE BY CHYIKWEI YAU WITH CODE BASED ON AN IMPLEMENTATION BY MATT HOFFMAN 3659
  - THE NEW SOLVER SAG IMPLEMENTS A STOCHASTIC AVERAGE GRADIENT DESCENT AND IS AVAILABLE IN BOTH LINEARMODELLOGISTICREGRESSION ANDLINEARMODELRIDGE THIS SOLVER IS VERY EFFICIENT FOR LARGE DATASETS BY DANNY SULLIVAN AND TOM DUPRE LA TOUR 4738
  - THE NEW SOLVER CDIMPLEMENTES A COORDINATE DESCENT IN DECOMPOSITIONNMF PREVIOUS SOLVER BASED ON PROJECTED GRADIENT IS STILL AVAILABLE SETTING NEW PARAMETER SOLVER TOPG BUT IS DEPRECATED AND WILL BE REMOVED IN 019 ALONG WITH DECOMPOSITIONPROJECTEDGRADIENTNMF AND PARAMETERS SPARSENESS ETA BETA ANDNLSMAXITER NEW PARAMETERS ALPHA ANDL1RATIO CONTROL L1 AND L2 REGULARIZATION AND SHUFFLE ADDS A SHUFFLING STEP IN THE CDSOLVER BY TOM DUPRE LA TOUR AND MATHIEU BLONDEL
- ENHANCEMENTS
- MANIFOLDTSNE NOW SUPPORTS APPROXIMATE OPTIMIZATION VIA THE BARNESHUT METHOD LEADING TO MUCH FASTER FITTING BY CHRISTOPHER ERICK MOODY 4025
  - CLUSTERMEANSHIFTMEANSHIFT NOW SUPPORTS PARALLEL EXECUTION AS IMPLEMENTED IN THE MEANSHIFT FUNCTION BY MARTINO SORBARO
  - NAIVEBAYESGAUSSIANNB NOW SUPPORTS FITTING WITH SAMPLEWEIGHT BY JAN HENDRIK METZEN
  - DUMMYDUMMYCLASSIFIER NOW SUPPORTS A PRIOR FITTING STRATEGY BY ARNAUD JOLY
  - ADDED AFITPREDICT METHOD FOR MIXTUREGMM AND SUBCLASSES BY CORY LORENZ

SCIKITLEARN USER GUIDE RELEASE 0213

- ADDED THE METRICSLABELRANKINGLOSS METRIC BY ARNAUD JOLY
  - ADDED THE METRICSCOHENKAPPASCORE METRIC
  - ADDED A WARMSTART CONSTRUCTOR PARAMETER TO THE BAGGING ENSEMBLE MODELS TO INCREASE THE SIZE OF THE ENSEMBLE BY TIM HEAD
  - ADDED OPTION TO USE MULTIOUTPUT REGRESSION METRICS WITHOUT AVERAGING BY KONSTANTIN SHMELKOV AND MICHAEL EICKENBERG
  - ADDEDSTRATIFY OPTION TOCROSSVALIDATIONTRAINTESTSPLIT FOR STRATIFIED SPLITTING BY MIROSLAV BATCHKAROV
  - THETREEEXPORTGRAPHVIZ FUNCTION NOW SUPPORTS AESTHETIC IMPROVEMENTS FOR TREE DECISIONTREECLASSIFIER ANDTREEDECISIONTREEREgressor INCLUDING OPTIONS FOR COLORING NODES BY THEIR MAJORITY CLASS OR IMPURITY SHOWING VARIABLE NAMES AND USING NODE PROPORTIONS INSTEAD OF RAW SAMPLE COUNTS BY TREVOR STEPHENS
  - IMPROVED SPEED OF NEWTONCG SOLVER INLINEARMODELLOGISTICREGRESSION BY AVOIDING LOSS COMPUTATION BY MATHIEU BLONDEL AND TOM DUPRE LA TOUR
  - THECLASSWEIGHTAUTO HEURISTIC IN CLASSIFIERS SUPPORTING CLASSWEIGHT WAS DEPRECATED AND REPLACED BY THE CLASSWEIGHTBALANCED OPTION WHICH HAS A SIMPLER FORMULA AND INTERPRETATION BY HANNA WALLACH AND ANDREAS MÜLLER
  - ADDCLASSWEIGHT PARAMETER TO AUTOMATICALLY WEIGHT SAMPLES BY CLASS FREQUENCY FOR LINEARMODEL PASSIVEAGGRESSIVECLASSIFIER BY TREVOR STEPHENS
  - ADDED BACKLINKS FROM THE API REFERENCE PAGES TO THE USER GUIDE BY ANDREAS MÜLLER
  - THELABELS PARAMETER TO SKLEARNMETRICSF1SCORE SKLEARNMETRICSF2SCORE SKLEARNMETRICSPRECISIONSCORE HAS BEEN EXTENDED IT IS NOW POSSIBLE TO IGNORE ONE OR MORE LABELS SUCH AS WHERE A MULTICLASS PROBLEM HAS A MAJORITY CLASS TO IGNORE BY JOEL NOTHMAN
  - ADDSAMPLEWEIGHT SUPPORT TOLINEARMODELRIDGECCLASSIFIER BY TREVOR STEPHENS
  - PROVIDE AN OPTION FOR SPARSE OUTPUT FROM SKLEARNMETRICSPAIRWISECOSINESIMILARITY BY JAIDEV DESHPANDE
  - ADDMINMAXSCALE TO PROVIDE A FUNCTION INTERFACE FOR MINMAXSCALER BY THOMAS UNTERTHINER
  - DUMPSVMLIGHTFILE NOW HANDLES MULTILABEL DATASETS BY CHIHWEI CHANG
  - RCV1 DATASET LOADER SKLEARNDATASETSFETCHRCV1 BY TOM DUPRE LA TOUR
  - THE “WISCONSIN BREAST CANCER” CLASSICAL TWOCCLASS CLASSIFICATION DATASET IS NOW INCLUDED IN SCIKITLEARN AVAILABLE WITHSKLEARNDATASETLOADBREASTCANCER
  - UPGRADED TO JOBLIB 093 TO BENEFIT FROM THE NEW AUTOMATIC BATCHING OF SHORT TASKS THIS MAKES IT POSSIBLE FOR SCIKITLEARN TO BENEFIT FROM PARALLELISM WHEN MANY VERY SHORT TASKS ARE EXECUTED IN PARALLEL FOR INSTANCE BY THE GRIDSEARCHGRIDSEARCHCV METAESTIMATOR WITH NJOBS 1 USED WITH A LARGE GRID OF PARAMETERS ON A SMALL DATASET BY VLAD NICULAE OLIVIER GRISEL AND LOIC ESTEVE
  - FOR MORE DETAILS ABOUT CHANGES IN JOBLIB 093 SEE THE RELEASE NOTES HTTPSGITHUBCOMJOBLIBJOBLIBBLOBBMASTERCHANGESRSTRELEASE093
  - IMPROVED SPEED 3 TIMES PER ITERATION OF DECOMPOSITIONDICTLEARNING WITH COORDINATE DESCENT METHOD FROM LINEARMODELLASSO BY ARTHUR MENSCH
  - PARALLEL PROCESSING THREADED FOR QUERIES OF NEAREST NEIGHBORS USING THE BALLTREE BY NIKOLAY MAYOROV
  - ALLOWDATASETSMAKEMULTILABELCLASSIFICATION TO OUTPUT A SPARSE Y BY KASHIF RASUL
- 117 PREVIOUS RELEASES 89

SCIKITLEARN USER GUIDE RELEASE 0213

- CLUSTERDBSCAN NOW ACCEPTS A SPARSE MATRIX OF PRECOMPUTED DISTANCES ALLOWING MEMORYEFFICIENT DISTANCE PRECOMPUTATION BY JOEL NOTHMAN
  - TREEDECISIONTREECLASSIFIER NOW EXPOSES AN APPLY METHOD FOR RETRIEVING THE LEAF INDICES SAMPLES ARE PREDICTED AS BY DANIEL GALVEZ AND GILLES LOUPPE
  - SPEED UP DECISION TREE REGRESSORS RANDOM FOREST REGRESSORS EXTRA TREES REGRESSORS AND GRADIENT BOOSTING ESTIMATORS BY COMPUTING A PROXY OF THE IMPURITY IMPROVEMENT DURING THE TREE GROWTH THE PROXY QUANTITY IS SUCH THAT THE SPLIT THAT MAXIMIZES THIS VALUE ALSO MAXIMIZES THE IMPURITY IMPROVEMENT BY ARNAUD JOLY JACOB SCHREIBER AND GILLES LOUPPE
  - SPEED UP TREE BASED METHODS BY REDUCING THE NUMBER OF COMPUTATIONS NEEDED WHEN COMPUTING THE IMPURITY MEASURE TAKING INTO ACCOUNT LINEAR RELATIONSHIP OF THE COMPUTED STATISTICS THE EFFECT IS PARTICULARLY VISIBLE WITH EXTRA TREES AND ON DATASETS WITH CATEGORICAL OR SPARSE FEATURES BY ARNAUD JOLY
  - ENSEMBLEGRADIENTBOOSTINGREGRESSOR ANDENSEMBLEGRADIENTBOOSTINGCLASSIFIER NOW EXPOSE AN APPLY METHOD FOR RETRIEVING THE LEAF INDICES EACH SAMPLE ENDS UP IN UNDER EACH TRY BY JACOB SCHREIBER
  - ADDSAMPLEWEIGHT SUPPORT TO LINEAR MODEL LINEAR REGRESSION BY SONNY HU 4881
  - ADD NITER WITHOUT PROGRESS TO MANIFOLD TSNE TO CONTROL THE STOPPING CRITERION BY SANTI VILLALBA 5186
  - ADDED OPTIONAL PARAMETER RANDOM STATE IN LINEAR MODEL RIDGE TO SET THE SEED OF THE PSEUDO RANDOM GENERATOR USED IN SAG SOLVER BY TOM DUPRE LA TOUR
  - ADDED OPTIONAL PARAMETER WARM START IN LINEAR MODEL LOGISTIC REGRESSION IF SET TO TRUE THE SOLVER SLBFGS NEWTONCG AND SAG WILL BE INITIALIZED WITH THE COEFFICIENTS COMPUTED IN THE PREVIOUS FIT BY TOM DUPRE LA TOUR
  - ADDED SAMPLEWEIGHT SUPPORT TO LINEAR MODEL LOGISTIC REGRESSION FOR THE LBFGS NEWTONCG AND SAG SOLVERS BY VALENTIN STOLBUNOV SUPPORT ADDED TO THE LIBLINEAR SOLVER BY MANOJ KUMAR
  - ADDED OPTIONAL PARAMETER PRESORT TO ENSEMBLE GRADIENT BOOSTING REGRESSOR AND ENSEMBLE GRADIENT BOOSTING CLASSIFIER KEEPING DEFAULT BEHAVIOR THE SAME THIS ALLOWS GRADIENT BOOSTERS TO TURN OFF PRESORTING WHEN BUILDING DEEP TREES OR USING SPARSE DATA BY JACOB SCHREIBER
  - ALTERED METRICS ROCCURVE TO DROP UNNECESSARY THRESHOLDS BY DEFAULT BY GRAHAM CLENAGHAN
  - ADDED FEATURE SELECTION SELECT FROM MODEL METATransformer WHICH CAN BE USED ALONG WITH ESTIMATORS THAT HAVE COEFF OR FEATURE IMPORTANCES ATTRIBUTE TO SELECT IMPORTANT FEATURES OF THE INPUT DATA BY MAHESHAKYA WIJEWARDENA JOEL NOTHMAN AND MANOJ KUMAR
  - ADDED METRICS PAIRWISE LAPLACIAN KERNEL BY CLYDE FARE
  - COVARIANCE GRAPH LASSO ALLOWS SEPARATE CONTROL OF THE CONVERGENCE CRITERION FOR THE ELASTIC NET SUBPROBLEM VIA THE ENET TOL PARAMETER
  - IMPROVED VERBOSITY IN DECOMPOSITION DICTIONARY LEARNING
  - ENSEMBLE RANDOM FOREST CLASSIFIER AND ENSEMBLE RANDOM FOREST REGRESSOR NO LONGER EXPLICITLY STORE THE SAMPLES USED IN BAGGING RESULTING IN A MUCH REDUCED MEMORY FOOTPRINT FOR STORING RANDOM FOREST MODELS
  - ADDED POSITIVE OPTION TO LINEAR MODEL LARS AND LINEAR MODEL LARS PATH TO FORCE COEFFICIENTS TO BE POSITIVE 5131
  - ADDED THE X NORMS SQUARED PARAMETER TO METRICS PAIRWISE EUCLIDEAN DISTANCES TO PROVIDE PRECOMPUTED SQUARED NORMS FOR X
  - ADDED THE FIT PREDICT METHOD TO PIPELINE PIPELINE
- 90 CHAPTER 1 WELCOME TO SCIKITLEARN



SCIKITLEARN USER GUIDE RELEASE 0213

- ADDED THE PREPROCESSINGMINMAXSCALE FUNCTION

BUG FIXES

- FIXED NONDETERMINISM IN DUMMYDUMMYCLASSIFIER WITH SPARSE MULTILABEL OUTPUT BY ANDREAS MÜLLER
  - FIXED THE OUTPUT SHAPE OF LINEARMODELRANSACREGRESSOR TONSAMPLES BY ANDREAS MÜLLER
  - FIXED BUG IN DECOMPOSITIONDICTLEARNING WHENNJOBS 0 BY ANDREAS MÜLLER
  - FIXED BUG WHERE GRIDSEARCHRANDOMIZEDSEARCHCV COULD CONSUME A LOT OF MEMORY FOR LARGE DISCRETE GRIDS BY JOEL NOTHMAN
  - FIXED BUG IN LINEARMODELLOGISTICREGRESSIONCV WHEREPENALTY WAS IGNORED IN THE FINAL FIT BY MANOJ KUMAR
  - FIXED BUG IN ENSEMBLEFORESTFORESTCLASSIFIER WHILE COMPUTING OOBSCORE AND X IS A SPARSESCMATRIX BY ANKUR ANKAN
  - ALL REGRESSORS NOW CONSISTENTLY HANDLE AND WARN WHEN GIVEN YTHAT IS OF SHAPE NSAMPLES 1 BY ANDREAS MÜLLER AND HENRY LIN 5431
  - FIX INCLUSTERKMEANS CLUSTER REASSIGNMENT FOR SPARSE INPUT BY LARS BUITINCK
  - FIXED A BUG IN LDALDA THAT COULD CAUSE ASYMMETRIC COVARIANCE MATRICES WHEN USING SHRINKAGE BY MARTIN BILLINGER
  - FIXEDCROSSVALIDATIONCROSSVALPREDICT FOR ESTIMATORS WITH SPARSE PREDICTIONS BY BUDDHA PRAKASH
  - FIXED THE PREDICTPROBA METHOD OFLINEARMODELLOGISTICREGRESSION TO USE SOFTMAX INSTEAD OF ONEVSREST NORMALIZATION BY MANOJ KUMAR 5182
  - FIXED THE PARTIALFIT METHOD OF LINEARMODELSGDCLASSIFIER WHEN CALLED WITH AVERAGETRUE BY ANDREW LAMB 5282
  - DATASET FETCHERS USE DIFFERENT FILENAMES UNDER PYTHON 2 AND PYTHON 3 TO AVOID PICKLING COMPATIBILITY ISSUES BY OLIVIER GRISEL 5355
  - FIXED A BUG IN NAIVEBAYESGAUSSIANNB WHICH CAUSED CLASSIFICATION RESULTS TO DEPEND ON SCALE BY JAKE VANDERPLAS
  - FIXED TEMPORARILY LINEARMODELRIDGE WHICH WAS INCORRECT WHEN FITTING THE INTERCEPT IN THE CASE OF SPARSE DATA THE FIX AUTOMATICALLY CHANGES THE SOLVER TO ‘SAG’ IN THIS CASE 5360 BY TOM DUPRE LA TOUR
  - FIXED A PERFORMANCE BUG IN DECOMPOSITIONRANDOMIZEDPCA ON DATA WITH A LARGE NUMBER OF FEATURES AND FEWER SAMPLES 4478 BY ANDREAS MÜLLER LOIC ESTEVE AND GIORGIO PATRINI
  - FIXED BUG IN CROSSDECOMPOSITIONPLS THAT YIELDED UNSTABLE AND PLATFORM DEPENDENT OUTPUT AND FAILED ONFITTRANSFORM BY ARTHUR MENSCH
  - FIXES TO THE BUNCH CLASS USED TO STORE DATASETS
  - FIXEDENSEMBLEPLOTPARTIALDEPENDENCE IGNORING THE PERCENTILES PARAMETER
  - PROVIDING A SET AS VOCABULARY IN COUNTVECTORIZER NO LONGER LEADS TO INCONSISTENT RESULTS WHEN PICKLING
  - FIXED THE CONDITIONS ON WHEN A PRECOMPUTED GRAM MATRIX NEEDS TO BE RECOMPUTED IN LINEARMODEL LINEARREGRESSION LINEARMODELORTHOGONALMATCHINGPURSUIT LINEARMODELLASSO ANDLINEARMODELELASTICNET
  - FIXED INCONSISTENT MEMORY LAYOUT IN THE COORDINATE DESCENT SOLVER THAT AFFECTED LINEARMODEL DICTIONARYLEARNING ANDCOVARIANCEGRAPHLASSO 5337 BY OLIVIER GRISEL
- 117 PREVIOUS RELEASES 91

SCIKITLEARN USER GUIDE RELEASE 0213

- MANIFOLDLOCALLYLINEAREMBEDDING NO LONGER IGNORES THE REG PARAMETER
  - NEAREST NEIGHBOR ESTIMATORS WITH CUSTOM DISTANCE METRICS CAN NOW BE PICKLED 4362
  - FIXED A BUG IN PIPELINEFEATUREUNION WHERE TRANSFORMERWEIGHTS WERE NOT PROPERLY HANDLED WHEN PERFORMING GRIDSEARCHES
  - FIXED A BUG IN LINEARMODELLOGISTICREGRESSION ANDLINEARMODEL LOGISTICREGRESSIONCV WHEN USING CLASSWEIGHTBALANCED ORCLASSWEIGHTAUTO
- BY TOM DUPRE LA TOUR
- FIXED BUG 5495 WHEN DOING OVR SVCDECISIONFUNCTIONSHAPE"OVR" FIXED BY ELVIS DOHMATOB
- API CHANGES SUMMARY
- ATTRIBUTE DATAMIN DATAMAX ANDDATARANGE INPREPROCESSINGMINMAXSCALER ARE DEPRE CATED AND WON'T BE AVAILABLE FROM 019 INSTEAD THE CLASS NOW EXPOSES DATAMIN DATAMAX AND DATARANGE BY GIORGIO PATRINI
  - ALL SCALER CLASSES NOW HAVE AN SCALE ATTRIBUTE THE FEATUREWISE RESCALING APPLIED BY THEIR TRANSFORM METHODS THE OLD ATTRIBUTE STD INPREPROCESSINGSTANDARDSCALER IS DEPRECATED AND SUPERSEDED BY SCALE IT WON'T BE AVAILABLE IN 019 BY GIORGIO PATRINI
  - SVMSVC ANDSVMNUSVC NOW HAVE AN DECISIONFUNCTIONSHAPE PARAMETER TO MAKE THEIR DECISION FUNCTION OF SHAPE NSAMPLES NCLASSES BY SETTINGDECISIONFUNCTIONSHAPEOVR THIS WILL BE THE DEFAULT BEHAVIOR STARTING IN 019 BY ANDREAS MÜLLER
  - PASSING 1D DATA ARRAYS AS INPUT TO ESTIMATORS IS NOW DEPRECATED AS IT CAUSED CONFUSION IN HOW THE ARRAY ELE MENTS SHOULD BE INTERPRETED AS FEATURES OR AS SAMPLES ALL DATA ARRAYS ARE NOW EXPECTED TO BE EXPLICITLY SHAPED NSAMPLES NFEATURES BY VIGHNESH BIRODKAR
  - LDALDA ANDQDAQDA HAVE BEEN MOVED TO DISCRIMINANTANALYSIS LINEARDISCRIMINANTANALYSIS AND DISCRIMINANTANALYSIS QUADRATICDISCRIMINANTANALYSIS
  - THESTORECOVARIANCE ANDTOL PARAMETERS HAVE BEEN MOVED FROM THE FIT METHOD TO THE CONSTRUCTOR IN DISCRIMINANTANALYSISLINEARDISCRIMINANTANALYSIS AND THESTORECOVARIANCES AND TOL PARAMETERS HAVE BEEN MOVED FROM THE FIT METHOD TO THE CONSTRUCTOR IN DISCRIMINANTANALYSIS QUADRATICDISCRIMINANTANALYSIS
  - MODELS INHERITING FROM LEARN TSELECTORMIXIN WILL NO LONGER SUPPORT THE TRANSFORM METHODS IE RAN DOMFORESTS GRADIENTBOOSTING LOGISTICREGRESSION DECISIONTREES SVMS AND SGD RELATED MODELS WRAP THESE MODELS AROUND THE METATransfomer FEATURESELECTIONSELECTFROMMODEL TO REMOVE FEATURES ACCORDING TO COEFS ORFEATUREIMPORTANCES WHICH ARE BELOW A CERTAIN THRESHOLD VALUE INSTEAD
  - CLUSTERKMEANS RERUNS CLUSTERASSIGNMENTS IN CASE OF NONCONVERGENCE TO ENSURE CONSISTENCY OF PREDICTX ANDLABELS BY VIGHNESH BIRODKAR
  - CLASSIFIER AND REGRESSOR MODELS ARE NOW TAGGED AS SUCH USING THE ESTIMATOR TYPE ATTRIBUTE
  - CROSSVALIDATION ITERATORS ALWAYS PROVIDE INDICES INTO TRAINING AND TEST SET NOT BOOLEAN MASKS
  - THEDECISIONFUNCTION ON ALL REGRESSORS WAS DEPRECATED AND WILL BE REMOVED IN 019 USE PREDICT INSTEAD
  - DATASETSLOADLFWPAIRS IS DEPRECATED AND WILL BE REMOVED IN 019 USE DATASETS FETCHLFWPAIRS INSTEAD
  - THE DEPRECATED HMM MODULE WAS REMOVED
  - THE DEPRECATED BOOTSTRAP CROSSVALIDATION ITERATOR WAS REMOVED
- 92 CHAPTER 1 WELCOME TO SCIKITLEARN

SCIKITLEARN USER GUIDE RELEASE 0213

- THE DEPRECATED WARD ANDWARDAGGLOMERATIVE CLASSES HAVE BEEN REMOVED USE CLUSTERING AGGLOMERATIVECLUSTERING INSTEAD
- CROSSVALIDATIONCHECKCV IS NOW A PUBLIC FUNCTION
- THE PROPERTY RESIDUES OFLINEARMODELLINEARREGRESSION IS DEPRECATED AND WILL BE REMOVED IN 019
- THE DEPRECATED NJOBS PARAMETER OF LINEARMODELLINEARREGRESSION HAS BEEN MOVED TO THE CONSTRUCTOR
- REMOVED DEPRECATED CLASSWEIGHT PARAMETER FROM LINEARMODELSGDCLASSIFIER 'SFIT METHOD USE THE CONSTRUCTION PARAMETER INSTEAD
- THE DEPRECATED SUPPORT FOR THE SEQUENCE OF SEQUENCES OR LIST OF LISTS MULTILABEL FORMAT WAS REMOVED TO CONVERT TO AND FROM THE SUPPORTED BINARY INDICATOR MATRIX FORMAT USE MULTILABELBINARIZER
- THE BEHAVIOR OF CALLING THE INVERSETRANSFORM METHOD OFPIPELINEPIPELINE WILL CHANGE IN 019 IT WILL NO LONGER RESHAPE ONEDIMENSIONAL INPUT TO TWODIMENSIONAL INPUT
- THE DEPRECATED ATTRIBUTES INDICATORMATRIX MULTILABEL ANDCLASSES OFPREPROCESSING LABELBINARIZER WERE REMOVED
- USINGGAMMA0 INSVMSVC ANDSVMSVR TO AUTOMATICALLY SET THE GAMMA TO 1 NFEATURES IS DEPRECATED AND WILL BE REMOVED IN 019 USE GAMMAAUTO INSTEAD

CODE CONTRIBUTORS

AARON SCHUMACHER ADITHYA GANESH AKITTY ALEXANDRE GRAMFORT ALEXEY GRIGOREV ALI BAHAREV ALLEN RIDDELL ANDO SAABAS ANDREAS MUELLER ANDREW LAMB ANISH SHAH ANKUR ANKAN ANTHONY ERLINGER ARI ROUVINEN ARNAUD JOLY ARNAUD RACHEZ ARTHUR MENSCH BANILO BARMALAYEY BENJAMINIRVING BOYUAN DENG BRETT NAUL BRIAN MCFEE BUDDHA PRAKASH CHI ZHANG CHIHWEI CHANG CHRISTOF ANGERMUELLER CHRISTOPH GOHLKE CHRISTOPHE BOURGUIGNAT CHRISTOPHER ERICK MOODY CHYIKWEI YAU CINDY SRIDHARAN CJ CAREY CLYDEFARE CORY LORENZ DAN BLANCHARD DANIEL GALVEZ DANIEL KRONOVET DANNY SULLIVAN DATA1010 DAVID DAVID D LOWE DAVID DOTSON DJIPEY DMITRY SPIKHALSKIY DONNE MARTIN DOUGAL J SUTHERLAND DOUGAL SUTHERLAND EDSON DUARTE EDUARDO CARO ERIC LARSON ERIC MARTIN ERICH SCHUBERT FERNANDO CARRILLO FRANK C ECKERT FRANK ZALKOW GAELE VAROQUAUX GANIEV IBRAIM GILLES LOUPPE GIORGIO PATRINI GIORGIOP GRAHAM CLENAGHAN GRYLLOS PROKOPIS GWULFS HENRY LIN HSUANTIENTIEN LIN IM MANUEL BAYER ISHANK GULATI JACK MARTIN JACOB SCHREIBER JAIDEV DESHPANDE JAKE VANDERPLAS JAN HENDRIK METZEN JEAN KOSSAIFI JEFFREY04 JEREMY JFRAJ JIALI MEI JOE JEVNIK JOEL NOTHMAN JOHN KIRKHAM JOHN WITTENAUER JOSEPH JOSHUA LOYAL JUNGKOOK PARK KAMALAKERDADI KASHIF RASUL KEITH GOODMAN KIAN HO KONSTANTIN SHMELKOV KYLER BROWN LARS BUITINCK LILIAN BESSON LOIC ESTEVE LOUIS TIAO MAHESHAKYA MAHESHAKYA WIJEWARDENA MANOJ KU MAR MARKTAB MARKTABNET MARTIN KU MARTIN SPACEK MARTINBPR MARTINOSORB MARYANMOREL MASAFUMI OYAMADA MATHIEU BLONDEL MATT KRUMP MATTI LYRA MAXIM KOLGANOV MBILLINGER MHG MICHAEL HEILMAN MICHAEL PATTERSON MIROSLAV BATCHKAROV NELLE VAROQUAUX NICOLAS NIKOLAY MAYOROV OLIVIER GRISEL OMER KATZ OSCAR NÁJERA PAULI VIRTANEN PETER FISCHER PETER PRETTENHOFER PHIL ROTH PIANOMANIA PRESTON PARRY RAGHAV RV ROB ZINKOV ROBERT LAYTON ROHAN RAMANATH SAKET CHOUDHARY SAM ZHANG SANTI SAURABHBANSOD SCLS19FR SEBASTIAN RASCHKA SEBAS TIAN SAEGER SHIVAN SORNARAJAH SIMONPL SINHRKS SKIPPER SEABOLD SONNY HU SSEG STEPHEN HOOVER STEVEN DE GRYZE STEVEN SEGUIN THEODORE VASILOUDIS THOMAS UNTERTHINER TIAGO FREITAS PEREIRA TIAN WANG TIM HEAD TIMOTHY HOPPER TOKOROTEN TOM DUPRÉ LA TOUR TREVOR STEPHENS VALENTIN STOLBUNOV VIGHNESH BIRODKAR VINAYAK MEHTA VINCENT VINCENT MICHEL VSTOLBUNOV WANGZ10 WEI XUE YUCHENG LOW YURY ZHAUNIAROVICH ZAC STEWART ZHAIPRO ZICHEN WANG

1179 VERSION 0161  
APRIL 14 2015  
117 PREVIOUS RELEASES 93

SCIKITLEARN USER GUIDE RELEASE 0213

CHANGELOG

BUG FIXES

- ALLOW INPUT DATA LARGER THAN BLOCKSIZE IN COVARIANCELEDOITWOLF BY ANDREAS MÜLLER
- FIX A BUG IN ISOTONICISOTONICREGRESSION DEDUPLICATION THAT CAUSED UNSTABLE RESULT IN CALIBRATIONCALIBRATEDCLASSIFIERCV BY JAN HENDRIK METZEN
- FIX SORTING OF LABELS IN FUNC PREPROCESSINGLABELBINARIZE BY MICHAEL HEILMAN
- FIX SEVERAL STABILITY AND CONVERGENCE ISSUES IN CROSSDECOMPOSITIONCCA AND CROSSDECOMPOSITIONPLSCANONICAL BY ANDREAS MÜLLER
- FIX A BUG IN CLUSTERKMEANS WHENPRECOMPUTEDDISTANCESFALSE ON FORTRANORDERED DATA
- FIX A SPEED REGRESSION IN ENSEMBLERANDOMFORESTCLASSIFIER 'SPREDICT ANDPREDICTPROBA BY ANDREAS MÜLLER
- FIX A REGRESSION WHERE UTILSSHUFFLE CONVERTED LISTS AND DATAFRAMES TO ARRAYS BY OLIVIER GRISEL

11710 VERSION 016

MARCH 26 2015

HIGHLIGHTS

- SPEED IMPROVEMENTS NOTABLY IN CLUSTERDBSCAN REDUCED MEMORY REQUIREMENTS BUGFIXES AND BETTER DEFAULT SETTINGS
- MULTINOMIAL LOGISTIC REGRESSION AND A PATH ALGORITHM IN LINEARMODELLOGISTICREGRESSIONCV
- OUTOF CORE LEARNING OF PCA VIA DECOMPOSITIONINCREMENTALPCA
- PROBABILITY CALLIBRATION OF CLASSIFIERS USING CALIBRATIONCALIBRATEDCLASSIFIERCV
- CLUSTERBIRCH CLUSTERING METHOD FOR LARGESCALE DATASETS
- SCALABLE APPROXIMATE NEAREST NEIGHBORS SEARCH WITH LOCALITYSENSITIVE HASHING FORESTS IN NEIGHBORS LSHFOREST
- IMPROVED ERROR MESSAGES AND BETTER VALIDATION WHEN USING MALFORMED INPUT DATA
- MORE ROBUST INTEGRATION WITH PANDAS DATAFRAMES

CHANGELOG

NEW FEATURES

- THE NEW NEIGHBORSLSHFOREST IMPLEMENTS LOCALITYSENSITIVE HASHING FOR APPROXIMATE NEAREST NEIGHBORS SEARCH BY MAHESHAKYA WIJEWARDENA
- ADDEDSVMLINEARSVR THIS CLASS USES THE LIBLINEAR IMPLEMENTATION OF SUPPORT VECTOR REGRESSION WHICH IS MUCH FASTER FOR LARGE SAMPLE SIZES THAN SVM SVR WITH LINEAR KERNEL BY FABIAN PEDREGOSA AND QIANG LUO
- INCREMENTAL FIT FOR GAUSSIANNB
- ADDEDSAMPLEWEIGHT SUPPORT TODUMMYDUMMYCLASSIFIER ANDDUMMYDUMMYREGRESSOR BY ARNAUD JOLY

ARNAUD JOLY

94 CHAPTER 1 WELCOME TO SCIKITLEARN

SCIKITLEARN USER GUIDE RELEASE 0213

- ADDED THE METRICSLABELRANKINGAVERAGEPRECISIONSCORE METRICS BY ARNAUD JOLY
  - ADD THEMETRICSCOVERAGEERROR METRICS BY ARNAUD JOLY
  - ADDEDLINEARMODELLOGISTICREGRESSIONCV BY MANOJ KUMAR FABIAN PEDREGOSA GAE VAROQUAUX AND ALEXANDRE GRAMFORT
  - ADDEDWARMSTART CONSTRUCTOR PARAMETER TO MAKE IT POSSIBLE FOR ANY TRAINED FOREST MODEL TO GROW ADDITIONAL TREES INCREMENTALLY BY LAURENT DIRER
  - ADDEDSAMPLEWEIGHT SUPPORT TO ENSEMBLEGRADIENTBOOSTINGCLASSIFIER ANDENSEMBLE GRADIENTBOOSTINGREGRESSOR BY PETER PRETTENHOFER
  - ADDEDDECOMPOSITIONINCREMENTALPCA AN IMPLEMENTATION OF THE PCA ALGORITHM THAT SUPPORTS OUT OFCORE LEARNING WITH A PARTIALFIT METHOD BY KYLE KASTNER
  - AVERAGED SGD FOR SGDCLASSIFIER ANDSGDREGRESSOR BY DANNY SULLIVAN
  - ADDEDCROSSVALPREDICT FUNCTION WHICH COMPUTES CROSSVALIDATED ESTIMATES BY LUIS PEDRO COELHO
  - ADDEDLINEARMODELTHEILSENREGRESSOR A ROBUST GENERALIZEDMEDIANBASED ESTIMATOR BY FLORIAN WILHELM
  - ADDEDMETRICSMEDIANABSOLUTEERROR A ROBUST METRIC BY GAE VAROQUAUX AND FLORIAN WILHELM
  - ADDCLUSTERBIRCH AN ONLINE CLUSTERING ALGORITHM BY MANOJ KUMAR ALEXANDRE GRAMFORT AND JOEL NOTHMAN
  - ADDED SHRINKAGE SUPPORT TO DISCRIMINANTANALYSISLINEARDISCRIMINANTANALYSIS USING TWO NEW SOLVERS BY CLEMENS BRUNNER AND MARTIN BILLINGER
  - ADDEDKERNELRIDGEKERNELRIDGE AN IMPLEMENTATION OF KERNELIZED RIDGE REGRESSION BY MATHIEU BLONDEL AND JAN HENDRIK METZEN
  - ALL SOLVERS IN LINEARMODELRIDGE NOW SUPPORT SAMPLEWEIGHT BY MATHIEU BLONDEL
  - ADDEDCROSSVALIDATIONPREDEFINEDSPLIT CROSSVALIDATION FOR FIXED USERPROVIDED CROSSVALIDATION FOLDS BY THOMAS UNTERTHINER
  - ADDEDCALIBRATIONCALIBRATEDCLASSIFIERCV AN APPROACH FOR CALIBRATING THE PREDICTED PROBABILITIES OF A CLASSIFIER BY ALEXANDRE GRAMFORT JAN HENDRIK METZEN MATHIEU BLONDEL AND BALAZS KEGL
  - ENHANCEMENTS
  - ADD OPTION RETURNDISTANCE INHIERARCHICALWARDTREE TO RETURN DISTANCES BETWEEN NODES FOR BOTH STRUCTURED AND UNSTRUCTURED VERSIONS OF THE ALGORITHM BY MATTEO VISCONTI DI OLEGGIO CASTELLO THE SAME OPTION WAS ADDED IN HIERARCHICALLINKAGETREE BY MANOJ KUMAR
  - ADD SUPPORT FOR SAMPLE WEIGHTS IN SCORER OBJECTS METRICS WITH SAMPLE WEIGHT SUPPORT WILL AUTOMATICALLY BENEFIT FROM IT BY NOEL DAWE AND VLAD NICULAE
  - ADDEDNEWTONCG ANDLBFGS SOLVER SUPPORT IN LINEARMODELLOGISTICREGRESSION BY MANOJ KUMAR
  - ADDSELECTIONRANDOM PARAMETER TO IMPLEMENT STOCHASTIC COORDINATE DESCENT FOR LINEARMODEL LASSO LINEARMODELELASTICNET AND RELATED BY MANOJ KUMAR
  - ADDSAMPLEWEIGHT PARAMETER TO METRICSJACCARDSIMILARITYSCORE ANDMETRICS LOGLOSS BY JATIN SHAH
  - SUPPORT SPARSE MULTILABEL INDICATOR REPRESENTATION IN PREPROCESSINGLABELBINARIZER AND MULTICLASSONEVSRESTCLASSIFIER BY HAMZEH ALSALHI WITH THANKS TO ROHIT SIVAPRASAD AS WELL AS EVALUATION METRICS BY JOEL NOTHMAN
- 117 PREVIOUS RELEASES 95

SCIKITLEARN USER GUIDE RELEASE 0213

- ADDSAMPLEWEIGHT PARAMETER TO METRICSJACCARDSIMILARITYSCORE BYJATIN SHAH
- ADD SUPPORT FOR MULTICLASS IN METRICSHINGELOSS ADDEDLABELSNONE AS OPTIONAL PARAMETER BY SAURABH JHA
- ADDSAMPLEWEIGHT PARAMETER TO METRICSHINGELOSS BYSAURABH JHA
- ADDMULTICLASSMULTINOMIAL OPTION INLINEARMODELLOGISTICREGRESSION TO IMPLEMENT A LOGISTIC REGRESSION SOLVER THAT MINIMIZES THE CROSSENTROPY OR MULTINOMIAL LOSS INSTEAD OF THE DEFAULT ONEVSREST SETTING SUPPORTS LBFGS ANDNEWTONCG SOLVERS BY LARS BUITINCK AND MANOJ KUMAR SOLVER OPTIONNEWTONCG BY SIMON WU
- DICTVECTORIZER CAN NOW PERFORM FITTRANSFORM ON AN ITERABLE IN A SINGLE PASS WHEN GIVING THE OPTION SORTFALSE BY DAN BLANCHARD
- GRIDSEARCHCV ANDRANDOMIZEDSEARCHCV CAN NOW BE CONFIGURED TO WORK WITH ESTIMATORS THAT MAY FAIL AND RAISE ERRORS ON INDIVIDUAL FOLDS THIS OPTION IS CONTROLLED BY THE ERRORSORESCORE PARAMETER THIS DOES NOT AFFECT ERRORS RAISED ON REFIT BY MICHAL ROMANIUK
- ADDDIGITS PARAMETER TO METRICSCCLASSIFICATIONREPORT TO ALLOW REPORT TO SHOW DIFFERENT PRECISION OF FLOATING POINT NUMBERS BY IAN GILMORE
- ADD A QUANTILE PREDICTION STRATEGY TO THE DUMMYDUMMYREGRESSOR BY AARON STAPLE
- ADDHANDLEUNKNOWN OPTION TOPREPROCESSINGONEHOTENCODER TO HANDLE UNKNOWN CATEGORICAL FEATURES MORE GRACEFULLY DURING TRANSFORM BY MANOJ KUMAR
- ADDED SUPPORT FOR SPARSE INPUT DATA TO DECISION TREES AND THEIR ENSEMBLES BY FARES HEDYATI AND ARNAUD JOLY
- OPTIMIZED CLUSTERAFFINITYPROPAGATION BY REDUCING THE NUMBER OF MEMORY ALLOCATIONS OF LARGE TEMPORARY DATASTRUCTURES BY ANTONY LEE
- PARELLIZATION OF THE COMPUTATION OF FEATURE IMPORTANCES IN RANDOM FOREST BY OLIVIER GRISEL AND ARNAUD JOLY
- ADDNITER ATTRIBUTE TO ESTIMATORS THAT ACCEPT A MAXITER ATTRIBUTE IN THEIR CONSTRUCTOR BY MANOJ KUMAR
- ADDED DECISION FUNCTION FOR MULTICLASSONEVSONECLASSIFIER BY RAGHAV RV AND KYLE BEAUCHAMP
- NEIGHBORSKNEIGHBORSGRAPH ANDRADIUSNEIGHBORSGRAPH SUPPORT NONEUCLIDEAN METRICS BY MANOJ KUMAR
- PARAMETER CONNECTIVITY INCLUSTERAGGLOMERATIVECLUSTERING AND FAMILY NOW ACCEPT CALLABLES THAT RETURN A CONNECTIVITY MATRIX BY MANOJ KUMAR
- SPARSE SUPPORT FOR PAIREDDISTANCES BY JOEL NOTHMAN
- CLUSTERDBSCAN NOW SUPPORTS SPARSE INPUT AND SAMPLE WEIGHTS AND HAS BEEN OPTIMIZED THE INNER LOOP HAS BEEN REWRITTEN IN CYTHON AND RADIUS NEIGHBORS QUERIES ARE NOW COMPUTED IN BATCH BY JOEL NOTHMAN AND LARS BUITINCK
- ADDCLASSWEIGHT PARAMETER TO AUTOMATICALLY WEIGHT SAMPLES BY CLASS FREQUENCY FOR ENSEMBLERANDOMFORESTCLASSIFIER TREEDECISIONTREECLASSIFIER ENSEMBLE EXTRATREESCLASSIFIER ANDTREEEXTRATREECLASSIFIER BY TREVOR STEPHENS
- GRIDSEARCHRANDOMIZEDSEARCHCV NOW DOES SAMPLING WITHOUT REPLACEMENT IF ALL PARAMETERS ARE GIVEN AS LISTS BY ANDREAS MÜLLER
- PARALLELIZED CALCULATION OF PAIRWISEDISTANCES IS NOW SUPPORTED FOR SCIPY METRICS AND CUSTOM CALLABLES BY JOEL NOTHMAN
- ALLOW THE FITTING AND SCORING OF ALL CLUSTERING ALGORITHMS IN PIPELINEPIPELINE BY ANDREAS MÜLLER
- MORE ROBUST SEEDING AND IMPROVED ERROR MESSAGES IN CLUSTERMEANSHIFT BY ANDREAS MÜLLER

SCIKITLEARN USER GUIDE RELEASE 0213

- MAKE THE STOPPING CRITERION FOR MIXTUREGMM MIXTUREDPGMM ANDMIXTUREVBGMM LESS DEPENDENT ON THE NUMBER OF SAMPLES BY THRESHOLDING THE AVERAGE LOGLIKELIHOOD CHANGE INSTEAD OF ITS SUM OVER ALL SAMPLES BY HERVÉ BREDIN
  - THE OUTCOME OF MANIFOLDSPECTRALEMBEDDING WAS MADE DETERMINISTIC BY FLIPPING THE SIGN OF EIGEN VECTORS BY HASIL SHARMA
  - SIGNIFICANT PERFORMANCE AND MEMORY USAGE IMPROVEMENTS IN PREPROCESSINGPOLYNOMIALFEATURES BY ERIC MARTIN
  - NUMERICAL STABILITY IMPROVEMENTS FOR PREPROCESSINGSTANDARDSCALER ANDPREPROCESSING SCALE BY NICOLAS GOIX
  - SVM SVC FITTED ON SPARSE INPUT NOW IMPLEMENTS DECISIONFUNCTION BY ROB ZINKOV AND ANDREAS MÜLLER
  - CROSSVALIDATIONTRAINTESTSPLIT NOW PRESERVES THE INPUT TYPE INSTEAD OF CONVERTING TO NUMPY ARRAYS
- DOCUMENTATION IMPROVEMENTS
- ADDED EXAMPLE OF USING FEATUREUNION FOR HETEROGENEOUS INPUT BY MATT TERRY
  - DOCUMENTATION ON SCORERS WAS IMPROVED TO HIGHLIGHT THE HANDLING OF LOSS FUNCTIONS BY MATT PICO
  - A DISCREPANCY BETWEEN LIBLINEAR OUTPUT AND SCIKITLEARN'S WRAPPERS IS NOW NOTED BY MANOJ KUMAR
  - IMPROVED DOCUMENTATION GENERATION EXAMPLES REFERRING TO A CLASS OR FUNCTION ARE NOW SHOWN IN A GALLERY ON THE CLASSFUNCTION'S API REFERENCE PAGE BY JOEL NOTHMAN
  - MORE EXPLICIT DOCUMENTATION OF SAMPLE GENERATORS AND OF DATA TRANSFORMATION BY JOEL NOTHMAN
  - SKLEARNNEIGHBORSBALLTREE ANDSKLEARNNEIGHBORSKDTree USED TO POINT TO EMPTY PAGES
  - STATING THAT THEY ARE ALIASES OF BINARYTREE THIS HAS BEEN FIXED TO SHOW THE CORRECT CLASS DOCS BY MANOJ KUMAR
  - ADDED SILHOUETTE PLOTS FOR ANALYSIS OF KMEANS CLUSTERING USING METRICSSILHOUETTESAMPLES AND METRICSSILHOUETTESCORE SEE SELECTING THE NUMBER OF CLUSTERS WITH SILHOUETTE ANALYSIS ON KMEANS CLUSTERING
- BUG FIXES
- METAESTIMATORS NOW SUPPORT DUCKTYPING FOR THE PRESENCE OF DECISIONFUNCTION PREDICTPROBA AND OTHER METHODS THIS FIXES BEHAVIOR OF GRIDSEARCHGRIDSEARCHCV GRIDSEARCHRANDOMIZEDSEARCHCV PIPELINEPIPELINE FEATURESELECTIONRFE FEATURESELECTIONRFE CV WHEN NESTED BY JOEL NOTHMAN
  - THESCORING ATTRIBUTE OF GRIDSEARCH AND CROSSVALIDATION METHODS IS NO LONGER IGNORED WHEN A GRIDSEARCHGRIDSEARCHCV IS GIVEN AS A BASE ESTIMATOR OR THE BASE ESTIMATOR DOESN'T HAVE PREDICT
  - THE FUNCTION HIERARCHICALWARDTREE NOW RETURNS THE CHILDREN IN THE SAME ORDER FOR BOTH THE STRUCTURED AND UNSTRUCTURED VERSIONS BY MATTEO VISCONTI DI OLEGGIO CASTELLO
  - FEATURESELECTIONRFE CV NOW CORRECTLY HANDLES CASES WHEN STEP IS NOT EQUAL TO 1 BY NIKOLAY MAYOROV
  - THEDECOMPOSITIONPCA NOW UNDOES WHITENING IN ITS INVERSETRANSFORM ALSO ITS COMPONENTS NOW ALWAYS HAVE UNIT LENGTH BY MICHAEL EICKENBERG
  - FIX INCOMPLETE DOWNLOAD OF THE DATASET WHEN DATASETSDOWNLOAD20NEWSGROUPS IS CALLED BY MANOJ KUMAR

SCIKITLEARN USER GUIDE RELEASE 0213

- VARIOUS FIXES TO THE GAUSSIAN PROCESSES SUBPACKAGE BY VINCENT DUBOURG AND JAN HENDRIK METZEN
  - CALLINGPARTIALFIT WITHCLASSWEIGHTAUTO THROWS AN APPROPRIATE ERROR MESSAGE AND SUGGESTS A WORK AROUND BY DANNY SULLIVAN
  - RBFSAMPLER WITHGAMMAG FORMERLY APPROXIMATED RBFKERNEL WITHGAMMAG2 THE DEFINITION OF GAMMA IS NOW CONSISTENT WHICH MAY SUBSTANTIALLY CHANGE YOUR RESULTS IF YOU USE A FIXED VALUE IF YOU CROSS VALIDATED OVER GAMMA IT PROBABLY DOESN'T MATTER TOO MUCH BY DOUGAL SUTHERLAND
  - PIPELINE OBJECT DELEGATE THE CLASSES ATTRIBUTE TO THE UNDERLYING ESTIMATOR IT ALLOWS FOR INSTANCE TO MAKE BAGGING OF A PIPELINE OBJECT BY ARNAUD JOLY
  - NEIGHBORSNEARESTCENTROID NOW USES THE MEDIAN AS THE CENTROID WHEN METRIC IS SET TO MANHATTAN IT WAS USING THE MEAN BEFORE BY MANOJ KUMAR
  - FIX NUMERICAL STABILITY ISSUES IN LINEARMODELSGDCLASSIFIER ANDLINEARMODEL SGDREGRESSOR BY CLIPPING LARGE GRADIENTS AND ENSURING THAT WEIGHT DECAY RESCALING IS ALWAYS POSITIVE FOR LARGE L2 REGULARIZATION AND LARGE LEARNING RATE VALUES BY OLIVIER GRISEL
  - WHENCOMPUTEFULLTREE IS SET TO "AUTO" THE FULL TREE IS BUILT WHEN NCLUSTERS IS HIGH AND IS EARLY STOPPED WHEN NCLUSTERS IS LOW WHILE THE BEHAVIOR SHOULD BE VICEVERSA IN CLUSTER AGGLOMERATIVECLUSTERING AND FRIENDS THIS HAS BEEN FIXED BY MANOJ KUMAR
  - FIX LAZY CENTERING OF DATA IN LINEARMODELENETPATH ANDLINEARMODELLASSOPATH IT WAS CENTERED AROUND ONE IT HAS BEEN CHANGED TO BE CENTERED AROUND THE ORIGIN BY MANOJ KUMAR
  - FIX HANDLING OF PRECOMPUTED AFFINITY MATRICES IN CLUSTERAGGLOMERATIVECLUSTERING WHEN USING CONNECTIVITY CONSTRAINTS BY CATHY DENG
  - CORRECTPARTIALFIT HANDLING OF CLASSPRIOR FORSKLEARNNAIVEBAYESMULTINOMIALNB AND SKLEARNNAIVEBAYESBERNOULLINB BY TREVOR STEPHENS
  - FIXED A CRASH IN METRICSPRECISIONRECALLFScoresSupport WHEN USING UNSORTED LABELS IN THE MULTILABEL SETTING BY ANDREAS MÜLLER
  - AVOID SKIPPING THE FIRST NEAREST NEIGHBOR IN THE METHODS RADIUSNEIGHBORS KNEIGHBORS KNEIGHBORSGRAPH ANDRADIUSNEIGHBORSGRAPH INSKLEARNNEIGHBORS NEARESTNEIGHBORS AND FAMILY WHEN THE QUERY DATA IS NOT THE SAME AS FIT DATA BY MANOJ KUMAR
  - FIX LOGDENSITY CALCULATION IN THE MIXTUREGMM WITH TIED COVARIANCE BY WILL DAWSON
  - FIXED A SCALING ERROR IN FEATURESELECTIONSELECTFDR WHERE A FACTOR NFEATURES WAS MISSING BY ANDREW TULLOCH
  - FIX ZERO DIVISION IN NEIGHBORSKNEIGHBORSREGRESSOR AND RELATED CLASSES WHEN USING DISTANCE WEIGHTING AND HAVING IDENTICAL DATA POINTS BY GARRETR
  - FIXED ROUND OFF ERRORS WITH NON POSITIVEDEFINITE COVARIANCE MATRICES IN GMM BY ALEXIS MIGNON
  - FIXED A ERROR IN THE COMPUTATION OF CONDITIONAL PROBABILITIES IN NAIVEBAYESBERNOULLINB BY HANNA WALLACH
  - MAKE THE METHOD RADIUSNEIGHBORS OFNEIGHBORSNEARESTNEIGHBORS RETURN THE SAMPLES LYING ON THE BOUNDARY FOR ALGORITHMBRUTE BY YAN YI
  - FLIP SIGN OF DUALCOEF OFSVMSVC TO MAKE IT CONSISTENT WITH THE DOCUMENTATION AND DECISIONFUNCTION BY ARTEM SOBOLEV
  - FIXED HANDLING OF TIES IN ISOTONICISOTONICREGRESSION WE NOW USE THE WEIGHTED AVERAGE OF TARGETS SECONDARY METHOD BY ANDREAS MÜLLER AND MICHAEL BOMMARITO
- 98 CHAPTER 1 WELCOME TO SCIKITLEARN



SCIKITLEARN USER GUIDE RELEASE 0213

API CHANGES SUMMARY

- GRIDSEARCHCV ANDCROSSVALSCORE AND OTHER METAESTIMATORS DON'T CONVERT PANDAS DATAFRAMES INTO ARRAYS ANY MORE ALLOWING DATAFRAME SPECIFIC OPERATIONS IN CUSTOM ESTIMATORS
  - MULTICLASSFITOVR MULTICLASSPREDICTOVR PREDICTPROBAOVR MULTICLASS FITOVO MULTICLASSPREDICTOVO MULTICLASSFITECOC ANDMULTICLASS PREDICTECOC ARE DEPRECATED USE THE UNDERLYING ESTIMATORS INSTEAD
  - NEAREST NEIGHBORS ESTIMATORS USED TO TAKE ARBITRARY KEYWORD ARGUMENTS AND PASS THESE TO THEIR DISTANCE METRIC THIS WILL NO LONGER BE SUPPORTED IN SCIKITLEARN 018 USE THE METRICPARAMS ARGUMENT INSTEAD
  - NJOBS PARAMETER OF THE FIT METHOD SHIFTED TO THE CONSTRUCTOR OF THE LINEARREGRESSION CLASS
  - THEPREDICTPROBA METHOD OFMULTICLASSONEVSRESTCLASSIFIER NOW RETURNS TWO PROBABILITIES PER SAMPLE IN THE MULTICLASS CASE THIS IS CONSISTENT WITH OTHER ESTIMATORS AND WITH THE METHOD'S DOCUMENTA TION BUT PREVIOUS VERSIONS ACCIDENTALLY RETURNED ONLY THE POSITIVE PROBABILITY FIXED BY WILL LAMOND AND LARS BUITINCK
  - CHANGE DEFAULT VALUE OF PRECOMPUTE IN ELASTICNET ANDLASSO TO FALSE SETTING PRECOMPUTE TO "AUTO" WAS FOUND TO BE SLOWER WHEN NSAMPLES NFEATURES SINCE THE COMPUTATION OF THE GRAM MATRIX IS COMPUTATIONALLY EXPENSIVE AND OUTWEIGHS THE BENEFIT OF FITTING THE GRAM FOR JUST ONE ALPHA PRECOMPUTEAUTO IS NOW DEPRECATED AND WILL BE REMOVED IN 018 BY MANOJ KUMAR
  - EXPOSEPOSITIVE OPTION INLINEARMODELENETPATH ANDLINEARMODELENETPATH WHICH CONSTRAINS COEFFICIENTS TO BE POSITIVE BY MANOJ KUMAR
  - USERS SHOULD NOW SUPPLY AN EXPLICIT AVERAGE PARAMETER TO SKLEARNMETRICSF1SCORE SKLEARN METRICSBETASCORE SKLEARNMETRICSRECALLSCORE ANDSKLEARNMETRICS PRECISIONSCORE WHEN PERFORMING MULTICLASS OR MULTILABEL IE NOT BINARY CLASSIFICATION BY JOEL NOTHMAN
  - SCORING PARAMETER FOR CROSS VALIDATION NOW ACCEPTS F1MICRO F1MACRO ORF1WEIGHTED F1 IS NOW FOR BINARY CLASSIFICATION ONLY SIMILAR CHANGES APPLY TO PRECISION ANDRECALL BY JOEL NOTHMAN
  - THEFITINTERCEPT NORMALIZE ANDRETURNMODELS PARAMETERS IN LINEARMODELENETPATH ANDLINEARMODELLASSOPATH HAVE BEEN REMOVED THEY WERE DEPRECATED SINCE 014
  - FROM NOW ONWARDS ALL ESTIMATORS WILL UNIFORMLY RAISE NOTFITTEDERROR UTILSVALIDATION NOTFITTEDERROR WHEN ANY OF THE PREDICT LIKE METHODS ARE CALLED BEFORE THE MODEL IS FIT BY RAGHAV RV
  - INPUT DATA VALIDATION WAS REFACTORED FOR MORE CONSISTENT INPUT VALIDATION THE CHECKARRAYS FUNCTION WAS REPLACED BY CHECKARRAY ANDCHECKXY BY ANDREAS MÜLLER
  - ALLOWXNONE IN THE METHODS RADIUSNEIGHBORS KNEIGHBORS KNEIGHBORSGRAPH AND RADIUSNEIGHBORSGRAPH INSKLEARNNEIGHBORSNEARESTNEIGHBORS AND FAMILY IF SET TO NONE THEN FOR EVERY SAMPLE THIS AVOIDS SETTING THE SAMPLE ITSELF AS THE FIRST NEAREST NEIGHBOR BY MANOJ KUMAR
  - ADD PARAMETER INCLUDESELF INNEIGHBORSKNEIGHBORSGRAPH ANDNEIGHBORS RADIUSNEIGHBORSGRAPH WHICH HAS TO BE EXPLICITLY SET BY THE USER IF SET TO TRUE THEN THE SAMPLE ITSELF IS CONSIDERED AS THE FIRST NEAREST NEIGHBOR
  - THRESH PARAMETER IS DEPRECATED IN FAVOR OF NEW TOL PARAMETER IN GMMDPGMM ANDVBGMM SEE ENHANCEMENTS SECTION FOR DETAILS BY HERVÉ BREDIN
  - ESTIMATORS WILL TREAT INPUT WITH DTYPE OBJECT AS NUMERIC WHEN POSSIBLE BY ANDREAS MÜLLER
  - ESTIMATORS NOW RAISE VALUEERROR CONSISTENTLY WHEN FITTED ON EMPTY DATA LESS THAN 1 SAMPLE OR LESS THAN 1 FEATURE FOR 2D INPUT BY OLIVIER GRISEL
- 117 PREVIOUS RELEASES 99

SCIKITLEARN USER GUIDE RELEASE 0213

- THESHUFFLE OPTION OF LINEARMODELSGDCLASSIFIER LINEARMODELSGDSREGRESSOR LINEARMODELPERCEPTRON LINEARMODELPASSIVEAGGRESSIVECLASSIFIER AND LINEARMODELPASSIVEAGGRESSIVEREGRESSOR NOW DEFAULTS TO TRUE
- CLUSTERDBSCAN NOW USES A DETERMINISTIC INITIALIZATION THE RANDOMSTATE PARAMETER IS DEPRECATED BY ERICH SCHUBERT

CODE CONTRIBUTORS

A FLAXMAN AARON SCHUMACHER AARON STAPLE ABHISHEK THAKUR AKSHAY AKSHAYAH3 ALDRIAN OBAJA ALEXANDER FABISCH ALEXANDRE GRAMFORT ALEXIS MIGNON ANDERS AAGAARD ANDREAS MUELLER ANDREAS VAN CRANENBURGH AN DREW TULLOCH ANDREW WALKER ANTONY LEE ARNAUD JOLY BANILO BARMALEYEXE BEN DAVIES BENEDIKT KOEHLER BHSU BORIS FELD BORJA AYERDI BOYUAN DENG BRENT PEDERSEN BRIAN WIGNALL BROOKE OSBORN CALVIN GILES CATHY DENG CELEO CGOHLKE CHEBEE7I CHRISTIAN STADESCHULDT CHRISTOF ANGERMUELLER CHYIKWEI YAU CJ CAREY CLEMENS BRUN NER DAIKI AMINAKA DAN BLANCHARD DANFRANKJ DANNY SULLIVAN DAVID FLETCHER DMITRIJS MILAJEVS DOUGAL J SUTHER LAND ERICH SCHUBERT FABIAN PEDREGOSA FLORIAN WILHELM FLOYDSOFT FÉLIXANTOINE FORTIN GAELE VAROQUAUX GARRETT GILLES LOUPPE GPASSINO GWULFS HAMPUS BENGTSSON HAMZE AL SALHI HANNA WALLACH HARRY MAVROFORAKIS HASIL SHARMA HELDER HERVE BREDIN HSIANGFU YU HUGUES SALAMIN IAN GILMORE ILAMBHARATHI KANNIAH IMRAN HAQUE ISMS JAKE VANDERPLAS JAN DLABAL JAN HENDRIK METZEN JATIN SHAH JAVIER LÓPEZ PEÑA JDCABALLERO JEAN KOSSAIFI JEFF HAMMERBACHER JOEL NOTHMAN JONATHAN HELMUS JOSEPH KAICHENG ZHANG KEVIN MARKHAM KYLE BEAUCHAMP KYLE KASTNER LAGACHERIE MATTHIEU LARS BUITINCK LAURENT DIRER LEEPEI LOIC ESTEVE LUIS PEDRO COELHO LUKAS MICHEL BACHER MAHESHAKYA MANOJ KUMAR MANUEL MARIO MICHAEL KRELL MARTIN MARTIN BILLINGER MARTIN KU MATEUSZ SUSIK MATTHIEU BLONDEL MATT PICO MATT TERRY MATTEO VISCONTI DOC MATTI LYRA MAX LINKE MEHDI CHERTI MICHAEL BOMMARITO MICHAEL EICKENBERG MICHAL ROMANIUK MLG MRSHU NELLE VAROQUAUX NICOLA MONTECCHIO NICOLAS NIKOLAY MAYOROV NOEL DAWE OKAL BILLY OLIVIER GRISEL ÓSCAR NÁJERA PAOLO PUGGIONI PETER PRETTENHOFER PRATAP VARDHAN PVNGUYEN QUEQICHAO RAFAEL CARRASCOSA RAGHAV R V RAHIEL KASIM RANDALL MASON ROB ZINKOV ROBERT BRADSHAW SAKET CHOUDHARY SAM NICHOLLS SAMUEL CHARRON SAURABH JHA SETHDANDRIDGE SINHRKS SNUDERL STEFAN OTTE STEFAN VAN DER WALT STEVE TJOA SWU SYLVAIN ZIMMER TEJESH95 TERRYCOJONES THOMAS DELTEIL THOMAS UN TERTHINER TOMAS KAZMAR TREVORSTEPHENS TTTTHOMASSSSS TZUMING KUO UGURCALISKAN UGURTHEMASTER VINAYAK MEHTA VINCENT DUBOURG VJACHESLAV MURASHKIN VLAD NICULAE WADAWSON WEI XUE WILL LAMOND WU JIANG XOL XINFAN MENG YAN YI YUCHIN

11711 VERSION 0152

SEPTEMBER 4 2014

BUG FIXES

- FIXED HANDLING OF THE PPARAMETER OF THE MINKOWSKI DISTANCE THAT WAS PREVIOUSLY IGNORED IN NEAREST NEIGHBORS MODELS BY NIKOLAY MAYOROV
- FIXED DUPLICATED ALPHAS IN LINEARMODELLASSOLARS WITH EARLY STOPPING ON 32 BIT PYTHON BY OLIVIER GRISEL AND FABIAN PEDREGOSA
- FIXED THE BUILD UNDER WINDOWS WHEN SCIKITLEARN IS BUILT WITH MSVC WHILE NUMPY IS BUILT WITH MINGW BY OLIVIER GRISEL AND FEDERICO VAGGI
- FIXED AN ARRAY INDEX OVERFLOW BUG IN THE COORDINATE DESCENT SOLVER BY GAELE VAROQUAUX
- BETTER HANDLING OF NUMPY 19 DEPRECATION WARNINGS BY GAELE VAROQUAUX
- REMOVED UNNECESSARY DATA COPY IN CLUSTERKMEANS BY GAELE VAROQUAUX
- EXPLICITLY CLOSE OPEN FILES TO AVOID RESOURCEWARNINGS UNDER PYTHON 3 BY CALVIN GILES

100 CHAPTER 1 WELCOME TO SCIKITLEARN

SCIKITLEARN USER GUIDE RELEASE 0213

- THE TRANSFORM OF DISCRIMINANT ANALYSIS LINEAR DISCRIMINANT ANALYSIS NOW PROJECTS THE INPUT ON THE MOST DISCRIMINANT DIRECTIONS BY MARTIN BILLINGER
- FIXED POTENTIAL OVERFLOW IN TREE SAFER EALLOC BY LARS BUITINCK
- PERFORMANCE OPTIMIZATION IN ISOTONIC ISOTONIC REGRESSION BY ROBERT BRADSHAW
- NOSE IS NO LONGER A RUNTIME DEPENDENCY TO IMPORT SKLEARN ONLY FOR RUNNING THE TESTS BY JOEL NOTHMAN
- MANY DOCUMENTATION AND WEBSITE FIXES BY JOEL NOTHMAN LARS BUITINCK MATT PICO AND OTHERS

11712 VERSION 0151

AUGUST 1 2014

BUG FIXES

- MADE CROSS VALIDATION CROSS VAL SCORE USE CROSS VALIDATION K FOLD INSTEAD OF CROSS VALIDATION STRATIFIED K FOLD ON MULTI OUTPUT CLASSIFICATION PROBLEMS BY NIKOLAY MAYOROV
- SUPPORT UNSEEN LABELS PREPROCESSING LABEL BINARYZER TO RESTORE THE DEFAULT BEHAVIOR OF 0141 FOR BACKWARD COMPATIBILITY BY HAMZEH ALSALHI
- FIXED THE CLUSTER K MEANS STOPPING CRITERION THAT PREVENTED EARLY CONVERGENCE DETECTION BY EDWARD RAFF AND GAELE VAROQUAUX
- FIXED THE BEHAVIOR OF MULTICLASS ONE VS ONE CLASSIFIER IN CASE OF TIES AT THE PER CLASS VOTE LEVEL BY COMPUTING THE CORRECT PER CLASS SUM OF PREDICTION SCORES BY ANDREAS MÜLLER
- MADE CROSS VALIDATION CROSS VAL SCORE AND GRID SEARCH GRID SEARCH CV ACCEPT PYTHON LISTS AS INPUT DATA THIS IS ESPECIALLY USEFUL FOR CROSS VALIDATION AND MODEL SELECTION OF TEXT PROCESSING PIPELINES BY ANDREAS MÜLLER
- FIXED DATA INPUT CHECKS OF MOST ESTIMATORS TO ACCEPT INPUT DATA THAT IMPLEMENTS THE NUMPY ARRAY PROTOCOL THIS IS THE CASE FOR PANDAS SERIES AND PANDAS DATAFRAME IN RECENT VERSIONS OF PANDAS BY GAELE VAROQUAUX
- FIXED A REGRESSION FOR LINEAR MODEL SGD CLASSIFIER WITH CLASS WEIGHT AUTO ON DATA WITH NON CONTIGUOUS LABELS BY OLIVIER GRISEL

11713 VERSION 015

JULY 15 2014

HIGHLIGHTS

- MANY SPEED AND MEMORY IMPROVEMENTS ALL ACROSS THE CODE
- HUGE SPEED AND MEMORY IMPROVEMENTS TO RANDOM FORESTS AND EXTRA TREES THAT ALSO BENEFIT BETTER FROM PARALLEL COMPUTING
- INCREMENTAL FIT TO BERNOLLI RBM
- ADDED CLUSTER AGGLOMERATIVE CLUSTERING FOR HIERARCHICAL AGGLOMERATIVE CLUSTERING WITH AVERAGE LINKAGE COMPLETE LINKAGE AND WARD STRATEGIES
- ADDED LINEAR MODEL RANSAC REGRESSOR FOR ROBUST REGRESSION MODELS

117 PREVIOUS RELEASES 101

SCIKITLEARN USER GUIDE RELEASE 0213

- ADDED DIMENSIONALITY REDUCTION WITH MANIFOLDTSNE WHICH CAN BE USED TO VISUALIZE HIGHDIMENSIONAL DATA

CHANGELOG

NEW FEATURES

- ADDEDENSEMBLEBAGGINGCLASSIFIER ANDENSEMBLEBAGGINGREGRESSOR METAESTIMATORS FOR ENSEMBLING ANY KIND OF BASE ESTIMATOR SEE THE BAGGING SECTION OF THE USER GUIDE FOR DETAILS AND EXAMPLES BY GILLES LOUPPE

- NEW UNSUPERVISED FEATURE SELECTION ALGORITHM FEATURESELECTIONVARIANCETHRESHOLD BY LARS BUITINCK

- ADDEDLINEARMODELTRANSACREGRESSOR METAESTIMATOR FOR THE ROBUST FITTING OF REGRESSION MODELS BY JOHANNES SCHÖNBERGER

- ADDEDCLUSTERAGGLOMERATIVECLUSTERING FOR HIERARCHICAL AGGLOMERATIVE CLUSTERING WITH AVERAGE LINKAGE COMPLETE LINKAGE AND WARD STRATEGIES BY NELLE VAROQUAUX AND GAELE VAROQUAUX

- SHORTHAND CONSTRUCTORS PIPELINEMAKEPIPELINE ANDPIPELINEMAKEUNION WERE ADDED BY LARS BUITINCK

- SHUFFLE OPTION FOR CROSSVALIDATIONSTRATIFIEDKFOLD BY JEFFREY BLACKBURNE
- INCREMENTAL LEARNING PARTIALFIT FOR GAUSSIAN NAIVE BAYES BY IMRAN HAQUE

- ADDEDPARTIALFIT TOBERNOULLIRBM BY DANNY SULLIVAN

- ADDEDLEARNINGCURVE UTILITY TO CHART PERFORMANCE WITH RESPECT TO TRAINING SIZE SEE PLOTTING LEARNING CURVES BY ALEXANDER FABISCH

- ADD POSITIVE OPTION IN LASSOCV ANDELASTICNETCV BY BRIAN WIGNALL AND ALEXANDRE GRAMFORT

- ADDEDLINEARMODELMULTITASKELASTICNETCV ANDLINEARMODELMULTITASKLASSOCV BY MANOJ KUMAR

- ADDEDMANIFOLDTSNE BY ALEXANDER FABISCH

ENHANCEMENTS

- ADD SPARSE INPUT SUPPORT TO ENSEMBLEADABOOSTCLASSIFIER ANDENSEMBLE

- ADABOOSTREGRESSOR METAESTIMATORS BY HAMZEHALSALHI

- MEMORY IMPROVEMENTS OF DECISION TREES BY ARNAUD JOLY

- DECISION TREES CAN NOW BE BUILT IN BESTFIRST MANNER BY USING MAXLEAFNODES AS THE STOPPING CRITERIA

- REFACTORED THE TREE CODE TO USE EITHER A STACK OR A PRIORITY QUEUE FOR TREE BUILDING BY PETER PRETTENHOFER AND GILLES LOUPPE

- DECISION TREES CAN NOW BE FITTED ON FORTRAN AND CSTYLE ARRAYS AND NONCONTINUOUS ARRAYS WITHOUT THE NEED TO MAKE A COPY IF THE INPUT ARRAY HAS A DIFFERENT DTYPE THAN NPFLOAT32 A FORTRAN STYLE COPY WILL BE MADE SINCE FORTRANSTYLE MEMORY LAYOUT HAS SPEED ADVANTAGES BY PETER PRETTENHOFER AND GILLES LOUPPE

- SPEED IMPROVEMENT OF REGRESSION TREES BY OPTIMIZING THE THE COMPUTATION OF THE MEAN SQUARE ERROR CRITERION

- THIS LEAD TO SPEED IMPROVEMENT OF THE TREE FOREST AND GRADIENT BOOSTING TREE MODULES BY ARNAUD JOLY

- THEIMGTOGRAPH ANDGRIDTOGRAPH FUNCTIONS IN SKLEARNFEATUREEXTRACTIONIMAGE NOW

- RETURNNNPNDARRAY INSTEAD OFNPMATRIX WHENRETURNNASNPNDARRAY SEE THE NOTES SECTION FOR

- MORE INFORMATION ON COMPATIBILITY

- 102 CHAPTER 1 WELCOME TO SCIKITLEARN

SCIKITLEARN USER GUIDE RELEASE 0213

- CHANGED THE INTERNAL STORAGE OF DECISION TREES TO USE A STRUCT ARRAY THIS FIXED SOME SMALL BUGS WHILE IMPROVING CODE AND PROVIDING A SMALL SPEED GAIN BY JOEL NOTHMAN
  - REDUCE MEMORY USAGE AND OVERHEAD WHEN FITTING AND PREDICTING WITH FORESTS OF RANDOMIZED TREES IN PARALLEL WITH NJOBS 1 BY LEVERAGING NEW THREADING BACKEND OF JOBLIB 08 AND RELEASING THE GIL IN THE TREE FITTING CYTHON CODE BY OLIVIER GRISEL AND GILLES LOUPPE
  - SPEED IMPROVEMENT OF THE SKLEARNENSEMBLEGRADIENTBOOSTING MODULE BY GILLES LOUPPE AND PETER PRETTENHOFER
  - VARIOUS ENHANCEMENTS TO THE SKLEARNENSEMBLEGRADIENTBOOSTING MODULE A WARMSTART ARGUMENT TO FIT ADDITIONAL TREES A MAXLEAFNODES ARGUMENT TO FIT GBM STYLE TREES A MONITOR FIT ARGUMENT TO INSPECT THE ESTIMATOR DURING TRAINING AND REFACTORING OF THE VERBOSE CODE BY PETER PRETTENHOFER
  - FASTER SKLEARNENSEMBLEEXTRA TREES BY CACHING FEATURE VALUES BY ARNAUD JOLY
  - FASTER DEPTHBASED TREE BUILDING ALGORITHM SUCH AS DECISION TREE RANDOM FOREST EXTRA TREES OR GRADIENT TREE BOOSTING WITH DEPTH BASED GROWING STRATEGY BY AVOIDING TRYING TO SPLIT ON FOUND CONSTANT FEATURES IN THE SAMPLE SUBSET BY ARNAUD JOLY
  - ADD MINWEIGHT FRACTION LEAF PREPRUNING PARAMETER TO TREEBASED METHODS THE MINIMUM WEIGHTED FRACTION OF THE INPUT SAMPLES REQUIRED TO BE AT A LEAF NODE BY NOEL DAWE
  - ADDED METRICS PAIRWISE DISTANCES ARGMINMIN BY PHILIPPE GERVAIS
  - ADDED PREDICT METHOD TO CLUSTER AFFINITY PROPAGATION AND CLUSTER MEAN SHIFT BY MATHIEU BLONDEL
  - VECTOR AND MATRIX MULTIPLICATIONS HAVE BEEN OPTIMISED THROUGHOUT THE LIBRARY BY DENIS ENGEMANN AND ALEXANDRE GRAMFORT IN PARTICULAR THEY SHOULD TAKE LESS MEMORY WITH OLDER NUMPY VERSIONS PRIOR TO 1.7.2
  - PRECISION RECALL AND ROC EXAMPLES NOW USE TRAIN TESTSPLIT AND HAVE MORE EXPLANATION OF WHY THESE METRICS ARE USEFUL BY KYLE KASTNER
  - THE TRAINING ALGORITHM FOR DECOMPOSITION NMF IS FASTER FOR SPARSE MATRICES AND HAS MUCH LOWER MEMORY COMPLEXITY MEANING IT WILL SCALE UP GRACEFULLY TO LARGE DATASETS BY LARS BUITINCK
  - ADDED SVD METHOD OPTION WITH DEFAULT VALUE TO "RANDOMIZED" TO DECOMPOSITION FACTOR ANALYSIS TO SAVE MEMORY AND SIGNIFICANTLY SPEEDUP COMPUTATION BY DENIS ENGEMANN AND ALEXANDRE GRAMFORT
  - CHANGED CROSS VALIDATION STRATIFIED K FOLD TO TRY AND PRESERVE AS MUCH OF THE ORIGINAL ORDERING OF SAMPLES AS POSSIBLE SO AS NOT TO HIDE OVERFITTING ON DATASETS WITH A NON NEGLIGIBLE LEVEL OF SAMPLES DEPENDENCY BY DANIEL NOURI AND OLIVIER GRISEL
  - ADD MULTI OUTPUT SUPPORT TO GAUSSIAN PROCESS GAUSSIAN PROCESS BY JOHN NOVAK
  - SUPPORT FOR PRECOMPUTED DISTANCE MATRICES IN NEAREST NEIGHBOR ESTIMATORS BY ROBERT LAYTON AND JOEL NOTHMAN
  - NORM COMPUTATIONS OPTIMIZED FOR NUMPY 1.6 AND LATER VERSIONS BY LARS BUITINCK IN PARTICULAR THE KMEANS ALGORITHM NO LONGER NEEDS A TEMPORARY DATA STRUCTURE THE SIZE OF ITS INPUT
  - DUMMY DUMMY CLASSIFIER CAN NOW BE USED TO PREDICT A CONSTANT OUTPUT VALUE BY MANOJ KUMAR
  - DUMMY DUMMY REGRESSOR HAS NOW A STRATEGY PARAMETER WHICH ALLOWS TO PREDICT THE MEAN THE MEDIAN OF THE TRAINING SET OR A CONSTANT OUTPUT VALUE BY MAHESHAKYA WIJEWARDENA
  - MULTILABEL CLASSIFICATION OUTPUT IN MULTILABEL INDICATOR FORMAT IS NOW SUPPORTED BY METRICS ROCAUC SCORE AND METRICS AVERAGE PRECISION SCORE BY ARNAUD JOLY
  - SIGNIFICANT PERFORMANCE IMPROVEMENTS MORE THAN 100X SPEEDUP FOR LARGE PROBLEMS IN ISOTONIC ISOTONIC REGRESSION BY ANDREW TULLOCH
  - SPEED AND MEMORY USAGE IMPROVEMENTS TO THE SGD ALGORITHM FOR LINEAR MODELS IT NOW USES THREADS NOT SEPARATE PROCESSES WHEN NJOBS 1 BY LARS BUITINCK
- 117 PREVIOUS RELEASES 103

SCIKITLEARN USER GUIDE RELEASE 0213

- GRID SEARCH AND CROSS VALIDATION ALLOW NANS IN THE INPUT ARRAYS SO THAT PREPROCESSORS SUCH AS PREPROCESSINGIMPUTER CAN BE TRAINED WITHIN THE CROSS VALIDATION LOOP AVOIDING POTENTIALLY SKEWED RESULTS
  - RIDGE REGRESSION CAN NOW DEAL WITH SAMPLE WEIGHTS IN FEATURE SPACE ONLY SAMPLE SPACE UNTIL THEN BY MICHAEL EICKENBERG BOTH SOLUTIONS ARE PROVIDED BY THE CHOLESKY SOLVER
  - SEVERAL CLASSIFICATION AND REGRESSION METRICS NOW SUPPORT WEIGHTED SAMPLES WITH THE NEW SAMPLEWEIGHT ARGUMENT METRICSACCURACYScore METRICSZEROONELOSS METRICSPRECISIONScore METRICSAVERAGEPRECISIONScore METRICS F1Score METRICSFbetaScore METRICSRECALLScore METRICSROCAUCScore METRICSEXPLAINEDVARIANCEScore METRICSMEANSQUAREDERROR METRICS MEANABSOLUTEERROR METRICSR2Score BY NOEL DAWE
  - SPEED UP OF THE SAMPLE GENERATOR DATASETSMAKEMULTILABELCLASSIFICATION BY JOEL NOTHMAN
- DOCUMENTATION IMPROVEMENTS
- THE WORKING WITH TEXT DATA TUTORIAL HAS NOW BEEN WORKED IN TO THE MAIN DOCUMENTATION’S TUTORIAL SECTION INCLUDES EXERCISES AND SKELETONS FOR TUTORIAL PRESENTATION ORIGINAL TUTORIAL CREATED BY SEVERAL AUTHORS INCLUDING OLIVIER GISEL LARS BUITINCK AND MANY OTHERS TUTORIAL INTEGRATION INTO THE SCIKITLEARN DOCUMENTATION BY JAQUES GROBLER
  - ADDED COMPUTATIONAL PERFORMANCE DOCUMENTATION DISCUSSION AND EXAMPLES OF PREDICTION LATENCY THROUGHPUT AND DIFFERENT FACTORS THAT HAVE INFLUENCE OVER SPEED ADDITIONAL TIPS FOR BUILDING FASTER MODELS AND CHOOSING A RELEVANT COMPROMISE BETWEEN SPEED AND PREDICTIVE POWER BY EUSTACHE DIEMERT
- BUG FIXES
- FIXED BUG IN DECOMPOSITIONMINIBATCHDICTIONARYLEARNING PARTIALFIT WAS NOT WORKING PROPERLY
  - FIXED BUG IN LINEARMODELSTOCHASTICGRADIENT L1RATIO WAS USED AS 10 L1RATIO
  - FIXED BUG IN MULTICLASSONEVSONECLASSIFIER WITH STRING LABELS
  - FIXED A BUG IN LASSOCV ANDELASTICNETCV THEY WOULD NOT PRECOMPUTE THE GRAM MATRIX WITH PRECOMPUTETRUE ORPRECOMPUTEAUTO ANDNSAMPLES NFEATURES BY MANOJ KUMAR
  - FIXED INCORRECT ESTIMATION OF THE DEGREES OF FREEDOM IN FEATURESELECTIONFREGRESSION WHEN VARIATES ARE NOT CENTERED BY VIRGILE FRITSCH
  - FIXED A RACE CONDITION IN PARALLEL PROCESSING WITH PREDISPATCH ALL FOR INSTANCE IN CROSSVALScore BY OLIVIER GISEL
  - RAISE ERROR IN CLUSTERFEATUREAGGLOMERATION ANDCLUSTERWARDAGGLOMERATION WHEN NO SAMPLES ARE GIVEN RATHER THAN RETURNING MEANINGLESS CLUSTERING
  - FIXED BUG IN GRADIENTBOOSTINGGRADIENTBOOSTINGREGRESSOR WITHLOSSHUBER GAMMA MIGHT HAVE NOT BEEN INITIALIZED
  - FIXED FEATURE IMPORTANCES AS COMPUTED WITH A FOREST OF RANDOMIZED TREES WHEN FIT WITH SAMPLEWEIGHT NONE ANDOR WITH BOOTSTRAPTRUE BY GILLES LOUPPE
- 104 CHAPTER 1 WELCOME TO SCIKITLEARN

SCIKITLEARN USER GUIDE RELEASE 0213

API CHANGES SUMMARY

- SKLEARNHMM IS DEPRECATED ITS REMOVAL IS PLANNED FOR THE 017 RELEASE
- USE OF COVARIANCEELLIPTICENVELOP HAS NOW BEEN REMOVED AFTER DEPRECATION PLEASE USE COVARIANCEELLIPTICENVELOPE INSTEAD
- CLUSTERWARD IS DEPRECATED USE CLUSTERAGGLOMERATIVECLUSTERING INSTEAD
- CLUSTERWARDCLUSTERING IS DEPRECATED USE
- CLUSTERAGGLOMERATIVECLUSTERING INSTEAD
- CROSSVALIDATIONBOOTSTRAP IS DEPRECATED CROSSVALIDATIONKFOLD OR CROSSVALIDATIONSHUFFLESPLIT ARE RECOMMENDED INSTEAD
- DIRECT SUPPORT FOR THE SEQUENCE OF SEQUENCES OR LIST OF LISTS MULTILABEL FORMAT IS DEPRECATED TO CONVERT TO AND FROM THE SUPPORTED BINARY INDICATOR MATRIX FORMAT USE MULTILABELBINARIZER BY JOEL NOTHMAN
- ADD SCORE METHOD TO PCA FOLLOWING THE MODEL OF PROBABILISTIC PCA AND DEPRECATE PROBABILISTICPCA MODEL WHOSE SCORE IMPLEMENTATION IS NOT CORRECT THE COMPUTATION NOW ALSO EXPLOITS THE MATRIX INVERSION LEMMA FOR FASTER COMPUTATION BY ALEXANDRE GRAMFORT
- THE SCORE METHOD OF FACTORANALYSIS NOW RETURNS THE AVERAGE LOGLIKELIHOOD OF THE SAMPLES USE SCORESAMPLES TO GET LOGLIKELIHOOD OF EACH SAMPLE BY ALEXANDRE GRAMFORT
- GENERATING BOOLEAN MASKS THE SETTING INDICESFALSE FROM CROSSVALIDATION GENERATORS IS DEPRECATED SUP PORT FOR MASKS WILL BE REMOVED IN 017 THE GENERATORS HAVE PRODUCED ARRAYS OF INDICES BY DEFAULT SINCE 010 BY JOEL NOTHMAN
- 1D ARRAYS CONTAINING STRINGS WITH DTYPEOBJECT AS USED IN PANDAS ARE NOW CONSIDERED VALID CLASSIFICATION TARGETS THIS FIXES A REGRESSION FROM VERSION 013 IN SOME CLASSIFIERS BY JOEL NOTHMAN
- FIX WRONG EXPLAINEDVARIANCERATIO ATTRIBUTE IN RANDOMIZEDPCA BY ALEXANDRE GRAMFORT
- FIT ALPHAS FOR EACH L1RATIO INSTEAD OF MEANL1RATIO INLINEARMODELELASTICNETCV AND LINEARMODELLASSOCV THIS CHANGES THE SHAPE OF ALPHAS FROMNALPHAS TONL1RATIO NALPHAS IF THEL1RATIO PROVIDED IS A 1D ARRAY LIKE OBJECT OF LENGTH GREATER THAN ONE BY MANOJ KUMAR
- FIXLINEARMODELELASTICNETCV ANDLINEARMODELLASSOCV WHEN FITTING INTERCEPT AND INPUT DATA IS SPARSE THE AUTOMATIC GRID OF ALPHAS WAS NOT COMPUTED CORRECTLY AND THE SCALING WITH NORMALIZE WAS WRONG BY MANOJ KUMAR
- FIX WRONG MAXIMAL NUMBER OF FEATURES DRAWN MAXFEATURES AT EACH SPLIT FOR DECISION TREES RANDOM FORESTS AND GRADIENT TREE BOOSTING PREVIOUSLY THE COUNT FOR THE NUMBER OF DRAWN FEATURES STARTED ONLY AFTER ONE NON CONSTANT FEATURES IN THE SPLIT THIS BUG FIX WILL AFFECT COMPUTATIONAL AND GENERALIZATION PERFORMANCE OF THOSE ALGORITHMS IN THE PRESENCE OF CONSTANT FEATURES TO GET BACK PREVIOUS GENERALIZATION PERFORMANCE YOU SHOULD MODIFY THE VALUE OF MAXFEATURES BY ARNAUD JOLY
- FIX WRONG MAXIMAL NUMBER OF FEATURES DRAWN MAXFEATURES AT EACH SPLIT FOR ENSEMBLE EXTRATREESCLASSIFIER ANDENSEMBLEEXTRATREESREGRESSOR PREVIOUSLY ONLY NON CONSTANT FEATURES IN THE SPLIT WAS COUNTED AS DRAWN NOW CONSTANT FEATURES ARE COUNTED AS DRAWN FURTHERMORE AT LEAST ONE FEATURE MUST BE NON CONSTANT IN ORDER TO MAKE A VALID SPLIT THIS BUG FIX WILL AFFECT COMPUTATIONAL AND GEN ERALIZATION PERFORMANCE OF EXTRA TREES IN THE PRESENCE OF CONSTANT FEATURES TO GET BACK PREVIOUS GENERALIZATION PERFORMANCE YOU SHOULD MODIFY THE VALUE OF MAXFEATURES BY ARNAUD JOLY
- FIXUTILSCOMPUTECLASSWEIGHT WHENCLASSWEIGHTAUTO PREVIOUSLY IT WAS BROKEN FOR INPUT OF NONINTEGER DTYPE AND THE WEIGHTED ARRAY THAT WAS RETURNED WAS WRONG BY MANOJ KUMAR
- FIXCROSSVALIDATIONBOOTSTRAP TO RETURNVALUEERROR WHENNTRAIN NTEST N BY RONALD PHLYPO

SCIKITLEARN USER GUIDE RELEASE 0213

PEOPLE

LIST OF CONTRIBUTORS FOR RELEASE 015 BY NUMBER OF COMMITS

- 312 OLIVIER GRISEL
- 275 LARS BUITINCK
- 221 GAEL VAROQUAUX
- 148 ARNAUD JOLY
- 134 JOHANNES SCHÖNBERGER
- 119 GILLES LOUPPE
- 113 JOEL NOTHMAN
- 111 ALEXANDRE GRAMFORT
- 95 JAQUES GROBLER
- 89 DENIS ENGEMANN
- 83 PETER PRETTENHOFER
- 83 ALEXANDER FABISCH
- 62 MATHIEU BLONDEL
- 60 EUSTACHE DIEMERT
- 60 NELLE VAROQUAUX
- 49 MICHAEL BOMMARITO
- 45 MANOJKUMARS
- 28 KYLE KASTNER
- 26 ANDREAS MUELLER
- 22 NOEL DAWE
- 21 MAHESHAKYA WIJEWARDENA
- 21 BROOKE OSBORN
- 21 HAMZEH ALSALHI
- 21 JAKE VANDERPLAS
- 21 PHILIPPE GERVAIS
- 19 BALA SUBRAHMANYAM VARANASI
- 12 RONALD PHLYPO
- 10 MIKHAIL KOROBV
- 8 THOMAS UNTERTHINER
- 8 JEFFREY BLACKBURNE
- 8 ELTERMANN
- 8 BWIGNALL
- 7 ANKIT AGRAWAL
- 7 CJ CAREY



SCIKITLEARN USER GUIDE RELEASE 0213

- 6 DANIEL NOURI
- 6 CHEN LIU
- 6 MICHAEL EICKENBERG
- 6 UGURTHEMASTER
- 5 AARON SCHUMACHER
- 5 BAPTISTE LAGARDE
- 5 RAJAT KHANDUJA
- 5 ROBERT MCGIBBON
- 5 SERGIO PASCUAL
- 4 ALEXIS METAIREAU
- 4 IGNACIO ROSSI
- 4 VIRGILE FRITSCH
- 4 SEBASTIAN SÄGER
- 4 ILAMBHARATHI KANNIAH
- 4 SDENTON4
- 4 ROBERT LAYTON
- 4 ALYSSA
- 4 AMOS WATERLAND
- 3 ANDREW TULLOCH
- 3 MURAD
- 3 STEVEN MAUDE
- 3 KAROL PYSNIAK
- 3 JACQUES KVAM
- 3 CGOHLKE
- 3 CJLIN
- 3 MICHAEL BECKER
- 3 HAMZEH
- 3 ERIC JACOBSEN
- 3 JOHN COLLINS
- 3 KAUSHIK94
- 3 ERWIN MARSI
- 2 CSYTRACY
- 2 LK
- 2 VLAD NICULAE
- 2 LAURENT DIRER
- 2 ERIK SHILTS

SCIKITLEARN USER GUIDE RELEASE 0213

- 2 RAUL GARRETA
- 2 YOSHIKI VÁZQUEZ BAEZA
- 2 YUNG SIANG LIAU
- 2 ABHISHEK THAKUR
- 2 JAMES YU
- 2 ROHIT SIVAPRASAD
- 2 ROLAND SZABO
- 2 AMORMACHINE
- 2 ALEXIS MIGNON
- 2 OSCAR CARLSSON
- 2 NANTAS NARDELLI
- 2 JESS010
- 2 KOWALSKI87
- 2 ANDREW CLEGG
- 2 FEDERICO VAGGI
- 2 SIMON FRID
- 2 FÉLIXANTOINE FORTIN
- 1 RALF GOMMERS
- 1 TAFT
- 1 RONAN AMICEL
- 1 RUPESH KUMAR SRIVASTAVA
- 1 RYAN WANG
- 1 SAMUEL CHARRON
- 1 SAMUEL STJEAN
- 1 FABIAN PEDREGOSA
- 1 SKIPPER SEABOLD
- 1 STEFAN WALK
- 1 STEFAN VAN DER WALT
- 1 STEPHAN HOYER
- 1 ALLEN RIDDELL
- 1 VALENTIN HAENEL
- 1 VIJAY RAMESH
- 1 WILL MYERS
- 1 YAROSLAV HALCHENKO
- 1 YONI BENMESHULAM
- 1 YURY V ZAYTSEV

SCIKITLEARN USER GUIDE RELEASE 0213

- 1 ADRINJALALI
- 1 AI8RAHIM
- 1 ALEMAGNANI
- 1 ALEX
- 1 BENJAMIN WILSON
- 1 CHALMERLOWE
- 1 DZIKIE DRO `ZD`ZE
- 1 JAMESTWEBBER
- 1 MATRIXORZ
- 1 POPO
- 1 SAMUELA
- 1 FRANÇOIS BOULOGNE
- 1 ALEXANDER MEASURE
- 1 ETHAN WHITE
- 1 GUILHERME TREIN
- 1 HENDRIK HEUER
- 1 IVICAJOVIC
- 1 JAN HENDRIK METZEN
- 1 JEAN MICHEL ROULY
- 1 EDUARDO ARIÑO DE LA RUBIA
- 1 JELLE ZIJLSTRA
- 1 EDDY L O JANSSON
- 1 DENIS
- 1 JOHN
- 1 JOHN SCHMIDT
- 1 JORGE CAÑARDO ALASTUEY
- 1 JOSEPH PERLA
- 1 JOSHUA VREDEVOOGD
- 1 JOSÉ RICARDO
- 1 JULIEN MIOTTE
- 1 KEMAL EREN
- 1 KENTA SATO
- 1 DAVID COURNAPEAU
- 1 KYLE KELLEY
- 1 DANIELE MEDRI
- 1 LAURENT LUCE

SCIKITLEARN USER GUIDE RELEASE 0213

- 1 LAURENT PIERRON
- 1 LUIS PEDRO COELHO
- 1 DANIELWEITZENFELD
- 1 CRAIG THOMPSON
- 1 CHYIKWEI YAU
- 1 MATTHEW BRETT
- 1 MATTHIAS FEURER
- 1 MAX LINKE
- 1 CHRIS FILO GORGOLEWSKI
- 1 CHARLES EARL
- 1 MICHAEL HANKE
- 1 MICHELE ORRÙ
- 1 BRYAN LUNT
- 1 BRIAN KEARNS
- 1 PAUL BUTLER
- 1 PAWEŁ MANDERA
- 1 PETER
- 1 ANDREW ASH
- 1 PIETRO ZAMBELLI
- 1 STAUBDA

11714 VERSION 014

AUGUST 7 2013

CHANGELOG

- MISSING VALUES WITH SPARSE AND DENSE MATRICES CAN BE IMPUTED WITH THE TRANSFORMER PREPROCESSING IMPUTER BY NICOLAS TRÉSEGNIÉ
  - THE CORE IMPLEMENTATION OF DECISIONS TREES HAS BEEN REWRITTEN FROM SCRATCH ALLOWING FOR FASTER TREE INDUCTION AND LOWER MEMORY CONSUMPTION IN ALL TREEBASED ESTIMATORS BY GILLES LOUPPE
  - ADDEDENSEMBLEADABOOSTCLASSIFIER ANDENSEMBLEADABOOSTREGRESSOR BY NOEL DAWE AND GILLES LOUPPE SEE THE ADABOOST SECTION OF THE USER GUIDE FOR DETAILS AND EXAMPLES
  - ADDEDGRIDSEARCHRANDOMIZEDSEARCHCV ANDGRIDSEARCHPARAMETERSAMPLER FOR RANDOM IZED HYPERPARAMETER OPTIMIZATION BY ANDREAS MÜLLER
  - ADDED BICLUSTERING ALGORITHMS SKLEARNCLUSTERBICLUSTERSPECTRALCOCLUSTERING AND SKLEARNCLUSTERBICLUSTERSPECTRALBICLUSTERING DATA GENERATION METHODS SKLEARN DATASETSMAKEBICLUSTERS ANDSKLEARNDATASETSMAKECHECKERBOARD AND SCORING MET RICS SKLEARNMETRICSCONSENSUSSCORE BY KEMAL EREN
  - ADDED RESTRICTED BOLTZMANN MACHINES NEURALNETWORKBERNOULLIRBM BY YANN DAUPHIN
- 110 CHAPTER 1 WELCOME TO SCIKITLEARN

SCIKITLEARN USER GUIDE RELEASE 0213

- PYTHON 3 SUPPORT BY JUSTIN VINCENT LARS BUITINCK SUBHODEEP MOITRA AND OLIVIER GRISEL ALL TESTS NOW PASS UNDER PYTHON 33

- ABILITY TO PASS ONE PENALTY ALPHA VALUE PER TARGET IN LINEARMODELRIDGE BY EICKENBERG AND MATHIEU BLONDEL

- FIXEDSKLEARNLINEARMODELSTOCHASTICGRADIENTPY L2 REGULARIZATION ISSUE MINOR PRACTICAL SIGNIFICANCE BY NORBERT CROMBACH AND MATHIEU BLONDEL

- ADDED AN INTERACTIVE VERSION OF ANDREAS MÜLLER'S MACHINE LEARNING CHEAT SHEET FOR SCIKITLEARN TO THE DOCUMENTATION SEE CHOOSING THE RIGHT ESTIMATOR BY JAQUES GROBLER

•GRIDSEARCHGRIDSEARCHCV ANDCROSSVALIDATIONCROSSVALSCORE NOW SUPPORT THE USE OF ADVANCED SCORING FUNCTION SUCH AS AREA UNDER THE ROC CURVE AND FBETA SCORES SEE THE SCORING PARAMETER DEFINING MODEL EVALUATION RULES FOR DETAILS BY ANDREAS MÜLLER AND LARS BUITINCK PASSING A FUNCTION FROM SKLEARNMETRICS ASSCOREFUNC IS DEPRECATED

- MULTILABEL CLASSIFICATION OUTPUT IS NOW SUPPORTED BY METRICSACCURACYScore METRICZEROONELOSS METRICSF1Score METRICSFBScore METRICS CLASSIFICATIONREPORT METRICSPRECISIONScore ANDMETRICSRECALLScore BY

ARNAUD JOLY

- TWO NEW METRICS METRICSHAMMINGLOSS ANDMETRICSJACCARDSIMILARITYSCORE ARE ADDED WITH MULTILABEL SUPPORT BY ARNAUD IOLY

• **SPEED AND MEMORY USAGE IMPROVEMENTS IN FEATUREEXTRACTIONTEXTCOUNTVECTORIZER AND FEATUREEXTRACTIONTEXTTFIDFVECTORIZER** BY JOCHEN WERSDÖRFER AND ROMAN SINAYEV

- THEMINDF PARAMETER IN FEATUREEXTRACTIONTEXTCOUNTVECTORIZER AND

FEATUREEXTRACTIONTEXTTFIDFVECTORIZER WHICH USED TO BE 2 HAS BEEN RESET TO 1 TO

AVOID UNPLEASANT SURPRISES EMPTY VOCABULARIES FOR NOVICE  
VALUE OF AT LEAST 2 IS STILL RECOMMENDED FOR PRACTICAL USE

- SVM LINEAR SVC LINEAR MODEL SGD CLASSIFIER AND LINEAR MODEL SGD REGRESSOR NOW

HAVE ASPARSIFY METHOD THAT CONVERTS THEIR COEF INTO A SPARSE MATRIX MEANING STORED MODELS TRAINED USING THESE ESTIMATORS CAN BE MADE MUCH MORE COMPACT

- LINEARMODELSGDCLASSIFIER NOW PRODUCES MULTICLASS PROBABILITY ESTIMATES WHEN TRAINED UNDER LOG LOSS OR MODIFIED HUBER LOSS

- [HYPERLINKS TO DOCUMENTATION IN EXAMPLE CODE ON THE WEBSITE BY MARTIN LUESSI](#)

• FIXED BUG IN PREPROCESSINGMINMAXSCALER CAUSING INCORRECT SCALING OF THE FEATURES FOR NONDEFAULT FEATURERANGE SETTINGS BY ANDREAS MÜLLER

- MAXFEATURES INTREEDECISIONTREECLASSIFIER TREEDECISIONTREEREgressor AND ALL

## DERIVED ENSEMBLE ESTIMATORS NOW SUPPORTS PERCENTAGE VALUES BY GILLES LOUPPE

- PERFORMANCE IMPROVEMENTS IN ISOTONIC ISOTONIC REGRESSION BY NELLE VAROQUAUX

• METRICS ACCURACY SCORE HAS AN OPTION NORMALIZE TO RETURN THE FRACTION OR THE NUMBER OF CORRECTLY CLASSIFIED SAMPLE BY ARNAUD JOLY

- ADDEDMETRICSLOGLOSS THAT COMPUTES LOG LOSS AKA CROSSENTROPY LOSS BY JOCHEN WERSDÖRFER AND LARS BUITINCK

- A BUG THAT CAUSED ENSEMBLEADABOOSTCLASSIFIER 'S TO OUTPUT INCORRECT PROBABILITIES HAS BEEN FIXED

- FEATURE SELECTORS NOW SHARE A MIXIN PROVIDING CONSISTENT TRANSFORM INVERSE TRANSFORM AND

# GETSUPPORT METHODS BY JOEL NOTHMAN

- A `FITTEDGRIDSEARCH` `GRIDSEARCHCV` OR `GRIDSEARCHRANDOMIZEDSEARCHCV` CAN NOW GENERALLY

BE PICKLED BY JOEL NOTHMAN

117 PREVIOUS RELEASES 111

SCIKITLEARN USER GUIDE RELEASE 0213

- REFACTORED AND VECTORIZED IMPLEMENTATION OF METRICSROCCURVE ANDMETRICS PRECISIONRECALLCURVE BY JOEL NOTHMAN
- THE NEW ESTIMATOR SKLEARNDECOMPOSITIONTRUNCATEDSVD PERFORMS DIMENSIONALITY REDUCTION USING SVD ON SPARSE MATRICES AND CAN BE USED FOR LATENT SEMANTIC ANALYSIS LSA BY LARS BUITINCK
- ADDED SELFCONTAINED EXAMPLE OF OUTFCORE LEARNING ON TEXT DATA OUTFCORE CLASSIFICATION OF TEXT DOCUMENTS BY EUSTACHE DIEMERT
- THE DEFAULT NUMBER OF COMPONENTS FOR SKLEARNDECOMPOSITIONRANDOMIZEDPCA IS NOW CORRECTLY DOCUMENTED TO BE NFEATURES THIS WAS THE DEFAULT BEHAVIOR SO PROGRAMS USING IT WILL CONTINUE TO WORK AS THEY DID
- SKLEARNCLUSTERKMEANS NOW FITS SEVERAL ORDERS OF MAGNITUDE FASTER ON SPARSE DATA THE SPEEDUP DEPENDS ON THE SPARSITY BY LARS BUITINCK
- REDUCE MEMORY FOOTPRINT OF FASTICA BY DENIS ENGEMANN AND ALEXANDRE GRAMFORT
- VERBOSE OUTPUT IN SKLEARNENSEMBLEGRADIENTBOOSTING NOW USES A COLUMN FORMAT AND PRINTS PROGRESS IN DECREASING FREQUENCY IT ALSO SHOWS THE REMAINING TIME BY PETER PRETTENHOFER
- SKLEARNENSEMBLEGRADIENTBOOSTING PROVIDES OUTFBAG IMPROVEMENT OOBIMPROVEMENT RATHER THAN THE OOB SCORE FOR MODEL SELECTION AN EXAMPLE THAT SHOWS HOW TO USE OOB ESTIMATES TO SELECT THE NUMBER OF TREES WAS ADDED BY PETER PRETTENHOFER
- MOST METRICS NOW SUPPORT STRING LABELS FOR MULTICLASS CLASSIFICATION BY ARNAUD JOLY AND LARS BUITINCK
- NEW ORTHOGONALMATCHINGPURSUITCV CLASS BY ALEXANDRE GRAMFORT AND VLAD NICULAE
- FIXED A BUG IN SKLEARNCOVARIANCEGRAPHLASSOCV THE 'ALPHAS' PARAMETER NOW WORKS AS EXPECTED WHEN GIVEN A LIST OF VALUES BY PHILIPPE GERVAIS
- FIXED AN IMPORTANT BUG IN SKLEARNCOVARIANCEGRAPHLASSOCV THAT PREVENTED ALL FOLDS PROVIDED BY A CV OBJECT TO BE USED ONLY THE FIRST 3 WERE USED WHEN PROVIDING A CV OBJECT EXECUTION TIME MAY THUS INCREASE SIGNIFICANTLY COMPARED TO THE PREVIOUS VERSION BUG RESULTS ARE CORRECT NOW BY PHILIPPE GERVAIS
- CROSSVALIDATIONCROSSVALSCORE AND THEGRIDSEARCH MODULE IS NOW TESTED WITH MULTI OUTPUT DATA BY ARNAUD JOLY
- DATASETSMAKEMULTILABELCLASSIFICATION CAN NOW RETURN THE OUTPUT IN LABEL INDICATOR MULTILABEL FORMAT BY ARNAUD JOLY
- KNEAREST NEIGHBORS NEIGHBORSKNEIGHBORSREGRESSOR ANDNEIGHBORS RADIUSNEIGHBORSREGRESSOR AND RADIUS NEIGHBORS NEIGHBORSRADIUSNEIGHBORSREGRESSOR ANDNEIGHBORSRADIUSNEIGHBORSCLASSIFIER SUPPORT MULTIOUTPUT DATA BY ARNAUD JOLY
- RANDOM STATE IN LIBSVMBASED ESTIMATORS SVM SVC ONECLASS SVM SVMSVR SVMNUSVR CAN NOW BE CONTROLLED THIS IS USEFUL TO ENSURE CONSISTENCY IN THE PROBABILITY ESTIMATES FOR THE CLASSIFIERS TRAINED WITHPROBABILITYTRUE BY VLAD NICULAE
- OUTFCORE LEARNING SUPPORT FOR DISCRETE NAIVE BAYES CLASSIFIERS SKLEARNNAIVEBAYES MULTINOMIALNB ANDSKLEARNNAIVEBAYESBERNOULLINB BY ADDING THE PARTIALFIT METHOD BY OLIVIER GRISEL
- NEW WEBSITE DESIGN AND NAVIGATION BY GILLES LOUPPE NELLE VAROQUAUX VINCENT MICHEL AND ANDREAS MÜLLER
- IMPROVED DOCUMENTATION ON MULTICLASS MULTILABEL AND MULTIOUTPUT CLASSIFICATION BY YANNICK SCHWARTZ AND ARNAUD JOLY
- BETTER INPUT AND ERROR HANDLING IN THE METRICS MODULE BY ARNAUD JOLY AND JOEL NOTHMAN
- SPEED OPTIMIZATION OF THE HMM MODULE BY MIKHAIL KOROBOV
- SIGNIFICANT SPEED IMPROVEMENTS FOR SKLEARNCLUSTERDBSCAN BY CLEVERLESS

SCIKITLEARN USER GUIDE RELEASE 0213

API CHANGES SUMMARY

- THEAUUCSCORE WAS RENAMED ROCAUCSCORE
- TESTING SCIKITLEARN WITH SKLEARNTEST IS DEPRECATED USE NOSETESTS SKLEARN FROM THE COMMAND LINE
- FEATURE IMPORTANCES IN TREEDECISIONTREECLASSIFIER TREEDECISIONTREEREgressor AND ALL DERIVED ENSEMBLE ESTIMATORS ARE NOW COMPUTED ON THE FLY WHEN ACCESSING THE FEATUREIMPORTANCES ATTRIBUTE SETTING COMPUTEIMPORTANCESTRUE IS NO LONGER REQUIRED BY GILLES LOUPPE
- LINEARMODELLASSOPATH ANDLINEARMODELENETPATH CAN RETURN ITS RESULTS IN THE SAME FORMAT AS THAT OFLINEARMODELLARSPATH THIS IS DONE BY SETTING THE RETURNMODELS PARAMETER TO FALSE BY JAQUES GROBLER AND ALEXANDRE GRAMFORT
- GRIDSEARCHITERGRID WAS RENAMED TO GRIDSEARCHPARAMETERGRID
- FIXED BUG IN KFOLD CAUSING IMPERFECT CLASS BALANCE IN SOME CASES BY ALEXANDRE GRAMFORT AND TADEJ JANEŽ
- SKLEARNNEIGHBORSBALLTREE HAS BEEN REFACTORED AND A SKLEARNNEIGHBORSKDTREE HAS BEEN ADDED WHICH SHARES THE SAME INTERFACE THE BALL TREE NOW WORKS WITH A WIDE VARIETY OF DISTANCE METRICS BOTH CLASSES HAVE MANY NEW METHODS INCLUDING SINGLETREE AND DUALTREE QUERIES BREADTHFIRST AND DEPTHFIRST SEARCHING AND MORE ADVANCED QUERIES SUCH AS KERNEL DENSITY ESTIMATION AND 2POINT CORRELATION FUNCTIONS BY JAKE VANDERPLAS
- SUPPORT FOR SCIPYSPATIALCKDTREE WITHIN NEIGHBORS QUERIES HAS BEEN REMOVED AND THE FUNCTIONALITY REPLACED WITH THE NEW KDTREE CLASS
- SKLEARNNEIGHBORSKERNELDENSITY HAS BEEN ADDED WHICH PERFORMS EFFICIENT KERNEL DENSITY ESTIMATION WITH A VARIETY OF KERNELS
- SKLEARNDECOMPOSITIONKERNELPCA NOW ALWAYS RETURNS OUTPUT WITH NCOMPONENTS COMPONENTS UNLESS THE NEW PARAMETER REMOVEZEROEIG IS SET TOTRUE THIS NEW BEHAVIOR IS CONSISTENT WITH THE WAY KERNEL PCA WAS ALWAYS DOCUMENTED PREVIOUSLY THE REMOVAL OF COMPONENTS WITH ZERO EIGENVALUES WAS TACITLY PERFORMED ON ALL DATA
- GCVMODEAUTO NO LONGER TRIES TO PERFORM SVD ON A DENSIFIED SPARSE MATRIX IN SKLEARN LINEARMODELRIDGECV
- SPARSE MATRIX SUPPORT IN SKLEARNDECOMPOSITIONRANDOMIZEDPCA IS NOW DEPRECATED IN FAVOR OF THE NEWTRUNCATEDSVD
- CROSSVALIDATIONKFOLD ANDCROSSVALIDATIONSTRATIFIEDKFOLD NOW ENFORCE NFOLDS 2 OTHERWISE A VALUEERROR IS RAISED BY OLIVIER GRISEL
- DATASETSLOADFILES 'SCHARSET ANDCHARSETERRORS PARAMETERS WERE RENAMED ENCODING AND DECODEERRORS
- ATTRIBUTE OOBSCORE INSKLEARNENSEMBLEGRADIENTBOOSTINGREGRESSOR AND SKLEARNENSEMBLEGRADIENTBOOSTINGCLASSIFIER IS DEPRECATED AND HAS BEEN REPLACED BY OOBIMPROVEMENT
- ATTRIBUTES IN ORTHOGONALMATCHINGPURSUIT HAVE BEEN DEPRECATED COPYX GRAM AND PRECOMPUTEGRAM RENAMED PRECOMPUTE FOR CONSISTENCY SEE 2224
- SKLEARNPREPROCESSINGSTANDARDSCALER NOW CONVERTS INTEGER INPUT TO FLOAT AND RAISES A WARNING PREVIOUSLY IT ROUNDED FOR DENSE INTEGER INPUT
- SKLEARNMULTICLASSONEVSRESTCLASSIFIER NOW HAS A DECISIONFUNCTION METHOD THIS WILL RETURN THE DISTANCE OF EACH SAMPLE FROM THE DECISION BOUNDARY FOR EACH CLASS AS LONG AS THE UNDERLYING ESTIMATORS IMPLEMENT THE DECISIONFUNCTION METHOD BY KYLE KASTNER
- BETTER INPUT VALIDATION WARNING ON UNEXPECTED SHAPES FOR Y

SCIKITLEARN USER GUIDE RELEASE 0213

PEOPLE

LIST OF CONTRIBUTORS FOR RELEASE 014 BY NUMBER OF COMMITS

- 277 GILLES LOUPPE
- 245 LARS BUITINCK
- 187 ANDREAS MUELLER
- 124 ARNAUD JOLY
- 112 JAQUES GROBLER
- 109 GAEL VAROQUAUX
- 107 OLIVIER GRISEL
- 102 NOEL DAWE
- 99 KEMAL EREN
- 79 JOEL NOTHMAN
- 75 JAKE VANDERPLAS
- 73 NELLE VAROQUAUX
- 71 VLAD NICULAE
- 65 PETER PRETTENHOFER
- 64 ALEXANDRE GRAMFORT
- 54 MATHIEU BLONDEL
- 38 NICOLAS TRÉSEGNIE
- 35 EUSTACHE
- 27 DENIS ENGEMANN
- 25 YANN N DAUPHIN
- 19 JUSTIN VINCENT
- 17 ROBERT LAYTON
- 15 DOUG COLEMAN
- 14 MICHAEL EICKENBERG
- 13 ROBERT MARCHMAN
- 11 FABIAN PEDREGOSA
- 11 PHILIPPE GERVAIS
- 10 JIM HOLMSTRÖM
- 10 TADEJ JANEŽ
- 10 SYHW
- 9 MIKHAIL KOROBOV
- 9 STEVEN DE GRYZE
- 8 SERGEYF
- 7 BEN ROOT

114 CHAPTER 1 WELCOME TO SCIKITLEARN



SCIKITLEARN USER GUIDE RELEASE 0213

- 7 HRISHIKESH HUILGOLKAR
- 6 KYLE KASTNER
- 6 MARTIN LUESSI
- 6 ROB SPEER
- 5 FEDERICO VAGGI
- 5 RAUL GARRETA
- 5 ROB ZINKOV
- 4 KEN GEIS
- 3 A FLAXMAN
- 3 DENTON COCKBURN
- 3 DOUGAL SUTHERLAND
- 3 IAN OZSVALD
- 3 JOHANNES SCHÖNBERGER
- 3 ROBERT MCGIBBON
- 3 ROMAN SINAYEV
- 3 SZABO ROLAND
- 2 DIEGO MOLLA
- 2 IMRAN HAQUE
- 2 JOCHEN WERSDÖRFER
- 2 SERGEY KARAYEV
- 2 YANNICK SCHWARTZ
- 2 JAMESTWEBBER
- 1 ABHIJEET KOLHE
- 1 ALEXANDER FABISCH
- 1 BASTIAAN VAN DEN BERG
- 1 BENJAMIN PETERSON
- 1 DANIEL VELKOV
- 1 FAZLUL SHAHRIAR
- 1 FELIX BROCKHERDE
- 1 FÉLIXANTOINE FORTIN
- 1 HARIKRISHNAN S
- 1 JACK HALE
- 1 JAKEMICK
- 1 JAMES MCDERMOTT
- 1 JOHN BENEDIKTSSON
- 1 JOHN ZWINCK

SCIKITLEARN USER GUIDE RELEASE 0213

- 1 JOSHUA VREDEVOOGD
- 1 JUSTIN PATI
- 1 KEVIN HUGHES
- 1 KYLE KELLEY
- 1 MATTHIAS EKMAN
- 1 MIROSLAV SHUBERNETSKIY
- 1 NAOKI ORII
- 1 NORBERT CROMBACH
- 1 RAFAEL CUNHA DE ALMEIDA
- 1 ROLANDO ESPINOZA LA FUENTE
- 1 SEAMUS ABSHERE
- 1 SERGEY FELDMAN
- 1 SERGIO MEDINA
- 1 STEFANO LATTARINI
- 1 STEVE KOCH
- 1 STURLA MOLDEN
- 1 THOMAS JAROSCH
- 1 YAROSLAV HALCHENKO

11715 VERSION 0131

FEBRUARY 23 2013

THE 0131 RELEASE ONLY FIXES SOME BUGS AND DOES NOT ADD ANY NEW FUNCTIONALITY

CHANGELOG

- FIXED A TESTING ERROR CAUSED BY THE FUNCTION CROSSVALIDATIONTRAINTESTSPLIT BEING INTERPRETED AS A TEST BY YAROSLAV HALCHENKO
- FIXED A BUG IN THE REASSIGNMENT OF SMALL CLUSTERS IN THE CLUSTERMINIBATCHKMEANS BY GAELE VAROQUAUX
- FIXED DEFAULT VALUE OF GAMMA IN DECOMPOSITIONKERNELPCA BY LARS BUITINCK
- UPDATED JOBLIB TO 070D BY GAELE VAROQUAUX
- FIXED SCALING OF THE DEVIANCE IN ENSEMBLEGRADIENTBOOSTINGCLASSIFIER BY PETER PRETTENHOFER
- BETTER TIEBREAKING IN MULTICLASSONEVSONECLASSIFIER BY ANDREAS MÜLLER
- OTHER SMALL IMPROVEMENTS TO TESTS AND DOCUMENTATION

116 CHAPTER 1 WELCOME TO SCIKITLEARN

SCIKITLEARN USER GUIDE RELEASE 0213

PEOPLE

LIST OF CONTRIBUTORS FOR RELEASE 0131 BY NUMBER OF COMMITS

- 16 LARS BUITINCK
- 12 ANDREAS MÜLLER
- 8 GAELE VAROQUAUX
- 5 ROBERT MARCHMAN
- 3 PETER PRETTENHOFER
- 2 HRISHIKESH HUILGOLKAR
- 1 BASTIAAN VAN DEN BERG
- 1 DIEGO MOLLA
- 1 GILLES LOUPPE
- 1 MATHIEU BLONDEL
- 1 NELLE VAROQUAUX
- 1 RAFAEL CUNHA DE ALMEIDA
- 1 ROLANDO ESPINOZA LA FUENTE
- 1 VLAD NICULAE
- 1 YAROSLAV HALCHENKO

11716 VERSION 013

JANUARY 21 2013

NEW ESTIMATOR CLASSES

- DUMMYDUMMYCLASSIFIER ANDDUMMYDUMMYREGRESSOR TWO DATAINDEPENDENT PREDICTORS BY MATHIEU BLONDEL USEFUL TO SANITYCHECK YOUR ESTIMATORS SEE DUMMY ESTIMATORS IN THE USER GUIDE MULTIOUTPUT SUPPORT ADDED BY ARNAUD JOLY
- DECOMPOSITIONFACTORANALYSIS A TRANSFORMER IMPLEMENTING THE CLASSICAL FACTOR ANALYSIS BY CHRIS TIAN OSENDORFER AND ALEXANDRE GRAMFORT SEE FACTOR ANALYSIS IN THE USER GUIDE
- FEATUREEXTRACTIONFEATUREHASHER A TRANSFORMER IMPLEMENTING THE “HASHING TRICK” FOR FAST LOWMEMORY FEATURE EXTRACTION FROM STRING FIELDS BY LARS BUITINCK AND FEATUREEXTRACTIONTEXT HASHINGVECTORIZER FOR TEXT DOCUMENTS BY OLIVIER GRISEL SEE FEATURE HASHING ANDVECTORIZING A LARGE TEXT CORPUS WITH THE HASHING TRICK FOR THE DOCUMENTATION AND SAMPLE USAGE
- PIPELINEFEATUREUNION A TRANSFORMER THAT CONCATENATES RESULTS OF SEVERAL OTHER TRANSFORMERS BY ANDREAS MÜLLER SEE FEATUREUNION COMPOSITE FEATURE SPACES IN THE USER GUIDE
- RANDOMPROJECTIONGAUSSIANRANDOMPROJECTION RANDOMPROJECTION SPARSERANDOMPROJECTION AND THE FUNCTION RANDOMPROJECTION
- JOHNSONLINDENSTRAUSSMINDIM THE FIRST TWO ARE TRANSFORMERS IMPLEMENTING GAUSSIAN AND SPARSE RANDOM PROJECTION MATRIX BY OLIVIER GRISEL AND ARNAUD JOLY SEE RANDOM PROJECTION IN THE USER GUIDE
- KERNELAPPROXIMATIONNYSTROEM A TRANSFORMER FOR APPROXIMATING ARBITRARY KERNELS BY ANDREAS MÜLLER SEE NYSTROEM METHOD FOR KERNEL APPROXIMATION IN THE USER GUIDE

117 PREVIOUS RELEASES 117

SCIKITLEARN USER GUIDE RELEASE 0213

- PREPROCESSINGONEHOTENCODER A TRANSFORMER THAT COMPUTES BINARY ENCODINGS OF CATEGORICAL FEATURES BY ANDREAS MÜLLER SEE ENCODING CATEGORICAL FEATURES IN THE USER GUIDE
- LINEARMODELPASSIVEAGGRESSIVECLASSIFIER AND LINEARMODELPASSIVEAGGRESSIVEREGRESSOR PREDICTORS IMPLEMENTING AN EFFICIENT STOCHASTIC OPTIMIZATION FOR LINEAR MODELS BY ROB ZINKOV AND MATHIEU BLONDEL SEE PASSIVE AGGRESSIVE ALGORITHMS IN THE USER GUIDE
- ENSEMBLERANDOMTREESEMBEDDING A TRANSFORMER FOR CREATING HIGHDIMENSIONAL SPARSE REPRESENTATIONS USING ENSEMBLES OF TOTALLY RANDOM TREES BY ANDREAS MÜLLER SEE TOTALLY RANDOM TREES EMBEDDING IN THE USER GUIDE
- MANIFOLDSPECTRALEMBEDDING AND FUNCTION MANIFOLDSPECTRALEMBEDDING IMPLEMENTING THE “LAPLACIAN EIGENMAPS” TRANSFORMATION FOR NONLINEAR DIMENSIONALITY REDUCTION BY WEI LI SEE SPECTRAL EMBEDDING IN THE USER GUIDE
- ISOTONICISOTONICREGRESSION BY FABIAN PEDREGOSA ALEXANDRE GRAMFORT AND NELLE VAROQUAUX CHANGELOG
- METRICSZEROONELOSS FORMERLYMETRICSZEROONE NOW HAS OPTION FOR NORMALIZED OUTPUT THAT REPORTS THE FRACTION OF MISCLASSIFICATIONS RATHER THAN THE RAW NUMBER OF MISCLASSIFICATIONS BY KYLE BEAUCHAMP
- TREEDECISIONTREECLASSIFIER AND ALL DERIVED ENSEMBLE MODELS NOW SUPPORT SAMPLE WEIGHTING BY NOEL DAWE AND GILLES LOUPPE
- SPEEDUP IMPROVEMENT WHEN USING BOOTSTRAP SAMPLES IN FORESTS OF RANDOMIZED TREES BY PETER PRETTENHOFER AND GILLES LOUPPE
- PARTIALDEPENDENCE PLOTS FOR GRADIENT TREE BOOSTING INENSEMBLE PARTIALDEPENDENCEPARTIALDEPENDENCE BY PETER PRETTENHOFER SEE SPHXGLRAUTOEXAMPLESENSEMBLEPLOTPARTIALDEPENDENCEPY FOR AN EXAMPLE
- THE TABLE OF CONTENTS ON THE WEBSITE HAS NOW BEEN MADE EXPANDABLE BY JAQUES GROBLER
- FEATURESELECTIONSELECTPERCENTILE NOW BREAKS TIES DETERMINISTICALLY INSTEAD OF RETURNING ALL EQUALLY RANKED FEATURES
- FEATURESELECTIONSELECTKBBEST ANDFEATURESELECTIONSELECTPERCENTILE ARE MORE NUMERICALLY STABLE SINCE THEY USE SCORES RATHER THAN PVALUES TO RANK RESULTS THIS MEANS THAT THEY MIGHT SOMETIMES SELECT DIFFERENT FEATURES THAN THEY DID PREVIOUSLY
- RIDGE REGRESSION AND RIDGE CLASSIFICATION FITTING WITH SPARSECG SOLVER NO LONGER HAS QUADRATIC MEMORY COMPLEXITY BY LARS BUITINCK AND FABIAN PEDREGOSA
- RIDGE REGRESSION AND RIDGE CLASSIFICATION NOW SUPPORT A NEW FAST SOLVER CALLED LSQR BY MATHIEU BLONDEL
- SPEED UP OF METRICSPRECISIONRECALLCURVE BY CONRAD LEE
- ADDED SUPPORT FOR READINGWRITING SVMLIGHT FILES WITH PAIRWISE PREFERENCE ATTRIBUTE QID IN SVMLIGHT FILE FORMAT INDATASETSDUMPSVMLIGHTFILE ANDDATASETSLOADSVMLIGHTFILE BY FABIAN PEDREGOSA
- FASTER AND MORE ROBUST METRICSCONFUSIONMATRIX ANDCLUSTERING PERFORMANCE EVALUATION BY WEI LI
- CROSSVALIDATIONCROSSVALSCORE NOW WORKS WITH PRECOMPUTED KERNELS AND AFFINITY MATRICES BY ANDREAS MÜLLER
- LARS ALGORITHM MADE MORE NUMERICALLY STABLE WITH HEURISTICS TO DROP REGRESSORS TOO CORRELATED AS WELL AS TO STOP THE PATH WHEN NUMERICAL NOISE BECOMES PREDOMINANT BY GAELE VAROQUAUX
- FASTER IMPLEMENTATION OF METRICSPRECISIONRECALLCURVE BY CONRAD LEE
- NEW KERNEL METRICSCHI2KERNEL BY ANDREAS MÜLLER OFTEN USED IN COMPUTER VISION APPLICATIONS

SCIKITLEARN USER GUIDE RELEASE 0213

- FIX OF LONGSTANDING BUG IN NAIVEBAYESBERNOULLINB FIXED BY SHAUN JACKMAN
  - IMPLEMENTED PREDICTPROBA INMULTICLASSONEVSRESTCLASSIFIER BY ANDREW WINTERMAN
  - IMPROVE CONSISTENCY IN GRADIENT BOOSTING ESTIMATORS ENSEMBLEGRADIENTBOOSTINGREGRESSOR AND ENSEMBLEGRADIENTBOOSTINGCLASSIFIER USE THE ESTIMATOR TREEDECISIONTREEREgressor INSTEAD OF THE TREETREETREE DATA STRUCTURE BY ARNAUD JOLY
  - FIXED A FLOATING POINT EXCEPTION IN THE DECISION TREES MODULE BY SEBERG
  - FIXMETRICSR2SCORE FAILS WHEN YTRUE HAS ONLY ONE CLASS BY WEI LI
  - ADD THE METRICSMEANABSOLUTEERROR FUNCTION WHICH COMPUTES THE MEAN ABSOLUTE ERROR THE METRICSMEANSQUAREDERROR METRICSMEANABSOLUTEERROR ANDMETRICSR2SCORE METRICS SUPPORT MULTIOUTPUT BY ARNAUD JOLY
  - FIXEDCLASSWEIGHT SUPPORT INSVMLINEARSVC ANDLINEARMODELLOGISTICREGRESSION BY ANDREAS MÜLLER THE MEANING OF CLASSWEIGHT WAS REVERSED AS ERRONEOUSLY HIGHER WEIGHT MEANT LESS POSITIVES OF A GIVEN CLASS IN EARLIER RELEASES
  - IMPROVE NARRATIVE DOCUMENTATION AND CONSISTENCY IN SKLEARNMETRICS FOR REGRESSION AND CLASSIFICATION METRICS BY ARNAUD JOLY
  - FIXED A BUG IN SKLEARNVMSVC WHEN USING CSRMATRICES WITH UNSORTED INDICES BY XINFAN MENG AND AN DREAS MÜLLER
  - MINIBATCHKMEANS ADD RANDOM REASSIGNMENT OF CLUSTER CENTERS WITH LITTLE OBSERVATIONS ATTACHED TO THEM BY GAELEVAROQUAUX
- API CHANGES SUMMARY
- RENAMED ALL OCCURRENCES OF NATOMS TONCOMPONENTS FOR CONSISTENCY THIS APPLIES TO DECOMPOSITIONDICTIONARYLEARNING DECOMPOSITION MINIBATCHDICTIONARYLEARNING DECOMPOSITIONDICTLEARNING DECOMPOSITION DICTLEARNINGONLINE
  - RENAMED ALL OCCURRENCES OF MAXITERS TOMAXITER FOR CONSISTENCY THIS APPLIES TO SEMISUPERVISEDLABELPROPAGATION ANDSEMISUPERVISEDLABELPROPAGATION LABELSPREADING
  - RENAMED ALL OCCURRENCES OF LEARNRATE TOLEARNINGRATE FOR CONSISTENCY IN ENSEMBLE BASEGRADIENTBOOSTING ANDENSEMBLEGRADIENTBOOSTINGREGRESSOR
  - THE MODULE SKLEARNLINEARMODELSPARSE IS GONE SPARSE MATRIX SUPPORT WAS ALREADY INTEGRATED INTO THE “REGULAR” LINEAR MODELS
  - SKLEARNMETRICSMEANSQUAREERROR WHICH INCORRECTLY RETURNED THE ACCUMULATED ERROR WAS REMOVED USE MEANSQUAREDERROR INSTEAD
  - PASSINGCLASSWEIGHT PARAMETERS TO FIT METHODS IS NO LONGER SUPPORTED PASS THEM TO ESTIMATOR CONSTRUCTORS INSTEAD
  - GMMS NO LONGER HAVE DECODE ANDRVS METHODS USE THE SCORE PREDICT ORSAMPLE METHODS INSTEAD
  - THESOLVER FIT OPTION IN RIDGE REGRESSION AND CLASSIFICATION IS NOW DEPRECATED AND WILL BE REMOVED IN V014 USE THE CONSTRUCTOR OPTION INSTEAD
  - FEATUREEXTRACTIONTEXTDICTVECTORIZER NOW RETURNS SPARSE MATRICES IN THE CSR FORMAT INSTEAD OF COO
  - RENAMED KINCROSSVALIDATIONKFOLD ANDCROSSVALIDATIONSTRATIFIEDKFOLD TO NFOLDS RENAMEDNBOOTSTRAPS TONITER INCROSSVALIDATIONBOOTSTRAP
- 117 PREVIOUS RELEASES 119

SCIKITLEARN USER GUIDE RELEASE 0213

- RENAMED ALL OCCURRENCES OF NITERATIONS TONITER FOR CONSISTENCY THIS APPLIES TO CROSSVALIDATIONSHUFFLESPLIT CROSSVALIDATIONSTRATIFIEDSHUFFLESPLIT UTILSRANDOMIZEDRANGEFINDER ANDUTILSRANDOMIZEDSVD
- REPLACED RHO INLINEARMODELELASTICNET ANDLINEARMODELSGDCLASSIFIER BY L1RATIO THERHO PARAMETER HAD DIFFERENT MEANINGS L1RATIO WAS INTRODUCED TO AVOID CONFUSION IT HAS THE SAME MEANING AS PREVIOUSLY RHO INLINEARMODELELASTICNET AND1RHO IN LINEARMODELSGDCLASSIFIER
- LINEARMODELLASSOLARS ANDLINEARMODELLARS NOW STORE A LIST OF PATHS IN THE CASE OF MULTIPLE TARGETS RATHER THAN AN ARRAY OF PATHS
- THE ATTRIBUTE GMM OFHMMGMMHMM WAS RENAMED TO GMM TO ADHERE MORE STRICTLY WITH THE API
- CLUSTERSPECTRALEMMBEDDING WAS MOVED TO MANIFOLDSPECTRALEMMBEDDING
- RENAMED EIGTOL INMANIFOLDSPECTRALEMMBEDDING CLUSTERSPECTRALCLUSTERING TO EIGENTOL RENAMEDMODE TOEIGENSOLVER
- RENAMED MODE INMANIFOLDSPECTRALEMMBEDDING ANDCLUSTERSPECTRALCLUSTERING TO EIGENSOLVER
- CLASSES ANDNCLASSES ATTRIBUTES OF TREEDECISIONTREECLASSIFIER AND ALL DERIVED ENSEMBLE MODELS ARE NOW FLAT IN CASE OF SINGLE OUTPUT PROBLEMS AND NESTED IN CASE OF MULTIOUTPUT PROBLEMS
- THEESTIMATORS ATTRIBUTE OF ENSEMBLEGRADIENTBOOSTINGGRADIENTBOOSTINGREGRESSOR ANDENSEMBLEGRADIENTBOOSTINGGRADIENTBOOSTINGCLASSIFIER IS NOW AN ARRAY OF CLASS'TREEDECISIONTREEREGRESSOR'
- RENAMED CHUNKSIZE TOBATCHSIZE INDECOMPOSITIONMINIBATCHDICTIONARYLEARNING ANDDECOMPOSITIONMINIBATCHSPARSEPCA FOR CONSISTENCY
- SVMSVC ANDSVMNUSVC NOW PROVIDE A CLASSES ATTRIBUTE AND SUPPORT ARBITRARY DTYPES FOR LABELS Y ALSO THE DTYPE RETURNED BY PREDICT NOW REFLECTS THE DTYPE OF YDURINGFIT USED TO BE NPFLOAT
- CHANGED DEFAULT TESTSIZE IN CROSSVALIDATIONTRAINTESTSPLIT TO NONE ADDED POSSIBILITY TO INFER TESTSIZE FROMTRAINSIZE INCROSSVALIDATIONSHUFFLESPLIT AND CROSSVALIDATIONSTRATIFIEDSHUFFLESPLIT
- RENAMED FUNCTION SKLEARNMETRICSZEROONE TOSKLEARNMETRICSZEROONELOSS BE AWARE THAT THE DEFAULT BEHAVIOR IN SKLEARNMETRICSZEROONELOSS IS DIFFERENT FROM SKLEARN METRICSZEROONE NORMALIZEFALSE IS CHANGED TO NORMALIZETRUE
- RENAMED FUNCTION METRICSZEROONESCORE TOMETRICSACCURACYSORE
- DATASETSMAKECIRCLES NOW HAS THE SAME NUMBER OF INNER AND OUTER POINTS
- IN THE NAIVE BAYES CLASSIFIERS THE CLASSPRIOR PARAMETER WAS MOVED FROM FIT TOINIT PEOPLE

LIST OF CONTRIBUTORS FOR RELEASE 013 BY NUMBER OF COMMITS

- 364 ANDREAS MÜLLER
- 143 ARNAUD JOLY
- 137 PETER PRETTENHOFER
- 131 GAELEVAROQUAUX
- 117 MATHIEU BLONDEL
- 108 LARS BUITINCK

120 CHAPTER 1 WELCOME TO SCIKITLEARN

SCIKITLEARN USER GUIDE RELEASE 0213

- 106 WEI LI
- 101 OLIVIER GRISEL
- 65 VLAD NICULAE
- 54 GILLES LOUPPE
- 40 JAQUES GROBLER
- 38 ALEXANDRE GRAMFORT
- 30 ROB ZINKOV
- 19 AYMERIC MASURELLE
- 18 ANDREW WINTERMAN
- 17 FABIAN PEDREGOSA
- 17 NELLE VAROQUAUX
- 16 CHRISTIAN OSENDORFER
- 14 DANIEL NOURI
- 13 VIRGILE FRITSCH
- 13 SYHW
- 12 SATRAJIT GHOSH
- 10 COREY LYNCH
- 10 KYLE BEAUCHAMP
- 9 BRIAN CHEUNG
- 9 IMMANUEL BAYER
- 9 MRSHU
- 8 CONRAD LEE
- 8 JAMES BERGSTRA
- 7 TADEJ JANEŽ
- 6 BRIAN CAJES
- 6 JAKE VANDERPLAS
- 6 MICHAEL
- 6 NOEL DAWE
- 6 TIAGO NUNES
- 6 COW
- 5 ANZE
- 5 SHIQIAO DU
- 4 CHRISTIAN JAUVIN
- 4 JACQUES KVAM
- 4 RICHARD T GUY
- 4 ROBERT LAYTON

SCIKITLEARN USER GUIDE RELEASE 0213

- 3 ALEXANDRE ABRAHAM
- 3 DOUG COLEMAN
- 3 SCOTT DICKERSON
- 2 APPROXIMATEIDENTITY
- 2 JOHN BENEDIKTSSON
- 2 MARK VERONDA
- 2 MATTI LYRA
- 2 MIKHAIL KOROBV
- 2 XINFAN MENG
- 1 ALEJANDRO WEINSTEIN
- 1 ALEXANDRE PASSOS
- 1 CHRISTOPH DEIL
- 1 EUGENE NIZHIBITSKY
- 1 KENNETH C ARNOLD
- 1 LUIS PEDRO COELHO
- 1 MIROSLAV BATCHKAROV
- 1 PAVEL
- 1 SEBASTIAN BERG
- 1 SHAUN JACKMAN
- 1 SUBHODEEP MOITRA
- 1 BOB
- 1 DENGEMANN
- 1 EMANUELE
- 1 X006

11717 VERSION 0121

OCTOBER 8 2012

THE 0121 RELEASE IS A BUGFIX RELEASE WITH NO ADDITIONAL FEATURES BUT IS INSTEAD A SET OF BUG FIXES  
CHANGELOG

- IMPROVED NUMERICAL STABILITY IN SPECTRAL EMBEDDING BY GAELE VAROQUAUX
- DOCTEST UNDER WINDOWS 64BIT BY GAELE VAROQUAUX
- DOCUMENTATION FIXES FOR ELASTIC NET BY ANDREAS MÜLLER AND ALEXANDRE GRAMFORT
- PROPER BEHAVIOR WITH FORTRANORDERED NUMPY ARRAYS BY GAELE VAROQUAUX
- MAKE GRIDSEARCHCV WORK WITH NONCSR SPARSE MATRIX BY LARS BUITINCK
- FIX PARALLEL COMPUTING IN MDS BY GAELE VAROQUAUX

122 CHAPTER 1 WELCOME TO SCIKITLEARN



SCIKITLEARN USER GUIDE RELEASE 0213

- FIX UNICODE SUPPORT IN COUNT VECTORIZER BY ANDREAS MÜLLER
- FIX MINCOVDET BREAKING WITH XSHAPE 3 1 BY VIRGILE FRITSCH
- FIX CLONE OF SGD OBJECTS BY PETER PRETTENHOFER
- STABILIZE GMM BY VIRGILE FRITSCH

PEOPLE

- 14 PETER PRETTENHOFER
- 12 GAEL VAROQUAUX
- 10 ANDREAS MÜLLER
- 5 LARS BUITINCK
- 3 VIRGILE FRITSCH
- 1 ALEXANDRE GRAMFORT
- 1 GILLES LOUPPE
- 1 MATHIEU BLONDEL

11718 VERSION 012

SEPTEMBER 4 2012

CHANGELOG

- VARIOUS SPEED IMPROVEMENTS OF THE DECISION TREES MODULE BY GILLES LOUPPE
  - ENSEMBLEGRADIENTBOOSTINGREGRESSOR AND ENSEMBLEGRADIENTBOOSTINGCLASSIFIER NOW SUPPORT FEATURE SUBSAMPLING VIA THE MAXFEATURES ARGUMENT BY PETER PRETTENHOFER
  - ADDED HUBER AND QUANTILE LOSS FUNCTIONS TO ENSEMBLEGRADIENTBOOSTINGREGRESSOR BY PETER PRETTENHOFER
  - DECISION TREES AND FORESTS OF RANDOMIZED TREES NOW SUPPORT MULTIOUTPUT CLASSIFICATION AND REGRESSION PROBLEMS BY GILLES LOUPPE
  - ADDED PREPROCESSING LABEL ENCODER A SIMPLE UTILITY CLASS TO NORMALIZE LABELS OR TRANSFORM NON NUMERICAL LABELS BY MATHIEU BLONDEL
  - ADDED THE EPSILON INSENSITIVE LOSS AND THE ABILITY TO MAKE PROBABILISTIC PREDICTIONS WITH THE MODIFIED HUBER LOSS INSTOCHASTIC GRADIENT DESCENT BY MATHIEU BLONDEL
  - ADDED MULTIDIMENSIONAL SCALING MDS BY NELLE VAROQUAUX
  - SVMLIGHT FILE FORMAT LOADER NOW DETECTS COMPRESSED GZIPBZIP2 FILES AND DECOMPRESSES THEM ON THE FLY BY LARS BUITINCK
  - SVMLIGHT FILE FORMAT SERIALIZER NOW PRESERVES DOUBLE PRECISION FLOATING POINT VALUES BY OLIVIER GRISEL
  - A COMMON TESTING FRAMEWORK FOR ALL ESTIMATORS WAS ADDED BY ANDREAS MÜLLER
  - UNDERSTANDABLE ERROR MESSAGES FOR ESTIMATORS THAT DO NOT ACCEPT SPARSE INPUT BY GAEL VAROQUAUX
  - SPEEDUPS IN HIERARCHICAL CLUSTERING BY GAEL VAROQUAUX IN PARTICULAR BUILDING THE TREE NOW SUPPORTS EARLY STOPPING THIS IS USEFUL WHEN THE NUMBER OF CLUSTERS IS NOT SMALL COMPARED TO THE NUMBER OF SAMPLES
- 117 PREVIOUS RELEASES 123

SCIKITLEARN USER GUIDE RELEASE 0213

- ADD MULTITASKLASSO AND MULTITASKELASTICNET FOR JOINT FEATURE SELECTION BY ALEXANDRE GRAMFORT
- ADDED METRICS AUROC SCORE AND METRICS AVERAGE PRECISION SCORE CONVENIENCE FUNCTIONS BY ANDREAS MÜLLER
- IMPROVED SPARSE MATRIX SUPPORT IN THE FEATURE SELECTION MODULE BY ANDREAS MÜLLER
- NEW WORD BOUNDARIES AWARE CHARACTER NGRAM ANALYZER FOR THE TEXT FEATURE EXTRACTION MODULE BY KERNC
- FIXED BUG IN SPECTRAL CLUSTERING THAT LED TO SINGLE POINT CLUSTERS BY ANDREAS MÜLLER
- INFEATUREEXTRACTIONTEXTCOUNTVECTORIZER ADDED AN OPTION TO IGNORE INFREQUENT WORDS MINDF BY ANDREAS MÜLLER
- ADD SUPPORT FOR MULTIPLE TARGETS IN SOME LINEAR MODELS ELASTICNET LASSO AND ORTHOGONAL MATCHING PURSUIT BY VLAD NICULAE AND ALEXANDRE GRAMFORT
- FIXES IN DECOMPOSITION PROBABILISTIC PCA SCORE FUNCTION BY WEI LI
- FIXED FEATURE IMPORTANCE COMPUTATION IN GRADIENT TREE BOOSTING

API CHANGES SUMMARY

- THE OLD SCIKITSLearn PACKAGE HAS DISAPPEARED ALL CODE SHOULD IMPORT FROM SKLEARN INSTEAD WHICH WAS INTRODUCED IN 09
  - IN METRICS ROC CURVE THE THRESHOLDS ARRAY IS NOW RETURNED WITH IT'S ORDER REVERSED IN ORDER TO KEEP IT CONSISTENT WITH THE ORDER OF THE RETURNED FPR AND TPR
  - IN HMM OBJECTS LIKE HMM GAUSSIAN HMM HMM MULTINOMIAL HMM ETC ALL PARAMETERS MUST BE PASSED TO THE OBJECT WHEN INITIALISING IT AND NOT THROUGH FIT NOW FIT WILL ONLY ACCEPT THE DATA AS AN INPUT PARAMETER
  - FOR ALL SVM CLASSES A FAULTY BEHAVIOR OF GAMMA WAS FIXED PREVIOUSLY THE DEFAULT GAMMA VALUE WAS ONLY COMPUTED THE FIRST TIME FIT WAS CALLED AND THEN STORED IT IS NOW RECALCULATED ON EVERY CALL TO FIT
  - ALL BASE CLASSES ARE NOW ABSTRACT META CLASSES SO THAT THEY CAN NOT BE INSTANTIATED
  - CLUSTER WARD TREE NOW ALSO RETURNS THE PARENT ARRAY THIS IS NECESSARY FOR EARLY STOPPING IN WHICH CASE THE TREE IS NOT COMPLETELY BUILT
  - INFEATUREEXTRACTIONTEXTCOUNTVECTORIZER THE PARAMETERS MINN AND MAXN WERE JOINED TO THE PARAMETER NGRAM RANGE TO ENABLE GRID SEARCHING BOTH AT ONCE
  - INFEATUREEXTRACTIONTEXTCOUNTVECTORIZER WORDS THAT APPEAR ONLY IN ONE DOCUMENT ARE NOW IGNORED BY DEFAULT TO REPRODUCE THE PREVIOUS BEHAVIOR SET MINDF1
  - FIXED API INCONSISTENCY LINEAR MODELS SGD CLASSIFIER PREDICT PROBA NOW RETURNS 2D ARRAY WHEN FIT ON TWO CLASSES
  - FIXED API INCONSISTENCY DISCRIMINANT ANALYSIS QUADRATIC DISCRIMINANT ANALYSIS DECISION FUNCTION AND DISCRIMINANT ANALYSIS LINEAR DISCRIMINANT ANALYSIS DECISION FUNCTION NOW RETURN 1D ARRAYS WHEN FIT ON TWO CLASSES
  - GRID OF ALPHAS USED FOR FITTING LINEAR MODEL LASSO CV AND LINEAR MODEL ELASTIC NET CV IS NOW STORED IN THE ATTRIBUTE ALPHAS RATHER THAN OVERRIDING THE INIT PARAMETER ALPHAS
  - LINEAR MODELS WHEN ALPHA IS ESTIMATED BY CROSS VALIDATION STORE THE ESTIMATED VALUE IN THE ALPHA ATTRIBUTE RATHER THAN JUST ALPHA OR BEST ALPHA
  - ENSEMBLE GRADIENT BOOSTING CLASSIFIER NOW SUPPORTS ENSEMBLE GRADIENT BOOSTING CLASSIFIER STAGED PREDICT PROBA AND ENSEMBLE GRADIENT BOOSTING CLASSIFIER STAGED PREDICT
- 124 CHAPTER 1 WELCOME TO SCIKITLEARN

SCIKITLEARN USER GUIDE RELEASE 0213

- SVMSPARSE SVC AND OTHER SPARSE SVM CLASSES ARE NOW DEPRECATED THE ALL CLASSES IN THE SUPPORT VECTOR MACHINES MODULE NOW AUTOMATICALLY SELECT THE SPARSE OR DENSE REPRESENTATION BASE ON THE INPUT
- ALL CLUSTERING ALGORITHMS NOW INTERPRET THE ARRAY X GIVEN TO FIT AS INPUT DATA IN PARTICULAR CLUSTER SPECTRALCLUSTERING AND CLUSTER AFFINITY PROPAGATION WHICH PREVIOUSLY EXPECTED AFFINITY MATRICES

- FOR CLUSTERING ALGORITHMS THAT TAKE THE DESIRED NUMBER OF CLUSTERS AS A PARAMETER THIS PARAMETER IS NOW CALLED NCLUSTERS

PEOPLE

- 267 ANDREAS MÜLLER
- 94 GILLES LOUPPE
- 89 GAELE VAROQUAUX
- 79 PETER PRETTENHOFER
- 60 MATHIEU BLONDEL
- 57 ALEXANDRE GRAMFORT
- 52 VLAD NICULAE
- 45 LARS BUITINCK
- 44 NELLE VAROQUAUX
- 37 JACQUES GROBLER
- 30 ALEXIS MIGNON
- 30 IMMANUEL BAYER
- 27 OLIVIER GRISEL
- 16 SUBHODEEP MOITRA
- 13 YANNICK SCHWARTZ
- 12 KERN C
- 11 VIRGILE FRITSCH
- 9 DANIEL DUCKWORTH
- 9 FABIAN PEDREGOSA
- 9 ROBERT LAYTON
- 8 JOHN BENEDIKTSSON
- 7 MARKO BURJEK
- 5 NICOLAS PINTO
- 4 ALEXANDRE ABRAHAM
- 4 JAKE VANDERPLAS
- 3 BRIAN HOLT
- 3 EDOUARD DUCHESNAY
- 3 FLORIAN HOENIG

SCIKITLEARN USER GUIDE RELEASE 0213

- 3 FLYINGIMMIDEV
- 2 FRANCOIS SAVARD
- 2 HANNES SCHULZ
- 2 PETER WELINDER
- 2 YAROSLAV HALCHENKO
- 2 WEI LI
- 1 ALEX COMPANIONI
- 1 BRANDYN A WHITE
- 1 BUSSONNIER MATTHIAS
- 1 CHARLESPIERRE ASTOLFI
- 1 DAN O’HUGINN
- 1 DAVID COURNAPEAU
- 1 KEITH GOODMAN
- 1 LUDWIG SCHWARDT
- 1 OLIVIER HERVIEU
- 1 SERGIO MEDINA
- 1 SHIQIAO DU
- 1 TIM SHEERMANCHASE
- 1 BUGUEN

11719 VERSION 011

MAY 7 2012

CHANGELOG

HIGHLIGHTS

- GRADIENT BOOSTED REGRESSION TREES GRADIENT TREE BOOSTING FOR CLASSIFICATION AND REGRESSION BY PETER PRETTEN  
HOFER AND SCOTT WHITE
- SIMPLE DICTBASED FEATURE LOADER WITH SUPPORT FOR CATEGORICAL VARIABLES FEATUREEXTRACTION  
DICTVECTORIZER BY LARS BUITINCK
- ADDED MATTHEWS CORRELATION COEFFICIENT METRICSMATTHEWSCORRCOEF AND ADDED MACRO AND MICRO AV  
ERAGE OPTIONS TO METRICSPRECISIONSCORE METRICSRECALLSCORE ANDMETRICSF1SCORE  
BY SATRAJIT GHOSH
- OUT OF BAG ESTIMATES OF GENERALIZATION ERROR FOR ENSEMBLE METHODS BY ANDREAS MÜLLER
- RANDOMIZED SPARSE LINEAR MODELS FOR FEATURE SELECTION BY ALEXANDRE GRAMFORT AND GAELEVAROQUAUX
- LABEL PROPAGATION FOR SEMISUPERVISED LEARNING BY CLAY WOOLAM NOTE THE SEMISUPERVISED API IS STILL WORK  
IN PROGRESS AND MAY CHANGE

126 CHAPTER 1 WELCOME TO SCIKITLEARN

SCIKITLEARN USER GUIDE RELEASE 0213

- ADDED BICAIC MODEL SELECTION TO CLASSICAL GAUSSIAN MIXTURE MODELS AND UNIFIED THE API WITH THE REMAINDER OF SCIKITLEARN BY BERTRAND THIRION
  - ADDEDSKLEARNCROSSVALIDATIONSTRATIFIEDSHUFFLESPLIT WHICH IS A SKLEARN CROSSVALIDATIONSHUFFLESPLIT WITH BALANCED SPLITS BY YANNICK SCHWARTZ
  - SKLEARNNEIGHBORSNEARESTCENTROID CLASSIFIER ADDED ALONG WITH A SHRINKTHRESHOLD PARAMETER WHICH IMPLEMENTS SHRUNKEN CENTROID CLASSIFICATION BY ROBERT LAYTON
- OTHER CHANGES
- MERGED DENSE AND SPARSE IMPLEMENTATIONS OF STOCHASTIC GRADIENT DESCENT MODULE AND EXPOSED UTILITY EXTENSION TYPES FOR SEQUENTIAL DATASETS SEQDATASET AND WEIGHT VECTORS WEIGHTVECTOR BY PETER PRETTENHOFER
  - ADDEDPARTIALFIT SUPPORT FOR ONLINEMINIBATCH LEARNING AND WARMSTART TO THE STOCHASTIC GRADIENT DESCENT MODULE BY MATHIEU BLONDEL
  - DENSE AND SPARSE IMPLEMENTATIONS OF SUPPORT VECTOR MACHINES CLASSES AND LINEARMODEL LOGISTICREGRESSION MERGED BY LARS BUITINCK
  - REGRESSORS CAN NOW BE USED AS BASE ESTIMATOR IN THE MULTICLASS AND MULTILABEL ALGORITHMS MODULE BY MATHIEU BLONDEL
  - ADDED NJOBS OPTION TO METRICSPAIRWISEPAIRWISEDISTANCES ANDMETRICSPAIRWISE PAIRWISEKERNELS FOR PARALLEL COMPUTATION BY MATHIEU BLONDEL
  - KMEANS CAN NOW BE RUN IN PARALLEL USING THE NJOBS ARGUMENT TO EITHER KMEANS ORKMEANS BY ROBERT LAYTON
  - IMPROVED CROSSVALIDATION EVALUATING ESTIMATOR PERFORMANCE ANDTUNING THE HYPERPARAMETERS OF AN ESTIMATORDOCUMENTATION AND INTRODUCED THE NEW CROSSVALIDATIONTRAINTESTSPLIT HELPER FUNCTION BY OLIVIER GRISEL
  - SVM SVC MEMBERS COEF ANDINTERCEPT CHANGED SIGN FOR CONSISTENCY WITH DECISIONFUNCTION
  - FORKERNELLINEAR COEF WAS FIXED IN THE ONEVONE CASE BY ANDREAS MÜLLER
  - PERFORMANCE IMPROVEMENTS TO EFFICIENT LEAVEONEOUT CROSSVALIDATED RIDGE REGRESSION ESP FOR THE NSAMPLES NFEATURES CASE IN LINEARMODEL RIDGE CV BY REUBEN FLETCHER COSTIN
  - REFACTORING AND SIMPLIFICATION OF THE TEXT FEATURE EXTRACTION API AND FIXED A BUG THAT CAUSED POSSIBLE NEGATIVE IDF BY OLIVIER GRISEL
  - BEAM PRUNING OPTION IN BASEHMM MODULE HAS BEEN REMOVED SINCE IT IS DIFFICULT TO CYTHONIZE IF YOU ARE INTERESTED IN CONTRIBUTING A CYTHON VERSION YOU CAN USE THE PYTHON VERSION IN THE GIT HISTORY AS A REFERENCE
  - CLASSES IN NEAREST NEIGHBORS NOW SUPPORT ARBITRARY MINKOWSKI METRIC FOR NEAREST NEIGHBORS SEARCHES THE METRIC CAN BE SPECIFIED BY ARGUMENT P
- API CHANGES SUMMARY
- COVARIANCEELLIPTICENVELOPE IS NOW DEPRECATED PLEASE USE COVARIANCEELLIPTICENVELOPE INSTEAD
  - NEIGHBORSCLASSIFIER ANDNEIGHBORSREGRESSOR ARE GONE IN THE MODULE NEAREST NEIGHBORS USE THE CLASSES KNEIGHBORSCLASSIFIER RADIUSNEIGHBORSCLASSIFIER KNEIGHBORSREGRESSOR ANDORRADIUSNEIGHBORSREGRESSOR INSTEAD
  - SPARSE CLASSES IN THE STOCHASTIC GRADIENT DESCENT MODULE ARE NOW DEPRECATED
- 117 PREVIOUS RELEASES 127

SCIKITLEARN USER GUIDE RELEASE 0213

- INMIXTUREGMM MIXTUREDPGMM ANDMIXTUREVBGMM PARAMETERS MUST BE PASSED TO AN OBJECT WHEN INITIALISING IT AND NOT THROUGH FIT NOWFIT WILL ONLY ACCEPT THE DATA AS AN INPUT PARAMETER
- METHODS RVS ANDDECODE INGMM MODULE ARE NOW DEPRECATED SAMPLE ANDSCORE ORPREDICT SHOULD BE USED INSTEAD
- ATTRIBUTESCORES ANDPVALUES IN UNIVARIATE FEATURE SELECTION OBJECTS ARE NOW DEPRECATED SCORES OR PVALUES SHOULD BE USED INSTEAD
- INLOGISTICREGRESSION LINEARSVC SVC ANDNUSVC THECLASSWEIGHT PARAMETER IS NOW AN INITIALIZATION PARAMETER NOT A PARAMETER TO FIT THIS MAKES GRID SEARCHES OVER THIS PARAMETER POSSIBLE
- LFWDATA IS NOW ALWAYS SHAPE NSAMPLES NFEATURES TO BE CONSISTENT WITH THE OLIVETTI FACES DATASET USE IMAGES ANDPAIRS ATTRIBUTE TO ACCESS THE NATURAL IMAGES SHAPES INSTEAD
- INSVMLINEARSVC THE MEANING OF THE MULTICLASS PARAMETER CHANGED OPTIONS NOW ARE OVR AND CRAMMERSINGER WITHOVR BEING THE DEFAULT THIS DOES NOT CHANGE THE DEFAULT BEHAVIOR BUT HOPEFULLY IS LESS CONFUSING
- CLASS FEATURESELECTIONTEXTVECTORIZER IS DEPRECATED AND REPLACED BY FEATURESELECTIONTEXTTFIDFVECTORIZER
- THE PREPROCESSOR ANALYZER NESTED STRUCTURE FOR TEXT FEATURE EXTRACTION HAS BEEN REMOVED ALL THOSE FEATURES ARE NOW DIRECTLY PASSED AS FLAT CONSTRUCTOR ARGUMENTS TO FEATURESELECTIONTEXTTFIDFVECTORIZER ANDFEATURESELECTIONTEXTCOUNTVECTORIZER IN PARTICULAR THE FOLLOWING PARAMETERS ARE NOW USED
- ANALYZER CAN BEWORD ORCHAR TO SWITCH THE DEFAULT ANALYSIS SCHEME OR USE A SPECIFIC PYTHON CALLABLE AS PREVIOUSLY
- TOKENIZER ANDPREPROCESSOR HAVE BEEN INTRODUCED TO MAKE IT STILL POSSIBLE TO CUSTOMIZE THOSE STEPS WITH THE NEW API
- INPUT EXPLICITLY CONTROL HOW TO INTERPRET THE SEQUENCE PASSED TO FIT ANDPREDICT FILENAMES FILE OBJECTS OR DIRECT BYTE OR UNICODE STRINGS
- CHARSET DECODING IS EXPLICIT AND STRICT BY DEFAULT
- THEVOCABULARY FITTED OR NOT IS NOW STORED IN THE VOCABULARY ATTRIBUTE TO BE CONSISTENT WITH THE PROJECT CONVENTIONS
- CLASS FEATURESELECTIONTEXTTFIDFVECTORIZER NOW DERIVES DIRECTLY FROM FEATURESELECTIONTEXTCOUNTVECTORIZER TO MAKE GRID SEARCH TRIVIAL
- METHODS RVS INBASEHMM MODULE ARE NOW DEPRECATED SAMPLE SHOULD BE USED INSTEAD
- BEAM PRUNING OPTION IN BASEHMM MODULE IS REMOVED SINCE IT IS DIFFICULT TO BE CYTHONIZED IF YOU ARE INTERESTED YOU CAN LOOK IN THE HISTORY CODES BY GIT
- THE SVMLIGHT FORMAT LOADER NOW SUPPORTS FILES WITH BOTH ZEROBASED AND ONEBASED COLUMN INDICES SINCE BOTH OCCUR “IN THE WILD”
- ARGUMENTS IN CLASS SHUFFLESPLIT ARE NOW CONSISTENT WITH STRATIFIEDSHUFFLESPLIT ARGUMENTS TESTFRACTION ANDTRAINFRACTION ARE DEPRECATED AND RENAMED TO TESTSIZE ANDTRAINSIZ
- AND CAN ACCEPT BOTH FLOAT ANDINT
- ARGUMENTS IN CLASS BOOTSTRAP ARE NOW CONSISTENT WITH STRATIFIEDSHUFFLESPLIT ARGUMENTS NTEST ANDNTRAIN ARE DEPRECATED AND RENAMED TO TESTSIZE ANDTRAINSIZ
- AND CAN ACCEPT BOTH FLOAT ANDINT
- ARGUMENT PADDED TO CLASSES IN NEAREST NEIGHBORS TO SPECIFY AN ARBITRARY MINKOWSKI METRIC FOR NEAREST NEIGHBORS SEARCHES

SCIKITLEARN USER GUIDE RELEASE 0213

PEOPLE

- 282 ANDREAS MÜLLER
  - 239 PETER PRETTENHOFER
  - 198 GAEL VAROQUAUX
  - 129 OLIVIER GRISEL
  - 114 MATHIEU BLONDEL
  - 103 CLAY WOOLAM
  - 96 LARS BUITINCK
  - 88 JAQUES GROBLER
  - 82 ALEXANDRE GRAMFORT
  - 50 BERTRAND THIRION
  - 42 ROBERT LAYTON
  - 28 FLYINGIMMIDEV
  - 26 JAKE VANDERPLAS
  - 26 SHIQIAO DU
  - 21 SATRAJIT GHOSH
  - 17 DAVID MAREK
  - 17 GILLES LOUPPE
  - 14 VLAD NICULAE
  - 11 YANNICK SCHWARTZ
  - 10 FABIAN PEDREGOSA
  - 9 FCOSTIN
  - 7 NICK WILSON
  - 5 ADRIEN GAIDON
  - 5 NICOLAS PINTO
  - 4 DAVID WARDEFARLEY
  - 5 NELLE VAROQUAUX
  - 5 EMMANUELLE GOUILLART
  - 3 JOONAS SILLANPÄÄ
  - 3 PAOLO LOSI
  - 2 CHARLES MCCARTHY
  - 2 ROY HYUNJIN HAN
  - 2 SCOTT WHITE
  - 2 IBAYER
  - 1 BRANDYN WHITE
  - 1 CARLOS SCHEIDEGGER
- 117 PREVIOUS RELEASES 129

SCIKITLEARN USER GUIDE RELEASE 0213

- 1 CLAIRE REVILLET
- 1 CONRAD LEE
- 1 EDOUARD DUCHESNAY
- 1 JAN HENDRIK METZEN
- 1 MENG XINFAN
- 1 ROB ZINKOV
- 1 SHIQIAO
- 1 UDI WEINSBERG
- 1 VIRGILE FRITSCH
- 1 XINFAN MENG
- 1 YAROSLAV HALCHENKO
- 1 JANSOE
- 1 LEON PALAFOX

11720 VERSION 010

JANUARY 11 2012

CHANGELOG

- PYTHON 25 COMPATIBILITY WAS DROPPED THE MINIMUM PYTHON VERSION NEEDED TO USE SCIKITLEARN IS NOW 26
- SPARSE INVERSE COVARIANCE ESTIMATION USING THE GRAPH LASSO WITH ASSOCIATED CROSSVALIDATED ESTIMATOR BY GAELEVAROQUAUX
- NEW TREE MODULE BY BRIAN HOLT PETER PRETTENHOFER SATRAJIT GHOSH AND GILLES LOUPPE THE MODULE COMES WITH COMPLETE DOCUMENTATION AND EXAMPLES
- FIXED A BUG IN THE RFE MODULE BY GILLES LOUPPE ISSUE 378
- FIXED A MEMORY LEAK IN SUPPORT VECTOR MACHINES MODULE BY BRIAN HOLT ISSUE 367
- FASTER TESTS BY FABIAN PEDREGOSA AND OTHERS
- SILHOUETTE COEFFICIENT CLUSTER ANALYSIS EVALUATION METRIC ADDED AS SKLEARNMETRICS
- SILHOUETTESCORE BY ROBERT LAYTON
- FIXED A BUG IN KMEANS IN THE HANDLING OF THE NINIT PARAMETER THE CLUSTERING ALGORITHM USED TO BE RUN NINIT TIMES BUT THE LAST SOLUTION WAS RETAINED INSTEAD OF THE BEST SOLUTION BY OLIVIER GRISEL
- MINOR REFACTORING IN STOCHASTIC GRADIENT DESCENT MODULE CONSOLIDATED DENSE AND SPARSE PREDICT METHODS ENHANCED TEST TIME PERFORMANCE BY CONVERTING MODEL PARAMETERS TO FORTRANSTYLE ARRAYS AFTER FITTING ONLY MULTI CLASS
- ADJUSTED MUTUAL INFORMATION METRIC ADDED AS SKLEARNMETRICSDJUSTEDMUTUALINFOSCORE
- BY ROBERT LAYTON
- MODELS LIKE SVCSVRLINEARSVCLOGISTICREGRESSION FROM LIBSVM LIBLINEAR NOW SUPPORT SCALING OF C REGULARIZATION PARAMETER BY THE NUMBER OF SAMPLES BY ALEXANDRE GRAMFORT
- NEW ENSEMBLE METHODS MODULE BY GILLES LOUPPE AND BRIAN HOLT THE MODULE COMES WITH THE RANDOM FOREST ALGORITHM AND THE EXTRATREES METHOD ALONG WITH DOCUMENTATION AND EXAMPLES

130 CHAPTER 1 WELCOME TO SCIKITLEARN



SCIKITLEARN USER GUIDE RELEASE 0213

- NOVELTY AND OUTLIER DETECTION OUTLIER AND NOVELTY DETECTION BY VIRGILE FRITSCH
- KERNEL APPROXIMATION A TRANSFORM IMPLEMENTING KERNEL APPROXIMATION FOR FAST SGD ON NONLINEAR KERNELS BY ANDREAS MÜLLER
- FIXED A BUG DUE TO ATOM SWAPPING IN ORTHOGONAL MATCHING PURSUIT OMP BY VLAD NICULAE
- SPARSE CODING WITH A PRECOMPUTED DICTIONARY BY VLAD NICULAE
- MINI BATCH KMEANS PERFORMANCE IMPROVEMENTS BY OLIVIER GRISEL
- KMEANS SUPPORT FOR SPARSE MATRICES BY MATHIEU BLONDEL
- IMPROVED DOCUMENTATION FOR DEVELOPERS AND FOR THE SKLEARNUTILS MODULE BY JAKE VANDERPLAS
- VECTORIZED 20NEWSGROUPS DATASET LOADER SKLEARNDATASETSFETCH20NEWSGROUPSVECTORIZED BY MATHIEU BLONDEL
- MULTICLASS AND MULTILABEL ALGORITHMS BY LARS BUITINCK
- UTILITIES FOR FAST COMPUTATION OF MEAN AND VARIANCE FOR SPARSE MATRICES BY MATHIEU BLONDEL
- MAKESKLEARNPREPROCESSINGSCALE ANDSKLEARNPREPROCESSINGSCALER WORK ON SPARSE MATRICES BY OLIVIER GRISEL
- FEATURE IMPORTANCES USING DECISION TREES ANDOR FOREST OF TREES BY GILLES LOUPPE
- PARALLEL IMPLEMENTATION OF FORESTS OF RANDOMIZED TREES BY GILLES LOUPPE
- SKLEARNCROSSVALIDATIONSHUFFLESPLIT CAN SUBSAMPLE THE TRAIN SETS AS WELL AS THE TEST SETS BY OLIVIER GRISEL
- ERRORS IN THE BUILD OF THE DOCUMENTATION FIXED BY ANDREAS MÜLLER

API CHANGES SUMMARY

HERE ARE THE CODE MIGRATION INSTRUCTIONS WHEN UPGRADING FROM SCIKITLEARN VERSION 09

- SOME ESTIMATORS THAT MAY OVERWRITE THEIR INPUTS TO SAVE MEMORY PREVIOUSLY HAD OVERWRITE PARAMETERS THESE HAVE BEEN REPLACED WITH COPY PARAMETERS WITH EXACTLY THE OPPOSITE MEANING THIS PARTICULARLY AFFECTS SOME OF THE ESTIMATORS IN LINEARMODEL THE DEFAULT BEHAVIOR IS STILL TO COPY EVERYTHING PASSED IN
  - THE SVMLIGHT DATASET LOADER SKLEARNDATASETSLOADSVMLIGHTFILE NO LONGER SUPPORTS LOADING TWO FILES AT ONCE USE LOADSVMLIGHTFILES INSTEAD ALSO THE UNUSED BUFFERMB PARAMETER IS GONE
  - SPARSE ESTIMATORS IN THE STOCHASTIC GRADIENT DESCENT MODULE USE DENSE PARAMETER VECTOR COEF INSTEAD OF SPARSECOEF THIS SIGNIFICANTLY IMPROVES TEST TIME PERFORMANCE
  - THE COVARIANCE ESTIMATION MODULE NOW HAS A ROBUST ESTIMATOR OF COVARIANCE THE MINIMUM COVARIANCE DETERMINANT ESTIMATOR
  - CLUSTER EVALUATION METRICS IN METRICSCLUSTER HAVE BEEN REFACTORED BUT THE CHANGES ARE BACKWARDS COMPATIBLE THEY HAVE BEEN MOVED TO THE METRICSCLUSTERSUPERVISED ALONG WITH METRICSCLUSTER UNSUPERVISED WHICH CONTAINS THE SILHOUETTE COEFFICIENT
  - THEPERMUTATIONTESTSCORE FUNCTION NOW BEHAVES THE SAME WAY AS CROSSVALSCORE IE USES THE MEAN SCORE ACROSS THE FOLDS
  - CROSS VALIDATION GENERATORS NOW USE INTEGER INDICES INDICESTRUE BY DEFAULT INSTEAD OF BOOLEAN MASKS THIS MAKE IT MORE INTUITIVE TO USE WITH SPARSE MATRIX DATA
- 117 PREVIOUS RELEASES 131

SCIKITLEARN USER GUIDE RELEASE 0213

- THE FUNCTIONS USED FOR SPARSE CODING SPARSEENCODE ANDSPARSEENCODEPARALLEL HAVE BEEN COMBINED INTO SKLEARNDECOMPOSITIONSPARSEENCODE AND THE SHAPES OF THE ARRAYS HAVE BEEN TRANSPOSED FOR CONSISTENCY WITH THE MATRIX FACTORIZATION SETTING AS OPPOSED TO THE REGRESSION SETTING
- FIXED AN OFFBYONE ERROR IN THE SVMLIGHTLIBSVM FILE FORMAT HANDLING FILES GENERATED USING SKLEARN DATASETS DUMPSVMLIGHTFILE SHOULD BE REGENERATED THEY SHOULD CONTINUE TO WORK BUT ACCIDENTALLY HAD ONE EXTRA COLUMN OF ZEROS PREPENDED
- BASEDICTIONARYLEARNING CLASS REPLACED BY SPARSECODINGMIXIN
- SKLEARNUTILSEXTMATHFASTSVD HAS BEEN RENAMED SKLEARNUTILSEXTMATHRANDOMIZEDSVD AND THE DEFAULT OVERSAMPLING IS NOW FIXED TO 10 ADDITIONAL RANDOM VECTORS INSTEAD OF DOUBLING THE NUMBER OF COMPONENTS TO EXTRACT THE NEW BEHAVIOR FOLLOWS THE REFERENCE PAPER

PEOPLE  
THE FOLLOWING PEOPLE CONTRIBUTED TO SCIKITLEARN SINCE LAST RELEASE

- 246 ANDREAS MÜLLER
- 242 OLIVIER GRISEL
- 220 GILLES LOUPPE
- 183 BRIAN HOLT
- 166 GAELE VAROQUAUX
- 144 LARS BUITINCK
- 73 VLAD NICULAE
- 65 PETER PRETTENHOFER
- 64 FABIAN PEDREGOSA
- 60 ROBERT LAYTON
- 55 MATHIEU BLONDEL
- 52 JAKE VANDERPLAS
- 44 NOEL DAWE
- 38 ALEXANDRE GRAMFORT
- 24 VIRGILE FRITSCH
- 23 SATRAJIT GHOSH
- 3 JAN HENDRIK METZEN
- 3 KENNETH C ARNOLD
- 3 SHIQIAO DU
- 3 TIM SHEERMANCHASE
- 3 YAROSLAV HALCHENKO
- 2 BALASUBRAHMANYAM VARANASI
- 2 DRAXUS
- 2 MICHAEL EICKENBERG
- 1 BOGDAN TRACH

SCIKITLEARN USER GUIDE RELEASE 0213

- 1 FÉLIXANTOINE FORTIN
- 1 JUAN MANUEL CAICEDO CARVAJAL
- 1 NELLE VAROQUAUX
- 1 NICOLAS PINTO
- 1 TIZIANO ZITO
- 1 XINFAN MENG

11721 VERSION 09

SEPTEMBER 21 2011

SCIKITLEARN 09 WAS RELEASED ON SEPTEMBER 2011 THREE MONTHS AFTER THE 08 RELEASE AND INCLUDES THE NEW MODULES MANIFOLD LEARNING THE DIRICHLET PROCESS AS WELL AS SEVERAL NEW ALGORITHMS AND DOCUMENTATION IMPROVEMENTS THIS RELEASE ALSO INCLUDES THE DICTIONARYLEARNING WORK DEVELOPED BY VLAD NICULAE AS PART OF THE GOOGLE SUMMER OF CODE PROGRAM

117 PREVIOUS RELEASES 133

SCIKITLEARN USER GUIDE RELEASE 0213

CHANGELOG

- NEW MANIFOLD LEARNING MODULE BY JAKE VANDERPLAS AND FABIAN PEDREGOSA
- NEW DIRICHLET PROCESS GAUSSIAN MIXTURE MODEL BY ALEXANDRE PASSOS

134 CHAPTER 1 WELCOME TO SCIKITLEARN

SCIKITLEARN USER GUIDE RELEASE 0213

- NEAREST NEIGHBORS MODULE REFACTORING BY JAKE VANDERPLAS GENERAL REFACTORING SUPPORT FOR SPARSE MATRICES IN INPUT SPEED AND DOCUMENTATION IMPROVEMENTS SEE THE NEXT SECTION FOR A FULL LIST OF API CHANGES
- IMPROVEMENTS ON THE FEATURE SELECTION MODULE BY GILLES LOUPPE REFACTORING OF THE RFE CLASSES DOCUMENTATION REWRITE INCREASED EFFICIENCY AND MINOR API CHANGES
- SPARSE PRINCIPAL COMPONENTS ANALYSIS SPARSEPCA AND MINIBATCHSPARSEPCA BY VLAD NICULAE GAELE VAROQUAUX AND ALEXANDRE GRAMFORT
- PRINTING AN ESTIMATOR NOW BEHAVES INDEPENDENTLY OF ARCHITECTURES AND PYTHON VERSION THANKS TO JEAN KOSSAIFI
- LOADER FOR LIBSVM SVMLIGHT FORMAT BY MATHIEU BLONDEL AND LARS BUITINCK
- DOCUMENTATION IMPROVEMENTS THUMBNAILS IN EXAMPLE GALLERY BY FABIAN PEDREGOSA
- IMPORTANT BUGFIXES IN SUPPORT VECTOR MACHINES MODULE SEGFAULTS BAD PERFORMANCE BY FABIAN PEDREGOSA
- ADDED MULTINOMIAL NAIVE BAYES AND BERNOULLI NAIVE BAYES BY LARS BUITINCK
- TEXT FEATURE EXTRACTION OPTIMIZATIONS BY LARS BUITINCK
- CHISQUARE FEATURE SELECTION FEATURESELECTION UNIVARIATE SELECTION CHI2 BY LARS BUITINCK
- GENERATED DATASETS MODULE REFACTORING BY GILLES LOUPPE
- MULTICLASS AND MULTILABEL ALGORITHMS BY MATHIEU BLONDEL
- BALL TREE REWRITE BY JAKE VANDERPLAS
- IMPLEMENTATION OF DBSCAN ALGORITHM BY ROBERT LAYTON
- KMEANS PREDICT AND TRANSFORM BY ROBERT LAYTON
- PREPROCESSING MODULE REFACTORING BY OLIVIER GRISEL
- FASTER MEAN SHIFT BY CONRAD LEE
- NEWBOOTSTRAP RANDOM PERMUTATIONS CROSSVALIDATION AKA SHUFFLE SPLIT AND VARIOUS OTHER IMPROVEMENTS IN CROSS VALIDATION SCHEMES BY OLIVIER GRISEL AND GAELE VAROQUAUX
- ADJUSTED RAND INDEX AND VMEASURE CLUSTERING EVALUATION METRICS BY OLIVIER GRISEL
- ADDED ORTHOGONAL MATCHING PURSUIT BY VLAD NICULAE
- ADDED 2DPATCH EXTRACTOR UTILITIES IN THE FEATURE EXTRACTION MODULE BY VLAD NICULAE
- IMPLEMENTATION OF LINEAR MODEL LASSO LARS CV CROSSVALIDATED LASSO SOLVER USING THE LARS ALGORITHM AND LINEAR MODEL LASSO LARS BIC AIC MODEL SELECTION IN LARS BY GAELE VAROQUAUX AND ALEXANDRE GRAMFORT
- SCALABILITY IMPROVEMENTS TO METRICS ROC CURVE BY OLIVIER HERVIEU
- DISTANCE HELPER FUNCTIONS METRICS PAIRWISE PAIRWISE DISTANCES AND METRICS PAIRWISE PAIRWISE KERNELS BY ROBERT LAYTON
- MINIBATCH KMEANS BY NELLE VAROQUAUX AND PETER PRETTENHOFER
- ML DATA UTILITIES BY PIETRO BERKES
- OLIVETTI FACES BY DAVID WARDEFARLEY

API CHANGES SUMMARY

HERE ARE THE CODE MIGRATION INSTRUCTIONS WHEN UPGRADING FROM SCIKITLEARN VERSION 0.8 TO 1.17 PREVIOUS RELEASES 135

SCIKITLEARN USER GUIDE RELEASE 0213

- THE SCIKITLEARN PACKAGE WAS RENAMED SKLEARN THERE IS STILL A SCIKITLEARN PACKAGE ALIAS FOR BACKWARD COMPATIBILITY

THIRDPARTY PROJECTS WITH A DEPENDENCY ON SCIKITLEARN 09 SHOULD UPGRADE THEIR CODEBASE FOR INSTANCE UNDER LINUX MACOSX JUST RUN MAKE A BACKUP FIRST

FIND NAME PY XARGS SED I S BSCIKITLEARN BSKLEARN

- ESTIMATORS NO LONGER ACCEPT MODEL PARAMETERS AS FIT ARGUMENTS INSTEAD ALL PARAMETERS MUST BE ONLY BE PASSED AS CONSTRUCTOR ARGUMENTS OR USING THE NOW PUBLIC SETPARAMS METHOD INHERITED FROM BASE BASEESTIMATOR

SOME ESTIMATORS CAN STILL ACCEPT KEYWORD ARGUMENTS ON THE FIT BUT THIS IS RESTRICTED TO DATA DEPENDENT VALUES EG A GRAM MATRIX OR AN AFFINITY MATRIX THAT ARE PRECOMPUTED FROM THE XDATA MATRIX

- THE CROSSVAL PACKAGE HAS BEEN RENAMED TO CROSSVALIDATION ALTHOUGH THERE IS ALSO A CROSSVAL PACKAGE ALIAS IN PLACE FOR BACKWARD COMPATIBILITY

THIRDPARTY PROJECTS WITH A DEPENDENCY ON SCIKITLEARN 09 SHOULD UPGRADE THEIR CODEBASE FOR INSTANCE UNDER LINUX MACOSX JUST RUN MAKE A BACKUP FIRST

FIND NAME PY XARGS SED I S BCROSSVAL BCROSSVALIDATION

- THE SCORE\_FUNC ARGUMENT OF THE SKLEARN CROSSVALIDATION CROSSVAL\_SCORE FUNCTION IS NOW EXPECTED TO ACCEPT YTEST AND Y\_PREDICTED AS ONLY ARGUMENTS FOR CLASSIFICATION AND REGRESSION TASKS OR XTEST FOR UNSUPERVISED ESTIMATORS

- GAMMA PARAMETER FOR SUPPORT VECTOR MACHINE ALGORITHMS IS SET TO 1 N\_FEATURES BY DEFAULT INSTEAD OF 1 N\_SAMPLES

- THE SKLEARN\_HMM HAS BEEN MARKED AS ORPHANED IT WILL BE REMOVED FROM SCIKITLEARN IN VERSION 0.11 UNLESS SOMEONE STEPS UP TO CONTRIBUTE DOCUMENTATION EXAMPLES AND FIX LURKING NUMERICAL STABILITY ISSUES

- SKLEARN\_NEIGHBORS HAS BEEN MADE INTO A SUBMODULE THE TWO PREVIOUSLY AVAILABLE ESTIMATORS NEIGHBORS\_CLASSIFIER AND NEIGHBORS\_REGRESSOR HAVE BEEN MARKED AS DEPRECATED THEIR FUNCTIONALITY HAS BEEN DIVIDED AMONG FIVE NEW CLASSES NEAREST\_NEIGHBORS FOR UNSUPERVISED NEIGHBORS SEARCHES

NEIGHBORS\_CLASSIFIER\_RADIUS\_NEIGHBORS\_CLASSIFIER FOR SUPERVISED CLASSIFICATION PROBLEMS

AND KNEIGHBORS\_REGRESSOR\_RADIUS\_NEIGHBORS\_REGRESSOR FOR SUPERVISED REGRESSION PROBLEMS

- SKLEARN\_BALL\_TREE BALL\_TREE HAS BEEN MOVED TO SKLEARN\_NEIGHBORS\_BALL\_TREE USING THE FORMER WILL GENERATE A WARNING

- SKLEARN\_LINEAR\_MODEL LARS AND RELATED CLASSES LASSO LARS LASSO\_LARS CV ETC HAVE BEEN RENAMED TO SKLEARN\_LINEAR\_MODEL LARS

- ALL DISTANCE METRICS AND KERNELS IN SKLEARN\_METRICS PAIRWISE NOW HAVE A Y PARAMETER WHICH BY DEFAULT IS NONE IF NOT GIVEN THE RESULT IS THE DISTANCE OR KERNEL SIMILARITY BETWEEN EACH SAMPLE IN Y IF GIVEN THE RESULT IS THE PAIRWISE DISTANCE OR KERNEL SIMILARITY BETWEEN SAMPLES IN X TO Y

- SKLEARN\_METRICS PAIRWISE\_L1\_DISTANCE IS NOW CALLED MANHATTAN\_DISTANCE AND BY DEFAULT RETURNS THE PAIRWISE DISTANCE FOR THE COMPONENT WISE DISTANCE SET THE PARAMETER SUM\_OVER\_FEATURES TO FALSE

BACKWARD COMPATIBILITY PACKAGE ALIASES AND OTHER DEPRECATED CLASSES AND FUNCTIONS WILL BE REMOVED IN VERSION 0.11

PEOPLE

38 PEOPLE CONTRIBUTED TO THIS RELEASE

- 387 VLAD NICULAE

136 CHAPTER 1 WELCOME TO SCIKITLEARN

SCIKITLEARN USER GUIDE RELEASE 0213

- 320 OLIVIER GRISEL
  - 192 LARS BUITINCK
  - 179 GAEL VAROQUAUX
  - 168 FABIAN PEDREGOSA INRIA PARIETAL TEAM
  - 127 JAKE VANDERPLAS
  - 120 MATHIEU BLONDEL
  - 85 ALEXANDRE PASSOS
  - 67 ALEXANDRE GRAMFORT
  - 57 PETER PRETTENHOFER
  - 56 GILLES LOUPPE
  - 42 ROBERT LAYTON
  - 38 NELLE VAROQUAUX
  - 32 JEAN KOSSAIFI
  - 30 CONRAD LEE
  - 22 PIETRO BERKES
  - 18 ANDY
  - 17 DAVID WARDEFARLEY
  - 12 BRIAN HOLT
  - 11 ROBERT
  - 8 AMIT AIDES
  - 8 VIRGILE FRITSCH
  - 7 YAROSLAV HALCHENKO
  - 6 SALVATORE MASECCHIA
  - 5 PAOLO LOSI
  - 4 VINCENT SCHUT
  - 3 ALEXIS METAIREAU
  - 3 BRYAN SILVERTHORN
  - 3 ANDREAS MÜLLER
  - 2 MINWOO JAKE LEE
  - 1 EMMANUELLE GOUILLART
  - 1 KEITH GOODMAN
  - 1 LUCAS WIMAN
  - 1 NICOLAS PINTO
  - 1 THOUI RAY JONES
  - 1 TIM SHEERMANCHASE
- 117 PREVIOUS RELEASES 137

SCIKITLEARN USER GUIDE RELEASE 0213

11722 VERSION 08

MAY 11 2011

SCIKITLEARN 08 WAS RELEASED ON MAY 2011 ONE MONTH AFTER THE FIRST “INTERNATIONAL” SCIKITLEARN CODING SPRINT AND IS MARKED BY THE INCLUSION OF IMPORTANT MODULES HIERARCHICAL CLUSTERING CROSS DECOMPOSITION NONNEGATIVE MATRIX FACTORIZATION NMF OR NNMF INITIAL SUPPORT FOR PYTHON 3 AND BY IMPORTANT ENHANCEMENTS AND BUG FIXES

CHANGELOG

SEVERAL NEW MODULES WERE INTRODUCED DURING THIS RELEASE

- NEW HIERARCHICAL CLUSTERING MODULE BY VINCENT MICHEL BERTRAND THIRION ALEXANDRE GRAMFORT AND GAELE VAROQUAUX

- KERNEL PCA IMPLEMENTATION BY MATHIEU BLONDEL

- LABELED FACES IN THE WILD BY OLIVIER GRISEL

- NEW CROSS DECOMPOSITION MODULE BY EDOUARD DUCHESNAY

- NONNEGATIVE MATRIX FACTORIZATION NMF OR NNMF MODULE VLAD NICULAE

- IMPLEMENTATION OF THE ORACLE APPROXIMATING SHRINKAGE ALGORITHM BY VIRGILE FRITSCH IN THE COVARIANCE ESTIMATION MODULE

SOME OTHER MODULES BENEFITED FROM SIGNIFICANT IMPROVEMENTS OR CLEANUPS

- INITIAL SUPPORT FOR PYTHON 3 BUILDS AND IMPORTS CLEANLY SOME MODULES ARE USABLE WHILE OTHERS HAVE FAILING TESTS BY FABIAN PEDREGOSA

- DECOMPOSITION PCA IS NOW USABLE FROM THE PIPELINE OBJECT BY OLIVIER GRISEL

- GUIDE HOW TO OPTIMIZE FOR SPEED BY OLIVIER GRISEL

- FIXES FOR MEMORY LEAKS IN LIBSVM BINDINGS 64BIT SAFER BALLTREE BY LARS BUITINCK

- BUG AND STYLE FIXING IN KMEANS ALGORITHM BY JAN SCHLÜTER

- ADD ATTRIBUTE CONVERGED TO GAUSSIAN MIXTURE MODELS BY VINCENT SCHUT

- IMPLEMENTED TRANSFORM PREDICT LOG PROBA IN DISCRIMINANT ANALYSIS

LINEAR DISCRIMINANT ANALYSIS BY MATHIEU BLONDEL

- REFACTORING IN THE SUPPORT VECTOR MACHINES MODULE AND BUG FIXES BY FABIAN PEDREGOSA GAELE VAROQUAUX AND AMIT AIDES

- REFACTORED SGD MODULE REMOVED CODE DUPLICATION BETTER VARIABLE NAMING ADDED INTERFACE FOR SAMPLE WEIGHT BY PETER PRETTENHOFER

- WRAPPED BALLTREE WITH CYTHON BY THOUIS RAY JONES

- ADDED FUNCTION SVML1 MINC BY PAOLO LOSI

- TYPOS DOC STYLE ETC BY YAROSLAV HALCHENKO GAELE VAROQUAUX OLIVIER GRISEL YANN MALET NICOLAS PINTO LARS BUITINCK AND FABIAN PEDREGOSA

PEOPLE

PEOPLE THAT MADE THIS RELEASE POSSIBLE PRECEDED BY NUMBER OF COMMITS

- 159 OLIVIER GRISEL

138 CHAPTER 1 WELCOME TO SCIKITLEARN



SCIKITLEARN USER GUIDE RELEASE 0213

- 96 GAELE VAROQUAUX
- 96 VLAD NICULAE
- 94 FABIAN PEDREGOSA
- 36 ALEXANDRE GRAMFORT
- 32 PAOLO LOSI
- 31 EDOUARD DUCHESNAY
- 30 MATHIEU BLONDEL
- 25 PETER PRETTENHOFER
- 22 NICOLAS PINTO
- 11 VIRGILE FRITSCH
- 7 LARS BUITINCK
- 6 VINCENT MICHEL
- 5 BERTRAND THIRION
- 4 THOUIS RAY JONES
- 4 VINCENT SCHUT
- 3 JAN SCHLÜTER
- 2 JULIEN MIOTTE
- 2 MATTHIEU PERROT
- 2 YANN MALET
- 2 YAROSLAV HALCHENKO
- 1 AMIT AIDES
- 1 ANDREAS MÜLLER
- 1 FETH AREZKI
- 1 MENG XINFAN

11723 VERSION 07

MARCH 2 2011

SCIKITLEARN 07 WAS RELEASED IN MARCH 2011 ROUGHLY THREE MONTHS AFTER THE 06 RELEASE THIS RELEASE IS MARKED BY THE SPEED IMPROVEMENTS IN EXISTING ALGORITHMS LIKE KNEAREST NEIGHBORS AND KMEANS ALGORITHM AND BY THE INCLUSION OF AN EFFICIENT ALGORITHM FOR COMPUTING THE RIDGE GENERALIZED CROSS VALIDATION SOLUTION UNLIKE THE PRECEDING RELEASE NO NEW MODULES WERE ADDED TO THIS RELEASE

CHANGELOG

- PERFORMANCE IMPROVEMENTS FOR GAUSSIAN MIXTURE MODEL SAMPLING JAN SCHLÜTER
- IMPLEMENTATION OF EFFICIENT LEAVEONEOUT CROSSVALIDATED RIDGE IN LINEARMODELRIDGECV MATHIEU BLONDEL

117 PREVIOUS RELEASES 139

SCIKITLEARN USER GUIDE RELEASE 0213

- BETTER HANDLING OF COLLINEARITY AND EARLY STOPPING IN LINEARMODELLARSPATH ALEXANDRE GRAMFORT AND FABIAN PEDREGOSA
- FIXES FOR LIBLINEAR ORDERING OF LABELS AND SIGN OF COEFFICIENTS DAN YAMINS PAOLO LOSI MATHIEU BLONDEL AND FABIAN PEDREGOSA
- PERFORMANCE IMPROVEMENTS FOR NEAREST NEIGHBORS ALGORITHM IN HIGHDIMENSIONAL SPACES FABIAN PEDREGOSA
- PERFORMANCE IMPROVEMENTS FOR CLUSTERKMEANS GAELE VAROQUAUX AND JAMES BERGSTR
- SANITY CHECKS FOR SVMBASED CLASSES MATHIEU BLONDEL
- REFACTORING OF NEIGHBORSNEIGHBORSCLASSIFIER ANDNEIGHBORSKNEIGHBORSGRAPH ADDED

DIFFERENT ALGORITHMS FOR THE KNEAREST NEIGHBOR SEARCH AND IMPLEMENTED A MORE STABLE ALGORITHM FOR FINDING BARYCENTER WEIGHTS ALSO ADDED SOME DEVELOPER DOCUMENTATION FOR THIS MODULE SEE NOTESNEIGHBORS FOR MORE INFORMATION FABIAN PEDREGOSA

- DOCUMENTATION IMPROVEMENTS ADDED PCARANDOMIZEDPCA ANDLINEARMODEL LOGISTICREGRESSION TO THE CLASS REFERENCE ALSO ADDED REFERENCES OF MATRICES USED FOR CLUSTERING AND OTHER FIXES GAELE VAROQUAUX FABIAN PEDREGOSA MATHIEU BLONDEL OLIVIER GRISEL VIRGILE FRITSCH EMMANUELLE GOUILLART

• BINDED DECISIONFUNCTION IN CLASSES THAT MAKE USE OF LIBLINEAR DENSE AND SPARSE VARIANTS LIKE SVM

- LINEARSVC ORLINEARMODELLOGISTICREGRESSION FABIAN PEDREGOSA
- PERFORMANCE AND API IMPROVEMENTS TO METRICSEUCLIDEANDDISTANCES AND TOPCA RANDOMIZEDPCA JAMES BERGSTR

- FIX COMPILATION ISSUES UNDER NETBSD KAMELEBEN HASSEN DEROUICHE
- ALLOW INPUT SEQUENCES OF DIFFERENT LENGTHS IN HMMGAUSSIANHMM RON WEISS
- FIX BUG IN AFFINITY PROPAGATION CAUSED BY INCORRECT INDEXING XINFAN MENG

PEOPLE THAT MADE THIS RELEASE POSSIBLE PRECEDED BY NUMBER OF COMMITS

- 85 FABIAN PEDREGOSA
- 67 MATHIEU BLONDEL
- 20 ALEXANDRE GRAMFORT
- 19 JAMES BERGSTR
- 14 DAN YAMINS
- 13 OLIVIER GRISEL
- 12 GAELE VAROQUAUX
- 4 EDOUARD DUCHESNAY
- 4 RON WEISS
- 2 SATRAJIT GHOSH
- 2 VINCENT DUBOURG
- 1 EMMANUELLE GOUILLART
- 1 KAMELEBEN HASSEN DEROUICHE
- 1 PAOLO LOSI

140 CHAPTER 1 WELCOME TO SCIKITLEARN

SCIKITLEARN USER GUIDE RELEASE 0213

- 1 VIRGILEFRITSCH
- 1 YAROSLAV HALCHENKO
- 1 XINFAN MENG

11724 VERSION 06  
DECEMBER 21 2010

SCIKITLEARN 06 WAS RELEASED ON DECEMBER 2010 IT IS MARKED BY THE INCLUSION OF SEVERAL NEW MODULES AND A GENERAL RENAMING OF OLD ONES IT IS ALSO MARKED BY THE INCLUSION OF NEW EXAMPLE INCLUDING APPLICATIONS TO REALWORLD DATASETS

CHANGELOG

- NEW STOCHASTIC GRADIENT DESCENT MODULE BY PETER PRETTENHOFER THE MODULE COMES WITH COMPLETE DOCUMENTATION AND EXAMPLES
- IMPROVED SVM MODULE MEMORY CONSUMPTION HAS BEEN REDUCED BY 50 HEURISTIC TO AUTOMATICALLY SET CLASS WEIGHTS POSSIBILITY TO ASSIGN WEIGHTS TO SAMPLES SEE SVM WEIGHTED SAMPLES FOR AN EXAMPLE
- NEW GAUSSIAN PROCESSES MODULE BY VINCENT DUBOURG THIS MODULE ALSO HAS GREAT DOCUMENTATION AND SOME VERY NEAT EXAMPLES SEE EXAMPLEGAUSSIANPROCESSPLOTGPREGRESSIONPY OR EXAM PLEGAUSSIANPROCESSPLOTGPPROBABILISTICCLASSIFICATIONAFTERREGRESSIONPY FOR A TASTE OF WHAT CAN BE DONE
- IT IS NOW POSSIBLE TO USE LIBLINEAR'S MULTICLASS SVC OPTION MULTICLASS IN SVMLINEARSVC
- NEW FEATURES AND PERFORMANCE IMPROVEMENTS OF TEXT FEATURE EXTRACTION
- IMPROVED SPARSE MATRIX SUPPORT BOTH IN MAIN CLASSES GRIDSEARCHGRIDSEARCHCV AS IN MODULES SKLEARN SVMSPARSE AND SKLEARN LINEARMODELSPARSE
- LOTS OF COOL NEW EXAMPLES AND A NEW SECTION THAT USES REALWORLD DATASETS WAS CREATED THESE INCLUDE FACES RECOGNITION EXAMPLE USING EIGENFACES AND SVMS SPECIES DISTRIBUTION MODELING LIBSVM GUI WIKIPEDIA PRINCIPAL EIGENVECTOR AND OTHERS
- FASTER LEAST ANGLE REGRESSION ALGORITHM IT IS NOW 2X FASTER THAN THE R VERSION ON WORST CASE AND UP TO 10X TIMES FASTER ON SOME CASES
- FASTER COORDINATE DESCENT ALGORITHM IN PARTICULAR THE FULL PATH VERSION OF LASSO LINEARMODEL LASSOPATH IS MORE THAN 200X TIMES FASTER THAN BEFORE
- IT IS NOW POSSIBLE TO GET PROBABILITY ESTIMATES FROM A LINEARMODELLOGISTICREGRESSION MODEL
- MODULE RENAMING THE GLM MODULE HAS BEEN RENAMED TO LINEARMODEL THE GMM MODULE HAS BEEN INCLUDED INTO THE MORE GENERAL MIXTURE MODEL AND THE SGD MODULE HAS BEEN INCLUDED IN LINEARMODEL
- LOTS OF BUG FIXES AND DOCUMENTATION IMPROVEMENTS

PEOPLE  
PEOPLE THAT MADE THIS RELEASE POSSIBLE PRECEDED BY NUMBER OF COMMITS

- 207 OLIVIER GRISEL
  - 167 FABIAN PEDREGOSA
  - 97 PETER PRETTENHOFER
  - 68 ALEXANDRE GRAMFORT
- 117 PREVIOUS RELEASES 141

SCIKITLEARN USER GUIDE RELEASE 0213

- 59 MATHIEU BLONDEL
- 55 GAELE VAROQUAUX
- 33 VINCENT DUBOURG
- 21 RON WEISS
- 9 BERTRAND THIRION
- 3 ALEXANDRE PASSOS
- 3 ANNELAURE FOUQUE
- 2 RONAN AMICEL
- 1 CHRISTIAN OSENDORFER

11725 VERSION 05

OCTOBER 11 2010

CHANGELOG

NEW CLASSES

- SUPPORT FOR SPARSE MATRICES IN SOME CLASSIFIERS OF MODULES SVM AND LINEAR MODEL SEESVM

SPARSE SVC SVM SPARSE SVR SVM SPARSE LINEAR SVC LINEAR MODEL SPARSE LASSO

LINEAR MODEL SPARSE ELASTIC NET

- NEW PIPELINE PIPELINE OBJECT TO COMPOSE DIFFERENT ESTIMATORS

- RECURSIVE FEATURE ELIMINATION ROUTINES IN MODULE FEATURE SELECTION

- ADDITION OF VARIOUS CLASSES CAPABLE OF CROSS VALIDATION IN THE LINEAR MODEL MODULE LINEAR MODEL

LASSO CV LINEAR MODEL ELASTIC NET CV ETC

- NEW MORE EFFICIENT LARS ALGORITHM IMPLEMENTATION THE LASSO VARIANT OF THE ALGORITHM IS ALSO IMPLEMENTED

SEE LINEAR MODEL LARS PATH LINEAR MODEL LARS AND LINEAR MODEL LASSO LARS

- NEW HIDDEN MARKOV MODELS MODULE SEE CLASSES HMM GAUSSIAN HMM HMM MULTINOMIAL HMM HMM

GMM HMM

- NEW MODULE FEATURE EXTRACTION SEE CLASS REFERENCE

- NEW FASTICA ALGORITHM IN MODULE SKLEARN FASTICA

DOCUMENTATION

- IMPROVED DOCUMENTATION FOR MANY MODULES NOW SEPARATING NARRATIVE DOCUMENTATION FROM THE CLASS REFERENCE

AS AN EXAMPLE SEE DOCUMENTATION FOR THE SVM MODULE AND THE COMPLETE CLASS REFERENCE

FIXES

- API CHANGES ADHERE VARIABLE NAMES TO PEP8 GIVE MORE MEANINGFUL NAMES

- FIXES FOR SVM MODULE TO RUN ON A SHARED MEMORY CONTEXT MULTIPROCESSING

- IT IS AGAIN POSSIBLE TO GENERATE LATEX AND THUS PDF FROM THE SPHINX DOCS

142 CHAPTER 1 WELCOME TO SCIKITLEARN

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES

- NEW EXAMPLES USING SOME OF THE MLCOMP DATASETS SPHXGLRAUTOEXAMPLESMLCOMPSPARSEDOCUMENTCLASSIFICATION PYSINCE REMOVED AND CLASSIFICATION OF TEXT DOCUMENTS USING SPARSE FEATURES
- MANY MORE EXAMPLES SEE HERE THE FULL LIST OF EXAMPLES

EXTERNAL DEPENDENCIES

- JOBLIB IS NOW A DEPENDENCY OF THIS PACKAGE ALTHOUGH IT IS SHIPPED WITH SKLEARNEXTERNALSJOLIB REMOVED MODULES

- MODULE ANN ARTIFICIAL NEURAL NETWORKS HAS BEEN REMOVED FROM THE DISTRIBUTION USERS WANTING THIS SORT OF ALGORITHMS SHOULD TAKE A LOOK INTO PYBRAIN

MISC

- NEW SPHINX THEME FOR THE WEB PAGE

AUTHORS

THE FOLLOWING IS A LIST OF AUTHORS FOR THIS RELEASE PRECEDED BY NUMBER OF COMMITS

- 262 FABIAN PEDREGOSA
- 240 GAELE VAROQUAUX
- 149 ALEXANDRE GRAMFORT
- 116 OLIVIER GRISEL
- 40 VINCENT MICHEL
- 38 RON WEISS
- 23 MATTHIEU PERROT
- 10 BERTRAND THIRION
- 7 YAROSLAV HALCHENKO
- 9 VIRGILEFRITSCH
- 6 EDOUARD DUCHESNAY
- 4 MATHIEU BLONDEL
- 1 ARIEL ROKEM
- 1 MATTHIEU BRUCHER

11726 VERSION 04

AUGUST 26 2010

117 PREVIOUS RELEASES 143

SCIKITLEARN USER GUIDE RELEASE 0213

CHANGELOG

MAJOR CHANGES IN THIS RELEASE INCLUDE

- COORDINATE DESCENT ALGORITHM LASSO ELASTICNET REFACTORING SPEED IMPROVEMENTS ROUGHLY 100X TIMES FASTER
- COORDINATE DESCENT REFACTORING AND BUG FIXING FOR CONSISTENCY WITH R'S PACKAGE GLMNET
- NEW METRICS MODULE
- NEW GMM MODULE CONTRIBUTED BY RON WEISS
- IMPLEMENTATION OF THE LARS ALGORITHM WITHOUT LASSO VARIANT FOR NOW
- FEATURESELECTION MODULE REDESIGN
- MIGRATION TO GIT AS VERSION CONTROL SYSTEM
- REMOVAL OF OBSOLETE ATTRSELECT MODULE
- RENAME OF PRIVATE COMPILED EXTENSIONS ADDED UNDERSCORE
- REMOVAL OF LEGACY UNMAINTAINED CODE
- DOCUMENTATION IMPROVEMENTS BOTH DOCSTRING AND RST
- IMPROVEMENT OF THE BUILD SYSTEM TO OPTIONALLY LINK WITH MKL ALSO PROVIDE A LITE BLAS IMPLEMENTATION IN CASE NO SYSTEMWIDE BLAS IS FOUND
- LOTS OF NEW EXAMPLES
- MANY MANY BUG FIXES

AUTHORS

THE COMMITTER LIST FOR THIS RELEASE IS THE FOLLOWING PRECEDED BY NUMBER OF COMMITS

- 143 FABIAN PEDREGOSA
- 35 ALEXANDRE GRAMFORT
- 34 OLIVIER GRISEL
- 11 GAEL VAROQUAUX
- 5 YAROSLAV HALCHENKO
- 2 VINCENT MICHEL
- 1 CHRIS FILO GORGOLEWSKI

11727 EARLIER VERSIONS

EARLIER VERSIONS INCLUDED CONTRIBUTIONS BY FRED MAILHOT DAVID COOKE DAVID HUARD DAVE MORRILL ED SCHOFIELD

TRAVIS OLIPHANT PEARU PETERSON

144 CHAPTER 1 WELCOME TO SCIKITLEARN

SCIKITLEARN USER GUIDE RELEASE 0213

118 ROADMAP

1181 PURPOSE OF THIS DOCUMENT

THIS DOCUMENT LIST GENERAL DIRECTIONS THAT CORE CONTRIBUTORS ARE INTERESTED TO SEE DEVELOPED IN SCIKITLEARN THE FACT THAT AN ITEM IS LISTED HERE IS IN NO WAY A PROMISE THAT IT WILL HAPPEN AS RESOURCES ARE LIMITED RATHER IT IS AN INDICATION THAT HELP IS WELCOMED ON THIS TOPIC

1182 STATEMENT OF PURPOSE SCIKITLEARN IN 2018

ELEVEN YEARS AFTER THE INCEPTION OF SCIKITLEARN MUCH HAS CHANGED IN THE WORLD OF MACHINE LEARNING KEY CHANGES INCLUDE

- COMPUTATIONAL TOOLS THE EXPLOITATION OF GPUS DISTRIBUTED PROGRAMMING FRAMEWORKS LIKE SCALASPARK ETC
- HIGHLEVEL PYTHON LIBRARIES FOR EXPERIMENTATION PROCESSING AND DATA MANAGEMENT JUPYTER NOTEBOOK CYTHON PANDAS DASK NUMBA

• CHANGES IN THE FOCUS OF MACHINE LEARNING RESEARCH ARTIFICIAL INTELLIGENCE APPLICATIONS WHERE INPUT STRUCTURE IS KEY WITH DEEP LEARNING REPRESENTATION LEARNING REINFORCEMENT LEARNING DOMAIN TRANSFER ETC

A MORE SUBTLE CHANGE OVER THE LAST DECADE IS THAT DUE TO CHANGING INTERESTS IN ML PHD STUDENTS IN MACHINE LEARNING ARE MORE LIKELY TO CONTRIBUTE TO PYTORCH DASK ETC THAN TO SCIKITLEARN SO OUR CONTRIBUTOR POOL IS VERY DIFFERENT TO A DECADE AGO

SCIKITLEARN REMAINS VERY POPULAR IN PRACTICE FOR TRYING OUT CANONICAL MACHINE LEARNING TECHNIQUES PARTICULARLY FOR APPLICATIONS IN EXPERIMENTAL SCIENCE AND IN DATA SCIENCE A LOT OF WHAT WE PROVIDE IS NOW VERY MATURE BUT IT CAN BE COSTLY TO MAINTAIN AND WE CANNOT THEREFORE INCLUDE ARBITRARY NEW IMPLEMENTATIONS YET SCIKITLEARN IS ALSO ESSENTIAL IN DEFINING AN API FRAMEWORK FOR THE DEVELOPMENT OF INTEROPERABLE MACHINE LEARNING COMPONENTS EXTERNAL TO THE CORE LIBRARY

THUS OUR MAIN GOALS IN THIS ERA ARE TO

- CONTINUE MAINTAINING A HIGHQUALITY WELLDOCUMENTED COLLECTION OF CANONICAL TOOLS FOR DATA PROCESSING AND MACHINE LEARNING WITHIN THE CURRENT SCOPE IE RECTANGULAR DATA LARGELY INVARIANT TO COLUMN AND ROW ORDER PREDICTING TARGETS WITH SIMPLE STRUCTURE
- IMPROVE THE EASE FOR USERS TO DEVELOP AND PUBLISH EXTERNAL COMPONENTS
- IMPROVE INTEROPERABILITY WITH MODERN DATA SCIENCE TOOLS EG PANDAS DASK AND INFRASTRUCTURES EG DISTRIBUTED PROCESSING

MANY OF THE MORE FINEGRAINED GOALS CAN BE FOUND UNDER THE API TAG ON THE ISSUE TRACKER

1183 ARCHITECTURAL GENERAL GOALS

THE LIST IS NUMBERED NOT AS AN INDICATION OF THE ORDER OF PRIORITY BUT TO MAKE REFERRING TO SPECIFIC POINTS EASIER PLEASE ADD NEW ENTRIES ONLY AT THE BOTTOM

1 EVERYTHING IN SCIKITLEARN SHOULD CONFORM TO OUR API CONTRACT

- PIPELINE ANDFEATUREUNION MODIFY THEIR INPUT PARAMETERS IN FIT FIXING THIS REQUIRES MAKING SURE

WE HAVE A GOOD GRASP OF THEIR USE CASES TO MAKE SURE ALL CURRENT FUNCTIONALITY IS MAINTAINED 8157 7382

2 IMPROVED HANDLING OF PANDAS DATAFRAMES AND SPARSEDATAFRAMES

- DOCUMENT CURRENT HANDLING
- COLUMN REORDERING ISSUE 7242

118 ROADMAP 145

SCIKITLEARN USER GUIDE RELEASE 0213

- AVOIDING UNNECESSARY CONVERSION TO NDARRAY
- RETURNING DATAFRAMES FROM TRANSFORMERS 5523
- GETTING DATAFRAMES FROM DATASET LOADERS
- SPARSE CURRENTLY NOT CONSIDERED

3 IMPROVED HANDLING OF CATEGORICAL FEATURES

- TREEBASED MODELS SHOULD BE ABLE TO HANDLE BOTH CONTINUOUS AND CATEGORICAL FEATURES 4899
- IN DATASET LOADERS
- AS GENERIC TRANSFORMERS TO BE USED WITH COLUMNTRANSFORMS EG ORDINAL ENCODING SUPERVISED BY CORRELATION WITH TARGET VARIABLE

4 IMPROVED HANDLING OF MISSING DATA

- MAKING SURE METAESTIMATORS ARE LENIENT TOWARDS MISSING DATA
- NONTRIVIAL IMPUTERS
- LEARNERS DIRECTLY HANDLING MISSING DATA
- AN AMPUTATION SAMPLE GENERATOR TO MAKE PARTS OF A DATASET GO MISSING
- HANDLING MIXTURES OF CATEGORICAL AND CONTINUOUS VARIABLES

5 PASSING AROUND INFORMATION THAT IS NOT X Y SAMPLE PROPERTIES

- WE NEED TO BE ABLE TO PASS SAMPLE WEIGHTS TO SCORERS IN CROSS VALIDATION
- WE SHOULD HAVE STANDARDGENERALISED WAYS OF PASSING SAMPLEWISE PROPERTIES AROUND IN METAESTIMATORS 4497 7646

6 PASSING AROUND INFORMATION THAT IS NOT X Y FEATURE PROPERTIES

- FEATURE NAMES OR DESCRIPTIONS SHOULD IDEALLY BE AVAILABLE TO FIT FOR EG 6425 6424
- PERFEATURE HANDLING EG “IS THIS A NOMINAL ORDINAL ENGLISH LANGUAGE TEXT” SHOULD ALSO NOT NEED TO BE PROVIDED TO ESTIMATOR CONSTRUCTORS IDEALLY BUT SHOULD BE AVAILABLE AS METADATA ALONGSIDE X 8480

7 PASSING AROUND INFORMATION THAT IS NOT X Y TARGET INFORMATION

- WE HAVE PROBLEMS GETTING THE FULL SET OF CLASSES TO ALL COMPONENTS WHEN THE DATA IS SPLITSAMPLED 6231 8100

- WE HAVE NO WAY TO HANDLE A MIXTURE OF CATEGORICAL AND CONTINUOUS TARGETS

8 MAKE IT EASIER FOR EXTERNAL USERS TO WRITE SCIKITLEARNCOMPATIBLE COMPONENTS

- MORE FLEXIBLE ESTIMATOR CHECKS THAT DO NOT SELECT BY ESTIMATOR NAME 6599 6715
- EXAMPLE OF HOW TO DEVELOP A METAESTIMATOR
- MORE SELFSUFFICIENT RUNNING OF SCIKITLEARNCONTRIB OR A SIMILAR RESOURCE

9 SUPPORT RESAMPLING AND SAMPLE REDUCTION

- ALLOW SUBSAMPLING OF MAJORITY CLASSES IN A PIPELINE 3855
- IMPLEMENT RANDOM FORESTS WITH RESAMPLING 8732

10 BETTER INTERFACES FOR INTERACTIVE DEVELOPMENT

- REPR AND HTML VISUALISATIONS OF ESTIMATORS 6323
- INCLUDE PLOTTING TOOLS NOT JUST AS EXAMPLES 9173

146 CHAPTER 1 WELCOME TO SCIKITLEARN



SCIKITLEARN USER GUIDE RELEASE 0213

11 IMPROVED TOOLS FOR MODEL DIAGNOSTICS AND BASIC INFERENCE

- ALTERNATIVE FEATURE IMPORTANCES IMPLEMENTATIONS EG METHODS OR WRAPPERS
- BETTER WAYS TO HANDLE VALIDATION SETS WHEN FITTING
- BETTER WAYS TO FIND THRESHOLDS CREATE DECISION RULES 8614

12 BETTER TOOLS FOR SELECTING HYPERPARAMETERS WITH TRANSDUCTIVE ESTIMATORS

- GRID SEARCH AND CROSS VALIDATION ARE NOT APPLICABLE TO MOST CLUSTERING TASKS STABILITYBASED SELECTION IS MORE RELEVANT

13 IMPROVED TRACKING OF FITTING

- VERBOSE IS NOT VERY FRIENDLY AND SHOULD USE A STANDARD LOGGING LIBRARY 6929
- CALLBACKS OR A SIMILAR SYSTEM WOULD FACILITATE LOGGING AND EARLY STOPPING

14 DISTRIBUTED PARALLELISM

- JOBLIB CAN NOW PLUG ONTO SEVERAL BACKENDS SOME OF THEM CAN DISTRIBUTE THE COMPUTATION ACROSS COMPUTERS
- HOWEVER WE WANT TO STAY HIGH LEVEL IN SCIKITLEARN

15 A WAY FORWARD FOR MORE OUT OF CORE

- DASK ENABLES EASY OUTOF CORE COMPUTATION WHILE THE DASK MODEL PROBABLY CANNOT BE ADAPTABLE TO ALL MACHINELEARNING ALGORITHMS MOST MACHINE LEARNING IS ON SMALLER DATA THAN ETL HENCE WE CAN MAYBE ADAPT TO VERY LARGE SCALE WHILE SUPPORTING ONLY A FRACTION OF THE PATTERNS

16 BETTER SUPPORT FOR MANUAL AND AUTOMATIC PIPELINE BUILDING

- EASIER WAY TO CONSTRUCT COMPLEX PIPELINES AND VALID SEARCH SPACES 7608 5082 8243
- PROVIDE SEARCH RANGES FOR COMMON ESTIMATORS
- CF SEARCHGRID

17 SUPPORT FOR WORKING WITH PRETRAINED MODELS

- ESTIMATOR “FREEZING” IN PARTICULAR RIGHT NOW IT’S IMPOSSIBLE TO CLONE A CALIBRATEDCLASSIFIERCV WITH PREFIT 8370 6451

18 BACKWARDSCOMPATIBLE DESERIALIZATION OF SOME ESTIMATORS

- CURRENTLY SERIALIZATION WITH PICKLE BREAKS ACROSS VERSIONS WHILE WE MAY NOT BE ABLE TO GET AROUND OTHER LIMITATIONS OF PICKLE RE SECURITY ETC IT WOULD BE GREAT TO OFFER CROSSVERSION SAFETY FROM VERSION 10 NOTE GAEL AND OLIVIER THINK THAT THIS CAN CAUSE HEAVY MAINTENANCE BURDEN AND WE SHOULD MANAGE THE TRADEOFFS A POSSIBLE ALTERNATIVE IS PRESENTED IN THE FOLLOWING POINT

19 DOCUMENTATION AND TOOLING FOR MODEL LIFECYCLE MANAGEMENT

- DOCUMENT GOOD PRACTICES FOR MODEL DEPLOYMENTS AND LIFECYCLE BEFORE DEPLOYING A MODEL SNAPSHOT THE CODE VERSIONS NUMPY SCIPY SCIKITLEARN CUSTOM CODE REPO THE TRAINING SCRIPT AND AN ALIAS ON HOW TO RETRIEVE HISTORICAL TRAINING DATA SNAPSHOT A COPY OF A SMALL VALIDATION SET SNAPSHOT OF THE PREDICTIONS PREDICTED PROBABILITIES FOR CLASSIFIERS ON THAT VALIDATION SET
- DOCUMENT AND TOOLS TO MAKE IT EASY TO MANAGE UPGRADE OF SCIKITLEARN VERSIONS
- TRY TO LOAD THE OLD PICKLE IF IT WORKS USE THE VALIDATION SET PREDICTION SNAPSHOT TO DETECT THAT THE SERIALIZED MODEL STILL BEHAVE THE SAME
- IF JOBLIBLOAD PICKLELOAD NOT WORK USE THE VERSIONED CONTROL TRAINING SCRIPT HISTORICAL TRAINING SET TO RETRAIN THE MODEL AND USE THE VALIDATION SET PREDICTION SNAPSHOT TO ASSERT THAT IT IS POSSIBLE TO RECOVER THE PREVIOUS PREDICTIVE PERFORMANCE IF THIS IS NOT THE CASE THERE IS PROBABLY A BUG IN SCIKITLEARN THAT NEEDS TO BE REPORTED

118 ROADMAP 147

SCIKITLEARN USER GUIDE RELEASE 0213

20 OPTIONAL IMPROVE SCIKITLEARN COMMON TESTS SUITE TO MAKE SURE THAT AT LEAST FOR FREQUENTLY USED MODELS HAVE STABLE PREDICTIONS ACROSSVERSIONS TO BE DISCUSSED

- EXTEND DOCUMENTATION TO MENTION HOW TO DEPLOY MODELS IN PYTHONFREE ENVIRONMENTS FOR INSTANCE ONNX AND USE THE ABOVE BEST PRACTICES TO ASSESS PREDICTIVE CONSISTENCY BETWEEN SCIKITLEARN AND ONNX PREDICTION FUNCTIONS ON VALIDATION SET

- DOCUMENT GOOD PRACTICES TO DETECT TEMPORAL DISTRIBUTION DRIFT FOR DEPLOYED MODEL AND GOOD PRACTICES FOR RETRAINING ON FRESH DATA WITHOUT CAUSING CATASTROPHIC PREDICTIVE PERFORMANCE REGRESSIONS

21 MORE DIDACTIC DOCUMENTATION

- MORE AND MORE OPTIONS HAVE BEEN ADDED TO SCIKITLEARN AS A RESULT THE DOCUMENTATION IS CROWDED WHICH MAKES IT HARD FOR BEGINNERS TO GET THE BIG PICTURE SOME WORK COULD BE DONE IN PRIORITIZING THE INFORMATION

1184 SUBPACKAGESPECIFIC GOALS

SKLEARNCLUSTER

- KMEANS VARIANTS FOR NONEUCLIDEAN DISTANCES IF WE CAN SHOW THESE HAVE BENEFITS BEYOND HIERARCHICAL CLUSTERING

SKLEARNENSEMBLE

- A STACKING IMPLEMENTATION

SKLEARNMODELSELECTION

- MULTIMETRIC SCORING IS SLOW 9326

- PERHAPS WE WANT TO BE ABLE TO GET BACK MORE THAN MULTIPLE METRICS

- THE HANDLING OF RANDOM STATES IN CV SPLITTERS IS A POOR DESIGN AND CONTRADICTS THE VALIDATION OF SIMILAR PARAMETERS IN ESTIMATORS

- EXPLOIT WARMSTARTING AND PATH ALGORITHMS SO THE BENEFITS OF ESTIMATORCV OBJECTS CAN BE ACCESSED VIA

GRIDSEARCHCV AND USED IN PIPELINES 1626

- CROSSVALIDATION SHOULD BE ABLE TO BE REPLACED BY OOB ESTIMATES WHENEVER A CROSSVALIDATION ITERATOR IS USED

- REDUNDANT COMPUTATIONS IN PIPELINES SHOULD BE AVOIDED RELATED TO POINT ABOVE CF DASKML

SKLEARNNEIGHBORS

- ABILITY TO SUBSTITUTE A CUSTOMAPPROXIMATEPRECOMPUTED NEAREST NEIGHBORS IMPLEMENTATION FOR OURS IN ALLMOST CONTEXTS THAT NEAREST NEIGHBORS ARE USED FOR LEARNING 10463

SKLEARNPIPELINE

- PERFORMANCE ISSUES WITH PIPELINEMEMORY

- SEE “EVERYTHING IN SCIKITLEARN SHOULD CONFORM TO OUR API CONTRACT” ABOVE

119 SCIKITLEARN GOVERNANCE AND DECISIONMAKING

THE PURPOSE OF THIS DOCUMENT IS TO FORMALIZE THE GOVERNANCE PROCESS USED BY THE SCIKITLEARN PROJECT TO CLARIFY HOW DECISIONS ARE MADE AND HOW THE VARIOUS ELEMENTS OF OUR COMMUNITY INTERACT THIS DOCUMENT ESTABLISHES A DECISION MAKING STRUCTURE THAT TAKES INTO ACCOUNT FEEDBACK FROM ALL MEMBERS OF THE COMMUNITY AND STRIVES TO FIND CONSENSUS WHILE AVOIDING ANY DEADLOCKS

148 CHAPTER 1 WELCOME TO SCIKITLEARN

SCIKITLEARN USER GUIDE RELEASE 0213

THIS IS A MERITOCRATIC CONSENSUSBASED COMMUNITY PROJECT ANYONE WITH AN INTEREST IN THE PROJECT CAN JOIN THE COMMUNITY CONTRIBUTE TO THE PROJECT DESIGN AND PARTICIPATE IN THE DECISION MAKING PROCESS THIS DOCUMENT DESCRIBES HOW THAT PARTICIPATION TAKES PLACE AND HOW TO SET ABOUT EARNING MERIT WITHIN THE PROJECT COMMUNITY

1191 ROLES AND RESPONSIBILITIES

CONTRIBUTORS

CONTRIBUTORS ARE COMMUNITY MEMBERS WHO CONTRIBUTE IN CONCRETE WAYS TO THE PROJECT ANYONE CAN BECOME A CONTRIBUTOR AND CONTRIBUTIONS CAN TAKE MANY FORMS - NOT ONLY CODE - AS DETAILED IN THE CONTRIBUTORS GUIDE

CORE DEVELOPERS

CORE DEVELOPERS ARE COMMUNITY MEMBERS WHO HAVE SHOWN THAT THEY ARE DEDICATED TO THE CONTINUED DEVELOPMENT OF THE PROJECT THROUGH ONGOING ENGAGEMENT WITH THE COMMUNITY THEY HAVE SHOWN THEY CAN BE TRUSTED TO MAINTAIN SCIKITLEARN WITH CARE BEING A CORE DEVELOPER ALLOWS CONTRIBUTORS TO MORE EASILY CARRY ON WITH THEIR PROJECT RELATED ACTIVITIES BY GIVING THEM DIRECT ACCESS TO THE PROJECT'S REPOSITORY AND IS REPRESENTED AS BEING AN ORGANIZATION MEMBER ON THE SCIKITLEARN GITHUB ORGANIZATION CORE DEVELOPERS ARE EXPECTED TO REVIEW CODE CONTRIBUTIONS CAN MERGE APPROVED PULL REQUESTS CAN CAST VOTES FOR AND AGAINST MERGING A PULLREQUEST AND CAN BE INVOLVED IN DECIDING MAJOR CHANGES TO THE API

NEW CORE DEVELOPERS CAN BE NOMINATED BY ANY EXISTING CORE DEVELOPERS ONCE THEY HAVE BEEN NOMINATED THERE WILL BE A VOTE BY THE CURRENT CORE DEVELOPERS VOTING ON NEW CORE DEVELOPERS IS ONE OF THE FEW ACTIVITIES THAT TAKES PLACE ON THE PROJECT'S PRIVATE MANAGEMENT LIST WHILE IT IS EXPECTED THAT MOST VOTES WILL BE UNANIMOUS A TWO THIRDS MAJORITY OF THE CAST VOTES IS ENOUGH THE VOTE NEEDS TO BE OPEN FOR AT LEAST 1 WEEK

CORE DEVELOPERS THAT HAVE NOT CONTRIBUTED TO THE PROJECT COMMITS OR GITHUB COMMENTS IN THE PAST 12 MONTHS WILL BE ASKED IF THEY WANT TO BECOME EMERITUS CORE DEVELOPERS AND RECAN THEIR COMMIT AND VOTING RIGHTS UNTIL THEY BECOME ACTIVE AGAIN THE LIST OF CORE DEVELOPERS ACTIVE AND EMERITUS WITH DATES AT WHICH THEY BECAME ACTIVE IS PUBLIC ON THE SCIKITLEARN WEBSITE

TECHNICAL COMMITTEE

THE TECHNICAL COMMITTEE TC MEMBERS ARE CORE DEVELOPERS WHO HAVE ADDITIONAL RESPONSIBILITIES TO ENSURE THE SMOOTH RUNNING OF THE PROJECT TC MEMBERS ARE EXPECTED TO PARTICIPATE IN STRATEGIC PLANNING AND APPROVE CHANGES TO THE GOVERNANCE MODEL THE PURPOSE OF THE TC IS TO ENSURE A SMOOTH PROGRESS FROM THE BIGPICTURE PERSPECTIVE INDEED CHANGES THAT IMPACT THE FULL PROJECT REQUIRE A SYNTHETIC ANALYSIS AND A CONSENSUS THAT IS BOTH EXPLICIT AND INFORMED IN CASES THAT THE CORE DEVELOPER COMMUNITY WHICH INCLUDES THE TC MEMBERS FAILS TO REACH SUCH A CONSENSUS IN THE REQUIRED TIME FRAME THE TC IS THE ENTITY TO RESOLVE THE ISSUE MEMBERSHIP OF THE TC IS BY NOMINATION BY A CORE DEVELOPER A NOMINATION WILL RESULT IN DISCUSSION WHICH CANNOT TAKE MORE THAN A MONTH AND THEN A VOTE BY THE CORE DEVELOPERS WHICH WILL STAY OPEN FOR A WEEK TC MEMBERSHIP VOTES ARE SUBJECT TO A TWO THIRDS MAJORITY OF ALL CAST VOTES AS WELL AS A SIMPLE MAJORITY APPROVAL OF ALL THE CURRENT TC MEMBERS TC MEMBERS WHO DO NOT ACTIVELY ENGAGE WITH THE TC DUTIES ARE EXPECTED TO RESIGN

THE INITIAL TECHNICAL COMMITTEE OF SCIKITLEARN CONSISTS OF ALEXANDRE GRAMFORT OLIVIER GRISEL ANDREAS MÜLLER JOEL NOTHMAN HANMIN QIN GAËL VAROQUAUX AND ROMAN YURCHAK

1192 DECISION MAKING PROCESS

DECISIONS ABOUT THE FUTURE OF THE PROJECT ARE MADE THROUGH DISCUSSION WITH ALL MEMBERS OF THE COMMUNITY ALL NON SENSITIVE PROJECT MANAGEMENT DISCUSSION TAKES PLACE ON THE PROJECT CONTRIBUTORS' MAILING LIST AND THE ISSUE TRACKER OCCASIONALLY SENSITIVE DISCUSSION OCCURS ON A PRIVATE LIST

119 SCIKITLEARN GOVERNANCE AND DECISIONMAKING 149

SCIKITLEARN USER GUIDE RELEASE 0213

SCIKITLEARN USES A “CONSENSUS SEEKING” PROCESS FOR MAKING DECISIONS THE GROUP TRIES TO FIND A RESOLUTION THAT HAS NO OPEN OBJECTIONS AMONG CORE DEVELOPERS AT ANY POINT DURING THE DISCUSSION ANY COREDEVELOPER CAN CALL FOR A VOTE WHICH WILL CONCLUDE ONE MONTH FROM THE CALL FOR THE VOTE ANY VOTE MUST BE BACKED BY A SLEP IF NO OPTION CAN GATHER TWO THIRDS OF THE VOTES CAST THE DECISION IS ESCALATED TO THE TC WHICH IN TURN WILL USE CONSENSUS SEEKING WITH THE FALLBACK OPTION OF A SIMPLE MAJORITY VOTE IF NO CONSENSUS CAN BE FOUND WITHIN A MONTH THIS IS WHAT WE HEREAFTER MAY REFER TO AS “THE DECISION MAKING PROCESS”

DECISIONS IN ADDITION TO ADDING CORE DEVELOPERS AND TC MEMBERSHIP AS ABOVE ARE MADE ACCORDING TO THE FOLLOWING RULES

- MINOR DOCUMENTATION CHANGES SUCH AS TYPO FIXES OR ADDITION CORRECTION OF A SENTENCE BUT NO CHANGE OF THE SCIKITLEARNORG LANDING PAGE OR THE “ABOUT” PAGE REQUIRES 1 BY A CORE DEVELOPER NO 1 BY A CORE DEVELOPER LAZY CONSENSUS HAPPENS ON THE ISSUE OR PULL REQUEST PAGE CORE DEVELOPERS ARE EXPECTED TO GIVE “REASONABLE TIME” TO OTHERS TO GIVE THEIR OPINION ON THE PULL REQUEST IF THEY’RE NOT CONFIDENT OTHERS WOULD AGREE
- CODE CHANGES AND MAJOR DOCUMENTATION CHANGES REQUIRE 1 BY TWO CORE DEVELOPERS NO 1 BY A CORE DEVELOPER LAZY CONSENSUS HAPPENS ON THE ISSUE OF PULLREQUEST PAGE
- CHANGES TO THE API PRINCIPLES AND CHANGES TO DEPENDENCIES OR SUPPORTED VERSIONS HAPPEN VIA A ENHANCEMENT PROPOSALS SLEPS AND FOLLOWS THE DECISIONMAKING PROCESS OUTLINED ABOVE
- CHANGES TO THE GOVERNANCE MODEL USE THE SAME DECISION PROCESS OUTLINED ABOVE

IF A VETO 1 VOTE IS CAST ON A LAZY CONSENSUS THE PROPOSER CAN APPEAL TO THE COMMUNITY AND CORE DEVELOPERS AND THE CHANGE CAN BE APPROVED OR REJECTED USING THE DECISION MAKING PROCEDURE OUTLINED ABOVE

1193 ENHANCEMENT PROPOSALS SLEPS

FOR ALL VOTES A PROPOSAL MUST HAVE BEEN MADE PUBLIC AND DISCUSSED BEFORE THE VOTE SUCH PROPOSAL MUST BE A CONSOLIDATED DOCUMENT IN THE FORM OF A ‘SCIKITLEARN ENHANCEMENT PROPOSAL’ SLEP RATHER THAN A LONG DISCUSSION ON AN ISSUE A SLEP MUST BE SUBMITTED AS A PULLREQUEST TO ENHANCEMENT PROPOSALS USING THE SLEP TEMPLATE

CHAPTER  
TWO  
SCIKITLEARN TUTORIALS  
21 AN INTRODUCTION TO MACHINE LEARNING WITH SCIKITLEARN  
SECTION CONTENTS  
IN THIS SECTION WE INTRODUCE THE MACHINE LEARNING VOCABULARY THAT WE USE THROUGHOUT SCIKITLEARN AND GIVE A SIMPLE LEARNING EXAMPLE  
21.1 MACHINE LEARNING THE PROBLEM SETTING  
IN GENERAL A LEARNING PROBLEM CONSIDERS A SET OF N SAMPLES OF DATA AND THEN TRIES TO PREDICT PROPERTIES OF UNKNOWN DATA IF EACH SAMPLE IS MORE THAN A SINGLE NUMBER AND FOR INSTANCE A MULTIDIMENSIONAL ENTRY AKA MULTIVARIATE DATA IT IS SAID TO HAVE SEVERAL ATTRIBUTES OR FEATURES  
LEARNING PROBLEMS FALL INTO A FEW CATEGORIES  
• SUPERVISED LEARNING IN WHICH THE DATA COMES WITH ADDITIONAL ATTRIBUTES THAT WE WANT TO PREDICT [CLICK HERE TO GO TO THE SCIKITLEARN SUPERVISED LEARNING PAGE](#) THIS PROBLEM CAN BE EITHER  
-CLASSIFICATION SAMPLES BELONG TO TWO OR MORE CLASSES AND WE WANT TO LEARN FROM ALREADY LABELED DATA HOW TO PREDICT THE CLASS OF UNLABELED DATA AN EXAMPLE OF A CLASSIFICATION PROBLEM WOULD BE HANDWRITTEN DIGIT RECOGNITION IN WHICH THE AIM IS TO ASSIGN EACH INPUT VECTOR TO ONE OF A FINITE NUMBER OF DISCRETE CATEGORIES ANOTHER WAY TO THINK OF CLASSIFICATION IS AS A DISCRETE AS OPPOSED TO CONTINUOUS FORM OF SUPERVISED LEARNING WHERE ONE HAS A LIMITED NUMBER OF CATEGORIES AND FOR EACH OF THE N SAMPLES PROVIDED ONE IS TO TRY TO LABEL THEM WITH THE CORRECT CATEGORY OR CLASS  
-REGRESSION IF THE DESIRED OUTPUT CONSISTS OF ONE OR MORE CONTINUOUS VARIABLES THEN THE TASK IS CALLED REGRESSION AN EXAMPLE OF A REGRESSION PROBLEM WOULD BE THE PREDICTION OF THE LENGTH OF A SALMON AS A FUNCTION OF ITS AGE AND WEIGHT  
• UNSUPERVISED LEARNING IN WHICH THE TRAINING DATA CONSISTS OF A SET OF INPUT VECTORS X WITHOUT ANY CORRESPONDING TARGET VALUES THE GOAL IN SUCH PROBLEMS MAY BE TO DISCOVER GROUPS OF SIMILAR EXAMPLES WITHIN THE DATA WHERE IT IS CALLED CLUSTERING OR TO DETERMINE THE DISTRIBUTION OF DATA WITHIN THE INPUT SPACE KNOWN AS DENSITY ESTIMATION OR TO PROJECT THE DATA FROM A HIGHDIMENSIONAL SPACE DOWN TO TWO OR THREE DIMENSIONS FOR THE PURPOSE OF VISUALIZATION [CLICK HERE TO GO TO THE SCIKITLEARN UNSUPERVISED LEARNING PAGE](#)

SCIKITLEARN USER GUIDE RELEASE 0213

TRAINING SET AND TESTING SET

MACHINE LEARNING IS ABOUT LEARNING SOME PROPERTIES OF A DATA SET AND THEN TESTING THOSE PROPERTIES AGAINST ANOTHER DATA SET A COMMON PRACTICE IN MACHINE LEARNING IS TO EVALUATE AN ALGORITHM BY SPLITTING A DATA SET INTO TWO WE CALL ONE OF THOSE SETS THE TRAINING SET ON WHICH WE LEARN SOME PROPERTIES WE CALL THE OTHER SET THE TESTING SET ON WHICH WE TEST THE LEARNED PROPERTIES

212 LOADING AN EXAMPLE DATASET

SCIKITLEARN COMES WITH A FEW STANDARD DATASETS FOR INSTANCE THE IRIS AND DIGITS DATASETS FOR CLASSIFICATION AND THE BOSTON HOUSE PRICES DATASET FOR REGRESSION

IN THE FOLLOWING WE START A PYTHON INTERPRETER FROM OUR SHELL AND THEN LOAD THE IRIS AND DIGITS DATASETS OUR NOTATIONAL CONVENTION IS THAT `python` DENOTES THE SHELL PROMPT WHILE `>` DENOTES THE PYTHON INTERPRETER PROMPT

```
python
from sklearn import datasets
iris = datasets.load_iris()
digits = datasets.load_digits()
A DATASET IS A DICTIONARYLIKE OBJECT THAT HOLDS ALL THE DATA AND SOME METADATA ABOUT THE DATA THIS DATA IS STORED IN THE DATA MEMBER WHICH IS A (NSAMPLES, NFEATURES) ARRAY IN THE CASE OF SUPERVISED PROBLEM ONE OR MORE RESPONSE VARIABLES ARE STORED IN THE TARGET MEMBER MORE DETAILS ON THE DIFFERENT DATASETS CAN BE FOUND IN THE DEDICATED SECTION
FOR INSTANCE IN THE CASE OF THE DIGITS DATASET digits.data GIVES ACCESS TO THE FEATURES THAT CAN BE USED TO CLASSIFY THE DIGITS SAMPLES
print(digits.data)
0 0 5 0 0 0
0 0 0 10 0 0
0 0 0 16 9 0

0 0 1 6 0 0
0 0 2 12 0 0
0 0 10 12 1 0
AND DIGIT TARGET GIVES THE GROUND TRUTH FOR THE DIGIT DATASET THAT IS THE NUMBER CORRESPONDING TO EACH DIGIT IMAGE THAT WE ARE TRYING TO LEARN
digit.target
array([0, 1, 2, 8, 9, 8])
SHAPE OF THE DATA ARRAYS
THE DATA IS ALWAYS A 2D ARRAY SHAPE (NSAMPLES, NFEATURES) ALTHOUGH THE ORIGINAL DATA MAY HAVE HAD A DIFFERENT SHAPE IN THE CASE OF THE DIGITS EACH ORIGINAL SAMPLE IS AN IMAGE OF SHAPE (8, 8) AND CAN BE ACCESSED USING
digits.images[0]
array([[0, 0, 5, 13, 9, 1, 0, 0],
       [0, 0, 13, 15, 10, 15, 5, 0],
       [0, 3, 15, 2, 0, 11, 8, 0],
       [0, 4, 12, 0, 0, 8, 8, 0],
       [0, 5, 8, 0, 0, 9, 8, 0],
       [0, 4, 11, 0, 1, 12, 7, 0],
       [0, 2, 14, 5, 10, 12, 0, 0],
       [0, 0, 6, 13, 10, 0, 0, 0]])
```

152 CHAPTER 2 SCIKITLEARN TUTORIALS

SCIKITLEARN USER GUIDE RELEASE 0213

THE SIMPLE EXAMPLE ON THIS DATASET ILLUSTRATES HOW STARTING FROM THE ORIGINAL PROBLEM ONE CAN SHAPE THE DATA FOR CONSUMPTION IN SCIKITLEARN

LOADING FROM EXTERNAL DATASETS

TO LOAD FROM AN EXTERNAL DATASET PLEASE REFER TO LOADING EXTERNAL DATASETS

213 LEARNING AND PREDICTING

IN THE CASE OF THE DIGITS DATASET THE TASK IS TO PREDICT GIVEN AN IMAGE WHICH DIGIT IT REPRESENTS WE ARE GIVEN SAMPLES OF EACH OF THE 10 POSSIBLE CLASSES THE DIGITS ZERO THROUGH NINE ON WHICH WE FIT AN ESTIMATOR TO BE ABLE TO PREDICT THE CLASSES TO WHICH UNSEEN SAMPLES BELONG

IN SCIKITLEARN AN ESTIMATOR FOR CLASSIFICATION IS A PYTHON OBJECT THAT IMPLEMENTS THE METHODS FIT\_X\_Y AND PREDICT

AN EXAMPLE OF AN ESTIMATOR IS THE CLASS SKLEARN.SVM.SVC WHICH IMPLEMENTS SUPPORT VECTOR CLASSIFICATION THE ESTIMATOR'S CONSTRUCTOR TAKES AS ARGUMENTS THE MODEL'S PARAMETERS

FOR NOW WE WILL CONSIDER THE ESTIMATOR AS A BLACK BOX

```
from sklearn import svm
clf = svm.SVC(gamma=0.001, C=100)
```

CHOOSING THE PARAMETERS OF THE MODEL

IN THIS EXAMPLE WE SET THE VALUE OF GAMMA MANUALLY TO FIND GOOD VALUES FOR THESE PARAMETERS WE CAN USE TOOLS SUCH AS GRID SEARCH AND CROSS VALIDATION

THE CLF FOR CLASSIFIER ESTIMATOR INSTANCE IS FIRST FITTED TO THE MODEL THAT IS IT MUST LEARN FROM THE MODEL THIS IS DONE BY PASSING OUR TRAINING SET TO THE FIT METHOD FOR THE TRAINING SET WE'LL USE ALL THE IMAGES FROM OUR DATASET EXCEPT FOR THE LAST IMAGE WHICH WE'LL RESERVE FOR OUR PREDICTING WE SELECT THE TRAINING SET WITH THE 1 PYTHON SYNTAX WHICH PRODUCES A NEW ARRAY THAT CONTAINS ALL BUT THE LAST ITEM FROM DIGITS\_DATA

```
clf.fit(digits_data[:-1], digit_target[:-1])
svcc1000, cachesize200, class_weight=None, coef0=0.0,
decision_function_shape='ovr', degree=3, gamma=0.001, kernel='rbf',
max_iter=1, probability=False, random_state=None, shrinking=True,
tol=0.001, verbose=False
```

NOW YOU CAN PREDICT NEW VALUES IN THIS CASE YOU'LL PREDICT USING THE LAST IMAGE FROM DIGITS\_DATA BY PREDICTING YOU'LL DETERMINE THE IMAGE FROM THE TRAINING SET THAT BEST MATCHES THE LAST IMAGE

```
clf.predict(digits_data[-1])
array(8)
```

21 AN INTRODUCTION TO MACHINE LEARNING WITH SCIKITLEARN 153

SCIKITLEARN USER GUIDE RELEASE 0213

THE CORRESPONDING IMAGE IS

AS YOU CAN SEE IT IS A CHALLENGING TASK AFTER ALL THE IMAGES

ARE OF POOR RESOLUTION DO YOU AGREE WITH THE CLASSIFIER

A COMPLETE EXAMPLE OF THIS CLASSIFICATION PROBLEM IS AVAILABLE AS AN EXAMPLE THAT YOU CAN RUN AND STUDY RECOGNIZING

HANDWRITTEN DIGITS

214 MODEL PERSISTENCE

IT IS POSSIBLE TO SAVE A MODEL IN SCIKITLEARN BY USING PYTHON’S BUILTIN PERSISTENCE MODEL PICKLE

```
FROM SKLEARN IMPORT SVM
FROM SKLEARN IMPORT DATASETS
CLF  SVMSVCGAMMASCALE
IRIS  DATASETSLOADIRIS
X Y  IRISDATA IRISTARGET
CLFFITX Y
SVCC10 CACHESIZE200 CLASSWEIGHTNONE COEF000
DECISIONFUNCTIONSHAPEOVR DEGREE3 GAMMASCALE KERNELRBF
MAXITER1 PROBABILITYFALSE RANDOMSTATENONE SHRINKINGTRUE
TOL0001 VERBOSEFALSE
IMPORT PICKLE
S  PICKLEDUMPSCLF
CLF2  PICKLELOADSS
CLF2PREDICTX01
ARRAY0
Y0
0
```

IN THE SPECIFIC CASE OF SCIKITLEARN IT MAY BE MORE INTERESTING TO USE JOBLIB’S REPLACEMENT FOR PICKLE JOBLIBDUMP

JOBLIBLOAD WHICH IS MORE EFFICIENT ON BIG DATA BUT IT CAN ONLY PICKLE TO THE DISK AND NOT TO A STRING

```
FROM JOBLIB IMPORT DUMP LOAD
DUMPCLF FILENAMEJOBLIB
```

LATER YOU CAN RELOAD THE PICKLED MODEL POSSIBLY IN ANOTHER PYTHON PROCESS WITH

```
CLF  LOADFILENAMEJOBLIB
```

NOTEJOBLIBDUMP ANDJOBLIBLOAD FUNCTIONS ALSO ACCEPT FILELIKE OBJECT INSTEAD OF FILENAMES MORE INFOR

MATION ON DATA PERSISTENCE WITH JOBLIB IS AVAILABLE HERE

NOTE THAT PICKLE HAS SOME SECURITY AND MAINTAINABILITY ISSUES PLEASE REFER TO SECTION MODEL PERSISTENCE FOR MORE

DETAILED INFORMATION ABOUT MODEL PERSISTENCE WITH SCIKITLEARN

154 CHAPTER 2 SCIKITLEARN TUTORIALS



SCIKITLEARN USER GUIDE RELEASE 0213

215 CONVENTIONS

SCIKITLEARN ESTIMATORS FOLLOW CERTAIN RULES TO MAKE THEIR BEHAVIOR MORE PREDICTIVE THESE ARE DESCRIBED IN MORE DETAIL IN THE GLOSSARY OF COMMON TERMS AND API ELEMENTS

TYPE CASTING

UNLESS OTHERWISE SPECIFIED INPUT WILL BE CAST TO FLOAT64

```
import numpy as np
from sklearn import randomprojection
rng = np.random.RandomState(0)
X = rng.randn(10, 2000)
X = X.astype(np.float32)
dtype = np.float32
transformer = randomprojection.GaussianRandomProjection()
X_new = transformer.fit_transform(X)
X_new.dtype = np.float64
```

IN THIS EXAMPLE X IS FLOAT32 WHICH IS CAST TO FLOAT64 BY FITTRANSFORMX

REGRESSION TARGETS ARE CAST TO FLOAT64 AND CLASSIFICATION TARGETS ARE MAINTAINED

```
from sklearn import datasets
from sklearn.svm import SVC
iris = datasets.load_iris()
clf = SVC(gamma=scale,
          decision_function_shape='ovr',
          degree=3,
          gamma=scale,
          kernel='rbf',
          max_iter=1,
          probability=False,
          random_state=None,
          shrinking=True,
          tol=0.001,
          verbose=False)
list(clf.predict(iris.data[3]))
0 0 0
clf.fit(iris.data, iris.target_names[iris.target])
svcc10 = cachesize=200, class_weight=None, coef=0.0,
decision_function_shape='ovr', degree=3, gamma=scale, kernel='rbf',
max_iter=1, probability=False, random_state=None, shrinking=True,
tol=0.001, verbose=False)
list(clf.predict(iris.data[3]))
setosa setosa setosa
```

HERE THE FIRST PREDICT RETURNS AN INTEGER ARRAY SINCE IRISTARGET AN INTEGER ARRAY WAS USED IN FIT THE SECOND PREDICT RETURNS A STRING ARRAY SINCE IRISTARGETNAMES WAS FOR FITTING

REFITTING AND UPDATING PARAMETERS

HYPERPARAMETERS OF AN ESTIMATOR CAN BE UPDATED AFTER IT HAS BEEN CONSTRUCTED VIA THE SETPARAMS METHOD CALLING FIT MORE THAN ONCE WILL OVERWRITE WHAT WAS LEARNED BY ANY PREVIOUS FIT

21 AN INTRODUCTION TO MACHINE LEARNING WITH SCIKITLEARN 155

SCIKITLEARN USER GUIDE RELEASE 0213

```
import numpy as np
from sklearn.datasets import load_iris
from sklearn.svm import SVC
X, y = load_iris(return_X_y=True)
clf = SVC(
    kernel='linear',
    svd_solver='eigh',
    cache_size=200,
    class_weight=None,
    coef0=0,
    decision_function_shape='ovr',
    degree=3,
    gamma='auto_deprecated',
    kernel='linear',
    max_iter=1,
    probability=False,
    random_state=None,
    shrinking=True,
    tol=0.001,
    verbose=False)
clf.predict(X)
array([0, 0, 0, 0])
```

CLFSETPARAMSKERNELRBF GAMMASCALEFITX Y  
SVCC10 CACHESIZE200 CLASSWEIGHTNONE COEF000  
DECISIONFUNCTIONSHAPEOVR DEGREE3 GAMMAAUTODEPRECATED  
KERNELLINEAR MAXITER1 PROBABILITYFALSE RANDOMSTATENONE  
SHRINKINGTRUE TOL0001 VERBOSEFALSE  
CLFPREDICTX5  
ARRAY0 0 0 0 0

CLFSETPARAMSKERNELRBF GAMMASCALEFITX Y  
SVCC10 CACHESIZE200 CLASSWEIGHTNONE COEF000  
DECISIONFUNCTIONSHAPEOVR DEGREE3 GAMMASCALE KERNELRBF  
MAXITER1 PROBABILITYFALSE RANDOMSTATENONE SHRINKINGTRUE  
TOL0001 VERBOSEFALSE  
CLFPREDICTX5  
ARRAY0 0 0 0 0

HERE THE DEFAULT KERNEL RBF IS FIRST CHANGED TO LINEAR VIA SVC SETPARAMS AFTER THE ESTIMATOR HAS BEEN  
CONSTRUCTED AND CHANGED BACK TO RBF TO REFIT THE ESTIMATOR AND TO MAKE A SECOND PREDICTION  
MULTICLASS VS MULTILABEL FITTING  
WHEN USING MULTICLASS CLASSIFIERS THE LEARNING AND PREDICTION TASK THAT IS PERFORMED IS DEPENDENT ON THE  
FORMAT OF THE TARGET DATA FIT UPON

```
from sklearn.svm import SVC
from sklearn.multiclass import OneVsRestClassifier
from sklearn.preprocessing import LabelBinarizer
X = [[1, 2, 2, 4, 4, 5, 3, 2, 3, 1],
     [0, 0, 1, 1, 2]]
clf = OneVsRestClassifier(SVC(gamma='scale'))
clf.fit(X, y)
clf.predict(X)
array([0, 1, 1, 2])
```

IN THE ABOVE CASE THE CLASSIFIER IS FIT ON A 1D ARRAY OF MULTICLASS LABELS AND THE PREDICT METHOD THEREFORE PROVIDES  
CORRESPONDING MULTICLASS PREDICTIONS IT IS ALSO POSSIBLE TO FIT UPON A 2D ARRAY OF BINARY LABEL INDICATORS

```
Y = LabelBinarizer(fit_transform=1)
Y.fit(X)
Y.predict(X)
array([[1, 0, 0],
       [0, 1, 0],
       [0, 0, 0],
       [0, 0, 0]])
```

HERE THE CLASSIFIER IS FIT ON A 2D BINARY LABEL REPRESENTATION OF Y USING THE LABELBINARIZER IN THIS CASE  
PREDICT RETURNS A 2D ARRAY REPRESENTING THE CORRESPONDING MULTILABEL PREDICTIONS

156 CHAPTER 2 SCIKITLEARN TUTORIALS

SCIKITLEARN USER GUIDE RELEASE 0213

NOTE THAT THE FOURTH AND FIFTH INSTANCES RETURNED ALL ZEROES INDICATING THAT THEY MATCHED NONE OF THE THREE LABELS FIT UPON WITH MULTILABEL OUTPUTS IT IS SIMILARLY POSSIBLE FOR AN INSTANCE TO BE ASSIGNED MULTIPLE LABELS

```
FROM SKLEARNPREPROCESSING IMPORT MULTILABELBINARIZER
Y 0 1 0 2 1 3 0 2 3 2 4
Y MULTILABELBINARIZERFITTRANSFORMY
CLASSIFFITX YPREDICTX
ARRAY1 1 0 0 0
1 0 1 0 0
0 1 0 1 0
1 0 1 0 0
1 0 1 0 0
```

IN THIS CASE THE CLASSIFIER IS FIT UPON INSTANCES EACH ASSIGNED MULTIPLE LABELS THE MULTILABELBINARIZER IS USED TO BINARIZE THE 2D ARRAY OF MULTILABELS TO FIT UPON AS A RESULT PREDICT RETURNS A 2D ARRAY WITH MULTIPLE PREDICTED LABELS FOR EACH INSTANCE

22 A TUTORIAL ON STATISTICALLEARNING FOR SCIENTIFIC DATA PROCESSING

STATISTICAL LEARNING

MACHINE LEARNING IS A TECHNIQUE WITH A GROWING IMPORTANCE AS THE SIZE OF THE DATASETS EXPERIMENTAL SCIENCES ARE FAC ING IS RAPIDLY GROWING PROBLEMS IT TACKLES RANGE FROM BUILDING A PREDICTION FUNCTION LINKING DIFFERENT OBSERVATIONS TO CLASSIFYING OBSERVATIONS OR LEARNING THE STRUCTURE IN AN UNLABELED DATASET

THIS TUTORIAL WILL EXPLORE STATISTICAL LEARNING THE USE OF MACHINE LEARNING TECHNIQUES WITH THE GOAL OF STATISTICAL INFERENCE DRAWING CONCLUSIONS ON THE DATA AT HAND

SCIKITLEARN IS A PYTHON MODULE INTEGRATING CLASSIC MACHINE LEARNING ALGORITHMS IN THE TIGHTLYKNIT WORLD OF SCIENTIFIC PYTHON PACKAGES NUMPY SCIPY MATPLOTLIB

221 STATISTICAL LEARNING THE SETTING AND THE ESTIMATOR OBJECT IN SCIKITLEARN

DATASETS

SCIKITLEARN DEALS WITH LEARNING INFORMATION FROM ONE OR MORE DATASETS THAT ARE REPRESENTED AS 2D ARRAYS THEY CAN BE UNDERSTOOD AS A LIST OF MULTIDIMENSIONAL OBSERVATIONS WE SAY THAT THE FIRST AXIS OF THESE ARRAYS IS THE SAMPLES AXIS WHILE THE SECOND IS THE FEATURES AXIS

A SIMPLE EXAMPLE SHIPPED WITH SCIKITLEARN IRIS DATASET

```
FROM SKLEARN IMPORT DATASETS
IRIS DATASETSLOADIRIS
DATA IRISDATA
DATASHAPE
150 4
```

IT IS MADE OF 150 OBSERVATIONS OF IRISES EACH DESCRIBED BY 4 FEATURES THEIR SEPAL AND PETAL LENGTH AND WIDTH AS DETAILED INIRISDESCR

22 A TUTORIAL ON STATISTICALLEARNING FOR SCIENTIFIC DATA PROCESSING 157

SCIKITLEARN USER GUIDE RELEASE 0213

WHEN THE DATA IS NOT INITIALLY IN THE NSAMPLES NFEATURES SHAPE IT NEEDS TO BE PREPROCESSED IN ORDER TO BE USED BY SCIKITLEARN

AN EXAMPLE OF RESHAPING DATA WOULD BE THE DIGITS DATASET

THE DIGITS DATASET IS MADE OF 1797 8X8 IMAGES OF HANDWRITTEN DIGITS

DIGITS DATASETSLOADDIGITS

DIGITSIMAGESSHAPE

1797 8 8

IMPORT MATPLOTLIBPYPLOT AS PLT

PLTIMSHOWDIGITSIMAGES1 CMAPPLTCMGRAYR

MATPLOTLIBIMAGEAXESIMAGE OBJECT AT

TO USE THIS DATASET WITH SCIKITLEARN WE TRANSFORM EACH 8X8 IMAGE INTO A FEATURE VECTOR OF LENGTH 64

DATA DIGITSIMAGESRESHAPEDIGITSIMAGESSHAPE0 1

ESTIMATORS OBJECTS

FITTING DATA THE MAIN API IMPLEMENTED BY SCIKITLEARN IS THAT OF THE ESTIMATOR AN ESTIMATOR IS ANY OBJECT THAT LEARNS FROM DATA IT MAY BE A CLASSIFICATION REGRESSION OR CLUSTERING ALGORITHM OR A TRANSFORMER THAT EXTRACTSFILTERS USEFUL FEATURES FROM RAW DATA

ALL ESTIMATOR OBJECTS EXPOSE A FIT METHOD THAT TAKES A DATASET USUALLY A 2D ARRAY

ESTIMATORFITDATA

ESTIMATOR PARAMETERS ALL THE PARAMETERS OF AN ESTIMATOR CAN BE SET WHEN IT IS INSTANTIATED OR BY MODIFYING THE CORRESPONDING ATTRIBUTE

ESTIMATOR ESTIMATORPARAM11 PARAM22

ESTIMATORPARAM1

1

ESTIMATED PARAMETERS WHEN DATA IS FITTED WITH AN ESTIMATOR PARAMETERS ARE ESTIMATED FROM THE DATA AT HAND ALL THE ESTIMATED PARAMETERS ARE ATTRIBUTES OF THE ESTIMATOR OBJECT ENDING BY AN UNDERSCORE

ESTIMATORESTIMATEDPARAM

158 CHAPTER 2 SCIKITLEARN TUTORIALS

SCIKITLEARN USER GUIDE RELEASE 0213

222 SUPERVISED LEARNING PREDICTING AN OUTPUT VARIABLE FROM HIGHDIMENSIONAL OBSERVATIONS

THE PROBLEM SOLVED IN SUPERVISED LEARNING

SUPERVISED LEARNING CONSISTS IN LEARNING THE LINK BETWEEN TWO DATASETS THE OBSERVED DATA  $X$  AND AN EXTERNAL VARIABLE  $Y$  THAT WE ARE TRYING TO PREDICT USUALLY CALLED “TARGET” OR “LABELS” MOST OFTEN  $Y$  IS A 1D ARRAY OF LENGTH  $N$  SAMPLES

ALL SUPERVISED ESTIMATORS IN SCIKITLEARN IMPLEMENT A  $\text{fit}(X, Y)$  METHOD TO FIT THE MODEL AND A  $\text{predict}(X)$  METHOD THAT GIVEN UNLABELED OBSERVATIONS  $X$  RETURNS THE PREDICTED LABELS  $Y$

VOCABULARY CLASSIFICATION AND REGRESSION

IF THE PREDICTION TASK IS TO CLASSIFY THE OBSERVATIONS IN A SET OF FINITE LABELS IN OTHER WORDS TO “NAME” THE OBJECTS OBSERVED THE TASK IS SAID TO BE A CLASSIFICATION TASK ON THE OTHER HAND IF THE GOAL IS TO PREDICT A CONTINUOUS TARGET VARIABLE IT IS SAID TO BE A REGRESSION TASK

WHEN DOING CLASSIFICATION IN SCIKITLEARN  $Y$  IS A VECTOR OF INTEGERS OR STRINGS

NOTE SEE THE INTRODUCTION TO MACHINE LEARNING WITH SCIKITLEARN TUTORIAL FOR A QUICK RUNTHROUGH ON THE BASIC MACHINE LEARNING VOCABULARY USED WITHIN SCIKITLEARN

NEAREST NEIGHBOR AND THE CURSE OF DIMENSIONALITY

CLASSIFYING IRISES

22 A TUTORIAL ON STATISTICAL LEARNING FOR SCIENTIFIC DATA PROCESSING 159

SCIKITLEARN USER GUIDE RELEASE 0213

THE IRIS DATASET IS A CLASSIFICATION TASK CONSISTING IN IDENTIFYING 3 DIFFERENT TYPES OF IRISES SETOSA VERSICOLOUR AND VIRGINICA FROM THEIR PETAL AND SEPAL LENGTH AND WIDTH

```
import numpy as np
from sklearn import datasets
iris = datasets.load_iris()
irisX = iris.data
irisY = iris.target
np.unique(irisY)
array([0, 1, 2])
```

KNEAREST NEIGHBORS CLASSIFIER

THE SIMPLEST POSSIBLE CLASSIFIER IS THE NEAREST NEIGHBOR GIVEN A NEW OBSERVATION XTEST FIND IN THE TRAINING SET IE THE DATA USED TO TRAIN THE ESTIMATOR THE OBSERVATION WITH THE CLOSEST FEATURE VECTOR PLEASE SEE THE NEAREST NEIGHBORS SECTION OF THE ONLINE SCIKITLEARN DOCUMENTATION FOR MORE INFORMATION ABOUT THIS TYPE OF CLASSIFIER

TRAINING SET AND TESTING SET

WHILE EXPERIMENTING WITH ANY LEARNING ALGORITHM IT IS IMPORTANT NOT TO TEST THE PREDICTION OF AN ESTIMATOR ON THE DATA USED TO FIT THE ESTIMATOR AS THIS WOULD NOT BE EVALUATING THE PERFORMANCE OF THE ESTIMATOR ON NEW DATA THIS IS WHY DATASETS ARE OFTEN SPLIT INTO TRAIN ANDTESTDATA

160 CHAPTER 2 SCIKITLEARN TUTORIALS

SCIKITLEARN USER GUIDE RELEASE 0213

KNN K NEAREST NEIGHBORS CLASSIFICATION EXAMPLE

SPLIT IRIS DATA IN TRAIN AND TEST DATA

A RANDOM PERMUTATION TO SPLIT THE DATA RANDOMLY

NPRANDOMSEED0

INDICES NPRANDOMPERMUTATIONLENIRISX

IRISXTRAIN IRISXINDICES10

IRISYTRAIN IRISYINDICES10

IRISXTEST IRISXINDICES10

IRISYTEST IRISYINDICES10

CREATE AND FIT A NEARESTNEIGHBOR CLASSIFIER

FROM SKLEARNNEIGHBORS IMPORT KNEIGHBORSCLASSIFIER

KNN KNEIGHBORSCLASSIFIER

KNNFITIRISXTRAIN IRISYTRAIN

KNEIGHBORSCLASSIFIERALGORITHM AUTO LEAF SIZE30 METRICMINKOWSKI

METRICPARAMSNONE NJOBSNONE NNEIGHBORS5 P2

WEIGHTSUNIFORM

KNNPREDICTIRISXTEST

ARRAY1 2 1 0 0 0 2 1 2 0

IRISYTEST

ARRAY1 1 1 0 0 0 2 1 2 0

THE CURSE OF DIMENSIONALITY

FOR AN ESTIMATOR TO BE EFFECTIVE YOU NEED THE DISTANCE BETWEEN NEIGHBORING POINTS TO BE LESS THAN SOME VALUE  $\epsilon$  WHICH DEPENDS ON THE PROBLEM IN ONE DIMENSION THIS REQUIRES ON AVERAGE  $\epsilon^{-1}$  POINTS IN THE CONTEXT OF THE ABOVE  $\epsilon$ NN

EXAMPLE IF THE DATA IS DESCRIBED BY JUST ONE FEATURE WITH VALUES RANGING FROM 0 TO 1 AND WITH  $\epsilon$ TRAINING OBSERVATIONS THEN NEW DATA WILL BE NO FURTHER AWAY THAN  $1/\epsilon$  THEREFORE THE NEAREST NEIGHBOR DECISION RULE WILL BE EFFICIENT AS SOON AS  $1/\epsilon$  IS SMALL COMPARED TO THE SCALE OF BETWEENCLASS FEATURE VARIATIONS

IF THE NUMBER OF FEATURES IS  $d$  YOU NOW REQUIRE  $\epsilon^{-1/d}$  POINTS LET’S SAY THAT WE REQUIRE 10 POINTS IN ONE DIMENSION NOW  $10^d$  POINTS ARE REQUIRED IN  $d$  DIMENSIONS TO PAVE THE  $01$ SPACE AS  $d$  BECOMES LARGE THE NUMBER OF TRAINING POINTS REQUIRED FOR A GOOD ESTIMATOR GROWS EXPONENTIALLY

FOR EXAMPLE IF EACH POINT IS JUST A SINGLE NUMBER 8 BYTES THEN AN EFFECTIVE  $\epsilon$ NN ESTIMATOR IN A PALTRY  $\epsilon \sim 20$  DIMENSIONS WOULD REQUIRE MORE TRAINING DATA THAN THE CURRENT ESTIMATED SIZE OF THE ENTIRE INTERNET  $\pm 1000$  EXABYTES OR  $2^{22}$  A TUTORIAL ON STATISTICALLEARNING FOR SCIENTIFIC DATA PROCESSING 161

SCIKITLEARN USER GUIDE RELEASE 0213  
SO  
THIS IS CALLED THE CURSE OF DIMENSIONALITY AND IS A CORE PROBLEM THAT MACHINE LEARNING ADDRESSES  
LINEAR MODEL FROM REGRESSION TO SPARSITY  
DIABETES DATASET  
THE DIABETES DATASET CONSISTS OF 10 PHYSIOLOGICAL VARIABLES AGE SEX WEIGHT BLOOD PRESSURE MEASURE ON 442  
PATIENTS AND AN INDICATION OF DISEASE PROGRESSION AFTER ONE YEAR  
DIABETES DATASETSLOADDIABETES  
DIABETESXTRAIN DIABETESDATA20  
DIABETESXTEST DIABETESDATA20  
DIABETESYTRAIN DIABETESTARGET20  
DIABETESYTEST DIABETESTARGET20  
THE TASK AT HAND IS TO PREDICT DISEASE PROGRESSION FROM PHYSIOLOGICAL VARIABLES  
LINEAR REGRESSION  
LINEARREGRESSION IN ITS SIMPLEST FORM FITS A LINEAR MODEL TO THE DATA SET BY ADJUSTING A SET  
OF PARAMETERS IN ORDER TO MAKE THE SUM OF THE SQUARED RESIDUALS OF THE MODEL AS SMALL AS POSSIBLE  
LINEAR MODELS □□□□  
•□ DATA  
•□ TARGET VARIABLE  
•□ COEFFICIENTS  
•□ OBSERVATION NOISE  
FROM SKLEARN IMPORT LINEARMODEL  
REGR LINEARMODELLINEARREGRESSION  
REGRFITDIABETESXTRAIN DIABETESYTRAIN  
  
LINEARREGRESSIONCOPYXTRUE FITINTERCEPTTRUE NJOBSNONE  
NORMALIZEFALSE  
PRINTREGRCOEF  
030349955 23763931533 51053060544 32773698041 81413170937  
49281458798 10284845219 18460648906 74351961675 7609517222  
162 CHAPTER 2 SCIKITLEARN TUTORIALS



SCIKITLEARN USER GUIDE RELEASE 0213  
THE MEAN SQUARE ERROR  
NPMEANREGRPREDICTDIABETESXTEST DIABETESYTEST 2  
  
200456760268  
EXPLAINED VARIANCE SCORE 1 IS PERFECT PREDICTION  
AND 0 MEANS THAT THERE IS NO LINEAR RELATIONSHIP  
BETWEEN X AND Y  
REGRSCOREDIABETESXTEST DIABETESYTEST  
05850753022690  
SHRINKAGE  
IF THERE ARE FEW DATA POINTS PER DIMENSION NOISE IN THE OBSERVATIONS INDUCES HIGH VARIANCE  
X NPC 5 1T  
Y 5 1  
TEST NPC 0 2T  
REGR LINEARMODELLINEARREGRESSION  
IMPORT MATPLOTLIBPYPLOT AS PLT  
PLTFigure  
NPRANDOMSEED0  
FOR INRANGE6  
THISX 1 NPRANDOMNORMALSIZE2 1 X  
REGRFITTHISX Y  
PLTPLOTTEST REGRPREDICTTEST  
PLTSCATTERTHISX Y S3  
A SOLUTION IN HIGHDIMENSIONAL STATISTICAL LEARNING IS TO SHRINK THE REGRESSION COEFFICIENTS TO ZERO ANY  
TWO RANDOMLY CHOSEN SET OF OBSERVATIONS ARE LIKELY TO BE UNCORRELATED THIS IS CALLED RIDGE REGRESSION  
22 A TUTORIAL ON STATISTICALLEARNING FOR SCIENTIFIC DATA PROCESSING 163

SCIKITLEARN USER GUIDE RELEASE 0213  
REGR LINEARMODELRIDGEALPHA1  
PLTFigure  
NPRandomSeed0  
FOR INRange6  
THISX 1 NPRandomNormalSize2 1 X  
REGRFITTHISX Y  
PLTPLOTTEST REGRPREDICTTEST  
PLTSCATTERTHISX Y S3  
THIS IS AN EXAMPLE OF BIASVARIANCE TRADEOFF THE LARGER THE RIDGE ALPHA PARAMETER THE HIGHER THE BIAS AND THE LOWER THE VARIANCE  
WE CAN CHOOSE ALPHA TO MINIMIZE LEFT OUT ERROR THIS TIME USING THE DIABETES DATASET RATHER THAN OUR SYNTHETIC DATA  
ALPHAS NPLOGSPACE4 1 6  
PRINTREGRSETPARAMSALPHAALPHA  
FITDIABETESXTRAIN DIABETESYTRAIN  
SCOREDIABETESXTEST DIABETESYTEST  
FOR ALPHAINALPHAS  
  
05851110683883 05852073015444 05854677540698  
05855512036503 05830717085554 057058999437  
NOTE CAPTURING IN THE FITTED PARAMETERS NOISE THAT PREVENTS THE MODEL TO GENERALIZE TO NEW DATA IS CALLED OVERFITTING  
THE BIAS INTRODUCED BY THE RIDGE REGRESSION IS CALLED A REGULARIZATION  
SPARSITY  
FITTING ONLY FEATURES 1 AND 2  
164 CHAPTER 2 SCIKITLEARN TUTORIALS

SCIKITLEARN USER GUIDE RELEASE 0213

NOTE A REPRESENTATION OF THE FULL DIABETES DATASET WOULD INVOLVE 11 DIMENSIONS 10 FEATURE DIMENSIONS AND ONE OF THE TARGET VARIABLE IT IS HARD TO DEVELOP AN INTUITION ON SUCH REPRESENTATION BUT IT MAY BE USEFUL TO KEEP IN MIND THAT IT WOULD BE A FAIRLY EMPTY SPACE

WE CAN SEE THAT ALTHOUGH FEATURE 2 HAS A STRONG COEFFICIENT ON THE FULL MODEL IT CONVEYS LITTLE INFORMATION ON YWHEN CONSIDERED WITH FEATURE 1

TO IMPROVE THE CONDITIONING OF THE PROBLEM IE MITIGATING THE THE CURSE OF DIMENSIONALITY IT WOULD BE INTERESTING TO SELECT ONLY THE INFORMATIVE FEATURES AND SET NONINFORMATIVE ONES LIKE FEATURE 2 TO 0 RIDGE REGRESSION WILL DECREASE THEIR CONTRIBUTION BUT NOT SET THEM TO ZERO ANOTHER PENALIZATION APPROACH CALLED LASSO LEAST ABSOLUTE SHRINKAGE AND SELECTION OPERATOR CAN SET SOME COEFFICIENTS TO ZERO SUCH METHODS ARE CALLED SPARSE METHOD AND SPARSITY CAN BE SEEN AS AN APPLICATION OF OCCAM’S RAZOR PREFER SIMPLER MODELS

REGR LINEARMODELLASSO

SCORES REGRSETPARAMSALPHAALPHA

FITDIABETESXTRAIN DIABETESYTRAIN

SCOREDIABETESXTEST DIABETESYTEST

FOR ALPHAINALPHAS

BESTALPHA ALPHASSCORESINDEXMAXSCORES

REGRALPHA BESTALPHA

REGRFITDIABETESXTRAIN DIABETESYTRAIN

LASSOALPHA0025118864315095794 COPYXTRUE FITINTERCEPTTRUE

MAXITER1000 NORMALIZEFALSE POSITIVEFALSE PRECOMPUTEFALSE

RANDOMSTATENONE SELECTIONCYCLIC TOL00001 WARMSTARTFALSE

PRINTREGRCOEF

0 21243764548 51719478111 31377959962 1608303982 0

18719554705 6938229038 50866011217 7184239008

22 A TUTORIAL ON STATISTICALLEARNING FOR SCIENTIFIC DATA PROCESSING 165

SCIKITLEARN USER GUIDE RELEASE 0213

DIFFERENT ALGORITHMS FOR THE SAME PROBLEM

DIFFERENT ALGORITHMS CAN BE USED TO SOLVE THE SAME MATHEMATICAL PROBLEM FOR INSTANCE THE LASSO OBJECT IN SCIKIT LEARN SOLVES THE LASSO REGRESSION PROBLEM USING A COORDINATE DESCENT METHOD THAT IS EFFICIENT ON LARGE DATASETS HOWEVER SCIKITLEARN ALSO PROVIDES THE LASSOLARS OBJECT USING THE LARS ALGORITHM WHICH IS VERY EFFICIENT FOR PROBLEMS IN WHICH THE WEIGHT VECTOR ESTIMATED IS VERY SPARSE IE PROBLEMS WITH VERY FEW OBSERVATIONS

CLASSIFICATION

FOR CLASSIFICATION AS IN THE LABELING IRIS TASK LINEAR REGRESSION IS NOT THE RIGHT APPROACH AS IT WILL GIVE TOO MUCH WEIGHT TO DATA FAR FROM THE DECISION FRONTIER A LINEAR APPROACH IS TO FIT A SIGMOID FUNCTION OR LOGISTIC FUNCTION

LOGITSIGMOID --OFFSET 1

1 EXP--OFFSET

LOG LINEARMODELLOGISTICREGRESSIONSOLVERLBFGS C1E5

MULTICLASSMULTINOMIAL

LOGFITIRISXTRAIN IRISYTRAIN

LOGISTICREGRESSIONC1000000 CLASSWEIGHTNONE DUALFALSE

FITINTERCEPTTRUE INTERCEPTSCALING1 L1RATIONONE MAXITER100

MULTICLASSMULTINOMIAL NJOBSNONE PENALTYL2 RANDOMSTATENONE

SOLVERLBFGS TOL00001 VERBOSE0 WARMSTARTFALSE

THIS IS KNOWN AS LOGISTICREGRESSION

166 CHAPTER 2 SCIKITLEARN TUTORIALS

SCIKITLEARN USER GUIDE RELEASE 0213

MULTICLASS CLASSIFICATION

IF YOU HAVE SEVERAL CLASSES TO PREDICT AN OPTION OFTEN USED IS TO FIT ONEVERSUSALL CLASSIFIERS AND THEN USE A VOTING HEURISTIC FOR THE FINAL DECISION

SHRINKAGE AND SPARSITY WITH LOGISTIC REGRESSION

THECPARAMETER CONTROLS THE AMOUNT OF REGULARIZATION IN THE LOGISTICREGRESSION OBJECT A LARGE VALUE FORCRESULTS IN LESS REGULARIZATION PENALTYL2 GIVES SHRINKAGE IE NONSPARSE COEFFICIENTS WHILE PENALTYL1 GIVES SPARSITY

EXERCISE

TRY CLASSIFYING THE DIGITS DATASET WITH NEAREST NEIGHBORS AND A LINEAR MODEL LEAVE OUT THE LAST 10 AND TEST PREDICTION PERFORMANCE ON THESE OBSERVATIONS

FROM SKLEARN IMPORT DATASETS NEIGHBORS LINEARMODEL

DIGITS DATASETSLOADDIGITS

XDIGITS DIGITSDATA DIGITSDATAMAX

YDIGITS DIGITSTARGET

SOLUTIONAUTOEXAMPLESEXERCISESPLOTDIGITSClassificationEXERCISEPY

SUPPORT VECTOR MACHINES SVMs

LINEAR SVMs

SUPPORT VECTOR MACHINES BELONG TO THE DISCRIMINANT MODEL FAMILY THEY TRY TO FIND A COMBINATION OF SAMPLES TO BUILD A PLANE MAXIMIZING THE MARGIN BETWEEN THE TWO CLASSES REGULARIZATION IS SET BY THE CPARAMETER A SMALL VALUE FOR C MEANS THE MARGIN IS CALCULATED USING MANY OR ALL OF THE OBSERVATIONS AROUND THE SEPARATING LINE MORE REGULARIZATION A LARGE VALUE FOR CMEANS THE MARGIN IS CALCULATED ON OBSERVATIONS CLOSE TO THE SEPARATING LINE LESS REGULARIZATION UNREGULARIZED SVM REGULARIZED SVM DEFAULT

22 A TUTORIAL ON STATISTICALLEARNING FOR SCIENTIFIC DATA PROCESSING 167

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLE

- PLOT DIFFERENT SVM CLASSIFIERS IN THE IRIS DATASET

SVMS CAN BE USED IN REGRESSION – SVR SUPPORT VECTOR REGRESSION- OR IN CLASSIFICATION – SVC SUPPORT VECTOR CLASSIFICATION

```
FROM SKLEARN IMPORT SVM
SVC SVMSVCKERNELLINEAR
SVCFITIRISXTRAIN IRISYTRAIN
SVCC10 CACHESIZE200 CLASSWEIGHTNONE COEF000
DECISIONFUNCTIONSHAPEOVR DEGREE3 GAMMAAUTODEPRECATED
KERNELLINEAR MAXITER1 PROBABILITYFALSE RANDOMSTATENONE
SHRINKINGTRUE TOL0001 VERBOSEFALSE
WARNING NORMALIZING DATA
```

FOR MANY ESTIMATORS INCLUDING THE SVMS HAVING DATASETS WITH UNIT STANDARD DEVIATION FOR EACH FEATURE IS IMPORTANT TO GET GOOD PREDICTION

USING KERNELS

CLASSES ARE NOT ALWAYS LINEARLY SEPARABLE IN FEATURE SPACE THE SOLUTION IS TO BUILD A DECISION FUNCTION THAT IS NOT LINEAR BUT MAY BE POLYNOMIAL INSTEAD THIS IS DONE USING THE KERNEL TRICK THAT CAN BE SEEN AS CREATING A DECISION ENERGY BY POSITIONING KERNELS ON OBSERVATIONS

LINEAR KERNEL POLYNOMIAL KERNEL

```
SVC SVMSVCKERNELLINEAR SVC SVMSVCKERNELPOLY
DEGREE3
DEGREE POLYNOMIAL DEGREE
```

168 CHAPTER 2 SCIKITLEARN TUTORIALS

SCIKITLEARN USER GUIDE RELEASE 0213  
RBF KERNEL RADIAL BASIS FUNCTION  
SVC SVMSVCKERNELRBF  
GAMMA INVERSE OF SIZE OF  
RADIAL KERNEL  
INTERACTIVE EXAMPLE  
SEE THE SVM GUI TO DOWNLOAD SVMGUIPY ADD DATA POINTS OF BOTH CLASSES WITH RIGHT AND LEFT BUTTON FIT THE  
MODEL AND CHANGE PARAMETERS AND DATA  
22 A TUTORIAL ON STATISTICALLEARNING FOR SCIENTIFIC DATA PROCESSING 169

```
SCIKITLEARN USER GUIDE RELEASE 0213
EXERCISE
TRY CLASSIFYING CLASSES 1 AND 2 FROM THE IRIS DATASET WITH SVMs WITH THE 2 FIRST FEATURES LEAVE OUT 10 OF EACH
CLASS AND TEST PREDICTION PERFORMANCE ON THESE OBSERVATIONS
WARNING THE CLASSES ARE ORDERED DO NOT LEAVE OUT THE LAST 10 YOU WOULD BE TESTING ON ONLY ONE CLASS
HINT YOU CAN USE THE DECISIONFUNCTION METHOD ON A GRID TO GET INTUITIONS
IRIS DATASETSLOADIRIS
X IRISDATA
Y IRISTARGET
X XY 0 2
Y YY 0
SOLUTIONAUTOEXAMPLESEXERCISESPLOTIRISEXERCISEPY
223 MODEL SELECTION CHOOSING ESTIMATORS AND THEIR PARAMETERS
SCORE AND CROSSVALIDATED SCORES
AS WE HAVE SEEN EVERY ESTIMATOR EXPOSES A SCORE METHOD THAT CAN JUDGE THE QUALITY OF THE FIT OR THE PREDICTION ON
NEW DATA BIGGER IS BETTER
FROM SKLEARN IMPORT DATASETS SVM
DIGITS DATASETSLOADDIGITS
XDIGITS DIGITSDATA
YDIGITS DIGITSTARGET
SVC SVMsvcc1 KERNELLINEAR
SVCFITXDIGITS100 YDIGITS100SCOREXDIGITS100 YDIGITS100
098
TO GET A BETTER MEASURE OF PREDICTION ACCURACY WHICH WE CAN USE AS A PROXY FOR GOODNESS OF FIT OF THE MODEL WE CAN
SUCCESSIVELY SPLIT THE DATA IN FOLDS THAT WE USE FOR TRAINING AND TESTING
IMPORT NUMPY AS NP
XFOLDS NPARRAYSPLITXDIGITS 3
YFOLDS NPARRAYSPLITXDIGITS 3
SCORES LIST
FOR KINRANGE3
WE USE LIST TO COPY IN ORDER TO POP LATER ON
XTRAIN LISTXFOLDS
XTEST XTRAINPOP
XTRAIN NPCONCATENATEXTRAIN
YTRAIN LISTYFOLDS
YTEST YTRAINPOP
YTRAIN NPCONCATENATEYTRAIN
SCORESAPPENDSVCFITXTRAIN YTRAINSCOREXTEST YTEST
PRINTSCORES
0934 0956 0939
THIS IS CALLED A KFOLD CROSSVALIDATION
170 CHAPTER 2 SCIKITLEARN TUTORIALS
```



SCIKITLEARN USER GUIDE RELEASE 0213

CROSSVALIDATION GENERATORS

SCIKITLEARN HAS A COLLECTION OF CLASSES WHICH CAN BE USED TO GENERATE LISTS OF TRAINTEST INDICES FOR POPULAR CROSS VALIDATION STRATEGIES

THEY EXPOSE A SPLIT METHOD WHICH ACCEPTS THE INPUT DATASET TO BE SPLIT AND YIELDS THE TRAINTEST SET INDICES FOR EACH ITERATION OF THE CHOSEN CROSSVALIDATION STRATEGY

THIS EXAMPLE SHOWS AN EXAMPLE USAGE OF THE SPLIT METHOD

```
FROM SKLEARNMODELSELECTION IMPORT KFOLD CROSSVALSCORE
X A A A B B C C C C C
KFOLD KFOLDNSPLITS5
FOR TRAININDICES TESTINDICES INKFOLDSPLITX
  PRINTTRAIN S TEST S TRAININDICES TESTINDICES
TRAIN 2 3 4 5 6 7 8 9 TEST 0 1
TRAIN 0 1 4 5 6 7 8 9 TEST 2 3
TRAIN 0 1 2 3 6 7 8 9 TEST 4 5
TRAIN 0 1 2 3 4 5 8 9 TEST 6 7
TRAIN 0 1 2 3 4 5 6 7 TEST 8 9
```

THE CROSSVALIDATION CAN THEN BE PERFORMED EASILY

```
SVCFITXDIGITSTRAIN YDIGITSTRAINSCOREXDIGITSTEST YDIGITSTEST
FOR TRAIN TEST INKFOLDSPLITXDIGITS
0963 0922 0963 0963 0930
```

THE CROSSVALIDATION SCORE CAN BE DIRECTLY CALCULATED USING THE CROSSVALSCORE HELPER GIVEN AN ESTIMATOR THE CROSSVALIDATION OBJECT AND THE INPUT DATASET THE CROSSVALSCORE SPLITS THE DATA REPEATEDLY INTO A TRAINING AND A TESTING SET TRAINS THE ESTIMATOR USING THE TRAINING SET AND COMPUTES THE SCORES BASED ON THE TESTING SET FOR EACH ITERATION OF CROSSVALIDATION

BY DEFAULT THE ESTIMATOR'S SCORE METHOD IS USED TO COMPUTE THE INDIVIDUAL SCORES

REFER THE METRICS MODULE TO LEARN MORE ON THE AVAILABLE SCORING METHODS

```
CROSSVALSCORESVC XDIGITS YDIGITS CVKFOLD NJOBS1
ARRAY096388889 092222222 09637883 09637883 093036212
```

NJOBS1 MEANS THAT THE COMPUTATION WILL BE DISPATCHED ON ALL THE CPUS OF THE COMPUTER

ALTERNATIVELY THE SCORING ARGUMENT CAN BE PROVIDED TO SPECIFY AN ALTERNATIVE SCORING METHOD

```
CROSSVALSCORESVC XDIGITS YDIGITS CVKFOLD
SCORINGPRECISIONMACRO
ARRAY096578289 092708922 096681476 096362897 093192644
```

CROSSVALIDATION GENERATORS

```
KFOLD NSPLITS SHUFFLE RAN
DOMSTATESTRATIFIEDKFOLD NSPLITS
SHUFFLE RANDOMSTATEGROUPKFOLD NSPLITS
SPLITS IT INTO K FOLDS TRAINS ON K1
AND THEN TESTS ON THE LEFTOUTSAME AS KFOLD BUT PRESERVES THE
CLASS DISTRIBUTION WITHIN EACH FOLDENSURES THAT THE SAME GROUP IS NOT IN
BOTH TESTING AND TRAINING SETS
```

22 A TUTORIAL ON STATISTICALLEARNING FOR SCIENTIFIC DATA PROCESSING 171

SCIKITLEARN USER GUIDE RELEASE 0213

SHUFFLESPLIT NSPLITS

TESTSIZE TRAINSIZE RAN

DOMSTATESTRATIFIEDSHUFFLESPLIT GROUPSHUFFLESPLIT

GENERATES TRAINTEST INDICES BASED

ON RANDOM PERMUTATIONSAME AS SHUFFLE SPLIT BUT PRESERVES THE

CLASS DISTRIBUTION WITHIN EACH ITERATIONENSURES THAT THE SAME GROUP IS NOT

IN BOTH TESTING AND TRAINING SETS

LEAVEONEGROUPOUT LEAVEPGROUPSOUT NGROUPS LEAVEONEOUT

TAKES A GROUP ARRAY TO GROUP OBSERVATIONS LEAVE P GROUPS OUT LEAVE ONE OBSERVATION OUT

LEAVEPOUT P PREDEFINEDSPLIT

LEAVE P OBSERVATIONS OUT GENERATES TRAINTEST INDICES BASED ON PREDEFINED SPLITS

EXERCISE

ON THE

DIGITS DATASET PLOT THE CROSSVALIDATION SCORE OF A SVC ESTIMATOR WITH AN LINEAR KERNEL AS A FUNCTION OF PARAMETER C

USE A LOGARITHMIC GRID OF POINTS FROM 1 TO 10

IMPORT NUMPY AS NP

FROM SKLEARNMODELSELECTION IMPORT CROSSVALSCORE

FROM SKLEARN IMPORT DATASETS SVM

DIGITS DATASETSLOADDIGITS

X DIGITSDATA

Y DIGITSTARGET

SVC SVMKERNELLINER

CS NPLOGSPACE10 0 10

172 CHAPTER 2 SCIKITLEARN TUTORIALS

SCIKITLEARN USER GUIDE RELEASE 0213  
SOLUTION CROSSVALIDATION ON DIGITS DATASET EXERCISE  
GRIDSEARCH AND CROSSVALIDATED ESTIMATORS  
GRIDSEARCH  
SCIKITLEARN PROVIDES AN OBJECT THAT GIVEN DATA COMPUTES THE SCORE DURING THE FIT OF AN ESTIMATOR ON A PARAMETER GRID AND  
CHOOSES THE PARAMETERS TO MAXIMIZE THE CROSSVALIDATION SCORE THIS OBJECT TAKES AN ESTIMATOR DURING THE CONSTRUCTION  
AND EXPOSES AN ESTIMATOR API  
FROM SKLEARNMODELSELECTION IMPORT GRIDSEARCHCV CROSSVALSCORE  
CS NPLOGSPACE6 1 10  
CLF GRIDSEARCHCVESTIMATORSVC PARAMGRIDDICTCCS  
NJOBS1  
CLFFITXDIGITS1000 YDIGITS1000  
GRIDSEARCHCVCVNONE  
CLFBESTSCORE  
0925  
CLFBESTESTIMATORC  
00077  
PREDICTION PERFORMANCE ON TEST SET IS NOT AS GOOD AS ON TRAIN SET  
CLFSCOREXDIGITS1000 YDIGITS1000  
0943  
BY DEFAULT THE GRIDSEARCHCV USES A 3FOLD CROSSVALIDATION HOWEVER IF IT DETECTS THAT A CLASSIFIER IS PASSED RATHER  
THAN A REGRESSOR IT USES A STRATIFIED 3FOLD THE DEFAULT WILL CHANGE TO A 5FOLD CROSSVALIDATION IN VERSION 022  
NESTED CROSSVALIDATION  
CROSSVALSCORECLF XDIGITS YDIGITS  
ARRAY0938 0963 0944  
TWO CROSSVALIDATION LOOPS ARE PERFORMED IN PARALLEL ONE BY THE GRIDSEARCHCV ESTIMATOR TO SET GAMMA AND THE  
OTHER ONE BY CROSSVALSCORE TO MEASURE THE PREDICTION PERFORMANCE OF THE ESTIMATOR THE RESULTING SCORES  
ARE UNBIASED ESTIMATES OF THE PREDICTION SCORE ON NEW DATA  
WARNING YOU CANNOT NEST OBJECTS WITH PARALLEL COMPUTING NJOBS DIFFERENT THAN 1  
CROSSVALIDATED ESTIMATORS  
CROSSVALIDATION TO SET A PARAMETER CAN BE DONE MORE EFFICIENTLY ON AN ALGORITHMBYALGORITHM BASIS THIS IS WHY FOR  
CERTAIN ESTIMATORS SCIKITLEARN EXPOSES CROSSVALIDATION EVALUATING ESTIMATOR PERFORMANCE ESTIMATORS THAT SET THEIR  
PARAMETER AUTOMATICALLY BY CROSSVALIDATION  
FROM SKLEARN IMPORT LINEARMODEL DATASETS  
LASSO LINEARMODELLASSOCVCV3  
DIABETES DATASETSLOADDIABETES  
XDIABETES DIABETESDATA  
YDIABETES DIABETESTARGET  
22 A TUTORIAL ON STATISTICALLEARNING FOR SCIENTIFIC DATA PROCESSING 173

SCIKITLEARN USER GUIDE RELEASE 0213  
LASSOFITXDIABETES YDIABETES  
LASSOCVALPHASNONE COPYXTRUE CV3 EPS0001 FITINTERCEPTTRUE  
MAXITER1000 NALPHAS100 NJOBSNONE NORMALIZEFALSE  
POSITIVEFALSE PRECOMPUTEAUTO RANDOMSTATENONE  
SELECTIONCYCLIC TOL00001 VERBOSEFALSE  
THE ESTIMATOR CHOSE AUTOMATICALLY ITS LAMBDA  
LASSOALPHA  
001229  
THESE ESTIMATORS ARE CALLED SIMILARLY TO THEIR COUNTERPARTS WITH 'CV' APPENDED TO THEIR NAME  
EXERCISE  
ON THE DIABETES DATASET FIND THE OPTIMAL REGULARIZATION PARAMETER ALPHA  
BONUS HOW MUCH CAN YOU TRUST THE SELECTION OF ALPHA  
FROM SKLEARN IMPORT DATASETS  
FROM SKLEARNLINEARMODEL IMPORT LASSOCV  
FROM SKLEARNLINEARMODEL IMPORT LASSO  
FROM SKLEARNMODELSELECTION IMPORT KFOLD  
FROM SKLEARNMODELSELECTION IMPORT GRIDSEARCHCV  
DIABETES DATASETSLOADDIABETES  
X DIABETESDATA150  
SOLUTION CROSSVALIDATION ON DIABETES DATASET EXERCISE  
224 UNSUPERVISED LEARNING SEEKING REPRESENTATIONS OF THE DATA  
CLUSTERING GROUPING OBSERVATIONS TOGETHER  
THE PROBLEM SOLVED IN CLUSTERING  
GIVEN THE IRIS DATASET IF WE KNEW THAT THERE WERE 3 TYPES OF IRIS BUT DID NOT HAVE ACCESS TO A TAXONOMIST TO LABEL  
THEM WE COULD TRY A CLUSTERING TASK SPLIT THE OBSERVATIONS INTO WELLSEPARATED GROUP CALLED CLUSTERS  
174 CHAPTER 2 SCIKITLEARN TUTORIALS

SCIKITLEARN USER GUIDE RELEASE 0213

KMEANS CLUSTERING

NOTE THAT THERE EXIST A LOT OF DIFFERENT CLUSTERING CRITERIA AND ASSOCIATED ALGORITHMS THE SIMPLEST CLUSTERING ALGORITHM

ISKMEANS

FROM SKLEARN IMPORT CLUSTER DATASETS

IRIS DATASETSLOADIRIS

XIRIS IRISDATA

YIRIS IRISTARGET

KMEANS CLUSTERKMEANSNCLUSTERS3

KMEANSFITXIRIS

KMEANSALGORITHMMAUTO COPYXTRUE INITKMEANS

PRINTKMEANSLABELS10

1 1 1 1 1 0 0 0 0 0 2 2 2 2 2

PRINTYIRIS10

0 0 0 0 0 1 1 1 1 1 2 2 2 2 2

WARNING THERE IS ABSOLUTELY NO GUARANTEE OF RECOVERING A GROUND TRUTH FIRST CHOOSING THE RIGHT NUMBER OF CLUSTERS IS HARD SECOND THE ALGORITHM IS SENSITIVE TO INITIALIZATION AND CAN FALL INTO LOCAL MINIMA ALTHOUGH SCIKIT LEARN EMPLOYS SEVERAL TRICKS TO MITIGATE THIS ISSUE

BAD INITIALIZATION 8 CLUSTERS GROUND TRUTH

DON'T OVERINTERPRET CLUSTERING RESULTS

APPLICATION EXAMPLE VECTOR QUANTIZATION

CLUSTERING IN GENERAL AND KMEANS IN PARTICULAR CAN BE SEEN AS A WAY OF CHOOSING A SMALL NUMBER OF EXEMPLARS TO COMPRESS THE INFORMATION THE PROBLEM IS SOMETIMES KNOWN AS VECTOR QUANTIZATION FOR INSTANCE THIS CAN BE USED

22 A TUTORIAL ON STATISTICALLEARNING FOR SCIENTIFIC DATA PROCESSING 175

SCIKITLEARN USER GUIDE RELEASE 0213

TO POSTERIZE AN IMAGE

IMPORT SCIPY AS SP

TRY

FACE SPFACEGRAY TRUE

EXCEPT ATTRIBUTEERROR

FROM SCIPY IMPORT MISC

FACE MISCFACEGRAY TRUE

X FACERESHAPE1 1 WE NEED AN NSAMPLE NFEATURE ARRAY

KMEANS CLUSTERKMEANSNCLUSTERS5 NINIT1

KMEANSFITX

KMEANSALGORITHMMAUTO COPYXTRUE INITKMEANS

VALUES KMEANSCLUSTERCENTERSSQUEEZE

LABELS KMEANSLABELS

FACECOMPRESSED NPCHOOSELABELS VALUES

FACECOMPRESSED SHAPE FACESHAPE

RAW IMAGE KMEANS QUANTIZATION EQUAL BINS IMAGE HISTOGRAM

HIERARCHICAL AGGLOMERATIVE CLUSTERING WARD

AHIERARCHICAL CLUSTERING METHOD IS A TYPE OF CLUSTER ANALYSIS THAT AIMS TO BUILD A HIERARCHY OF CLUSTERS IN GENERAL THE VARIOUS APPROACHES OF THIS TECHNIQUE ARE EITHER

•AGGLOMERATIVE BOTTOMUP APPROACHES EACH OBSERVATION STARTS IN ITS OWN CLUSTER AND CLUSTERS ARE ITERATIVELY MERGED IN SUCH A WAY TO MINIMIZE A LINKAGE CRITERION THIS APPROACH IS PARTICULARLY INTERESTING WHEN THE CLUSTERS OF INTEREST ARE MADE OF ONLY A FEW OBSERVATIONS WHEN THE NUMBER OF CLUSTERS IS LARGE IT IS MUCH MORE COMPUTATIONALLY EFFICIENT THAN KMEANS

•DIVISIVE TOPDOWN APPROACHES ALL OBSERVATIONS START IN ONE CLUSTER WHICH IS ITERATIVELY SPLIT AS ONE MOVES DOWN THE HIERARCHY FOR ESTIMATING LARGE NUMBERS OF CLUSTERS THIS APPROACH IS BOTH SLOW DUE TO ALL OBSERVATIONS STARTING AS ONE CLUSTER WHICH IT SPLITS RECURSIVELY AND STATISTICALLY ILLPOSED

CONNECTIVITYCONSTRAINED CLUSTERING

WITH AGGLOMERATIVE CLUSTERING IT IS POSSIBLE TO SPECIFY WHICH SAMPLES CAN BE CLUSTERED TOGETHER BY GIVING A CONNECTIVITY GRAPH GRAPHS IN SCIKITLEARN ARE REPRESENTED BY THEIR ADJACENCY MATRIX OFTEN A SPARSE MATRIX IS USED THIS CAN BE USEFUL FOR INSTANCE TO RETRIEVE CONNECTED REGIONS SOMETIMES ALSO REFERRED TO AS CONNECTED COMPONENTS WHEN

176 CHAPTER 2 SCIKITLEARN TUTORIALS

SCIKITLEARN USER GUIDE RELEASE 0213

CLUSTERING AN IMAGE

```
FROM SCIPYNDIMAGEFILTERS IMPORT GAUSSIANSFILTER
IMPORT MATPLOTLIBPYPLOT AS PLT
IMPORT SKIMAGE
FROM SKIMAGEDATA IMPORT COINS
FROM SKIMAGETRANSFORM IMPORT RESCALE
FROM SKLEARNFEATUREEXTRACTIONIMAGE IMPORT GRIDTOGRAPH
FROM SKLEARNCLUSTER IMPORT AGGLOMERATIVECLUSTERING
THESE WERE INTRODUCED IN SKIMAGE014
IFLOOSEVERSIONSKIMAGEVERSION 014
RESCALEPARAMS ANTIALIASING FALSE MULTICHANNEL FALSE
ELSE
RESCALEPARAMS
```

GENERATE DATA

```
ORIGCOINS COINS
RESIZE IT TO 20 OF THE ORIGINAL SIZE TO SPEED UP THE PROCESSING
APPLYING A GAUSSIAN FILTER FOR SMOOTHING PRIOR TO DOWNSCALING
REDUCES ALIASING ARTIFACTS
SMOOTHENEDCOINS GAUSSIANSFILTERORIGCOINS SIGMA2
FEATURE AGGLOMERATION
```

WE HAVE SEEN THAT SPARSITY COULD BE USED TO MITIGATE THE CURSE OF DIMENSIONALITY IEAN INSUFFICIENT AMOUNT OF OBSERVATIONS COMPARED TO THE NUMBER OF FEATURES ANOTHER APPROACH IS TO MERGE TOGETHER SIMILAR FEATURES FEATURE AGGLOMERATION THIS APPROACH CAN BE IMPLEMENTED BY CLUSTERING IN THE FEATURE DIRECTION IN OTHER WORDS CLUSTERING

22 A TUTORIAL ON STATISTICALLEARNING FOR SCIENTIFIC DATA PROCESSING 177

SCIKITLEARN USER GUIDE RELEASE 0213

THE TRANSPOSED DATA

DIGITS DATASETSLOADDIGITS

IMAGES DIGITSIMAGES

X NPRESHAPEIMAGES LENIMAGES 1

CONNECTIVITY GRIDTOGRAPH IMAGES0SHAPE

AGGLO CLUSTERFEATUREAGGLOMERATIONCONNECTIVITYCONNECTIVITY

NCLUSTERS32

AGGLOFITX

FEATUREAGGLOMERATIONAFFINITYEUCLIDEAN COMPUTEFULLTREEAUTO

XREDUCED AGGLOTRANSFORMX

XAPPROX AGGLOINVERSETRANSFORMXREDUCED

IMAGESAPPROX NPRESHAPEXAPPROX IMAGESSHAPE

TRANSFORM ANDINVERSETRANSFORM METHODS

SOME ESTIMATORS EXPOSE A TRANSFORM METHOD FOR INSTANCE TO REDUCE THE DIMENSIONALITY OF THE DATASET

DECOMPOSITIONS FROM A SIGNAL TO COMPONENTS AND LOADINGS

COMPONENTS AND LOADINGS

IF X IS OUR MULTIVARIATE DATA THEN THE PROBLEM THAT WE ARE TRYING TO SOLVE IS TO REWRITE IT ON A DIFFERENT OBSERVATIONAL

BASIS WE WANT TO LEARN LOADINGS L AND A SET OF COMPONENTS C SUCH THAT  $X = L C$  DIFFERENT CRITERIA EXIST TO CHOOSE

THE COMPONENTS

PRINCIPAL COMPONENT ANALYSIS PCA

PRINCIPAL COMPONENT ANALYSIS PCA SELECTS THE SUCCESSIVE COMPONENTS THAT EXPLAIN THE MAXIMUM VARIANCE IN THE

SIGNAL

178 CHAPTER 2 SCIKITLEARN TUTORIALS



SCIKITLEARN USER GUIDE RELEASE 0213

THE POINT CLOUD SPANNED BY THE OBSERVATIONS ABOVE IS VERY FLAT IN ONE DIRECTION ONE OF THE THREE UNIVARIATE FEATURES CAN ALMOST BE EXACTLY COMPUTED USING THE OTHER TWO PCA FINDS THE DIRECTIONS IN WHICH THE DATA IS NOT FLAT WHEN USED TO TRANSFORM DATA PCA CAN REDUCE THE DIMENSIONALITY OF THE DATA BY PROJECTING ON A PRINCIPAL SUBSPACE

CREATE A SIGNAL WITH ONLY 2 USEFUL DIMENSIONS

X1 NPRANDOMNORMALSIZE100

X2 NPRANDOMNORMALSIZE100

X3 X1 X2

X NPCX1 X2 X3

FROM SKLEARN IMPORT DECOMPOSITION

PCA DECOMPOSITIONPCA

PCAFITX

PCACOPYTRUE ITERATEDPOWERAUTO NCOMPONENTSNONE RANDOMSTATENONE

SVDSOLVERAUTO TOL00 WHITENFALSE

PRINTPCAEXPLAINEDVARIANCE

218565811E00 119346747E00 843026679E32

AS WE CAN SEE ONLY THE 2 FIRST COMPONENTS ARE USEFUL

PCANCOMPONENTS 2

XREDUCED PCAFITTRANSFORMX

XREDUCEDSHAPE

100 2

INDEPENDENT COMPONENT ANALYSIS ICA

INDEPENDENT COMPONENT ANALYSIS ICA SELECTS COMPONENTS SO THAT THE DISTRIBUTION OF THEIR LOADINGS CARRIES A MAXIMUM AMOUNT OF INDEPENDENT INFORMATION IT IS ABLE TO RECOVER NONGAUSSIAN INDEPENDENT SIGNALS

22 A TUTORIAL ON STATISTICALLARNING FOR SCIENTIFIC DATA PROCESSING 179

```
SCIKITLEARN USER GUIDE RELEASE 0213
GENERATE SAMPLE DATA
IMPORT NUMPY AS NP
FROM SCIPY IMPORT SIGNAL
TIME  NPLinspace0 10 2000
S1  NPSIN2 TIME SIGNAL 1  SINUSOIDAL SIGNAL
S2  NPSIGNNPSIN3 TIME  SIGNAL 2  SQUARE SIGNAL
S3  SIGNALSAWTOOTH2 NPPITIME SIGNAL 3 SAW TOOTH SIGNAL
S  NPCS1 S2 S3
S  02 NPRANDOMNORMALSIZESSHAPE  ADD NOISE
S  SSTDAXIS0  STANDARDIZE DATA
MIX DATA
A  NPARRAY1 1 1 05 2 1 15 1 2  MIXING MATRIX
X  NPDOTS AT  GENERATE OBSERVATIONS
COMPUTE ICA
ICA  DECOMPOSITIONFASTICA
S  ICAFITTRANSFORMX  GET THE ESTIMATED SOURCES
A  ICAMIXINGT
NPALLCLOSEX NPDOTS A  ICAMEAN
TRUE
225 PUTTING IT ALL TOGETHER
180 CHAPTER 2 SCIKITLEARN TUTORIALS
```

```
SCIKITLEARN USER GUIDE RELEASE 0213
PIPELINING
WE HAVE SEEN THAT SOME ESTIMATORS CAN TRANSFORM DATA AND THAT SOME ESTIMATORS CAN PREDICT VARIABLES WE CAN ALSO
CREATE COMBINED ESTIMATORS
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOTT AS PLT
IMPORT PANDAS AS PD
FROM SKLEARN IMPORT DATASETS
FROM SKLEARNDECOMPOSITION IMPORT PCA
FROM SKLEARNLINEARMODEL IMPORT SGDCLASSIFIER
FROM SKLEARNPIPELINE IMPORT PIPELINE
FROM SKLEARNMODELSELECTION IMPORT GRIDSEARCHCV
DEFINE A PIPELINE TO SEARCH FOR THE BEST COMBINATION OF PCA TRUNCATION
AND CLASSIFIER REGULARIZATION
LOGISTIC SGDCLASSIFIERLOSSLOG PENALTYL2 EARLYSTOPPING TRUE
MAXITER10000 TOL1E5 RANDOMSTATE0
PCA PCA
PIPE PIPELINESTEPSPCA PCA LOGISTIC LOGISTIC
DIGITS DATASETSLOADDIGITS
XDIGITS DIGITS DATA
YDIGITS DIGITS TARGET
PARAMETERS OF PIPELINES CAN BE SET USING " SEPARATED PARAMETER NAMES
PARAMGRID
PCANCOMPONENTS 5 20 30 40 50 64
LOGISTICALPHA NPLOGSPACE4 4 5

SEARCH GRIDSEARCHCVPIPE PARAMGRID IID FALSE CV5
SEARCHFITXDIGITS YDIGITS
22 A TUTORIAL ON STATISTICALLEARNING FOR SCIENTIFIC DATA PROCESSING 181
```

SCIKITLEARN USER GUIDE RELEASE 0213  
PRINTBEST PARAMETER CV SCORE 03F SEARCHBESTSCORE  
PRINTSEARCHBESTPARAMS  
PLOT THE PCA SPECTRUM  
PCAFITXDIGITS  
FIG AX0 AX1 PLTSUBPLOTSNROWS2 SHAREX TRUE FIGSIZE6 6  
AX0PLOTPCAEXPLAINEDVARIANCERATIO LINEWIDTH2  
AX0SETYLABELPCA EXPLAINED VARIANCE  
AX0AXVLINERSEARCHBESTESTIMATORNAMEDSTEPSPCANCOMPONENTS  
LINESTYLE LABELNCOMPONENTS CHOSEN  
FACE RECOGNITION WITH EIGENFACES  
THE DATASET USED IN THIS EXAMPLE IS A PREPROCESSED EXCERPT OF THE “LABELED FACES IN THE WILD” ALSO KNOWN AS LFW  
HTTPVISWWWCSUMASSEDULFWLFWFUNNELEDTGZ 233MB

FACES RECOGNITION EXAMPLE USING EIGENFACES AND SVMs

THE DATASET USED IN THIS EXAMPLE IS A PREPROCESSED EXCERPT OF THE  
LABELED FACES IN THE WILD AKA LFW  
HTTPVISWWWCSUMASSEDULFWLFWFUNNELEDTGZ 233MB  
LFW HTTPVISWWWCSUMASSEDULFW  
EXPECTED RESULTS FOR THE TOP 5 MOST REPRESENTED PEOPLE IN THE DATASET

PRECISION RECALL F1SCORE SUPPORT

ARIEL SHARON 067 092 077 13  
COLIN POWELL 075 078 076 60  
DONALD RUMSFELD 078 067 072 27  
GEORGE W BUSH 086 086 086 146  
GERHARD SCHROEDER 076 076 076 25  
HUGO CHAVEZ 067 067 067 15  
TONY BLAIR 081 069 075 36  
AVG TOTAL 080 080 080 322

FROM TIME IMPORT TIME  
IMPORT LOGGING  
IMPORT MATPLOTLIBPYPLOT AS PLT  
FROM SKLEARNMODELSELECTION IMPORT TRAINTESTSPLIT  
FROM SKLEARNMODELSELECTION IMPORT GRIDSEARCHCV  
FROM SKLEARNDATASETS IMPORT FETCHLFWPEOPLE  
FROM SKLEARNMETRICS IMPORT CLASSIFICATIONREPORT  
FROM SKLEARNMETRICS IMPORT CONFUSIONMATRIX  
182 CHAPTER 2 SCIKITLEARN TUTORIALS

```
SCIKITLEARN USER GUIDE RELEASE 0213
FROM SKLEARNDECOMPOSITION IMPORT PCA
FROM SKLEARN SVM IMPORT SVC
PRINTDOC
    DISPLAY PROGRESS LOGS ON STDOUT
LOGGINGBASICCONFIGLEVELLOGGINGINFO FORMAT ASCTIMES MESSAGES

    DOWNLOAD THE DATA IF NOT ALREADY ON DISK AND LOAD IT AS NUMPY ARRAYS
LFWPEOPLE FETCHLFWPEOPLEMINFACESPERPERSON70 RESIZE04
INTROSPECT THE IMAGES ARRAYS TO FIND THE SHAPES FOR PLOTTING
NSAMPLES H W LFWPEOPLEIMAGESSHAPE
    FOR MACHINE LEARNING WE USE THE 2 DATA DIRECTLY AS RELATIVE PIXEL
    POSITIONS INFO IS IGNORED BY THIS MODEL
X LFWPEOPLEDATA
NFEATURES XSHAPE1
    THE LABEL TO PREDICT IS THE ID OF THE PERSON
Y LFWPEOPLETARGET
TARGETNAMES LFWPEOPLETARGETNAMES
NCLASSES TARGETNAMESSHAPE0
PRINTTOTAL DATASET SIZE
PRINTNSAMPLES D NSAMPLES
PRINTNFEATURES D NFEATURES
PRINTNCLASSES D NCLASSES

    SPLIT INTO A TRAINING SET AND A TEST SET USING A STRATIFIED K FOLD
    SPLIT INTO A TRAINING AND TESTING SET
XTRAIN XTEST YTRAIN YTEST TRAINTESTSPLIT
X Y TESTSIZE025 RANDOMSTATE42

    COMPUTE A PCA EIGENFACES ON THE FACE DATASET TREATED AS UNLABELED
    DATASET UNSUPERVISED FEATURE EXTRACTION DIMENSIONALITY REDUCTION
NCOMPONENTS 150
PRINTEXTRACTING THE TOP DEIGENFACES FROM DFACES
NCOMPONENTS XTRAINSHAPE0
T0 TIME
PCA PCANCOMPONENTSNCOMPONENTS SVDSOLVERRANDOMIZED
WHITENTRUEFITXTRAIN
PRINTDONE IN 03FS TIME T0
EIGENFACES PCACOMPONENTSRESHAPENCOMPONENTS H W
PRINTPROJECTING THE INPUT DATA ON THE EIGENFACES ORTHONORMAL BASIS
T0 TIME
22 A TUTORIAL ON STATISTICALLEARNING FOR SCIENTIFIC DATA PROCESSING 183
```

```
SCIKITLEARN USER GUIDE RELEASE 0213
XTRAINPCA PCATransformXTrain
XTESTPCA PCATransformXTest
PRINTDONE IN 03FS TIME T0

TRAIN A SVM CLASSIFICATION MODEL
PRINTFITTING THE CLASSIFIER TO THE TRAINING SET
T0 TIME
PARAMGRID C 1E3 5E3 1E4 5E4 1E5
GAMMA 00001 00005 0001 0005 001 01
CLF GRIDSEARCHCVSVCKernelRBF CLASSWEIGHTBALANCED
PARAMGRID CV5 IID FALSE
CLF CLFFITXTRAINPCA YTRAIN
PRINTDONE IN 03FS TIME T0
PRINTBEST ESTIMATOR FOUND BY GRID SEARCH
PRINTCLFBESTESTIMATOR

QUANTITATIVE EVALUATION OF THE MODEL QUALITY ON THE TEST SET
PRINTPREDICTING PEOPLES NAMES ON THE TEST SET
T0 TIME
YPRED CLFPREDICTXTESTPCA
PRINTDONE IN 03FS TIME T0
PRINTCLASSIFICATIONREPORTYTEST YPRED TARGETNAMESTARGETNAMES
PRINTCONFUSIONMATRIXTEST YPRED LABELSRANGENCLASSES

QUALITATIVE EVALUATION OF THE PREDICTIONS USING MATPLOTLIB
DEFPlotGalleryImages titles H W NROW3 NCOL4
HELPER FUNCTION TO PLOT A GALLERY OF PORTRAITS
PLTFigureFigSize18 NCOL 24 NROW
PLTSubplotsAdjustBottom0 LEFT01 RIGHT99 TOP90 HSPACE35
FORIINRANGENROW NCOL
PLTSubPlotNROW NCOL I 1
PLTImShowImagesiReshapeH W CMAPPLTCMGRAY
PLTTITLEtitlesI SIZE12
PLTXTICKS
PLTYTICKS
PLOT THE RESULT OF THE PREDICTION ON A PORTION OF THE TEST SET
DEFTITLEYPRED YTEST TARGETNAMES I
PREDNAME TARGETNAMESYPREDIRSPIT 11
TRUENAME TARGETNAMESYTESTIRSPIT 11
RETURNPREDICTED SNTRUE S PREDNAME TRUENAME
PREDICTIONTITLES TITLEYPRED YTEST TARGETNAMES I
FORIINRANGEYPREDSHAPE0
PLOTGALLERYXTEST PREDICTIONTITLES H W
184 CHAPTER 2 SCIKITLEARN TUTORIALS
```

SCIKITLEARN USER GUIDE RELEASE 0213  
PLOT THE GALLERY OF THE MOST SIGNIFICATIVE EIGENFACES  
EIGENFACETITLES EIGENFACE D IFORIINRANGEEIGENFACESSHAPE0  
PLOTGALLERYEIGENFACES EIGENFACETITLES H W  
PLTSHOW  
PREDICTION EIGENFACES  
EXPECTED RESULTS FOR THE TOP 5 MOST REPRESENTED PEOPLE IN THE DATASET  
PRECISION RECALL F1SCORE SUPPORT  
GERHARDSCHROEDER 091 075 082 28  
DONALDRUMSFELD 084 082 083 33  
TONYBLAIR 065 082 073 34  
COLINPOWELL 078 088 083 58  
GEORGEWBUSH 093 086 090 129  
AVG TOTAL 086 084 085 282  
OPEN PROBLEM STOCK MARKET STRUCTURE  
CAN WE PREDICT THE VARIATION IN STOCK PRICES FOR GOOGLE OVER A GIVEN TIME FRAME  
LEARNING A GRAPH STRUCTURE  
226 FINDING HELP  
THE PROJECT MAILING LIST  
IF YOU ENCOUNTER A BUG WITH SCIKITLEARN OR SOMETHING THAT NEEDS CLARIFICATION IN THE DOCSTRING OR THE ONLINE  
DOCUMENTATION PLEASE FEEL FREE TO ASK ON THE MAILING LIST  
22 A TUTORIAL ON STATISTICALLEARNING FOR SCIENTIFIC DATA PROCESSING 185

SCIKITLEARN USER GUIDE RELEASE 0213

QA COMMUNITIES WITH MACHINE LEARNING PRACTITIONERS

QUORACOM QUORA HAS A TOPIC FOR MACHINE LEARNING RELATED QUESTIONS THAT ALSO FEATURES SOME INTERESTING DISCUSSIONS [HTTPSWWWQUORACOMTOPICMACHINELEARNING](https://www.quora.com/topic/machine-learning)

STACK EXCHANGE THE STACK EXCHANGE FAMILY OF SITES HOSTS MULTIPLE SUBDOMAINS FOR MACHINE LEARNING QUESTIONS

- 'AN EXCELLENT FREE ONLINE COURSE FOR MACHINE LEARNING TAUGHT BY PROFESSOR ANDREW NG OF STANFORD' [HTTPSWWW COURSEARAORGLearnMACHINELEARNING](https://www.coursera.org/learn/machine-learning)

- 'ANOTHER EXCELLENT FREE ONLINE COURSE THAT TAKES A MORE GENERAL APPROACH TO ARTIFICIAL INTELLIGENCE' [HTTPSWWWUDACITYCOMCOURSEINTROTOARTIFICIALINTELLIGENCE-CS271](https://www.udacity.com/course/intro-to-artificial-intelligence-cs271)

23 WORKING WITH TEXT DATA

THE GOAL OF THIS GUIDE IS TO EXPLORE SOME OF THE MAIN SCIKITLEARN TOOLS ON A SINGLE PRACTICAL TASK ANALYZING A COLLECTION OF TEXT DOCUMENTS NEWSGROUPS POSTS ON TWENTY DIFFERENT TOPICS

IN THIS SECTION WE WILL SEE HOW TO

- LOAD THE FILE CONTENTS AND THE CATEGORIES
- EXTRACT FEATURE VECTORS SUITABLE FOR MACHINE LEARNING
- TRAIN A LINEAR MODEL TO PERFORM CATEGORIZATION
- USE A GRID SEARCH STRATEGY TO FIND A GOOD CONFIGURATION OF BOTH THE FEATURE EXTRACTION COMPONENTS AND THE CLASSIFIER

231 TUTORIAL SETUP

TO GET STARTED WITH THIS TUTORIAL YOU MUST FIRST INSTALL SCIKITLEARN AND ALL OF ITS REQUIRED DEPENDENCIES

PLEASE REFER TO THE INSTALLATION INSTRUCTIONS PAGE FOR MORE INFORMATION AND FOR SYSTEMSPECIFIC INSTRUCTIONS

THE SOURCE OF THIS TUTORIAL CAN BE FOUND WITHIN YOUR SCIKITLEARN FOLDER

SCIKITLEARNDOCTUTORIALTEXTANALYTICS

THE SOURCE CAN ALSO BE FOUND ON GITHUB

THE TUTORIAL FOLDER SHOULD CONTAIN THE FOLLOWING SUBFOLDERS

- RST FILES THE SOURCE OF THE TUTORIAL DOCUMENT WRITTEN WITH SPHINX
- DATA FOLDER TO PUT THE DATASETS USED DURING THE TUTORIAL
- SKELETONS SAMPLE INCOMPLETE SCRIPTS FOR THE EXERCISES
- SOLUTIONS SOLUTIONS OF THE EXERCISES

YOU CAN ALREADY COPY THE SKELETONS INTO A NEW FOLDER SOMEWHERE ON YOUR HARDDRIVE NAMED SKLEARNTUTWORKSPACE WHERE YOU WILL EDIT YOUR OWN FILES FOR THE EXERCISES WHILE KEEPING THE ORIGINAL SKELETONS INTACT

CP R SKELETONS WORKDIRECTORYSKLEARNTUTWORKSPACE

186 CHAPTER 2 SCIKITLEARN TUTORIALS



SCIKITLEARN USER GUIDE RELEASE 0213

MACHINE LEARNING ALGORITHMS NEED DATA GO TO EACH TUTORIALHOMEDATA SUBFOLDER AND RUN THE  
FETCHDATAPY SCRIPT FROM THERE AFTER HAVING READ THEM FIRST

FOR INSTANCE

CD TUTORIALHOMEDATALANGUAGES  
LESS FETCHDATAPY  
PYTHON FETCHDATAPY

232 LOADING THE 20 NEWSGROUPS DATASET

THE DATASET IS CALLED “TWENTY NEWSGROUPS” HERE IS THE OFFICIAL DESCRIPTION QUOTED FROM THE WEBSITE

THE 20 NEWSGROUPS DATA SET IS A COLLECTION OF APPROXIMATELY 20000 NEWSGROUP DOCUMENTS PARTITIONED  
NEARLY EVENLY ACROSS 20 DIFFERENT NEWSGROUPS TO THE BEST OF OUR KNOWLEDGE IT WAS ORIGINALLY COLLECTED  
BY KEN LANG PROBABLY FOR HIS PAPER “NEWSWEEDER LEARNING TO FILTER NETNEWS” THOUGH HE DOES NOT EXPLIC  
ITLY MENTION THIS COLLECTION THE 20 NEWSGROUPS COLLECTION HAS BECOME A POPULAR DATA SET FOR EXPERIMENTS  
IN TEXT APPLICATIONS OF MACHINE LEARNING TECHNIQUES SUCH AS TEXT CLASSIFICATION AND TEXT CLUSTERING  
IN THE FOLLOWING WE WILL USE THE BUILTIN DATASET LOADER FOR 20 NEWSGROUPS FROM SCIKITLEARN ALTERNATIVELY IT IS POSSIBLE  
TO DOWNLOAD THE DATASET MANUALLY FROM THE WEBSITE AND USE THE SKLEARNDATASETSLOADFILES FUNCTION BY  
POINTING IT TO THE 20NEWSBYDATETRAIN SUBFOLDER OF THE UNCOMPRESSED ARCHIVE FOLDER  
IN ORDER TO GET FASTER EXECUTION TIMES FOR THIS FIRST EXAMPLE WE WILL WORK ON A PARTIAL DATASET WITH ONLY 4 CATEGORIES OUT  
OF THE 20 AVAILABLE IN THE DATASET

CATEGORIES ALTATHEISM SOCRELIGIONCHRISTIAN  
COMPGRAPHICS SCIMED

WE CAN NOW LOAD THE LIST OF FILES MATCHING THOSE CATEGORIES AS FOLLOWS

FROM SKLEARNDATASETS IMPORT FETCH20NEWSGROUPS  
TWENTYTRAIN FETCH20NEWSGROUPSSUBSETTRAIN  
CATEGORIESCATEGORIES SHUFFLE TRUE RANDOMSTATE42

THE RETURNED DATASET IS A SCIKITLEARN “BUNCH” A SIMPLE HOLDER OBJECT WITH FIELDS THAT CAN BE BOTH ACCESSED  
AS PYTHONDICT KEYS OROBJECT ATTRIBUTES FOR CONVENIENCE FOR INSTANCE THE TARGETNAMES HOLDS THE LIST OF THE  
REQUESTED CATEGORY NAMES

TWENTYTRAINTARGETNAMES  
ALTATHEISM COMPGRAPHICS SCIMED SOCRELIGIONCHRISTIAN

THE FILES THEMSELVES ARE LOADED IN MEMORY IN THE DATA ATTRIBUTE FOR REFERENCE THE FILENAMES ARE ALSO AVAILABLE

LENTWENTYTRAINDATA  
2257  
LENTWENTYTRAINFILENAMES  
2257

LET’S PRINT THE FIRST LINES OF THE FIRST LOADED FILE

PRINTNJOINTWENTYTRAINDATA0SPLIT N3  
FROM SD345CITYACUK MICHAEL COLLIER  
SUBJECT CONVERTING IMAGES TO HP LASERJET III  
NNTPPOSTINGHOST HAMPTON  
PRINTTWENTYTRAINTARGETNAMESTWENTYTRAINTARGETO  
COMPGRAPHICS

23 WORKING WITH TEXT DATA 187

SCIKITLEARN USER GUIDE RELEASE 0213

SUPERVISED LEARNING ALGORITHMS WILL REQUIRE A CATEGORY LABEL FOR EACH DOCUMENT IN THE TRAINING SET IN THIS CASE THE CATEGORY IS THE NAME OF THE NEWSGROUP WHICH ALSO HAPPENS TO BE THE NAME OF THE FOLDER HOLDING THE INDIVIDUAL DOCUMENTS FOR SPEED AND SPACE EFFICIENCY REASONS SCIKITLEARN LOADS THE TARGET ATTRIBUTE AS AN ARRAY OF INTEGERS THAT CORRESPONDS TO THE INDEX OF THE CATEGORY NAME IN THE TARGETNAMES LIST THE CATEGORY INTEGER ID OF EACH SAMPLE IS STORED IN THE TARGET ATTRIBUTE

TWENTYTRAIN TARGET10  
ARRAY1 1 3 3 3 3 3 2 2 2

IT IS POSSIBLE TO GET BACK THE CATEGORY NAMES AS FOLLOWS

FOR TINTWENTYTRAIN TARGET10  
PRINT TWENTYTRAIN TARGETNAMEST

COMPGRAPHICS  
COMPGRAPHICS  
SOCRELIGIONCHRISTIAN  
SOCRELIGIONCHRISTIAN  
SOCRELIGIONCHRISTIAN  
SOCRELIGIONCHRISTIAN  
SOCRELIGIONCHRISTIAN  
SCIMED  
SCIMED  
SCIMED

YOU MIGHT HAVE NOTICED THAT THE SAMPLES WERE SHUFFLED RANDOMLY WHEN WE CALLED FETCH20NEWSGROUPS SHUFFLETRUE RANDOMSTATE42 THIS IS USEFUL IF YOU WISH TO SELECT ONLY A SUBSET OF SAMPLES TO QUICKLY TRAIN A MODEL AND GET A FIRST IDEA OF THE RESULTS BEFORE RETRAINING ON THE COMPLETE DATASET LATER

233 EXTRACTING FEATURES FROM TEXT FILES

IN ORDER TO PERFORM MACHINE LEARNING ON TEXT DOCUMENTS WE FIRST NEED TO TURN THE TEXT CONTENT INTO NUMERICAL FEATURE VECTORS

BAGS OF WORDS

THE MOST INTUITIVE WAY TO DO SO IS TO USE A BAGS OF WORDS REPRESENTATION

1 ASSIGN A FIXED INTEGER ID TO EACH WORD OCCURRING IN ANY DOCUMENT OF THE TRAINING SET FOR INSTANCE BY BUILDING A DICTIONARY FROM WORDS TO INTEGER INDICES

2 FOR EACH DOCUMENT I COUNT THE NUMBER OF OCCURRENCES OF EACH WORD AND STORE IT IN  $X_{ij}$  AS THE VALUE OF  $FEATURE_j$  WHERE  $j$  IS THE INDEX OF WORD  $w_i$  IN THE DICTIONARY

THE BAGS OF WORDS REPRESENTATION IMPLIES THAT  $N_{FEATURES}$  IS THE NUMBER OF DISTINCT WORDS IN THE CORPUS THIS NUMBER IS TYPICALLY LARGER THAN 100000

IF  $N_{SAMPLES} = 10000$  STORING  $X$  AS A NUMPY ARRAY OF TYPE FLOAT32 WOULD REQUIRE  $10000 \times 100000 \times 4$  BYTES 4GB IN RAM WHICH IS BARELY MANAGEABLE ON TODAY'S COMPUTERS

FORTUNATELY MOST VALUES IN  $X$  WILL BE ZEROS SINCE FOR A GIVEN DOCUMENT LESS THAN A FEW THOUSAND DISTINCT WORDS WILL BE USED FOR THIS REASON WE SAY THAT BAGS OF WORDS ARE TYPICALLY HIGH-DIMENSIONAL SPARSE DATASETS WE CAN SAVE A LOT OF MEMORY BY ONLY STORING THE NONZERO PARTS OF THE FEATURE VECTORS IN MEMORY

188 CHAPTER 2 SCIKITLEARN TUTORIALS

SCIKITLEARN USER GUIDE RELEASE 0213

SCIPYSPARSE MATRICES ARE DATA STRUCTURES THAT DO EXACTLY THIS AND SCIKITLEARN HAS BUILTIN SUPPORT FOR THESE STRUCTURES

TOKENIZING TEXT WITH SCIKITLEARN

TEXT PREPROCESSING TOKENIZING AND FILTERING OF STOPWORDS ARE ALL INCLUDED IN COUNTVECTORIZER WHICH BUILDS A DICTIONARY OF FEATURES AND TRANSFORMS DOCUMENTS TO FEATURE VECTORS

```
from sklearn.feature_extraction.text import CountVectorizer
countvect = CountVectorizer()
xtraincounts = countvect.fit_transform(traindata)
xtraincountsshape
```

2257 35788

COUNTVECTORIZER SUPPORTS COUNTS OF NGRAMS OF WORDS OR CONSECUTIVE CHARACTERS ONCE FITTED THE VECTORIZER HAS BUILT A DICTIONARY OF FEATURE INDICES

```
countvect.vocabulary.get('algorithm')
4690
```

THE INDEX VALUE OF A WORD IN THE VOCABULARY IS LINKED TO ITS FREQUENCY IN THE WHOLE TRAINING CORPUS

FROM OCCURRENCES TO FREQUENCIES

OCCURRENCE COUNT IS A GOOD START BUT THERE IS AN ISSUE LONGER DOCUMENTS WILL HAVE HIGHER AVERAGE COUNT VALUES THAN SHORTER DOCUMENTS EVEN THOUGH THEY MIGHT TALK ABOUT THE SAME TOPICS

TO AVOID THESE POTENTIAL DISCREPANCIES IT SUFFICES TO DIVIDE THE NUMBER OF OCCURRENCES OF EACH WORD IN A DOCUMENT BY THE TOTAL NUMBER OF WORDS IN THE DOCUMENT THESE NEW FEATURES ARE CALLED TF FOR TERM FREQUENCIES

ANOTHER REFINEMENT ON TOP OF TF IS TO DOWNSCALE WEIGHTS FOR WORDS THAT OCCUR IN MANY DOCUMENTS IN THE CORPUS AND ARE THEREFORE LESS INFORMATIVE THAN THOSE THAT OCCUR ONLY IN A SMALLER PORTION OF THE CORPUS

THIS DOWNSCALING IS CALLED TF-IDF FOR “TERM FREQUENCY TIMES INVERSE DOCUMENT FREQUENCY”

BOTH TF AND TF-IDF CAN BE COMPUTED AS FOLLOWS USING TfidfTransformer

```
from sklearn.feature_extraction.text import TfidfTransformer
tftransformer = TfidfTransformer(use_idf=False)
xtraintf = tftransformer.transform(xtraincounts)
xtraintfshape
```

2257 35788

IN THE ABOVE EXAMPLE CODE WE FIRSTLY USE THE FIT METHOD TO FIT OUR ESTIMATOR TO THE DATA AND SECONDLY THE TRANSFORM METHOD TO TRANSFORM OUR COUNTMATRIX TO A TFIDF REPRESENTATION THESE TWO STEPS CAN BE COMBINED TO ACHIEVE THE SAME END RESULT FASTER BY SKIPPING REDUNDANT PROCESSING THIS IS DONE THROUGH USING THE FITTRANSFORM METHOD AS SHOWN BELOW AND AS MENTIONED IN THE NOTE IN THE PREVIOUS SECTION

```
TfidfTransformer()
TfidfTransformer()
xtraintfidf = TfidfTransformer().fit_transform(xtraincounts)
xtraintfidfshape
```

2257 35788

23 WORKING WITH TEXT DATA 189

SCIKITLEARN USER GUIDE RELEASE 0213

234 TRAINING A CLASSIFIER

NOW THAT WE HAVE OUR FEATURES WE CAN TRAIN A CLASSIFIER TO TRY TO PREDICT THE CATEGORY OF A POST LET’S START WITH A NAÏVE BAYES CLASSIFIER WHICH PROVIDES A NICE BASELINE FOR THIS TASK SCIKITLEARN INCLUDES SEVERAL VARIANTS OF THIS CLASSIFIER THE ONE MOST SUITABLE FOR WORD COUNTS IS THE MULTINOMIAL VARIANT

```
FROM SKLEARNNAIVEBAYES IMPORT MULTINOMIALNB
CLF MULTINOMIALNBFITXTRAINTFIDF TWENTYTRAINTARGET
TO TRY TO PREDICT THE OUTCOME ON A NEW DOCUMENT WE NEED TO EXTRACT THE FEATURES USING ALMOST THE SAME FEATURE EXTRACT
ING CHAIN AS BEFORE THE DIFFERENCE IS THAT WE CALL TRANSFORM INSTEAD OFFITTRANSFORM ON THE TRANSFORMERS
SINCE THEY HAVE ALREADY BEEN FIT TO THE TRAINING SET
DOCSNEW GOD IS LOVE OPENGL ON THE GPU IS FAST
XNEWCOUNTS COUNTVECTTRANSFORMDOCSNEW
XNEWTFIDF TFIDFTRANSFORMERTRANSFORMMXNEWCOUNTS
PREDICTED CLFPREDICTXNEWTFIDF
FOR DOC CATEGORY INZIPDOCSNEW PREDICTED
PRINTRS DOC TWENTYTRAINTARGETNAMECATEGORY
```

GOD IS LOVE SOCRELIGIONCHRISTIAN

OPENGL ON THE GPU IS FAST COMPGRAPHICS

235 BUILDING A PIPELINE

IN ORDER TO MAKE THE VECTORIZER TRANSFORMER CLASSIFIER EASIER TO WORK WITH SCIKITLEARN PROVIDES A PIPELINE CLASS THAT BEHAVES LIKE A COMPOUND CLASSIFIER

```
FROM SKLEARNPIPELINE IMPORT PIPELINE
TEXTCLF PIPELINE
VECT COUNTVECTORIZER
TFIDF TFIDFTRANSFORMER
CLF MULTINOMIALNB
```

THE NAMES VECT TFIDF ANDCLF CLASSIFIER ARE ARBITRARY WE WILL USE THEM TO PERFORM GRID SEARCH FOR SUITABLE HYPERPARAMETERS BELOW WE CAN NOW TRAIN THE MODEL WITH A SINGLE COMMAND

```
TEXTCLFFITTWENTYTRAINDATA TWENTYTRAINTARGET
PIPELINE
```

236 EVALUATION OF THE PERFORMANCE ON THE TEST SET

EVALUATING THE PREDICTIVE ACCURACY OF THE MODEL IS EQUALLY EASY

```
IMPORT NUMPY AS NP
TWENTYTEST FETCH20NEWSGROUPSSUBSETTEST
CATEGORIESCATEGORIES SHUFFLE TRUE RANDOMSTATE42
DOCSTEST TWENTYTESTDATA
PREDICTED TEXTCLFPREDICTDOCSTEST
NPMEANPREDICTED TWENTYTESTTARGET
08348
190 CHAPTER 2 SCIKITLEARN TUTORIALS
```

SCIKITLEARN USER GUIDE RELEASE 0213

WE ACHIEVED 835 ACCURACY LET’S SEE IF WE CAN DO BETTER WITH A LINEAR SUPPORT VECTOR MACHINE SVM WHICH IS WIDELY REGARDED AS ONE OF THE BEST TEXT CLASSIFICATION ALGORITHMS ALTHOUGH IT’S ALSO A BIT SLOWER THAN NAÏVE BAYES WE CAN CHANGE THE LEARNER BY SIMPLY PLUGGING A DIFFERENT CLASSIFIER OBJECT INTO OUR PIPELINE

```
FROM SKLEARNLINEARMODEL IMPORT SGDCLASSIFIER
TEXTCLF PIPELINE
VECT COUNTVECTORIZER
TFIDF TFIDFTRANSFORMER
CLF SGDCLASSIFIERLOSSHINGE PENALTYL2
ALPHA1E3 RANDOMSTATE42
MAXITER5 TOL NONE
```

TEXTCLFFITTWENTYTRAINDATA TWENTYTRAINTARGET PIPELINE

```
PREDICTED TEXTCLFPREDICTDOCSTEST
NPMEANPREDICTED TWENTYTESTTARGET
09101
```

WE ACHIEVED 913 ACCURACY USING THE SVM SCIKITLEARN PROVIDES FURTHER UTILITIES FOR MORE DETAILED PERFORMANCE ANALYSIS OF THE RESULTS

```
FROM SKLEARN IMPORT METRICS
PRINTMETRICSCCLASSIFICATIONREPORTTWENTYTESTTARGET PREDICTED
TARGETNAMESTWENTYTESTTARGETNAMES
```

PRECISION RECALL F1SCORE SUPPORT

```
ALTATHEISM 095 080 087 319
COMPGRAPHICS 087 098 092 389
SCIMED 094 089 091 396
SOCRELIGIONCHRISTIAN 090 095 093 398
ACCURACY 091 1502
MACRO AVG 091 091 091 1502
WEIGHTED AVG 091 091 091 1502
METRICSCONFUSIONMATRIXTWENTYTESTTARGET PREDICTED
ARRAY256 11 16 36
4 380 3 2
5 35 353 3
5 11 4 378
```

AS EXPECTED THE CONFUSION MATRIX SHOWS THAT POSTS FROM THE NEWSGROUPS ON ATHEISM AND CHRISTIANITY ARE MORE OFTEN CONFUSED FOR ONE ANOTHER THAN WITH COMPUTER GRAPHICS

237 PARAMETER TUNING USING GRID SEARCH

WE’VE ALREADY ENCOUNTERED SOME PARAMETERS SUCH AS USEIDF IN THETFIDFTRANSFORMER CLASSIFIERS TEND TO HAVE MANY PARAMETERS AS WELL EG MULTINOMIALNB INCLUDES A SMOOTHING PARAMETER ALPHA ANDSGDCLASSIFIER HAS A PENALTY PARAMETER ALPHA AND CONFIGURABLE LOSS AND PENALTY TERMS IN THE OBJECTIVE FUNCTION SEE THE MODULE DOCUMENTATION OR USE THE PYTHON HELP FUNCTION TO GET A DESCRIPTION OF THESE

INSTEAD OF TWEAKING THE PARAMETERS OF THE VARIOUS COMPONENTS OF THE CHAIN IT IS POSSIBLE TO RUN AN EXHAUSTIVE SEARCH OF THE BEST PARAMETERS ON A GRID OF POSSIBLE VALUES WE TRY OUT ALL CLASSIFIERS ON EITHER WORDS OR BIGRAMS WITH OR WITHOUT IDF AND WITH A PENALTY PARAMETER OF EITHER 001 OR 0001 FOR THE LINEAR SVM

23 WORKING WITH TEXT DATA 191

```
SCIKITLEARN USER GUIDE RELEASE 0213
FROM SKLEARNMODELSELECTION IMPORT GRIDSEARCHCV
PARAMETERS
VECTNGRAMRANGE 1 1 1 2
TFIDFUSEIDF TRUEFALSE
CLFALPHA 1E2 1E3
```

```
OBVIOUSLY SUCH AN EXHAUSTIVE SEARCH CAN BE EXPENSIVE IF WE HAVE MULTIPLE CPU CORES AT OUR DISPOSAL WE CAN TELL
THE GRID SEARCHER TO TRY THESE EIGHT PARAMETER COMBINATIONS IN PARALLEL WITH THE NJOBS PARAMETER IF WE GIVE THIS
PARAMETER A VALUE OF 1 GRID SEARCH WILL DETECT HOW MANY CORES ARE INSTALLED AND USE THEM ALL
GSCLF GRIDSEARCHCVTEXTCLF PARAMETERS CV5 IID FALSE NJOBS1
THE GRID SEARCH INSTANCE BEHAVES LIKE A NORMAL SCIKITLEARN MODEL LET’S PERFORM THE SEARCH ON A SMALLER SUBSET
OF THE TRAINING DATA TO SPEED UP THE COMPUTATION
GSCLF GSCLFFITTWENTYTRAINDATA400 TWENTYTRAINTARGET400
THE RESULT OF CALLING FIT ON AGRIDSEARCHCV OBJECT IS A CLASSIFIER THAT WE CAN USE TO PREDICT
TWENTYTRAINTARGETNAMESGSCLFPREDICTGOD IS LOVE0
SOCRELIGIONCHRISTIAN
THE OBJECT’S BESTSCORE ANDBESTPARAMS ATTRIBUTES STORE THE BEST MEAN SCORE AND THE PARAMETERS SETTING
CORRESPONDING TO THAT SCORE
GSCLFBESTSCORE
09
FOR PARAMNAME INSERTEDPARAMETERSKEYS
PRINTSR PARAMNAME GSCLFBESTPARAMSPARAMNAME
```

```
CLFALPHA 0001
TFIDFUSEIDF TRUE
VECTNGRAMRANGE 1 1
A MORE DETAILED SUMMARY OF THE SEARCH IS AVAILABLE AT GSCLFCVRESULTS
THECVRESULTS PARAMETER CAN BE EASILY IMPORTED INTO PANDAS AS A DATAFRAME FOR FURTHER INSPECTION
EXERCISES
TO DO THE EXERCISES COPY THE CONTENT OF THE ‘SKELETONS’ FOLDER AS A NEW FOLDER NAMED ‘WORKSPACE’
CP R SKELETONS WORKSPACE
YOU CAN THEN EDIT THE CONTENT OF THE WORKSPACE WITHOUT FEAR OF LOSING THE ORIGINAL EXERCISE INSTRUCTIONS
THEN FIRE AN IPYTHON SHELL AND RUN THE WORKINPROGRESS SCRIPT WITH
1 RUN WORKSPACEEXERCISEXXSCRIPTPY ARG1 ARG2 ARG3
IF AN EXCEPTION IS TRIGGERED USE DEBUG TO FIREUP A POST MORTEM IPDB SESSION
REFINE THE IMPLEMENTATION AND ITERATE UNTIL THE EXERCISE IS SOLVED
FOR EACH EXERCISE THE SKELETON FILE PROVIDES ALL THE NECESSARY IMPORT STATEMENTS BOILERPLATE CODE TO LOAD THE
DATA AND SAMPLE CODE TO EVALUATE THE PREDICTIVE ACCURACY OF THE MODEL
192 CHAPTER 2 SCIKITLEARN TUTORIALS
```

SCIKITLEARN USER GUIDE RELEASE 0213

238 EXERCISE 1 LANGUAGE IDENTIFICATION

- WRITE A TEXT CLASSIFICATION PIPELINE USING A CUSTOM PREPROCESSOR AND CHARNGRAMANALYZER USING DATA FROM WIKIPEDIA ARTICLES AS TRAINING SET
- EVALUATE THE PERFORMANCE ON SOME HELD OUT TEST SET

IPYTHON COMMAND LINE

RUN WORKSPACEEXERCISE01LANGUAGETRAINMODEL PY DATALANGUAGESPARAGRAPHS

239 EXERCISE 2 SENTIMENT ANALYSIS ON MOVIE REVIEWS

- WRITE A TEXT CLASSIFICATION PIPELINE TO CLASSIFY MOVIE REVIEWS AS EITHER POSITIVE OR NEGATIVE
- FIND A GOOD SET OF PARAMETERS USING GRID SEARCH
- EVALUATE THE PERFORMANCE ON A HELD OUT TEST SET

IPYTHON COMMAND LINE

RUN WORKSPACEEXERCISE02SENTIMENTPY DATAMOVIEREVIEWSTXTSENTOKEN

2310 EXERCISE 3 CLI TEXT CLASSIFICATION UTILITY

USING THE RESULTS OF THE PREVIOUS EXERCISES AND THE CPICKLE MODULE OF THE STANDARD LIBRARY WRITE A COMMAND LINE UTILITY THAT DETECTS THE LANGUAGE OF SOME TEXT PROVIDED ON STDIN AND ESTIMATE THE POLARITY POSITIVE OR NEGATIVE IF THE TEXT IS WRITTEN IN ENGLISH

BONUS POINT IF THE UTILITY IS ABLE TO GIVE A CONFIDENCE LEVEL FOR ITS PREDICTIONS

2311 WHERE TO FROM HERE

HERE ARE A FEW SUGGESTIONS TO HELP FURTHER YOUR SCIKITLEARN INTUITION UPON THE COMPLETION OF THIS TUTORIAL

- TRY PLAYING AROUND WITH THE ANALYZER ANDTOKEN NORMALISATION UNDERCOUNTVECTORIZER
- IF YOU DON'T HAVE LABELS TRY USING CLUSTERING ON YOUR PROBLEM
- IF YOU HAVE MULTIPLE LABELS PER DOCUMENT EG CATEGORIES HAVE A LOOK AT THE MULTICLASS AND MULTILABEL SECTION
- TRY USING TRUNCATED SVD FOR LATENT SEMANTIC ANALYSIS
- HAVE A LOOK AT USING OUTOF CORE CLASSIFICATION TO LEARN FROM DATA THAT WOULD NOT FIT INTO THE COMPUTER MAIN MEMORY
- HAVE A LOOK AT THE HASHING VECTORIZER AS A MEMORY EFFICIENT ALTERNATIVE TO COUNTVECTORIZER

24 CHOOSING THE RIGHT ESTIMATOR

OFTEN THE HARDEST PART OF SOLVING A MACHINE LEARNING PROBLEM CAN BE FINDING THE RIGHT ESTIMATOR FOR THE JOB

DIFFERENT ESTIMATORS ARE BETTER SUITED FOR DIFFERENT TYPES OF DATA AND DIFFERENT PROBLEMS

THE FLOWCHART BELOW IS DESIGNED TO GIVE USERS A BIT OF A ROUGH GUIDE ON HOW TO APPROACH PROBLEMS WITH REGARD TO WHICH ESTIMATORS TO TRY ON YOUR DATA

24 CHOOSING THE RIGHT ESTIMATOR 193

SCIKITLEARN USER GUIDE RELEASE 0213  
CLICK ON ANY ESTIMATOR IN THE CHART BELOW TO SEE ITS DOCUMENTATION  
25 EXTERNAL RESOURCES VIDEOS AND TALKS  
FOR WRITTEN TUTORIALS SEE THE TUTORIAL SECTION OF THE DOCUMENTATION  
251 NEW TO SCIENTIFIC PYTHON  
FOR THOSE THAT ARE STILL NEW TO THE SCIENTIFIC PYTHON ECOSYSTEM WE HIGHLY RECOMMEND THE PYTHON SCIENTIFIC LECTURE NOTES THIS WILL HELP YOU FIND YOUR FOOTING A BIT AND WILL DEFINITELY IMPROVE YOUR SCIKITLEARN EXPERIENCE A BASIC UNDERSTANDING OF NUMPY ARRAYS IS RECOMMENDED TO MAKE THE MOST OF SCIKITLEARN  
252 EXTERNAL TUTORIALS  
THERE ARE SEVERAL ONLINE TUTORIALS AVAILABLE WHICH ARE GEARED TOWARD SPECIFIC SUBJECT AREAS

- MACHINE LEARNING FOR NEUROIMAGING IN PYTHON
- MACHINE LEARNING FOR ASTRONOMICAL DATA ANALYSIS

253 VIDEOS

- AN INTRODUCTION TO SCIKITLEARN PART I AND PART II AT SCIPY 2013 BY GAELE VAROQUAUX JAKE VANDERPLAS AND OLIVIER GRISEL NOTEBOOKS ON GITHUB
- INTRODUCTION TO SCIKITLEARN BY GAELE VAROQUAUX AT ICML 2010

A THREE MINUTE VIDEO FROM A VERY EARLY STAGE OF SCIKITLEARN EXPLAINING THE BASIC IDEA AND APPROACH WE ARE FOLLOWING

- INTRODUCTION TO STATISTICAL LEARNING WITH SCIKITLEARN BY GAELE VAROQUAUX AT SCIPY 2011

AN EXTENSIVE TUTORIAL CONSISTING OF FOUR SESSIONS OF ONE HOUR THE TUTORIAL COVERS THE BASICS OF MACHINE LEARNING MANY ALGORITHMS AND HOW TO APPLY THEM USING SCIKITLEARN THE MATERIAL CORRESPONDING IS NOW IN THE SCIKITLEARN DOCUMENTATION SECTION A TUTORIAL ON STATISTICAL LEARNING FOR SCIENTIFIC DATA PROCESSING

- STATISTICAL LEARNING FOR TEXT CLASSIFICATION WITH SCIKITLEARN AND NLTK AND SLIDES BY OLIVIER GRISEL AT PYCON 2011

THIRTY MINUTE INTRODUCTION TO TEXT CLASSIFICATION EXPLAINS HOW TO USE NLTK AND SCIKITLEARN TO SOLVE REALWORLD TEXT CLASSIFICATION TASKS AND COMPARES AGAINST CLOUDBASED SOLUTIONS

- INTRODUCTION TO INTERACTIVE PREDICTIVE ANALYTICS IN PYTHON WITH SCIKITLEARN BY OLIVIER GRISEL AT PYCON 2012

3HOURS LONG INTRODUCTION TO PREDICTION TASKS USING SCIKITLEARN

- SCIKITLEARN MACHINE LEARNING IN PYTHON BY JAKE VANDERPLAS AT THE 2012 PYDATA WORKSHOP AT GOOGLE

INTERACTIVE DEMONSTRATION OF SOME SCIKITLEARN FEATURES 75 MINUTES

- SCIKITLEARN TUTORIAL BY JAKE VANDERPLAS AT PYDATA NYC 2012

PRESENTATION USING THE ONLINE TUTORIAL 45 MINUTES

194 CHAPTER 2 SCIKITLEARN TUTORIALS



SCIKITLEARN USER GUIDE RELEASE 0213

NOTE DOCTEST MODE

THE CODEEXAMPLES IN THE ABOVE TUTORIALS ARE WRITTEN IN A PYTHONCONSOLE FORMAT IF YOU WISH TO EASILY EXECUTE THESE EXAMPLES IN IPYTHON USE DOCTESTMODE

IN THE IPYTHONCONSOLE YOU CAN THEN SIMPLY COPY AND PASTE THE EXAMPLES DIRECTLY INTO IPYTHON WITHOUT HAVING TO WORRY ABOUT REMOVING THE MANUALLY

25 EXTERNAL RESOURCES VIDEOS AND TALKS 195



CHAPTER  
THREE  
USER GUIDE  
31 SUPERVISED LEARNING  
311 GENERALIZED LINEAR MODELS  
THE FOLLOWING ARE A SET OF METHODS INTENDED FOR REGRESSION IN WHICH THE TARGET VALUE IS EXPECTED TO BE A LINEAR COMBI  
NATION OF THE FEATURES IN MATHEMATICAL NOTATION IF  $\hat{y}$  IS THE PREDICTED VALUE  
$$\hat{y} = \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n$$
  
ACROSS THE MODULE WE DESIGNATE THE VECTOR  $\beta = [\beta_0 \ \beta_1 \ \dots \ \beta_n]^T$  AS COEF AND  $\beta_0$  AS INTERCEPT  
TO PERFORM CLASSIFICATION WITH GENERALIZED LINEAR MODELS SEE LOGISTIC REGRESSION  
ORDINARY LEAST SQUARES  
LINEARREGRESSION FITS A LINEAR MODEL WITH COEFFICIENTS  $\beta = [\beta_0 \ \beta_1 \ \dots \ \beta_n]^T$  TO MINIMIZE THE RESIDUAL SUM OF SQUARES  
BETWEEN THE OBSERVED TARGETS IN THE DATASET AND THE TARGETS PREDICTED BY THE LINEAR APPROXIMATION MATHEMATICALLY IT  
SOLVES A PROBLEM OF THE FORM  
MIN  
$$\sum_{i=1}^n (y_i - \hat{y}_i)^2$$
  
2  
LINEARREGRESSION WILL TAKE IN ITS FIT METHOD ARRAYS X Y AND WILL STORE THE COEFFICIENTS  $\beta$  OF THE LINEAR MODEL  
IN ITS COEF MEMBER  
197

SCIKITLEARN USER GUIDE RELEASE 0213  
FROM SKLEARN IMPORT LINEARMODEL  
REG LINEARMODELLINEARREGRESSION  
REGFIT0 0 1 1 2 2 0 1 2

LINEARREGRESSIONCOPYXTRUE FITINTERCEPTTRUE NJOBSNONE  
NORMALIZEFALSE  
REGCOEF  
ARRAY05 05

THE COEFFICIENT ESTIMATES FOR ORDINARY LEAST SQUARES RELY ON THE INDEPENDENCE OF THE FEATURES WHEN FEATURES ARE CORRELATED AND THE COLUMNS OF THE DESIGN MATRIX  $X$  HAVE AN APPROXIMATE LINEAR DEPENDENCE THE DESIGN MATRIX BECOMES CLOSE TO SINGULAR AND AS A RESULT THE LEASTSQUARES ESTIMATE BECOMES HIGHLY SENSITIVE TO RANDOM ERRORS IN THE OBSERVED TARGET PRODUCING A LARGE VARIANCE THIS SITUATION OF MULTICOLLINEARITY CAN ARISE FOR EXAMPLE WHEN DATA ARE COLLECTED WITHOUT AN EXPERIMENTAL DESIGN

EXAMPLES  
•LINEAR REGRESSION EXAMPLE

ORDINARY LEAST SQUARES COMPLEXITY

THE LEAST SQUARES SOLUTION IS COMPUTED USING THE SINGULAR VALUE DECOMPOSITION OF  $X$  IF  $X$  IS A MATRIX OF SHAPE  $N \times P$  WHERE  $N$  IS THE NUMBER OF SAMPLES AND  $P$  IS THE NUMBER OF FEATURES THIS METHOD HAS A COST OF  $O(N^2P)$  ASSUMING THAT  $N \geq P$

RIDGE REGRESSION

RIDGE REGRESSION ADDRESSES SOME OF THE PROBLEMS OF ORDINARY LEAST SQUARES BY IMPOSING A PENALTY ON THE SIZE OF THE COEFFICIENTS THE RIDGE COEFFICIENTS MINIMIZE A PENALIZED RESIDUAL SUM OF SQUARES

MIN  
$$\sum_{i=1}^N (y_i - \hat{y}_i)^2 + \lambda \sum_{j=1}^P \beta_j^2$$

THE COMPLEXITY PARAMETER  $\lambda \geq 0$  CONTROLS THE AMOUNT OF SHRINKAGE THE LARGER THE VALUE OF  $\lambda$  THE GREATER THE AMOUNT OF SHRINKAGE AND THUS THE COEFFICIENTS BECOME MORE ROBUST TO COLLINEARITY

AS WITH OTHER LINEAR MODELS RIDGE WILL TAKE IN ITS FIT METHOD ARRAYS  $X$   $Y$  AND WILL STORE THE COEFFICIENTS  $\beta$  OF THE LINEAR MODEL IN ITS COEF MEMBER

SCIKITLEARN USER GUIDE RELEASE 0213

```
FROM SKLEARN IMPORT LINEARMODEL
REG LINEARMODELRIDGEALPHA5
REGFIT0 0 0 0 1 1 0 1 1
RIDGEALPHA05 COPYXTRUE FITINTERCEPTTRUE MAXITERNONE
NORMALIZEFALSE RANDOMSTATENONE SOLVERAUTO TOL0001
REGCOEF
ARRAY034545455 034545455
REGINTERCEPT
013636
EXAMPLES
•PLOT RIDGE COEFFICIENTS AS A FUNCTION OF THE REGULARIZATION
•CLASSIFICATION OF TEXT DOCUMENTS USING SPARSE FEATURES
RIDGE COMPLEXITY
THIS METHOD HAS THE SAME ORDER OF COMPLEXITY AS ORDINARY LEAST SQUARES
SETTING THE REGULARIZATION PARAMETER GENERALIZED CROSSVALIDATION
RIDGECV IMPLEMENTS RIDGE REGRESSION WITH BUILTIN CROSSVALIDATION OF THE ALPHA PARAMETER THE OBJECT WORKS IN
THE SAME WAY AS GRIDSEARCHCV EXCEPT THAT IT DEFAULTS TO GENERALIZED CROSSVALIDATION GCV AN EFFICIENT FORM OF
LEAVEONEOUT CROSSVALIDATION
IMPORT NUMPY AS NP
FROM SKLEARN IMPORT LINEARMODEL
REG LINEARMODELRIDGECVALPHASNPLOGSPACE6 6 13
REGFIT0 0 0 0 1 1 0 1 1
RIDGECVALPHASARRAY1E06 1E05 1E04 1E03 1E02 1E01 1E00 1E01
1E02 1E03 1E04 1E05 1E06
CVNONE FITINTERCEPTTRUE GCVMODENONE NORMALIZEFALSE
SCORINGNONE STORECVVALUESFALSE
REGALPHA
001
SPECIFYING THE VALUE OF THE CVATTRIBUTE WILL TRIGGER THE USE OF CROSSVALIDATION WITH GRIDSEARCHCV FOR EXAMPLE
CV10 FOR 10FOLD CROSSVALIDATION RATHER THAN GENERALIZED CROSSVALIDATION
REFERENCES
• “NOTES ON REGULARIZED LEAST SQUARES” RIFKIN LIPPERT TECHNICAL REPORT COURSE SLIDES
LASSO
THELASSO IS A LINEAR MODEL THAT ESTIMATES SPARSE COEFFICIENTS IT IS USEFUL IN SOME CONTEXTS DUE TO ITS TENDENCY TO
PREFER SOLUTIONS WITH FEWER NONZERO COEFFICIENTS EFFECTIVELY REDUCING THE NUMBER OF FEATURES UPON WHICH THE GIVEN
SOLUTION IS DEPENDENT FOR THIS REASON LASSO AND ITS VARIANTS ARE FUNDAMENTAL TO THE FIELD OF COMPRESSED SENSING
31 SUPERVISED LEARNING 199
```

SCIKITLEARN USER GUIDE RELEASE 0213

UNDER CERTAIN CONDITIONS IT CAN RECOVER THE EXACT SET OF NONZERO COEFFICIENTS SEE COMPRESSIVE SENSING TOMOGRAPHY RECONSTRUCTION WITH L1 PRIOR LASSO

MATHEMATICALLY IT CONSISTS OF A LINEAR MODEL WITH AN ADDED REGULARIZATION TERM THE OBJECTIVE FUNCTION TO MINIMIZE IS

$$\min_{\beta} \frac{1}{2} \sum_{i=1}^n \text{SAMPLES}_i - \frac{\lambda}{2} \|\beta\|_1$$

THE LASSO ESTIMATE THUS SOLVES THE MINIMIZATION OF THE LEASTSQUARES PENALTY WITH  $\lambda$  ADDED WHERE  $\lambda$  IS A CONSTANT AND  $\|\cdot\|_1$  IS THE  $\ell_1$  NORM OF THE COEFFICIENT VECTOR

THE IMPLEMENTATION IN THE CLASS LASSO USES COORDINATE DESCENT AS THE ALGORITHM TO FIT THE COEFFICIENTS SEE LEAST ANGLE REGRESSION FOR ANOTHER IMPLEMENTATION

```
FROM SKLEARN IMPORT LINEARMODEL
REG LINEARMODELLASSOALPHA01
REGFIT0 0 1 1 0 1
LASSOALPHA01 COPYXTRUE FITINTERCEPTTRUE MAXITER1000
NORMALIZEFALSE POSITIVEFALSE PRECOMPUTEFALSE RANDOMSTATENONE
SELECTIONCYCLIC TOL00001 WARMSTARTFALSE
REGPREDICT1 1
ARRAY08
```

THE FUNCTION LASSOPATH IS USEFUL FOR LOWERLEVEL TASKS AS IT COMPUTES THE COEFFICIENTS ALONG THE FULL PATH OF POSSIBLE VALUES

EXAMPLES

- LASSO AND ELASTIC NET FOR SPARSE SIGNALS
- COMPRESSIVE SENSING TOMOGRAPHY RECONSTRUCTION WITH L1 PRIOR LASSO

NOTE FEATURE SELECTION WITH LASSO

AS THE LASSO REGRESSION YIELDS SPARSE MODELS IT CAN THUS BE USED TO PERFORM FEATURE SELECTION AS DETAILED IN L1BASED FEATURE SELECTION

THE FOLLOWING TWO REFERENCES EXPLAIN THE ITERATIONS USED IN THE COORDINATE DESCENT SOLVER OF SCIKITLEARN AS WELL AS THE DUALITY GAP COMPUTATION USED FOR CONVERGENCE CONTROL

REFERENCES

- “REGULARIZATION PATH FOR GENERALIZED LINEAR MODELS BY COORDINATE DESCENT” FRIEDMAN HASTIE TIBSHIRANI J STAT SOFTW 2010 PAPER
- “AN INTERIORPOINT METHOD FOR LARGESCALE L1REGULARIZED LEAST SQUARES” S J KIM K KOH M LUSTIG S BOYD AND D GORINEVSKY IN IEEE JOURNAL OF SELECTED TOPICS IN SIGNAL PROCESSING 2007 PAPER

SETTING REGULARIZATION PARAMETER

THEALPHA PARAMETER CONTROLS THE DEGREE OF SPARSITY OF THE ESTIMATED COEFFICIENTS

200 CHAPTER 3 USER GUIDE

USING CROSSVALIDATION

SCIKITLEARN EXPOSES OBJECTS THAT SET THE LASSO ALPHA PARAMETER BY CROSSVALIDATION LASSOCV ANDLASSOLARSCV  
LASSOLARSCV IS BASED ON THE LEAST ANGLE REGRESSION ALGORITHM EXPLAINED BELOW

FOR HIGHDIMENSIONAL DATASETS WITH MANY COLLINEAR FEATURES LASSOCV IS MOST OFTEN PREFERABLE HOWEVER  
LASSOLARSCV HAS THE ADVANTAGE OF EXPLORING MORE RELEVANT VALUES OF ALPHA PARAMETER AND IF THE NUMBER OF  
SAMPLES IS VERY SMALL COMPARED TO THE NUMBER OF FEATURES IT IS OFTEN FASTER THAN LASSOCV

INFORMATIONCRITERIA BASED MODEL SELECTION

ALTERNATIVELY THE ESTIMATOR LASSOLARSIC PROPOSES TO USE THE AKAIKE INFORMATION CRITERION AIC AND THE BAYES  
INFORMATION CRITERION BIC IT IS A COMPUTATIONALLY CHEAPER ALTERNATIVE TO FIND THE OPTIMAL VALUE OF ALPHA AS THE REGU  
LARIZATION PATH IS COMPUTED ONLY ONCE INSTEAD OF K1 TIMES WHEN USING KFOLD CROSSVALIDATION HOWEVER SUCH CRITERIA  
NEEDS A PROPER ESTIMATION OF THE DEGREES OF FREEDOM OF THE SOLUTION ARE DERIVED FOR LARGE SAMPLES ASYMPTOTIC RESULTS  
AND ASSUME THE MODEL IS CORRECT IE THAT THE DATA ARE ACTUALLY GENERATED BY THIS MODEL THEY ALSO TEND TO BREAK WHEN  
THE PROBLEM IS BADLY CONDITIONED MORE FEATURES THAN SAMPLES

EXAMPLES

•LASSO MODEL SELECTION CROSSVALIDATION AIC BIC

SCIKITLEARN USER GUIDE RELEASE 0213

COMPARISON WITH THE REGULARIZATION PARAMETER OF SVM

THE EQUIVALENCE BETWEEN ALPHA AND THE REGULARIZATION PARAMETER OF SVM IS GIVEN BY  $\alpha = 1/C$  OR  $\alpha = 1/(n \cdot C)$  DEPENDING ON THE ESTIMATOR AND THE EXACT OBJECTIVE FUNCTION OPTIMIZED BY THE MODEL

MULTITASK LASSO

THE MULTITASK LASSO IS A LINEAR MODEL THAT ESTIMATES SPARSE COEFFICIENTS FOR MULTIPLE REGRESSION PROBLEMS JOINTLY.  $Y$  IS A 2D ARRAY OF SHAPE  $(n \text{ samples}, n \text{ tasks})$ . THE CONSTRAINT IS THAT THE SELECTED FEATURES ARE THE SAME FOR ALL THE REGRESSION PROBLEMS, ALSO CALLED TASKS.

THE FOLLOWING FIGURE COMPARES THE LOCATION OF THE NONZERO ENTRIES IN THE COEFFICIENT MATRIX  $W$  OBTAINED WITH A SIMPLE LASSO OR A MULTITASK LASSO. THE LASSO ESTIMATES YIELD SCATTERED NONZEROS, WHILE THE NONZEROS OF THE MULTITASK LASSO ARE FULL COLUMNS.

FITTING A TIMESERIES MODEL IMPOSING THAT ANY ACTIVE FEATURE BE ACTIVE AT ALL TIMES

EXAMPLES

- JOINT FEATURE SELECTION WITH MULTITASK LASSO

MATHEMATICALLY, IT CONSISTS OF A LINEAR MODEL TRAINED WITH A MIXED  $\ell_1/\ell_2$  NORM FOR REGULARIZATION. THE OBJECTIVE FUNCTION

202 CHAPTER 3 USER GUIDE



TO MINIMIZE IS

MIN

$\frac{1}{2}$

$\frac{1}{2} \sum_{i=1}^n \text{SAMPLES}_{ii} - \frac{1}{2}$

$\text{FRO}^2$

WHERE FRO INDICATES THE FROBENIUS NORM

$\text{FRO} = \sqrt{\sum_{i,j} \Sigma_{ij}^2}$

$\frac{1}{2} \sum_{i,j} \Sigma_{ij}^2$

$\frac{1}{2}$

AND  $\frac{1}{2}$  READS

$\frac{1}{2} \sum_{i,j} \Sigma_{ij}^2$

$\frac{1}{2} \sum_{i,j} \Sigma_{ij}^2$

$\frac{1}{2}$

THE IMPLEMENTATION IN THE CLASS MULTITASKLASSO USES COORDINATE DESCENT AS THE ALGORITHM TO FIT THE COEFFICIENTS ELASTICNET

ELASTICNET IS A LINEAR REGRESSION MODEL TRAINED WITH BOTH  $\ell_1$  AND  $\ell_2$  NORM REGULARIZATION OF THE COEFFICIENTS THIS COMBINATION ALLOWS FOR LEARNING A SPARSE MODEL WHERE FEW OF THE WEIGHTS ARE NONZERO LIKE LASSO WHILE STILL MAINTAINING THE REGULARIZATION PROPERTIES OF RIDGE WE CONTROL THE CONVEX COMBINATION OF  $\ell_1$  AND  $\ell_2$  USING THE  $\ell_1$  RATIO PARAMETER

ELASTICNET IS USEFUL WHEN THERE ARE MULTIPLE FEATURES WHICH ARE CORRELATED WITH ONE ANOTHER LASSO IS LIKELY TO PICK ONE OF THESE AT RANDOM WHILE ELASTICNET IS LIKELY TO PICK BOTH A PRACTICAL ADVANTAGE OF TRADING OFF BETWEEN LASSO AND RIDGE IS THAT IT ALLOWS ELASTICNET TO INHERIT SOME OF RIDGE'S STABILITY UNDER ROTATION

THE OBJECTIVE FUNCTION TO MINIMIZE IS IN THIS CASE

MIN

$\frac{1}{2}$

$\frac{1}{2} \sum_{i=1}^n \text{SAMPLES}_{ii} - \frac{1}{2}$

$\frac{1}{2} \sum_{i,j} \Sigma_{ij}^2$

$\frac{1}{2}$

THE CLASS ELASTICNETCV CAN BE USED TO SET THE PARAMETERS ALPHA AND  $\ell_1$  RATIO BY CROSSVALIDATION

EXAMPLES

- LASSO AND ELASTIC NET FOR SPARSE SIGNALS
- LASSO AND ELASTIC NET

THE FOLLOWING TWO REFERENCES EXPLAIN THE ITERATIONS USED IN THE COORDINATE DESCENT SOLVER OF SCIKITLEARN AS WELL AS THE DUALITY GAP COMPUTATION USED FOR CONVERGENCE CONTROL

REFERENCES

- “REGULARIZATION PATH FOR GENERALIZED LINEAR MODELS BY COORDINATE DESCENT” FRIEDMAN HASTIE TIBSHIRANI J STAT SOFTW 2010 PAPER
- “AN INTERIORPOINT METHOD FOR LARGESCALE L1REGULARIZED LEAST SQUARES” S J KIM K KOH M LUSTIG S BOYD AND D GORINEVSKY IN IEEE JOURNAL OF SELECTED TOPICS IN SIGNAL PROCESSING 2007 PAPER

MULTITASK ELASTICNET

THEMULTITASKELASTICNET IS AN ELASTICNET MODEL THAT ESTIMATES SPARSE COEFFICIENTS FOR MULTIPLE REGRESSION PROBLEMS JOINTLY YIS A 2D ARRAY OF SHAPE NSAMPLES NTASKS THE CONSTRAINT IS THAT THE SELECTED FEATURES ARE THE SAME FOR ALL THE REGRESSION PROBLEMS ALSO CALLED TASKS MATHEMATICALLY IT CONSISTS OF A LINEAR MODEL TRAINED WITH A MIXED  $\ell_1/\ell_2$ NORM AND  $\ell_2$ NORM FOR REGULARIZATION THE OBJECTIVE FUNCTION TO MINIMIZE IS

$$\min_{\beta} \frac{1}{2} \sum_{i=1}^N \text{SAMPLES}_{i,:}^2 - \sum_{j=1}^N \text{FROM}_{:,j}^2 \beta_j^2$$

THE IMPLEMENTATION IN THE CLASS MULTITASKELASTICNET USES COORDINATE DESCENT AS THE ALGORITHM TO FIT THE COEFFICIENTS

THE CLASSMULTITASKELASTICNETCV CAN BE USED TO SET THE PARAMETERS ALPHA ANDL1RATIO BY CROSS VALIDATION

LEAST ANGLE REGRESSION

LEASTANGLE REGRESSION LARS IS A REGRESSION ALGORITHM FOR HIGHDIMENSIONAL DATA DEVELOPED BY BRADLEY EFRON TREVOR HASTIE IAIN JOHNSTONE AND ROBERT TIBSHIRANI LARS IS SIMILAR TO FORWARD STEPWISE REGRESSION AT EACH STEP IT FINDS THE FEATURE MOST CORRELATED WITH THE TARGET WHEN THERE ARE MULTIPLE FEATURES HAVING EQUAL CORRELATION INSTEAD OF CONTINUING ALONG THE SAME FEATURE IT PROCEEDS IN A DIRECTION EQUIANGULAR BETWEEN THE FEATURES THE ADVANTAGES OF LARS ARE

- IT IS NUMERICALLY EFFICIENT IN CONTEXTS WHERE THE NUMBER OF FEATURES IS SIGNIFICANTLY GREATER THAN THE NUMBER OF SAMPLES
- IT IS COMPUTATIONALLY JUST AS FAST AS FORWARD SELECTION AND HAS THE SAME ORDER OF COMPLEXITY AS ORDINARY LEAST SQUARES
- IT PRODUCES A FULL PIECEWISE LINEAR SOLUTION PATH WHICH IS USEFUL IN CROSSVALIDATION OR SIMILAR ATTEMPTS TO TUNE THE MODEL
- IF TWO FEATURES ARE ALMOST EQUALLY CORRELATED WITH THE TARGET THEN THEIR COEFFICIENTS SHOULD INCREASE AT APPROXIMATELY THE SAME RATE THE ALGORITHM THUS BEHAVES AS INTUITION WOULD EXPECT AND ALSO IS MORE STABLE

SCIKITLEARN USER GUIDE RELEASE 0213

- IT IS EASILY MODIFIED TO PRODUCE SOLUTIONS FOR OTHER ESTIMATORS LIKE THE LASSO

THE DISADVANTAGES OF THE LARS METHOD INCLUDE

- BECAUSE LARS IS BASED UPON AN ITERATIVE REFITTING OF THE RESIDUALS IT WOULD APPEAR TO BE ESPECIALLY SENSITIVE TO THE EFFECTS OF NOISE THIS PROBLEM IS DISCUSSED IN DETAIL BY WEISBERG IN THE DISCUSSION SECTION OF THE EFRON ET AL 2004 ANNALS OF STATISTICS ARTICLE

THE LARS MODEL CAN BE USED USING ESTIMATOR LARS OR ITS LOWLEVEL IMPLEMENTATION LARSPATH OR LARSPATHGRAM

LARS LASSO

LASSOLARS IS A LASSO MODEL IMPLEMENTED USING THE LARS ALGORITHM AND UNLIKE THE IMPLEMENTATION BASED ON COORDINATE DESCENT THIS YIELDS THE EXACT SOLUTION WHICH IS PIECEWISE LINEAR AS A FUNCTION OF THE NORM OF ITS COEFFICIENTS

```
FROM SKLEARN IMPORT LINEARMODEL
REG LINEARMODELLASSOLARSALPHA1
REGFIT0 0 1 1 0 1
LASSOLARSALPHA01 COPYXTRUE EPS FITINTERCEPTTRUE
FITPATHTRUE MAXITER500 NORMALIZETRUE POSITIVEFALSE
PRECOMPUTEAUTO VERBOSEFALSE
REGCOEF
ARRAY0717157 0
```

EXAMPLES

- LASSO PATH USING LARS

THE LARS ALGORITHM PROVIDES THE FULL PATH OF THE COEFFICIENTS ALONG THE REGULARIZATION PARAMETER ALMOST FOR FREE THUS A COMMON OPERATION IS TO RETRIEVE THE PATH WITH ONE OF THE FUNCTIONS LARSPATH ORLARSPATHGRAM

MATHEMATICAL FORMULATION

THE ALGORITHM IS SIMILAR TO FORWARD STEPWISE REGRESSION BUT INSTEAD OF INCLUDING FEATURES AT EACH STEP THE ESTIMATED COEFFICIENTS ARE INCREASED IN A DIRECTION EQUIANGULAR TO EACH ONE’S CORRELATIONS WITH THE RESIDUAL

31 SUPERVISED LEARNING 205

INSTEAD OF GIVING A VECTOR RESULT THE LARS SOLUTION CONSISTS OF A CURVE DENOTING THE SOLUTION FOR EACH VALUE OF THE  $\ell_1$  NORM OF THE PARAMETER VECTOR THE FULL COEFFICIENTS PATH IS STORED IN THE ARRAY COEFPATH WHICH HAS SIZE NFEATURES MAXFEATURES1 THE FIRST COLUMN IS ALWAYS ZERO

REFERENCES

- ORIGINAL ALGORITHM IS DETAILED IN THE PAPER LEAST ANGLE REGRESSION BY HASTIE ET AL

ORTHOGONAL MATCHING PURSUIT OMP

ORTHOGONALMATCHINGPURSUIT ANDORTHOGONALMP IMPLEMENTS THE OMP ALGORITHM FOR APPROXIMATING THE FIT OF A LINEAR MODEL WITH CONSTRAINTS IMPOSED ON THE NUMBER OF NONZERO COEFFICIENTS IE THE  $\ell_0$  PSEUDONORM

BEING A FORWARD FEATURE SELECTION METHOD LIKE LEAST ANGLE REGRESSION ORTHOGONAL MATCHING PURSUIT CAN APPROXIMATE THE OPTIMUM SOLUTION VECTOR WITH A FIXED NUMBER OF NONZERO ELEMENTS

ARG MIN

$$\|w\|_0 - \|w\|_2^2$$

SUBJECT TO  $0 \leq \text{NONZEROCOEF}$

ALTERNATIVELY ORTHOGONAL MATCHING PURSUIT CAN TARGET A SPECIFIC ERROR INSTEAD OF A SPECIFIC NUMBER OF NONZERO COEFFICIENTS THIS CAN BE EXPRESSED AS

ARG MIN

$$\|w\|_0 \text{ SUBJECT TO } \|w\|_2^2$$

$$2 \leq \text{TOL}$$

OMP IS BASED ON A GREEDY ALGORITHM THAT INCLUDES AT EACH STEP THE ATOM MOST HIGHLY CORRELATED WITH THE CURRENT RESIDUAL IT IS SIMILAR TO THE SIMPLER MATCHING PURSUIT MP METHOD BUT BETTER IN THAT AT EACH ITERATION THE RESIDUAL IS RECOMPUTED USING AN ORTHOGONAL PROJECTION ON THE SPACE OF THE PREVIOUSLY CHOSEN DICTIONARY ELEMENTS

EXAMPLES

- ORTHOGONAL MATCHING PURSUIT

REFERENCES

- [HTTPSWWWCSSTECHNIONACILRONRUBINPUBLICATIONSKSVDOMPV2PDF](https://www.cs.technion.ac.il/~ronrubin/publications/sk_svd_omp_v2.pdf)
- MATCHING PURSUITS WITH TIMEFREQUENCY DICTIONARIES S G MALLAT Z ZHANG

BAYESIAN REGRESSION

BAYESIAN REGRESSION TECHNIQUES CAN BE USED TO INCLUDE REGULARIZATION PARAMETERS IN THE ESTIMATION PROCEDURE THE REGULARIZATION PARAMETER IS NOT SET IN A HARD SENSE BUT TUNED TO THE DATA AT HAND

THIS CAN BE DONE BY INTRODUCING UNINFORMATIVE PRIORS OVER THE HYPER PARAMETERS OF THE MODEL THE  $\ell_2$  REGULARIZATION USED IN RIDGE REGRESSION IS EQUIVALENT TO FINDING A MAXIMUM A POSTERIORI ESTIMATION UNDER A GAUSSIAN PRIOR OVER THE COEFFICIENTS WITH PRECISION  $\lambda^{-1}$  INSTEAD OF SETTING LAMBDA MANUALLY IT IS POSSIBLE TO TREAT IT AS A RANDOM VARIABLE TO BE ESTIMATED FROM THE DATA TO OBTAIN A FULLY PROBABILISTIC MODEL THE OUTPUT  $\hat{w}$  IS ASSUMED TO BE GAUSSIAN DISTRIBUTED AROUND  $\hat{w}$

SCIKITLEARN USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

WHERE  $\beta$  IS AGAIN TREATED AS A RANDOM VARIABLE THAT IS TO BE ESTIMATED FROM THE DATA  
THE ADVANTAGES OF BAYESIAN REGRESSION ARE

- IT ADAPTS TO THE DATA AT HAND
- IT CAN BE USED TO INCLUDE REGULARIZATION PARAMETERS IN THE ESTIMATION PROCEDURE

THE DISADVANTAGES OF BAYESIAN REGRESSION INCLUDE

- INFERENCE OF THE MODEL CAN BE TIME CONSUMING

REFERENCES

- A GOOD INTRODUCTION TO BAYESIAN METHODS IS GIVEN IN C BISHOP PATTERN RECOGNITION AND MACHINE LEARNING
- ORIGINAL ALGORITHM IS DETAILED IN THE BOOK BAYESIAN LEARNING FOR NEURAL NETWORKS BY RAD FORD M NEAL

BAYESIAN RIDGE REGRESSION

BAYESIANRIDGE ESTIMATES A PROBABILISTIC MODEL OF THE REGRESSION PROBLEM AS DESCRIBED ABOVE THE PRIOR FOR THE COEFFICIENT  $\beta$  IS GIVEN BY A SPHERICAL GAUSSIAN

$\beta \sim \mathcal{N}(0, \lambda^{-1} I)$

THE PRIORS OVER  $\beta$  AND  $\lambda$  ARE CHOSEN TO BE GAMMA DISTRIBUTIONS THE CONJUGATE PRIOR FOR THE PRECISION OF THE GAUSSIAN THE RESULTING MODEL IS CALLED BAYESIAN RIDGE REGRESSION AND IS SIMILAR TO THE CLASSICAL RIDGE

THE PARAMETERS  $\alpha$  AND  $\lambda$  ARE ESTIMATED JOINTLY DURING THE FIT OF THE MODEL THE REGULARIZATION PARAMETERS  $\alpha$  AND  $\lambda$  BEING ESTIMATED BY MAXIMIZING THE LOG MARGINAL LIKELIHOOD THE SCIKITLEARN IMPLEMENTATION IS BASED ON THE ALGORITHM DESCRIBED IN APPENDIX A OF TIPPING 2001 WHERE THE UPDATE OF THE PARAMETERS  $\alpha$  AND  $\lambda$  IS DONE AS SUGGESTED IN MACKAY 1992

THE REMAINING HYPERPARAMETERS ARE THE PARAMETERS  $\alpha_1$  AND  $\alpha_2$  OF THE GAMMA PRIORS OVER  $\alpha$  AND  $\lambda$  THESE ARE USUALLY CHOSEN TO BE NONINFORMATIVE BY DEFAULT  $\alpha_1$  AND  $\alpha_2$  10–6

BAYESIAN RIDGE REGRESSION IS USED FOR REGRESSION

SCIKITLEARN USER GUIDE RELEASE 0213

FROM SKLEARN IMPORT LINEARMODEL

X 0 0 1 1 2 2 3 3

Y 0 1 2 3

REG LINEARMODELBAYESIANRIDGE

REGFITX Y

BAYESIANRIDGEALPHA11E06 ALPHA21E06 COMPUTESCOREFALSE COPYXTRUE

FITINTERCEPTTRUE LAMBDA11E06 LAMBDA21E06 NITER300

NORMALIZEFALSE TOL0001 VERBOSEFALSE

AFTER BEING FITTED THE MODEL CAN THEN BE USED TO PREDICT NEW VALUES

REGPREDICT1 0

ARRAY050000013

THE COEFFICIENTS OF THE MODEL CAN BE ACCESSED

REGCOEF

ARRAY049999993 049999993

DUE TO THE BAYESIAN FRAMEWORK THE WEIGHTS FOUND ARE SLIGHTLY DIFFERENT TO THE ONES FOUND BY ORDINARY LEAST SQUARES

HOWEVER BAYESIAN RIDGE REGRESSION IS MORE ROBUST TO ILLPOSED PROBLEMS

EXAMPLES

- BAYESIAN RIDGE REGRESSION

REFERENCES

- SECTION 33 IN CHRISTOPHER M BISHOP PATTERN RECOGNITION AND MACHINE LEARNING 2006
- DAVID J C MACKAY BAYESIAN INTERPOLATION 1992
- MICHAEL E TIPPING SPARSE BAYESIAN LEARNING AND THE RELEVANCE VECTOR MACHINE 2001

AUTOMATIC RELEVANCE DETERMINATION ARD

ARDREGRESSION IS VERY SIMILAR TO BAYESIAN RIDGE REGRESSION BUT CAN LEAD TO SPARSER COEFFICIENTS 12

ARDREGRESSION POSES A DIFFERENT PRIOR OVER BY DROPPING THE ASSUMPTION OF THE GAUSSIAN BEING SPHERICAL

INSTEAD THE DISTRIBUTION OVER IS ASSUMED TO BE AN AXISPARALLEL ELLIPTICAL GAUSSIAN DISTRIBUTION

THIS MEANS EACH COEFFICIENT IS DRAWN FROM A GAUSSIAN DISTRIBUTION CENTERED ON ZERO AND WITH A PRECISION

0 -1

WITH DIAG 1

IN CONTRAST TO BAYESIAN RIDGE REGRESSION EACH COORDINATE OF HAS ITS OWN STANDARD DEVIATION THE PRIOR OVER ALL

IS CHOSEN TO BE THE SAME GAMMA DISTRIBUTION GIVEN BY HYPERPARAMETERS 1AND2

ARD IS ALSO KNOWN IN THE LITERATURE AS SPARSE BAYESIAN LEARNING ANDRELEVANCE VECTOR MACHINE34

1CHRISTOPHER M BISHOP PATTERN RECOGNITION AND MACHINE LEARNING CHAPTER 721

2DAVID WIPF AND SRIKANTAN NAGARAJAN A NEW VIEW OF AUTOMATIC RELEVANCE DETERMINATION

3MICHAEL E TIPPING SPARSE BAYESIAN LEARNING AND THE RELEVANCE VECTOR MACHINE

4TRISTAN FLETCHER RELEVANCE VECTOR MACHINES EXPLAINED

208 CHAPTER 3 USER GUIDE

EXAMPLES

•AUTOMATIC RELEVANCE DETERMINATION REGRESSION AND

REFERENCES

LOGISTIC REGRESSION

LOGISTIC REGRESSION DESPITE ITS NAME IS A LINEAR MODEL FOR CLASSIFICATION RATHER THAN REGRESSION LOGISTIC REGRESSION IS ALSO KNOWN IN THE LITERATURE AS LOGIT REGRESSION MAXIMUMENTROPY CLASSIFICATION MAXENT OR THE LOGLINEAR CLASSIFIER IN THIS MODEL THE PROBABILITIES DESCRIBING THE POSSIBLE OUTCOMES OF A SINGLE TRIAL ARE MODELED USING A LOGISTIC FUNCTION LOGISTIC REGRESSION IS IMPLEMENTED IN LOGISTICREGRESSION THIS IMPLEMENTATION CAN FIT BINARY ONEVSREST OR MULTINOMIAL LOGISTIC REGRESSION WITH OPTIONAL  $\ell_1/\ell_2$ OR ELASTICNET REGULARIZATION NOTE REGULARIZATION IS APPLIED BY DEFAULT WHICH IS COMMON IN MACHINE LEARNING BUT NOT IN STATISTICS ANOTHER ADVANTAGE OF REGULARIZATION IS THAT IT IMPROVES NUMERICAL STABILITY NO REGULARIZATION AMOUNTS TO SETTING C TO A VERY HIGH VALUE

AS AN OPTIMIZATION PROBLEM BINARY CLASS  $\ell_2$ PENALIZED LOGISTIC REGRESSION MINIMIZES THE FOLLOWING COST FUNCTION

$$\begin{aligned} \min_{\mathbf{w}} & \frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_{i=1}^n \log \exp(-\mathbf{w}^T \mathbf{x}_i) \\ & \text{subject to } \mathbf{w}^T \mathbf{x}_i \geq 1 \end{aligned}$$

SIMILARLY  $\ell_1$ REGULARIZED LOGISTIC REGRESSION SOLVES THE FOLLOWING OPTIMIZATION PROBLEM

$$\begin{aligned} \min_{\mathbf{w}} & \frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_{i=1}^n \log \exp(-\mathbf{w}^T \mathbf{x}_i) \\ & \text{subject to } \mathbf{w}^T \mathbf{x}_i \geq 1 \end{aligned}$$

ELASTICNET REGULARIZATION IS A COMBINATION OF  $\ell_1$ AND $\ell_2$  AND MINIMIZES THE FOLLOWING COST FUNCTION

$$\begin{aligned} \min_{\mathbf{w}} & \frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_{i=1}^n \log \exp(-\mathbf{w}^T \mathbf{x}_i) \\ & \text{subject to } \mathbf{w}^T \mathbf{x}_i \geq 1 \end{aligned}$$

SCIKITLEARN USER GUIDE RELEASE 0213

WHERE  $\alpha$  CONTROLS THE STRENGTH OF  $\ell_1$  REGULARIZATION VS  $\ell_2$  REGULARIZATION IT CORRESPONDS TO THE  $L1RATIO$  PARAMETER  
NOTE THAT IN THIS NOTATION IT'S ASSUMED THAT THE TARGET  $y$  TAKES VALUES IN THE SET  $\{-1,1\}$  AT TRIAL  $i$  WE CAN ALSO SEE THAT  
ELASTICNET IS EQUIVALENT TO  $\ell_1$  WHEN  $\alpha = 1$  AND EQUIVALENT TO  $\ell_2$  WHEN  $\alpha = 0$   
THE SOLVERS IMPLEMENTED IN THE CLASS `LOGISTICREGRESSION` ARE "LIBLINEAR" "NEWTONCG" "LBFGS" "SAG" AND  
"SAGA"

THE SOLVER "LIBLINEAR" USES A COORDINATE DESCENT CD ALGORITHM AND RELIES ON THE EXCELLENT C LIBLINEAR LIBRARY  
WHICH IS SHIPPED WITH SCIKITLEARN HOWEVER THE CD ALGORITHM IMPLEMENTED IN LIBLINEAR CANNOT LEARN A TRUE MULTI  
NOMIAL MULTICLASS MODEL INSTEAD THE OPTIMIZATION PROBLEM IS DECOMPOSED IN A "ONEVSREST" FASHION SO SEPARATE  
BINARY CLASSIFIERS ARE TRAINED FOR ALL CLASSES THIS HAPPENS UNDER THE HOOD SO LOGISTICREGRESSION INSTANCES US  
ING THIS SOLVER BEHAVE AS MULTICLASS CLASSIFIERS FOR  $\ell_1$  REGULARIZATION SKLEARN SVML1 MINC ALLOWS TO CALCULATE  
THE LOWER BOUND FOR C IN ORDER TO GET A NON "NULL" ALL FEATURE WEIGHTS TO ZERO MODEL  
THE "LBFGS" "SAG" AND "NEWTONCG" SOLVERS ONLY SUPPORT  $\ell_2$  REGULARIZATION OR NO REGULARIZATION AND ARE FOUND TO  
CONVERGE FASTER FOR SOME HIGH DIMENSIONAL DATA SETTING MULTICLASS TO "MULTINOMIAL" WITH THESE SOLVERS LEARNS  
A TRUE MULTINOMIAL LOGISTIC REGRESSION MODEL<sup>5</sup> WHICH MEANS THAT ITS PROBABILITY ESTIMATES SHOULD BE BETTER CALIBRATED  
THAN THE DEFAULT "ONEVSREST" SETTING

THE "SAG" SOLVER USES STOCHASTIC AVERAGE GRADIENT DESCENT<sup>6</sup> IT IS FASTER THAN OTHER SOLVERS FOR LARGE DATASETS WHEN  
BOTH THE NUMBER OF SAMPLES AND THE NUMBER OF FEATURES ARE LARGE  
THE "SAGA" SOLVER<sup>7</sup> IS A VARIANT OF "SAG" THAT ALSO SUPPORTS THE NONSMOOTH PENALTY  $L1$  THIS IS THERE  
FORE THE SOLVER OF CHOICE FOR SPARSE MULTINOMIAL LOGISTIC REGRESSION IT IS ALSO THE ONLY SOLVER THAT SUPPORTS  
PENALTY ELASTICNET

THE "LBFGS" IS AN OPTIMIZATION ALGORITHM THAT APPROXIMATES THE BROYDEN-FLETCHER-GOLDFARB-SHANNO ALGORITHM<sup>8</sup>  
WHICH BELONGS TO QUASINEWTON METHODS THE "LBFGS" SOLVER IS RECOMMENDED FOR USE FOR SMALL DATASETS BUT FOR  
LARGER DATASETS ITS PERFORMANCE SUFFERS<sup>9</sup>  
THE FOLLOWING TABLE SUMMARIZES THE PENALTIES SUPPORTED BY EACH SOLVER

SOLVERS	PENALTIES	'LIBLINEAR'	'LBFGS'	'NEWTONCG'	'SAG'	'SAGA'
MULTINOMIAL $L2$ PENALTY	NO	YES	YES	YES	YES	YES
OVR $L2$ PENALTY	YES	YES	YES	YES	YES	YES
MULTINOMIAL $L1$ PENALTY	NO	NO	NO	NO	NO	YES
OVR $L1$ PENALTY	YES	NO	NO	NO	NO	YES
ELASTICNET	NO	NO	NO	NO	YES	
NO PENALTY 'NONE'	NO	YES	YES	YES	YES	YES

BEHAVIORS

PENALIZE THE INTERCEPT	BAD	YES	NO	NO	NO	NO
FASTER FOR LARGE DATASETS	NO	NO	NO	YES	YES	
ROBUST TO UNSCALED DATASETS	YES	YES	YES	NO	NO	

THE "LBFGS" SOLVER IS USED BY DEFAULT FOR ITS ROBUSTNESS FOR LARGE DATASETS THE "SAGA" SOLVER IS USUALLY FASTER FOR LARGE  
DATASET YOU MAY ALSO CONSIDER USING `SGDCLASSIFIER` WITH 'LOG' LOSS WHICH MIGHT BE EVEN FASTER BUT REQUIRES MORE  
TUNING

5CHRISTOPHER M BISHOP PATTERN RECOGNITION AND MACHINE LEARNING CHAPTER 434  
6MARK SCHMIDT NICOLAS LE ROUX AND FRANCIS BACH MINIMIZING FINITE SUMS WITH THE STOCHASTIC AVERAGE GRADIENT  
7AARON DEFAZIO FRANCIS BACH SIMON LACOSTEJULIEN SAGA A FAST INCREMENTAL GRADIENT METHOD WITH SUPPORT FOR NONSTRONGLY CONVEX  
COMPOSITE OBJECTIVES  
8[HTTPS://ENWIKIPEDIA.ORG/WIKI/BROYDENE28093FLETCHERE28093GOLDFARBE28093SHANNOALGORITHM](https://en.wikipedia.org/wiki/Broyden-Fletcher-Goldfarb-Shanno_algorithm)  
9"PERFORMANCE EVALUATION OF LBFGS VS OTHER SOLVERS"  
210 CHAPTER 3 USER GUIDE



SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES

- L1 PENALTY AND SPARSITY IN LOGISTIC REGRESSION
- REGULARIZATION PATH OF L1 LOGISTIC REGRESSION
- PLOT MULTINOMIAL AND ONEVSREST LOGISTIC REGRESSION
- MULTICLASS SPARSE LOGISITIC REGRESSION ON NEWGROUPS20
- MNIST CLASSIFICATION USING MULTINOMIAL LOGISTIC L1

DIFFERENCES FROM LIBLINEAR

THERE MIGHT BE A DIFFERENCE IN THE SCORES OBTAINED BETWEEN LOGISTICREGRESSION WITHSOLVERLIBLINEAR ORLINEARSVC AND THE EXTERNAL LIBLINEAR LIBRARY DIRECTLY WHEN FITINTERCEPTFALSE AND THE FIT COEF OR THE DATA TO BE PREDICTED ARE ZEROES THIS IS BECAUSE FOR THE SAMPLES WITH DECISIONFUNCTION ZERO LOGISTICREGRESSION ANDLINEARSVC PREDICT THE NEGATIVE CLASS WHILE LIBLINEAR PREDICTS THE POSITIVE CLASS NOTE THAT A MODEL WITH FITINTERCEPTFALSE AND HAVING MANY SAMPLES WITH DECISIONFUNCTION ZERO IS LIKELY TO BE A UNDERFIT BAD MODEL AND YOU ARE ADVISED TO SET FITINTERCEPTTRUE AND INCREASE THE INTERCEPTSCALING

NOTE FEATURE SELECTION WITH SPARSE LOGISTIC REGRESSION

A LOGISTIC REGRESSION WITH  $l_1$ PENALTY YIELDS SPARSE MODELS AND CAN THUS BE USED TO PERFORM FEATURE SELECTION AS DETAILED IN L1BASED FEATURE SELECTION

LOGISTICREGRESSIONCV IMPLEMENTS LOGISTIC REGRESSION WITH BUILTIN CROSSVALIDATION SUPPORT TO FIND THE OPTIMALCANDL1RATIO PARAMETERS ACCORDING TO THE SCORING ATTRIBUTE THE “NEWTONCG” “SAG” “SAGA” AND “LBFGS” SOLVERS ARE FOUND TO BE FASTER FOR HIGHDIMENSIONAL DENSE DATA DUE TO WARMSTARTING SEE GLOSSARY

REFERENCES

STOCHASTIC GRADIENT DESCENT SGD

STOCHASTIC GRADIENT DESCENT IS A SIMPLE YET VERY EFFICIENT APPROACH TO FIT LINEAR MODELS IT IS PARTICULARLY USEFUL WHEN THE NUMBER OF SAMPLES AND THE NUMBER OF FEATURES IS VERY LARGE THE PARTIALFIT METHOD ALLOWS ONLINEOUTOF CORE LEARNING

THE CLASSES SGDCLASSIFIER ANDSGDREGRESSOR PROVIDE FUNCTIONALITY TO FIT LINEAR MODELS FOR CLASSIFICATION AND REGRESSION USING DIFFERENT CONVEX LOSS FUNCTIONS AND DIFFERENT PENALTIES EG WITH LOSSLOG SGDCLASSIFIER FITS A LOGISTIC REGRESSION MODEL WHILE WITH LOSSHINGE IT FITS A LINEAR SUPPORT VECTOR MACHINE SVM

REFERENCES

- STOCHASTIC GRADIENT DESCENT

PERCEPTRON

THEPERCEPTRON IS ANOTHER SIMPLE CLASSIFICATION ALGORITHM SUITABLE FOR LARGE SCALE LEARNING BY DEFAULT

- IT DOES NOT REQUIRE A LEARNING RATE
- IT IS NOT REGULARIZED PENALIZED
- IT UPDATES ITS MODEL ONLY ON MISTAKES

THE LAST CHARACTERISTIC IMPLIES THAT THE PERCEPTRON IS SLIGHTLY FASTER TO TRAIN THAN SGD WITH THE HINGE LOSS AND THAT THE RESULTING MODELS ARE SPARSER

PASSIVE AGGRESSIVE ALGORITHMS

THE PASSIVEAGGRESSIVE ALGORITHMS ARE A FAMILY OF ALGORITHMS FOR LARGESCALE LEARNING THEY ARE SIMILAR TO THE PERCEPTRON IN THAT THEY DO NOT REQUIRE A LEARNING RATE HOWEVER CONTRARY TO THE PERCEPTRON THEY INCLUDE A REGULARIZATION PARAMETERC

FOR CLASSIFICATION PASSIVEAGGRESSIVECLASSIFIER CAN BE USED WITH LOSSHINGE PAI OR LOSSSQUAREDHINGE PAII FOR REGRESSION PASSIVEAGGRESSIVEREGRESSOR CAN BE USED WITH LOSSEPSILONINSENSITIVE PAI ORLOSSSQUAREDEPSILONINSENSITIVE PAII

REFERENCES

- “ONLINE PASSIVEAGGRESSIVE ALGORITHMS” K CRAMMER O DEKEL J KESHAT S SHALEVSHWARTZ Y SINGER

JMLR 7 2006

ROBUSTNESS REGRESSION OUTLIERS AND MODELING ERRORS

ROBUST REGRESSION AIMS TO FIT A REGRESSION MODEL IN THE PRESENCE OF CORRUPT DATA EITHER OUTLIERS OR ERROR IN THE MODEL DIFFERENT SCENARIO AND USEFUL CONCEPTS

THERE ARE DIFFERENT THINGS TO KEEP IN MIND WHEN DEALING WITH DATA CORRUPTED BY OUTLIERS

SCIKITLEARN USER GUIDE RELEASE 0213

- OUTLIERS IN X OR IN Y

OUTLIERS IN THE Y DIRECTION OUTLIERS IN THE X DIRECTION

- FRACTION OF OUTLIERS VERSUS AMPLITUDE OF ERROR

THE NUMBER OF OUTLYING POINTS MATTERS BUT ALSO HOW MUCH THEY ARE OUTLIERS

SMALL OUTLIERS LARGE OUTLIERS

AN IMPORTANT NOTION OF ROBUST FITTING IS THAT OF BREAKDOWN POINT THE FRACTION OF DATA THAT CAN BE OUTLYING FOR THE FIT TO START MISSING THE INLYING DATA

NOTE THAT IN GENERAL ROBUST FITTING IN HIGHDIMENSIONAL SETTING LARGE NFEATURES IS VERY HARD THE ROBUST MODELS HERE WILL PROBABLY NOT WORK IN THESE SETTINGS

TRADEOFFS WHICH ESTIMATOR

SCIKITLEARN PROVIDES 3 ROBUST REGRESSION ESTIMATORS RANSAC THEIL SEN ANDHUBERREGRESSOR

- HUBERREGRESSOR SHOULD BE FASTER THAN RANSAC ANDTHEIL SEN UNLESS THE NUMBER OF SAMPLES ARE VERY LARGE IE NSAMPLES NFEATURES THIS IS BECAUSE RANSAC ANDTHEIL SEN FIT ON SMALLER SUBSETS OF THE DATA HOWEVER BOTH THEIL SEN ANDRANSAC ARE UNLIKELY TO BE AS ROBUST AS HUBERREGRESSOR FOR THE DEFAULT PARAMETERS

31 SUPERVISED LEARNING 213

SCIKITLEARN USER GUIDE RELEASE 0213

- RANSAC IS FASTER THAN THEIL SEN AND SCALES MUCH BETTER WITH THE NUMBER OF SAMPLES
- RANSAC WILL DEAL BETTER WITH LARGE OUTLIERS IN THE Y DIRECTION MOST COMMON SITUATION
- THEIL SEN WILL COPE BETTER WITH MEDIUMSIZE OUTLIERS IN THE X DIRECTION BUT THIS PROPERTY WILL DISAPPEAR IN HIGHDIMENSIONAL SETTINGS

WHEN IN DOUBT USE RANSAC

RANSAC RANDOM SAMPLE CONSENSUS

RANSAC RANDOM SAMPLE CONSENSUS FITS A MODEL FROM RANDOM SUBSETS OF INLIERS FROM THE COMPLETE DATA SET

RANSAC IS A NONDETERMINISTIC ALGORITHM PRODUCING ONLY A REASONABLE RESULT WITH A CERTAIN PROBABILITY WHICH IS DEPENDENT ON THE NUMBER OF ITERATIONS SEE MAXTRIALS PARAMETER IT IS TYPICALLY USED FOR LINEAR AND NONLINEAR REGRESSION PROBLEMS AND IS ESPECIALLY POPULAR IN THE FIELD OF PHOTOGRAMMETRIC COMPUTER VISION

THE ALGORITHM SPLITS THE COMPLETE INPUT SAMPLE DATA INTO A SET OF INLIERS WHICH MAY BE SUBJECT TO NOISE AND OUTLIERS WHICH ARE EG CAUSED BY ERRONEOUS MEASUREMENTS OR INVALID HYPOTHESES ABOUT THE DATA THE RESULTING MODEL IS THEN ESTIMATED ONLY FROM THE DETERMINED INLIERS

DETAILS OF THE ALGORITHM

EACH ITERATION PERFORMS THE FOLLOWING STEPS

1 SELECTMINSAMPLES RANDOM SAMPLES FROM THE ORIGINAL DATA AND CHECK WHETHER THE SET OF DATA IS VALID SEE ISDATAVALID

2 FIT A MODEL TO THE RANDOM SUBSET BASEESTIMATORFIT AND CHECK WHETHER THE ESTIMATED MODEL IS VALID SEEISMODELVALID

3 CLASSIFY ALL DATA AS INLIERS OR OUTLIERS BY CALCULATING THE RESIDUALS TO THE ESTIMATED MODEL BASEESTIMATOR PREDICTX Y ALL DATA SAMPLES WITH ABSOLUTE RESIDUALS SMALLER THAN THE RESIDUALTHRESHOLD ARE CONSIDERED AS INLIERS

4 SAVE FITTED MODEL AS BEST MODEL IF NUMBER OF INLIER SAMPLES IS MAXIMAL IN CASE THE CURRENT ESTIMATED MODEL HAS THE SAME NUMBER OF INLIERS IT IS ONLY CONSIDERED AS THE BEST MODEL IF IT HAS BETTER SCORE

SCIKITLEARN USER GUIDE RELEASE 0213

THESE STEPS ARE PERFORMED EITHER A MAXIMUM NUMBER OF TIMES MAXTRIALS OR UNTIL ONE OF THE SPECIAL STOP CRITERIA ARE MET SEE STOPNINLIERS ANDSTOPSCORE THE FINAL MODEL IS ESTIMATED USING ALL INLIER SAMPLES CONSENSUS SET OF THE PREVIOUSLY DETERMINED BEST MODEL

THEISDATAVALID ANDISMODELVALID FUNCTIONS ALLOW TO IDENTIFY AND REJECT DEGENERATE COMBINATIONS OF RANDOM SUBSAMPLES IF THE ESTIMATED MODEL IS NOT NEEDED FOR IDENTIFYING DEGENERATE CASES ISDATAVALID SHOULD BE USED AS IT IS CALLED PRIOR TO FITTING THE MODEL AND THUS LEADING TO BETTER COMPUTATIONAL PERFORMANCE

- EXAMPLES
- ROBUST LINEAR MODEL ESTIMATION USING RANSAC
  - ROBUST LINEAR ESTIMATOR FITTING

REFERENCES

- HTTPSENWIKIPEDIAORGWIKIRANSAC
- “RANDOM SAMPLE CONSENSUS A PARADIGM FOR MODEL FITTING WITH APPLICATIONS TO IMAGE ANALYSIS AND AUTOMATED CARTOGRAPHY” MARTIN A FISCHLER AND ROBERT C BOLLES SRI INTERNATIONAL 1981
- “PERFORMANCE EVALUATION OF RANSAC FAMILY” SUNGLOK CHOI TAEMIN KIM AND WONPIL YU BMVC 2009

THEILSEN ESTIMATOR GENERALIZEDMEDIANBASED ESTIMATOR

THETHEILSENREGRESSOR ESTIMATOR USES A GENERALIZATION OF THE MEDIAN IN MULTIPLE DIMENSIONS IT IS THUS ROBUST TO MULTIVARIATE OUTLIERS NOTE HOWEVER THAT THE ROBUSTNESS OF THE ESTIMATOR DECREASES QUICKLY WITH THE DIMENSIONALITY OF THE PROBLEM IT LOSES ITS ROBUSTNESS PROPERTIES AND BECOMES NO BETTER THAN AN ORDINARY LEAST SQUARES IN HIGH DIMENSION

- EXAMPLES
- THEILSEN REGRESSION
  - ROBUST LINEAR ESTIMATOR FITTING

REFERENCES

- HTTPSENWIKIPEDIAORGWIKITHEILE28093SENESTIMATOR

THEORETICAL CONSIDERATIONS

THEILSENREGRESSOR IS COMPARABLE TO THE ORDINARY LEAST SQUARES OLS IN TERMS OF ASYMPTOTIC EFFICIENCY AND AS AN UNBIASED ESTIMATOR IN CONTRAST TO OLS THEILSEN IS A NONPARAMETRIC METHOD WHICH MEANS IT MAKES NO ASSUMPTION ABOUT THE UNDERLYING DISTRIBUTION OF THE DATA SINCE THEILSEN IS A MEDIANBASED ESTIMATOR IT IS MORE ROBUST AGAINST CORRUPTED DATA AKA OUTLIERS IN UNIVARIATE SETTING THEILSEN HAS A BREAKDOWN POINT OF ABOUT 293 IN CASE OF A SIMPLE LINEAR REGRESSION WHICH MEANS THAT IT CAN TOLERATE ARBITRARY CORRUPTED DATA OF UP TO 293

THE IMPLEMENTATION OF THEILSENREGRESSOR IN SCIKITLEARN FOLLOWS A GENERALIZATION TO A MULTIVARIATE LINEAR REGRESSION MODEL10USING THE SPATIAL MEDIAN WHICH IS A GENERALIZATION OF THE MEDIAN TO MULTIPLE DIMENSIONS11 IN TERMS OF TIME AND SPACE COMPLEXITY THEILSEN SCALES ACCORDING TO

□SAMPLES  
□SUBSAMPLES

WHICH MAKES IT INFEASIBLE TO BE APPLIED EXHAUSTIVELY TO PROBLEMS WITH A LARGE NUMBER OF SAMPLES AND FEATURES THEREFORE THE MAGNITUDE OF A SUBPOPULATION CAN BE CHOSEN TO LIMIT THE TIME AND SPACE COMPLEXITY BY CONSIDERING ONLY A RANDOM SUBSET OF ALL POSSIBLE COMBINATIONS

EXAMPLES  
•THEILSEN REGRESSION

REFERENCES  
HUBER REGRESSION

THEHUBERREGRESSOR IS DIFFERENT TO RIDGE BECAUSE IT APPLIES A LINEAR LOSS TO SAMPLES THAT ARE CLASSIFIED AS OUTLIERS A SAMPLE IS CLASSIFIED AS AN INLIER IF THE ABSOLUTE ERROR OF THAT SAMPLE IS LESSER THAN A CERTAIN THRESHOLD IT DIFFERS FROM THEILSENREGRESSOR ANDRANSACREGRESSOR BECAUSE IT DOES NOT IGNORE THE EFFECT OF THE OUTLIERS BUT GIVES A LESSER WEIGHT TO THEM

THE LOSS FUNCTION THAT HUBERREGRESSOR MINIMIZES IS GIVEN BY  
MIN

$$\sum_{i=1}^n \min \left( \frac{1}{2} \epsilon^2, \epsilon |r_i| \right)$$

□□22  
10XIN DANG HANXIANG PENG XUEQIN WANG AND HEPING ZHANG THEILSEN ESTIMATORS IN A MULTIPLE LINEAR REGRESSION MODEL  
11  
20 KÄRKKÄINEN AND S ÄYRÄMÖ ON COMPUTATION OF SPATIAL MEDIAN FOR ROBUST DATA MINING  
216 CHAPTER 3 USER GUIDE

WHERE

000

02 IF 0

200-20 OTHERWISE

IT IS ADVISED TO SET THE PARAMETER EPSILON TO 135 TO ACHIEVE 95 STATISTICAL EFFICIENCY

NOTES

THEHUBERREGRESSOR DIFFERS FROM USING SGDREGRESSOR WITH LOSS SET TO HUBER IN THE FOLLOWING WAYS

- HUBERREGRESSOR IS SCALING INVARIANT ONCE EPSILON IS SET SCALING XANDYDOWN OR UP BY DIFFERENT VALUES WOULD PRODUCE THE SAME ROBUSTNESS TO OUTLIERS AS BEFORE AS COMPARED TO SGDREGRESSOR WHEREEPSILON HAS TO BE SET AGAIN WHEN XANDYARE SCALED
- HUBERREGRESSOR SHOULD BE MORE EFFICIENT TO USE ON DATA WITH SMALL NUMBER OF SAMPLES WHILE SGDREGRESSOR NEEDS A NUMBER OF PASSES ON THE TRAINING DATA TO PRODUCE THE SAME ROBUSTNESS

EXAMPLES

- HUBERREGRESSOR VS RIDGE ON DATASET WITH STRONG OUTLIERS

REFERENCES

- PETER J HUBER ELVEZIO M RONCHETTI ROBUST STATISTICS CONCOMITANT SCALE ESTIMATES PG 172

NOTE THAT THIS ESTIMATOR IS DIFFERENT FROM THE R IMPLEMENTATION OF ROBUST REGRESSION HTTPWWWATSUCLAEDUSTATR DAERREGHTM BECAUSE THE R IMPLEMENTATION DOES A WEIGHTED LEAST SQUARES IMPLEMENTATION WITH WEIGHTS GIVEN TO EACH SAMPLE ON THE BASIS OF HOW MUCH THE RESIDUAL IS GREATER THAN A CERTAIN THRESHOLD

POLYNOMIAL REGRESSION EXTENDING LINEAR MODELS WITH BASIS FUNCTIONS

ONE COMMON PATTERN WITHIN MACHINE LEARNING IS TO USE LINEAR MODELS TRAINED ON NONLINEAR FUNCTIONS OF THE DATA THIS APPROACH MAINTAINS THE GENERALLY FAST PERFORMANCE OF LINEAR METHODS WHILE ALLOWING THEM TO FIT A MUCH WIDER RANGE OF DATA

SCIKITLEARN USER GUIDE RELEASE 0213

FOR EXAMPLE A SIMPLE LINEAR REGRESSION CAN BE EXTENDED BY CONSTRUCTING POLYNOMIAL FEATURES FROM THE COEFFICIENTS IN THE STANDARD LINEAR REGRESSION CASE YOU MIGHT HAVE A MODEL THAT LOOKS LIKE THIS FOR TWODIMENSIONAL DATA

$$y = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

IF WE WANT TO FIT A PARABOLOID TO THE DATA INSTEAD OF A PLANE WE CAN COMBINE THE FEATURES IN SECONDORDER POLYNOMIALS SO THAT THE MODEL LOOKS LIKE THIS

$$y = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2 + \theta_5 x_1 x_2$$

THE SOMETIMES SURPRISING OBSERVATION IS THAT THIS IS STILL A LINEAR MODEL TO SEE THIS IMAGINE CREATING A NEW SET OF FEATURES

$$z_1 = x_1^2, z_2 = x_1 x_2, z_3 = x_2^2$$

WITH THIS RELABELING OF THE DATA OUR PROBLEM CAN BE WRITTEN

$$y = \theta_0 + \theta_1 z_1 + \theta_2 z_2 + \theta_3 z_3 + \theta_4 x_1 + \theta_5 x_2$$

WE SEE THAT THE RESULTING POLYNOMIAL REGRESSION IS IN THE SAME CLASS OF LINEAR MODELS WE CONSIDERED ABOVE IE THE MODEL IS LINEAR IN  $z$  AND CAN BE SOLVED BY THE SAME TECHNIQUES BY CONSIDERING LINEAR FITS WITHIN A HIGHERDIMENSIONAL SPACE BUILT WITH THESE BASIS FUNCTIONS THE MODEL HAS THE FLEXIBILITY TO FIT A MUCH BROADER RANGE OF DATA

HERE IS AN EXAMPLE OF APPLYING THIS IDEA TO ONEDIMENSIONAL DATA USING POLYNOMIAL FEATURES OF VARYING DEGREES THIS FIGURE IS CREATED USING THE POLYNOMIALFEATURES TRANSFORMER WHICH TRANSFORMS AN INPUT DATA MATRIX INTO A NEW DATA MATRIX OF A GIVEN DEGREE IT CAN BE USED AS FOLLOWS

FROM SKLEARNPREPROCESSING IMPORT POLYNOMIALFEATURES

IMPORT NUMPY AS NP

X = NP.RANGE(6).RESHAPE(3, 2)

X

ARRAY0 1

2 3

4 5

POLY = POLYNOMIALFEATURES(DEGREE=2)

POLYFIT=POLY.TRANSFORM(X)

ARRAY1 0 1 0 0 1

1 2 3 4 6 9

1 4 5 16 20 25



SCIKITLEARN USER GUIDE RELEASE 0213

THE FEATURES OF X HAVE BEEN TRANSFORMED FROM  $\begin{bmatrix} x_1 & x_2 \end{bmatrix}$  TO  $\begin{bmatrix} x_1 & x_2 & x_1^2 & x_2^2 \end{bmatrix}$

1  $\begin{bmatrix} x_1 & x_2 \end{bmatrix}$

2 AND CAN NOW BE USED WITHIN ANY LINEAR MODEL

THIS SORT OF PREPROCESSING CAN BE STREAMLINED WITH THE PIPELINE TOOLS A SINGLE OBJECT REPRESENTING A SIMPLE POLYNOMIAL REGRESSION CAN BE CREATED AND USED AS FOLLOWS

```
FROM SKLEARNPREPROCESSING IMPORT POLYNOMIALFEATURES
FROM SKLEARNLINEARMODEL IMPORT LINEARREGRESSION
FROM SKLEARNPIPELINE IMPORT PIPELINE
IMPORT NUMPY AS NP
MODEL PIPELINEPOLY(POLYNOMIALFEATURESDEGREE=3,
LINEAR LINEARREGRESSIONFITINTERCEPT=False)
FIT TO AN ORDER=3 POLYNOMIAL DATA
X = NP.RANGE(5)
Y = 3 - 2 * X + X**2 + X**3
MODEL = MODEL.FIT(X, NP.NEWAXIS(Y))
MODEL.NAMED_STEPS[LINEAR].COEF
ARRAY([3, 2, 1, 1])
```

THE LINEAR MODEL TRAINED ON POLYNOMIAL FEATURES IS ABLE TO EXACTLY RECOVER THE INPUT POLYNOMIAL COEFFICIENTS

IN SOME CASES IT’S NOT NECESSARY TO INCLUDE HIGHER POWERS OF ANY SINGLE FEATURE BUT ONLY THE SO CALLED INTERACTION FEATURES THAT MULTIPLY TOGETHER AT MOST  $\begin{bmatrix} 2 \end{bmatrix}$  DISTINCT FEATURES THESE CAN BE GOTTEN FROM POLYNOMIALFEATURES WITH THE SETTING INTERACTION ONLY=True

FOR EXAMPLE WHEN DEALING WITH BOOLEAN FEATURES  $\begin{bmatrix} x_1 & x_2 \end{bmatrix}$

$\begin{bmatrix} x_1 & x_2 \end{bmatrix}$  FOR ALL  $\begin{bmatrix} x_1 & x_2 \end{bmatrix}$  AND IS THEREFORE USELESS BUT  $\begin{bmatrix} x_1 & x_2 \end{bmatrix}$  REPRESENTS THE CONJUNCTION OF TWO BOOLEANS THIS WAY WE CAN SOLVE THE XOR PROBLEM WITH A LINEAR CLASSIFIER

```
FROM SKLEARNLINEARMODEL IMPORT PERCEPTRON
FROM SKLEARNPREPROCESSING IMPORT POLYNOMIALFEATURES
IMPORT NUMPY AS NP
X = NP.ARRAY([0, 0, 1, 1, 0, 1, 1])
Y = X[0] * X[1]
Y
ARRAY([0, 1, 1, 0])
X = POLYNOMIALFEATURES.INTERACTION_ONLY=True.FIT_TRANSFORM(X, AS_TYPE='int')
X
ARRAY([1, 0, 0, 0, 1, 0, 1, 0, 1, 1, 0, 0, 1, 1, 1, 1])
CLF = PERCEPTRON.FIT_INTERCEPT=False, MAX_ITER=10, TOL=None,
SHUFFLE=False).FIT(X, Y)
AND THE CLASSIFIER “PREDICTIONS” ARE PERFECT
CLF.PREDICT(X)
ARRAY([0, 1, 1, 0])
CLF.SCORE(X, Y)
10
```

312 LINEAR AND QUADRATIC DISCRIMINANT ANALYSIS

LINEAR DISCRIMINANT ANALYSIS DISCRIMINANTANALYSIS LINEAR DISCRIMINANT ANALYSIS AND QUADRATIC DISCRIMINANT ANALYSIS DISCRIMINANTANALYSIS QUADRATIC DISCRIMINANT ANALYSIS ARE TWO CLASSIC CLASSIFIERS WITH AS THEIR NAMES SUGGEST A LINEAR AND A QUADRATIC DECISION SURFACE RESPECTIVELY

31 SUPERVISED LEARNING 219

SCIKITLEARN USER GUIDE RELEASE 0213

THESE CLASSIFIERS ARE ATTRACTIVE BECAUSE THEY HAVE CLOSEDFORM SOLUTIONS THAT CAN BE EASILY COMPUTED ARE INHERENTLY MULTICLASS HAVE PROVEN TO WORK WELL IN PRACTICE AND HAVE NO HYPERPARAMETERS TO TUNE

THE PLOT SHOWS DECISION BOUNDARIES FOR LINEAR DISCRIMINANT ANALYSIS AND QUADRATIC DISCRIMINANT ANALYSIS THE BOTTOM ROW DEMONSTRATES THAT LINEAR DISCRIMINANT ANALYSIS CAN ONLY LEARN LINEAR BOUNDARIES WHILE QUADRATIC DISCRIMINANT ANALYSIS CAN LEARN QUADRATIC BOUNDARIES AND IS THEREFORE MORE FLEXIBLE

EXAMPLES

LINEAR AND QUADRATIC DISCRIMINANT ANALYSIS WITH COVARIANCE ELLIPSOID COMPARISON OF LDA AND QDA ON SYNTHETIC DATA

DIMENSIONALITY REDUCTION USING LINEAR DISCRIMINANT ANALYSIS

DISCRIMINANTANALYSISLINEARDISCRIMINANTANALYSIS CAN BE USED TO PERFORM SUPERVISED DIMENSIONALITY REDUCTION BY PROJECTING THE INPUT DATA TO A LINEAR SUBSPACE CONSISTING OF THE DIRECTIONS WHICH MAXIMIZE THE SEPARATION BETWEEN CLASSES IN A PRECISE SENSE DISCUSSED IN THE MATHEMATICS SECTION BELOW THE DIMENSION OF THE OUTPUT IS NECESSARILY LESS THAN THE NUMBER OF CLASSES SO THIS IS IN GENERAL A RATHER STRONG DIMENSIONALITY REDUCTION AND ONLY MAKES SENSE IN A MULTICLASS SETTING

220 CHAPTER 3 USER GUIDE

THIS IS IMPLEMENTED IN DISCRIMINANTANALYSISLINEARDISCRIMINANTANALYSISITRANSFORM THE DESIRED DIMENSIONALITY CAN BE SET USING THE NCOMPONENTS CONSTRUCTOR PARAMETER THIS PARAMETER HAS NO INFLUENCE ONDISCRIMINANTANALYSISLINEARDISCRIMINANTANALYSISFIT ORDISCRIMINANTANALYSISLINEARDISCRIMINANTANALYSISPREDICT

EXAMPLES  
COMPARISON OF LDA AND PCA 2D PROJECTION OF IRIS DATASET COMPARISON OF LDA AND PCA FOR DIMENSIONALITY REDUCTION OF THE IRIS DATASET  
MATHEMATICAL FORMULATION OF THE LDA AND QDA CLASSIFIERS

BOTH LDA AND QDA CAN BE DERIVED FROM SIMPLE PROBABILISTIC MODELS WHICH MODEL THE CLASS CONDITIONAL DISTRIBUTION OF THE DATA  $p(\mathbf{x} | c)$  FOR EACH CLASS  $c$ . PREDICTIONS CAN THEN BE OBTAINED BY USING BAYES' RULE

$$p(c | \mathbf{x}) \propto p(\mathbf{x} | c) p(c)$$

AND WE SELECT THE CLASS  $c$  WHICH MAXIMIZES THIS CONDITIONAL PROBABILITY  
MORE SPECIFICALLY FOR LINEAR AND QUADRATIC DISCRIMINANT ANALYSIS  $p(\mathbf{x} | c)$  IS MODELED AS A MULTIVARIATE GAUSSIAN DISTRIBUTION WITH DENSITY

$$p(\mathbf{x} | c) = \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1}(\mathbf{x} - \mu)\right)$$

WHERE  $D$  IS THE NUMBER OF FEATURES  
TO USE THIS MODEL AS A CLASSIFIER WE JUST NEED TO ESTIMATE FROM THE TRAINING DATA THE CLASS PRIORS  $p(c)$  BY THE PROPORTION OF INSTANCES OF CLASS  $c$  THE CLASS MEANS  $\mu$  BY THE EMPIRICAL SAMPLE CLASS MEANS AND THE COVARIANCE MATRICES EITHER BY THE EMPIRICAL SAMPLE CLASS COVARIANCE MATRICES OR BY A REGULARIZED ESTIMATOR SEE THE SECTION ON SHRINKAGE BELOW

IN THE CASE OF LDA THE GAUSSIANS FOR EACH CLASS ARE ASSUMED TO SHARE THE SAME COVARIANCE MATRIX  $\Sigma$  FOR ALL  $c$ . THIS LEADS TO LINEAR DECISION SURFACES WHICH CAN BE SEEN BY COMPARING THE LOGPROBABILITY RATIOS  $\log \frac{p(\mathbf{x} | c_1) p(c_1)}{p(\mathbf{x} | c_2) p(c_2)}$

$$\log \frac{p(\mathbf{x} | c_1) p(c_1)}{p(\mathbf{x} | c_2) p(c_2)} = \log \frac{p(c_1)}{p(c_2)} - \frac{1}{2}(\mathbf{x} - \mu_1)^T \Sigma^{-1}(\mathbf{x} - \mu_1) + \frac{1}{2}(\mathbf{x} - \mu_2)^T \Sigma^{-1}(\mathbf{x} - \mu_2)$$

IN THE CASE OF QDA THERE ARE NO ASSUMPTIONS ON THE COVARIANCE MATRICES  $\Sigma_c$  OF THE GAUSSIANS LEADING TO QUADRATIC DECISION SURFACES SEE3FOR MORE DETAILS  
NOTE RELATION WITH GAUSSIAN NAIVE BAYES  
IF IN THE QDA MODEL ONE ASSUMES THAT THE COVARIANCE MATRICES ARE DIAGONAL THEN THE INPUTS ARE ASSUMED TO BE CONDITIONALLY INDEPENDENT IN EACH CLASS AND THE RESULTING CLASSIFIER IS EQUIVALENT TO THE GAUSSIAN NAIVE BAYES CLASSIFIER  
NAIVEBAYESGAUSSIANNB  
3"THE ELEMENTS OF STATISTICAL LEARNING" HASTIE T TIBSHIRANI R FRIEDMAN J SECTION 43 P106119 2008  
31 SUPERVISED LEARNING 221

MATHEMATICAL FORMULATION OF LDA DIMENSIONALITY REDUCTION

TO UNDERSTAND THE USE OF LDA IN DIMENSIONALITY REDUCTION IT IS USEFUL TO START WITH A GEOMETRIC REFORMULATION OF THE LDA CLASSIFICATION RULE EXPLAINED ABOVE WE WRITE  $\mu_k$  FOR THE TOTAL NUMBER OF TARGET CLASSES SINCE IN LDA WE ASSUME THAT ALL CLASSES HAVE THE SAME ESTIMATED COVARIANCE  $\Sigma$  WE CAN RESCALE THE DATA SO THAT THIS COVARIANCE IS THE IDENTITY  $\Sigma = I$  WITH  $\mu_k$

THEN ONE CAN SHOW THAT TO CLASSIFY A DATA POINT AFTER SCALING IS EQUIVALENT TO FINDING THE ESTIMATED CLASS MEAN  $\mu_k$  WHICH

IS CLOSEST TO THE DATA POINT IN THE EUCLIDEAN DISTANCE BUT THIS CAN BE DONE JUST AS WELL AFTER PROJECTING ON THE  $(k-1)$  AFFINE SUBSPACE  $\pi_k$  GENERATED BY ALL THE  $\mu_k$

FOR ALL CLASSES THIS SHOWS THAT IMPLICIT IN THE LDA CLASSIFIER THERE IS A DIMENSIONALITY REDUCTION BY LINEAR PROJECTION ONTO A  $(k-1)$  DIMENSIONAL SPACE

WE CAN REDUCE THE DIMENSION EVEN MORE TO A CHOSEN  $k$  BY PROJECTING ONTO THE LINEAR SUBSPACE  $\pi_k$  WHICH MAXIMIZES THE VARIANCE OF THE  $\mu_k$

AFTER PROJECTION IN EFFECT WE ARE DOING A FORM OF PCA FOR THE TRANSFORMED CLASS MEANS  $\mu_k$

THIS  $k$  CORRESPONDS TO THE NCOMPONENTS PARAMETER USED IN THE DISCRIMINANTANALYSIS LINEARDISCRIMINANTANALYSIS TRANSFORM METHOD SEE 3 FOR MORE DETAILS

SHRINKAGE

SHRINKAGE IS A TOOL TO IMPROVE ESTIMATION OF COVARIANCE MATRICES IN SITUATIONS WHERE THE NUMBER OF TRAINING SAMPLES IS SMALL COMPARED TO THE NUMBER OF FEATURES IN THIS SCENARIO THE EMPIRICAL SAMPLE COVARIANCE IS A POOR ESTIMATOR SHRINKAGE LDA CAN BE USED BY SETTING THE SHRINKAGE PARAMETER OF THE DISCRIMINANTANALYSIS

LINEARDISCRIMINANTANALYSIS CLASS TO 'AUTO' THIS AUTOMATICALLY DETERMINES THE OPTIMAL SHRINKAGE PARAMETER IN AN ANALYTIC WAY FOLLOWING THE LEMMA INTRODUCED BY LEDOIT AND WOLF 4 NOTE THAT CURRENTLY SHRINKAGE ONLY WORKS WHEN SETTING THE SOLVER PARAMETER TO 'LSQR' OR 'EIGEN'

THE SHRINKAGE PARAMETER CAN ALSO BE MANUALLY SET BETWEEN 0 AND 1 IN PARTICULAR A VALUE OF 0 CORRESPONDS TO NO SHRINKAGE WHICH MEANS THE EMPIRICAL COVARIANCE MATRIX WILL BE USED AND A VALUE OF 1 CORRESPONDS TO COMPLETE SHRINKAGE WHICH MEANS THAT THE DIAGONAL MATRIX OF VARIANCES WILL BE USED AS AN ESTIMATE FOR THE COVARIANCE MATRIX SETTING THIS PARAMETER TO A VALUE BETWEEN THESE TWO EXTREMA WILL ESTIMATE A SHRUNK VERSION OF THE COVARIANCE MATRIX 4 LEDOIT O WOLF M HONEY I SHRUNK THE SAMPLE COVARIANCE MATRIX THE JOURNAL OF PORTFOLIO MANAGEMENT 30 4 1101 19 2004 222 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

ESTIMATION ALGORITHMS

THE DEFAULT SOLVER IS ‘SVD’ IT CAN PERFORM BOTH CLASSIFICATION AND TRANSFORM AND IT DOES NOT RELY ON THE CALCULATION OF THE COVARIANCE MATRIX THIS CAN BE AN ADVANTAGE IN SITUATIONS WHERE THE NUMBER OF FEATURES IS LARGE HOWEVER THE ‘SVD’ SOLVER CANNOT BE USED WITH SHRINKAGE

THE ‘LSQR’ SOLVER IS AN EFFICIENT ALGORITHM THAT ONLY WORKS FOR CLASSIFICATION IT SUPPORTS SHRINKAGE

THE ‘EIGEN’ SOLVER IS BASED ON THE OPTIMIZATION OF THE BETWEEN CLASS SCATTER TO WITHIN CLASS SCATTER RATIO IT CAN BE USED FOR BOTH CLASSIFICATION AND TRANSFORM AND IT SUPPORTS SHRINKAGE HOWEVER THE ‘EIGEN’ SOLVER NEEDS TO COMPUTE THE COVARIANCE MATRIX SO IT MIGHT NOT BE SUITABLE FOR SITUATIONS WITH A HIGH NUMBER OF FEATURES

EXAMPLES

NORMAL AND SHRINKAGE LINEAR DISCRIMINANT ANALYSIS FOR CLASSIFICATION COMPARISON OF LDA CLASSIFIERS WITH AND WITHOUT SHRINKAGE

REFERENCES

313 KERNEL RIDGE REGRESSION

KERNEL RIDGE REGRESSION KRR M2012 COMBINES RIDGE REGRESSION LINEAR LEAST SQUARES WITH L2NORM REGULARIZATION WITH THE KERNEL TRICK IT THUS LEARNS A LINEAR FUNCTION IN THE SPACE INDUCED BY THE RESPECTIVE KERNEL AND THE DATA FOR NONLINEAR KERNELS THIS CORRESPONDS TO A NONLINEAR FUNCTION IN THE ORIGINAL SPACE

THE FORM OF THE MODEL LEARNED BY KERNELRIDGE IS IDENTICAL TO SUPPORT VECTOR REGRESSION SVR HOWEVER DIFFERENT LOSS FUNCTIONS ARE USED KRR USES SQUARED ERROR LOSS WHILE SUPPORT VECTOR REGRESSION USES  $\epsilon$ INSENSITIVE LOSS BOTH COMBINED WITH L2 REGULARIZATION IN CONTRAST TO SVR FITTINGKERNELRIDGE CAN BE DONE IN CLOSEDFORM AND IS TYPICALLY FASTER FOR MEDIUMSIZED DATASETS ON THE OTHER HAND THE LEARNED MODEL IS NONSPARSE AND THUS SLOWER THAN SVR WHICH LEARNS A SPARSE MODEL FOR  $\epsilon = 0$  AT PREDICTIONTIME

THE FOLLOWING FIGURE COMPARES KERNELRIDGE ANDSVR ON AN ARTIFICIAL DATASET WHICH CONSISTS OF A SINUSOIDAL TARGET FUNCTION AND STRONG NOISE ADDED TO EVERY FIFTH DATAPOINT THE LEARNED MODEL OF KERNELRIDGE ANDSVR IS PLOTTED WHERE BOTH COMPLEXITYREGULARIZATION AND BANDWIDTH OF THE RBF KERNEL HAVE BEEN OPTIMIZED USING GRIDSEARCH THE LEARNED FUNCTIONS ARE VERY SIMILAR HOWEVER FITTING KERNELRIDGE IS APPROX SEVEN TIMES FASTER THAN FITTING SVR BOTH WITH GRIDSEARCH HOWEVER PREDICTION OF 100000 TARGET VALUES IS MORE THAN THREE TIMES FASTER WITH SVR SINCE IT HAS LEARNED A SPARSE MODEL USING ONLY APPROX 13 OF THE 100 TRAINING DATAPOINTS AS SUPPORT VECTORS THE NEXT FIGURE COMPARES THE TIME FOR FITTING AND PREDICTION OF KERNELRIDGE ANDSVR FOR DIFFERENT SIZES OF THE TRAINING SET FITTING KERNELRIDGE IS FASTER THAN SVR FOR MEDIUMSIZED TRAINING SETS LESS THAN 1000 SAMPLES HOWEVER FOR LARGER TRAINING SETS SVR SCALES BETTER WITH REGARD TO PREDICTION TIME SVR IS FASTER THAN KERNELRIDGE FOR ALL SIZES OF THE TRAINING SET BECAUSE OF THE LEARNED SPARSE SOLUTION NOTE THAT THE DEGREE OF SPARSITY AND THUS THE PREDICTION TIME DEPENDS ON THE PARAMETERS  $\epsilon$ AND $\gamma$ OF THESVR $\epsilon = 0$ WOULD CORRESPOND TO A DENSE MODEL

REFERENCES

314 SUPPORT VECTOR MACHINES

SUPPORT VECTOR MACHINES SVMs ARE A SET OF SUPERVISED LEARNING METHODS USED FOR CLASSIFICATION REGRESSION AND OUTLIERS DETECTION

31 SUPERVISED LEARNING 223





SCIKITLEARN USER GUIDE RELEASE 0213

THE ADVANTAGES OF SUPPORT VECTOR MACHINES ARE

- EFFECTIVE IN HIGH DIMENSIONAL SPACES
- STILL EFFECTIVE IN CASES WHERE NUMBER OF DIMENSIONS IS GREATER THAN THE NUMBER OF SAMPLES
- USES A SUBSET OF TRAINING POINTS IN THE DECISION FUNCTION CALLED SUPPORT VECTORS SO IT IS ALSO MEMORY EFFICIENT
- VERSATILE DIFFERENT KERNEL FUNCTIONS CAN BE SPECIFIED FOR THE DECISION FUNCTION COMMON KERNELS ARE PROVIDED BUT IT IS ALSO POSSIBLE TO SPECIFY CUSTOM KERNELS

THE DISADVANTAGES OF SUPPORT VECTOR MACHINES INCLUDE

- IF THE NUMBER OF FEATURES IS MUCH GREATER THAN THE NUMBER OF SAMPLES AVOID OVERFITTING IN CHOOSING KERNEL FUNCTIONS AND REGULARIZATION TERM IS CRUCIAL
- SVMs DO NOT DIRECTLY PROVIDE PROBABILITY ESTIMATES THESE ARE CALCULATED USING AN EXPENSIVE FIVEFOLD CROSS VALIDATION SEE SCORES AND PROBABILITIES BELOW

THE SUPPORT VECTOR MACHINES IN SCIKITLEARN SUPPORT BOTH DENSE NUMPYNDARRAY AND CONVERTIBLE TO THAT BY NUMPY ASARRAY AND SPARSE ANY SCIPYSPARSE SAMPLE VECTORS AS INPUT HOWEVER TO USE AN SVM TO MAKE PREDICTIONS FOR SPARSE DATA IT MUST HAVE BEEN FIT ON SUCH DATA FOR OPTIMAL PERFORMANCE USE CORDERED NUMPYNDARRAY DENSE ORSCIPYSPARSECSRMATRIX SPARSE WITH DTTYPEFLOAT64

CLASSIFICATION

SVCNUSVC ANDLINEARSVC ARE CLASSES CAPABLE OF PERFORMING MULTICLASS CLASSIFICATION ON A DATASET



SCIKITLEARN USER GUIDE RELEASE 0213

SVC ANDNUSVC ARE SIMILAR METHODS BUT ACCEPT SLIGHTLY DIFFERENT SETS OF PARAMETERS AND HAVE DIFFERENT MATHEMATICAL FORMULATIONS SEE SECTION MATHEMATICAL FORMULATION ON THE OTHER HAND LINEARSVC IS ANOTHER IMPLEMENTATION OF SUPPORT VECTOR CLASSIFICATION FOR THE CASE OF A LINEAR KERNEL NOTE THAT LINEARSVC DOES NOT ACCEPT KEYWORD KERNEL AS THIS IS ASSUMED TO BE LINEAR IT ALSO LACKS SOME OF THE MEMBERS OF SVC ANDNUSVC LKESUPPORT AS OTHER CLASSIFIERS SVCNUSVC ANDLINEARSVC TAKE AS INPUT TWO ARRAYS AN ARRAY X OF SIZE NSAMPLES NFEATURES HOLDING THE TRAINING SAMPLES AND AN ARRAY Y OF CLASS LABELS STRINGS OR INTEGERS SIZE NSAMPLES

```
FROM SKLEARN IMPORT SVM
X 0 0 1 1
Y 0 1
CLF SVM SVC GAMMA SCALE
CLFFITX Y
SVCC10 CACHESIZE200 CLASSWEIGHTNONE COEF000
DECISIONFUNCTIONSHAPEOVR DEGREE3 GAMMA SCALE KERNELRBF
MAXITER1 PROBABILITYFALSE RANDOMSTATENONE SHRINKINGTRUE
TOL0001 VERBOSEFALSE
AFTER BEING FITTED THE MODEL CAN THEN BE USED TO PREDICT NEW VALUES
CLFPREDICT2 2
ARRAY1
SVMS DECISION FUNCTION DEPENDS ON SOME SUBSET OF THE TRAINING DATA CALLED THE SUPPORT VECTORS SOME PROPERTIES OF THESE SUPPORT VECTORS CAN BE FOUND IN MEMBERS SUPPORTVECTORS SUPPORT ANDNSUPPORT
GET SUPPORT VECTORS
CLFSUPPORTVECTORS
ARRAY0 0
1 1
GET INDICES OF SUPPORT VECTORS
CLFSUPPORT
ARRAY0 1
GET NUMBER OF SUPPORT VECTORS FOR EACH CLASS
CLFNSUPPORT
ARRAY1 1
MULTICLASS CLASSIFICATION
SVC ANDNUSVC IMPLEMENT THE “ONEAGAINSTONE” APPROACH KNERR ET AL 1990 FOR MULTI CLASS CLASSIFICATION IFNCLASS IS THE NUMBER OF CLASSES THEN NCLASS NCLASS 1 2 CLASSIFIERS ARE CONSTRUCTED AND EACH ONE TRAINS DATA FROM TWO CLASSES TO PROVIDE A CONSISTENT INTERFACE WITH OTHER CLASSIFIERS THE DECISIONFUNCTIONSHAPE OPTION ALLOWS TO MONOTICALLY TRANSFORM THE RESULTS OF THE “ONEAGAINSTONE” CLASSIFIERS TO A DECISION FUNCTION OF SHAPE NSAMPLES NCLASSES
X 0 1 2 3
Y 0 1 2 3
CLF SVM SVC GAMMA SCALE DECISIONFUNCTIONSHAPEOVO
CLFFITX Y
SVCC10 CACHESIZE200 CLASSWEIGHTNONE COEF000
DECISIONFUNCTIONSHAPEOVO DEGREE3 GAMMA SCALE KERNELRBF
MAXITER1 PROBABILITYFALSE RANDOMSTATENONE SHRINKINGTRUE
TOL0001 VERBOSEFALSE
DEC CLFDECISIONFUNCTION1
DEC SHAPE1 4 CLASSES 4 32 6
6
31 SUPERVISED LEARNING 227
```

CLFDECISIONFUNCTIONSHAPE OVR

DEC CLFDECISIONFUNCTION1

DECSHAPE1 4 CLASSES

4

ON THE OTHER HAND LINEARSVC IMPLEMENTS “ONEVSTHEREST” MULTICLASS STRATEGY THUS TRAINING NCLASS MODELS IF THERE ARE ONLY TWO CLASSES ONLY ONE MODEL IS TRAINED

LINCLF SVMLINEARSVC

LINCLFFITX Y

LINEARSVCC10 CLASSWEIGHTNONE DUALTRUE FITINTERCEPTTRUE

INTERCEPTSCALING1 LOSSSSQUAREDHINGE MAXITER1000

MULTICLASSOVR PENALTYL2 RANDOMSTATENONE TOL00001

VERBOSE0

DEC LINCLFDECISIONFUNCTION1

DECSHAPE1

4

SEEMATHEMATICAL FORMULATION FOR A COMPLETE DESCRIPTION OF THE DECISION FUNCTION

NOTE THAT THE LINEARSVC ALSO IMPLEMENTS AN ALTERNATIVE MULTICLASS STRATEGY THE SOCALLED MULTICLASS SVM FORMULATED BY CRAMMER AND SINGER BY USING THE OPTION MULTICLASSCRAMMERSINGER THIS METHOD IS CONSISTENT WHICH IS NOT TRUE FOR ONEVSREST CLASSIFICATION IN PRACTICE ONEVSREST CLASSIFICATION IS USUALLY PREFERRED SINCE THE RESULTS ARE MOSTLY SIMILAR BUT THE RUNTIME IS SIGNIFICANTLY LESS

FOR “ONEVSREST” LINEARSVC THE ATTRIBUTES COEF ANDINTERCEPT HAVE THE SHAPE NCLASS

NFEATURES ANDNCLASS RESPECTIVELY EACH ROW OF THE COEFFICIENTS CORRESPONDS TO ONE OF THE NCLASS

MANY “ONEVSREST” CLASSIFIERS AND SIMILAR FOR THE INTERCEPTS IN THE ORDER OF THE “ONE” CLASS

IN THE CASE OF “ONEVSONE” SVC THE LAYOUT OF THE ATTRIBUTES IS A LITTLE MORE INVOLVED IN THE CASE OF HAVING

A LINEAR KERNEL THE ATTRIBUTES COEF ANDINTERCEPT HAVE THE SHAPE NCLASS NCLASS 1

2 NFEATURES ANDNCLASS NCLASS 1 2 RESPECTIVELY THIS IS SIMILAR TO THE LAYOUT FOR

LINEARSVC DESCRIBED ABOVE WITH EACH ROW NOW CORRESPONDING TO A BINARY CLASSIFIER THE ORDER FOR CLASSES 0 TO N IS

“0 VS 1” “0 VS 2” “0 VS N” “1 VS 2” “1 VS 3” “1 VS N” “N1 VS N”

THE SHAPE OF DUALCOEF ISNCLASS1 NSV WITH A SOMEWHAT HARD TO GRASP LAYOUT THE COLUMNS CORRE

SPOND TO THE SUPPORT VECTORS INVOLVED IN ANY OF THE NCLASS NCLASS 1 2 “ONEVSONE” CLASSIFIERS

EACH OF THE SUPPORT VECTORS IS USED IN NCLASS 1 CLASSIFIERS THE NCLASS 1 ENTRIES IN EACH ROW CORRESPOND

TO THE DUAL COEFFICIENTS FOR THESE CLASSIFIERS

THIS MIGHT BE MADE MORE CLEAR BY AN EXAMPLE

CONSIDER A THREE CLASS PROBLEM WITH CLASS 0 HAVING THREE SUPPORT VECTORS

0

0

0

AND CLASS 1 AND 2 HAVING TWO

SUPPORT VECTORS

1

1AND

2

RESPECTIVELY FOR EACH SUPPORT VECTOR

THERE ARE TWO DUAL COEFFICIENTS LET’S CALL

THE COEFFICIENT OF SUPPORT VECTOR

IN THE CLASSIFIER BETWEEN CLASSES AND

THENDUALCOEF LOOKS LIKE THIS

0

01

COEFFICIENTS FOR SVS OF CLASS 0

1

01

02

2

01

02

0

10

COEFFICIENTS FOR SVS OF CLASS 1

1

10

12

0

20

COEFFICIENTS FOR SVS OF CLASS 2

1

20

21

SCIKITLEARN USER GUIDE RELEASE 0213

SCORES AND PROBABILITIES

THEDECISIONFUNCTION METHOD OF SVC AND NUSVC GIVES PERCLASS SCORES FOR EACH SAMPLE OR A SINGLE SCORE PER SAMPLE IN THE BINARY CASE WHEN THE CONSTRUCTOR OPTION PROBABILITY IS SET TO TRUE. CLASS MEMBERSHIP PROBABILITY ESTIMATES FROM THE METHODS PREDICTPROBA AND PREDICTLOGPROBA ARE ENABLED IN THE BINARY CASE. THE PROBABILITIES ARE CALIBRATED USING PLATT SCALING LOGISTIC REGRESSION ON THE SVM'S SCORES FIT BY AN ADDITIONAL CROSSVALIDATION ON THE TRAINING DATA IN THE MULTICLASS CASE. THIS IS EXTENDED AS PER WU ET AL 2004. NEEDLESS TO SAY THE CROSSVALIDATION INVOLVED IN PLATT SCALING IS AN EXPENSIVE OPERATION FOR LARGE DATASETS. IN ADDITION THE PROBABILITY ESTIMATES MAY BE INCONSISTENT WITH THE SCORES IN THE SENSE THAT THE "ARGMAX" OF THE SCORES MAY NOT BE THE ARGMAX OF THE PROBABILITIES. EG IN BINARY CLASSIFICATION A SAMPLE MAY BE LABELED BY PREDICT AS BELONGING TO A CLASS THAT HAS PROBABILITY  $\frac{1}{2}$  ACCORDING TO PREDICTPROBA. PLATT'S METHOD IS ALSO KNOWN TO HAVE THEORETICAL ISSUES IF CONFIDENCE SCORES ARE REQUIRED BUT THESE DO NOT HAVE TO BE PROBABILITIES THEN IT IS ADVISABLE TO SET PROBABILITYFALSE AND USEDECISIONFUNCTION INSTEAD OF PREDICTPROBA.

REFERENCES

- WU LIN AND WENG "PROBABILITY ESTIMATES FOR MULTICLASS CLASSIFICATION BY PAIRWISE COUPLING" JMLR 5975 1005 2004
- PLATT "PROBABILISTIC OUTPUTS FOR SVMs AND COMPARISONS TO REGULARIZED LIKELIHOOD METHODS"

UNBALANCED PROBLEMS

IN PROBLEMS WHERE IT IS DESIRED TO GIVE MORE IMPORTANCE TO CERTAIN CLASSES OR CERTAIN INDIVIDUAL SAMPLES KEYWORDS CLASSWEIGHT AND SAMPLEWEIGHT CAN BE USED.

SVC BUT NOT NUSVC IMPLEMENT A KEYWORD CLASSWEIGHT IN THE FIT METHOD. IT'S A DICTIONARY OF THE FORM CLASSLABEL: VALUE WHERE VALUE IS A FLOATING POINT NUMBER > 0 THAT SETS THE PARAMETER COF CLASS.

CLASSLABEL: TOCVALUE

SCIKITLEARN USER GUIDE RELEASE 0213  
SVCNUSVC SVRNUSVR ANDONECLASSSVM IMPLEMENT ALSO WEIGHTS FOR INDIVIDUAL SAMPLES IN METHOD FIT  
THROUGH KEYWORD SAMPLEWEIGHT SIMILAR TO CLASSWEIGHT THESE SET THE PARAMETER CFOR THE ITH EXAMPLE TO  
CSAMPLEWEIGHTI  
EXAMPLES

- PLOT DIFFERENT SVM CLASSIFIERS IN THE IRIS DATASET
- SVM MAXIMUM MARGIN SEPARATING HYPERPLANE
- SVM SEPARATING HYPERPLANE FOR UNBALANCED CLASSES
- SVMANOVA SVM WITH UNIVARIATE FEATURE SELECTION
- NONLINEAR SVM
- SVM WEIGHTED SAMPLES

REGRESSION  
THE METHOD OF SUPPORT VECTOR CLASSIFICATION CAN BE EXTENDED TO SOLVE REGRESSION PROBLEMS THIS METHOD IS CALLED  
SUPPORT VECTOR REGRESSION

THE MODEL PRODUCED BY SUPPORT VECTOR CLASSIFICATION AS DESCRIBED ABOVE DEPENDS ONLY ON A SUBSET OF THE TRAINING  
DATA BECAUSE THE COST FUNCTION FOR BUILDING THE MODEL DOES NOT CARE ABOUT TRAINING POINTS THAT LIE BEYOND THE MARGIN  
ANALOGOUSLY THE MODEL PRODUCED BY SUPPORT VECTOR REGRESSION DEPENDS ONLY ON A SUBSET OF THE TRAINING DATA BECAUSE  
THE COST FUNCTION FOR BUILDING THE MODEL IGNORES ANY TRAINING DATA CLOSE TO THE MODEL PREDICTION  
THERE ARE THREE DIFFERENT IMPLEMENTATIONS OF SUPPORT VECTOR REGRESSION SVRNUSVR ANDLINEARSVR  
LINEARSVR PROVIDES A FASTER IMPLEMENTATION THAN SVR BUT ONLY CONSIDERS LINEAR KERNELS WHILE NUSVR IMPLEMENTS  
A SLIGHTLY DIFFERENT FORMULATION THAN SVR ANDLINEARSVR SEE IMPLEMENTATION DETAILS FOR FURTHER DETAILS  
AS WITH CLASSIFICATION CLASSES THE FIT METHOD WILL TAKE AS ARGUMENT VECTORS X Y ONLY THAT IN THIS CASE Y IS EXPECTED TO  
HAVE FLOATING POINT VALUES INSTEAD OF INTEGER VALUES

```
FROM SKLEARN IMPORT SVM  
X 0 0 2 2
```

230 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

Y 05 25  
CLF SVM SVR  
CLFFITX Y  
SVRC10 CACHESIZE200 COEF000 DEGREE3 EPSILON01  
GAMMAAUTODEPRECATED KERNELRBF MAXITER1 SHRINKINGTRUE  
TOL0001 VERBOSEFALSE  
CLFPREDICT1 1  
ARRAY15

EXAMPLES

•SUPPORT VECTOR REGRESSION SVR USING LINEAR AND NONLINEAR KERNELS

DENSITY ESTIMATION NOVELTY DETECTION

THE CLASSONECLASSSVM IMPLEMENTS A ONECLASS SVM WHICH IS USED IN OUTLIER DETECTION

SEENOVELTY AND OUTLIER DETECTION FOR THE DESCRIPTION AND USAGE OF ONECLASSSVM

COMPLEXITY

SUPPORT VECTOR MACHINES ARE POWERFUL TOOLS BUT THEIR COMPUTE AND STORAGE REQUIREMENTS INCREASE RAPIDLY WITH THE NUMBER OF TRAINING VECTORS THE CORE OF AN SVM IS A QUADRATIC PROGRAMMING PROBLEM QP SEPARATING SUPPORT VECTORS FROM THE REST OF THE TRAINING DATA THE QP SOLVER USED BY THIS LIBSVMBASED IMPLEMENTATION SCALES BETWEEN

$n^2$

AND  $n^3$

DEPENDING ON HOW EFFICIENTLY THE LIBSVM CACHE IS USED IN

PRACTICE DATASET DEPENDENT IF THE DATA IS VERY SPARSE SHOULD BE REPLACED BY THE AVERAGE NUMBER OF NON ZERO FEATURES IN A SAMPLE VECTOR

ALSO NOTE THAT FOR THE LINEAR CASE THE ALGORITHM USED IN LINEARSVC BY THE LIBLINEAR IMPLEMENTATION IS MUCH MORE EFFICIENT THAN ITS LIBSVMBASED SVC COUNTERPART AND CAN SCALE ALMOST LINEARLY TO MILLIONS OF SAMPLES ANDOR FEATURES

TIPS ON PRACTICAL USE

•AVOIDING DATA COPY FOR SVCSVRNUSVC ANDNUSVR IF THE DATA PASSED TO CERTAIN METHODS IS NOT CORDERED CONTIGUOUS AND DOUBLE PRECISION IT WILL BE COPIED BEFORE CALLING THE UNDERLYING C IMPLEMENTATION YOU CAN

CHECK WHETHER A GIVEN NUMPY ARRAY IS CCONTIGUOUS BY INSPECTING ITS FLAGS ATTRIBUTE

FOR LINEARSVC AND LOGISTIC REGRESSION ANY INPUT PASSED AS A NUMPY ARRAY WILL BE COPIED AND CONVERTED TO THE LIBLINEAR INTERNAL SPARSE DATA REPRESENTATION DOUBLE PRECISION FLOATS AND INT32 INDICES OF NONZERO COMPONENTS IF YOU WANT TO FIT A LARGESCALE LINEAR CLASSIFIER WITHOUT COPYING A DENSE NUMPY CCONTIGUOUS DOUBLE PRECISION ARRAY AS INPUT WE SUGGEST TO USE THE SGDCLASSIFIER CLASS INSTEAD THE OBJECTIVE FUNCTION CAN BE CONFIGURED TO BE ALMOST THE SAME AS THE LINEARSVC MODEL

•KERNEL CACHE SIZE FOR SVCSVRNUSVC ANDNUSVR THE SIZE OF THE KERNEL CACHE HAS A STRONG IMPACT ON RUN TIMES FOR LARGER PROBLEMS IF YOU HAVE ENOUGH RAM AVAILABLE IT IS RECOMMENDED TO SET CACHESIZE TO A HIGHER VALUE THAN THE DEFAULT OF 200MB SUCH AS 500MB OR 1000MB

•SETTING C IS1BY DEFAULT AND IT'S A REASONABLE DEFAULT CHOICE IF YOU HAVE A LOT OF NOISY OBSERVATIONS YOU SHOULD DECREASE IT IT CORRESPONDS TO REGULARIZE MORE THE ESTIMATION

LINEARSVC AND LINEAR SVR ARE LESS SENSITIVE TO C WHEN IT BECOMES LARGE AND PREDICTION RESULTS STOP IMPROVING AFTER A CERTAIN THRESHOLD MEANWHILE LARGER C VALUES WILL TAKE MORE TIME TO TRAIN SOMETIMES UP TO 10 TIMES LONGER AS SHOWN BY FAN ET AL 2008

SCIKITLEARN USER GUIDE RELEASE 0213

• SUPPORT VECTOR MACHINE ALGORITHMS ARE NOT SCALE INVARIANT SO IT IS HIGHLY RECOMMENDED TO SCALE YOUR DATA FOR EXAMPLE SCALE EACH ATTRIBUTE ON THE INPUT VECTOR X TO 01 OR 11 OR STANDARDIZE IT TO HAVE MEAN 0 AND VARIANCE 1 NOTE THAT THE SAME SCALING MUST BE APPLIED TO THE TEST VECTOR TO OBTAIN MEANINGFUL RESULTS SEE SECTION PREPROCESSING DATA FOR MORE DETAILS ON SCALING AND NORMALIZATION

• PARAMETER NUINNUSVC ONECLASSSVM NUSVR APPROXIMATES THE FRACTION OF TRAINING ERRORS AND SUPPORT VECTORS

• INSVC IF DATA FOR CLASSIFICATION ARE UNBALANCED EG MANY POSITIVE AND FEW NEGATIVE SET CLASSWEIGHTBALANCED ANDOR TRY DIFFERENT PENALTY PARAMETERS C

•RANDOMNESS OF THE UNDERLYING IMPLEMENTATIONS THE UNDERLYING IMPLEMENTATIONS OF SVC ANDNUSVC USE A RANDOM NUMBER GENERATOR ONLY TO SHUFFLE THE DATA FOR PROBABILITY ESTIMATION WHEN PROBABILITY IS SET TO TRUE THIS RANDOMNESS CAN BE CONTROLLED WITH THE RANDOMSTATE PARAMETER IF PROBABILITY IS SET TOFALSE THESE ESTIMATORS ARE NOT RANDOM AND RANDOMSTATE HAS NO EFFECT ON THE RESULTS THE UNDERLYING ONECLASSSVM IMPLEMENTATION IS SIMILAR TO THE ONES OF SVC ANDNUSVC AS NO PROBABILITY ESTIMATION IS PROVIDED FOR ONECLASSSVM IT IS NOT RANDOM

THE UNDERLYING LINEARSVC IMPLEMENTATION USES A RANDOM NUMBER GENERATOR TO SELECT FEATURES WHEN FITTING THE MODEL WITH A DUAL COORDINATE DESCENT IE WHEN DUAL IS SET TOTRUE IT IS THUS NOT UNCOMMON TO HAVE SLIGHTLY DIFFERENT RESULTS FOR THE SAME INPUT DATA IF THAT HAPPENS TRY WITH A SMALLER TOL PARAMETER THIS RANDOMNESS CAN ALSO BE CONTROLLED WITH THE RANDOMSTATE PARAMETER WHEN DUAL IS SET TOFALSE THE UNDERLYING IMPLEMENTATION OF LINEARSVC IS NOT RANDOM AND RANDOMSTATE HAS NO EFFECT ON THE RESULTS

• USING L1 PENALIZATION AS PROVIDED BY LINEARSVCLOSSL2 PENALTYL1 DUALFALSE

YIELDS A SPARSE SOLUTION IE ONLY A SUBSET OF FEATURE WEIGHTS IS DIFFERENT FROM ZERO AND CONTRIBUTE TO THE DECISION FUNCTION INCREASING CYIELDS A MORE COMPLEX MODEL MORE FEATURE ARE SELECTED THE CVALUE THAT YIELDS A “NULL” MODEL ALL WEIGHTS EQUAL TO ZERO CAN BE CALCULATED USING L1MINC

REFERENCES

• FAN RONGEN ET AL “LIBLINEAR A LIBRARY FOR LARGE LINEAR CLASSIFICATION” JOURNAL OF MACHINE LEARNING RESEARCH 9AUG 2008 18711874

KERNEL FUNCTIONS

THEKERNEL FUNCTION CAN BE ANY OF THE FOLLOWING

- LINEAR<math>\langle \phi \rangle</math>
- POLYNOMIAL  $\langle \phi \rangle^2$ IS SPECIFIED BY KEYWORD DEGREE  $\phi$ BYCOEFO
- RBF  $\exp(-\frac{\|\phi - \phi'\|^2}{2\sigma^2})$ IS SPECIFIED BY KEYWORD GAMMA MUST BE GREATER THAN 0
- SIGMOID  $\tanh \langle \phi \rangle$  WHERE  $\phi$ IS SPECIFIED BY COEFO

DIFFERENT KERNELS ARE SPECIFIED BY KEYWORD KERNEL AT INITIALIZATION

LINEARSVC SVMKVCKERNELLINEAR

LINEARSVCKERNEL

LINEAR

RBFSVC SVMKVCKERNELRBF

RBFSVCKERNEL

RBF

SCIKITLEARN USER GUIDE RELEASE 0213

CUSTOM KERNELS

YOU CAN DEFINE YOUR OWN KERNELS BY EITHER GIVING THE KERNEL AS A PYTHON FUNCTION OR BY PRECOMPUTING THE GRAM MATRIX  
CLASSIFIERS WITH CUSTOM KERNELS BEHAVE THE SAME WAY AS ANY OTHER CLASSIFIERS EXCEPT THAT

- FIELDSUPPORTVECTORS IS NOW EMPTY ONLY INDICES OF SUPPORT VECTORS ARE STORED IN SUPPORT
- A REFERENCE AND NOT A COPY OF THE FIRST ARGUMENT IN THE FIT METHOD IS STORED FOR FUTURE REFERENCE IF THAT  
ARRAY CHANGES BETWEEN THE USE OF FIT ANDPREDICT YOU WILL HAVE UNEXPECTED RESULTS  
USING PYTHON FUNCTIONS AS KERNELS

YOU CAN ALSO USE YOUR OWN DEFINED KERNELS BY PASSING A FUNCTION TO THE KEYWORD KERNEL IN THE CONSTRUCTOR  
YOUR KERNEL MUST TAKE AS ARGUMENTS TWO MATRICES OF SHAPE NSAMPLES1 NFEATURES NSAMPLES2  
NFEATURES AND RETURN A KERNEL MATRIX OF SHAPE NSAMPLES1 NSAMPLES2

THE FOLLOWING CODE DEFINES A LINEAR KERNEL AND CREATES A CLASSIFIER INSTANCE THAT WILL USE THAT KERNEL

```
import numpy as np
from sklearn import svm
def mykernel(x, y):
    return np.dot(x, y)
```

clf = svm.SVC(kernel=mykernel)
examples

- SVM WITH CUSTOM KERNEL  
USING THE GRAM MATRIX

setkernelprecomputed and pass the gram matrix instead of x in the fit method at the moment the kernel  
values between alltraining vectors and the test vectors must be provided

```
import numpy as np
from sklearn import svm
x = np.array([0, 1, 1])
y = [0, 1]
```

```
clf = svm.SVC(kernel=precomputed)
# linear kernel computation
gram = np.dot(x, x.T)
clf.fit(gram)
svcc10 = svm.SVC(kernel='linear', cache_size=200, class_weight=None, coef0=0.0,
decision_function_shape='ovr', degree=3, gamma='auto', deprecated_kernel=precomputed, max_iter=1, probability=False,
random_state=None, shrinking=True, tol=0.001, verbose=False)
predict on training examples
clf.predict(gram)
array([0, 1])
```

PARAMETERS OF THE RBF KERNEL

WHEN TRAINING AN SVM WITH THE RADIAL BASIS FUNCTION RBF KERNEL TWO PARAMETERS MUST BE CONSIDERED  $C$  AND  $\gamma$ . THE PARAMETER  $C$  COMMON TO ALL SVM KERNELS TRADES OFF MISCLASSIFICATION OF TRAINING EXAMPLES AGAINST SIMPLICITY OF THE DECISION SURFACE. A LOW  $C$  MAKES THE DECISION SURFACE SMOOTH WHILE A HIGH  $C$  AIMS AT CLASSIFYING ALL TRAINING EXAMPLES CORRECTLY.  $\gamma$  DEFINES HOW MUCH INFLUENCE A SINGLE TRAINING EXAMPLE HAS. THE LARGER  $\gamma$  IS THE CLOSER OTHER EXAMPLES MUST BE TO BE AFFECTED.

PROPER CHOICE OF  $C$  AND  $\gamma$  IS CRITICAL TO THE SVM'S PERFORMANCE. ONE IS ADVISED TO USE SKLEARN `ModelSelectionGridSearchCV` WITH `C` AND  $\gamma$  SPACED EXPONENTIALLY FAR APART TO CHOOSE GOOD VALUES.

EXAMPLES

•RBF SVM PARAMETERS

MATHEMATICAL FORMULATION

A SUPPORT VECTOR MACHINE CONSTRUCTS A HYPERPLANE OR SET OF HYPERPLANES IN A HIGH OR INFINITE DIMENSIONAL SPACE WHICH CAN BE USED FOR CLASSIFICATION, REGRESSION OR OTHER TASKS. INTUITIVELY A GOOD SEPARATION IS ACHIEVED BY THE HYPERPLANE THAT HAS THE LARGEST DISTANCE TO THE NEAREST TRAINING DATA POINTS OF ANY CLASS. SO CALLED FUNCTIONAL MARGIN. SINCE IN GENERAL THE LARGER THE MARGIN THE LOWER THE GENERALIZATION ERROR OF THE CLASSIFIER.



SVC

GIVEN TRAINING VECTORS  $x_i \in \mathbb{R}^l$   $i = 1 \dots n$  IN TWO CLASSES AND A VECTOR  $x \in \mathbb{R}^l$  SVC SOLVES THE FOLLOWING PRIMAL PROBLEM

MIN  
$$\frac{1}{2} \sum_{i=1}^n \xi_i^2$$
  
SUBJECT TO  $\xi_i \geq 1 - y_i w^T x_i$   
 $\xi_i \geq 0 \quad i = 1 \dots n$   
ITS DUAL IS

MIN  
$$\frac{1}{2} \sum_{i=1}^n \alpha_i^2$$
  
SUBJECT TO  $\alpha_i \geq 0$   
 $0 \leq \alpha_i \leq 1 \quad i = 1 \dots n$   
WHERE  $w$  IS THE VECTOR OF ALL ONES  $w = (1, \dots, 1)^T$  IS AN  $n$  BY  $n$  POSITIVE SEMIDEFINITE MATRIX  $\alpha_i \alpha_j =$   
 $\sum_{k=1}^n x_{ik} x_{jk}$  WHERE  $x_{ik} x_{jk}$  IS THE KERNEL HERE TRAINING VECTORS ARE IMPLICITLY MAPPED INTO A  
HIGHER MAYBE INFINITE DIMENSIONAL SPACE BY THE FUNCTION  $\phi$   
THE DECISION FUNCTION IS

$$\text{SGN} \left( \sum_{i=1}^n \alpha_i x_i^T x \right)$$
  
NOTE WHILE SVM MODELS DERIVED FROM LIBSVM AND LIBLINEAR USE C AS REGULARIZATION PARAMETER MOST OTHER ESTIMATORS  
USE ALPHA THE EXACT EQUIVALENCE BETWEEN THE AMOUNT OF REGULARIZATION OF TWO MODELS DEPENDS ON THE EXACT OBJECTIVE  
FUNCTION OPTIMIZED BY THE MODEL FOR EXAMPLE WHEN THE ESTIMATOR USED IS SKLEARNLINEARMODEL RIDGE  
REGRESSION THE RELATION BETWEEN THEM IS GIVEN AS  $C = \frac{1}{\alpha}$

$\alpha_i$   $h(x_i, x)$   
THIS PARAMETERS CAN BE ACCESSED THROUGH THE MEMBERS DUALCOEF WHICH HOLDS THE PRODUCT  $\alpha_i \alpha_j$   
SUPPORTVECTORS WHICH HOLDS THE SUPPORT VECTORS AND INTERCEPT WHICH HOLDS THE INDEPENDENT TERM  $w_0$

REFERENCES

- “AUTOMATIC CAPACITY TUNING OF VERY LARGE VCDIMENSION CLASSIFIERS” I GUYON B BOSER V VAPNIK  
ADVANCES IN NEURAL INFORMATION PROCESSING 1993
- “SUPPORTVECTOR NETWORKS” C CORTES V VAPNIK MACHINE LEARNING 20 273297 1995

NUSVC

WE INTRODUCE A NEW PARAMETER  $\nu$  WHICH CONTROLS THE NUMBER OF SUPPORT VECTORS AND TRAINING ERRORS THE PARAMETER  
 $\nu \in [0, 1]$  IS AN UPPER BOUND ON THE FRACTION OF TRAINING ERRORS AND A LOWER BOUND OF THE FRACTION OF SUPPORT VECTORS  
IT CAN BE SHOWN THAT THE  $\nu$  SVC FORMULATION IS A REPARAMETERIZATION OF THE  $C$  SVC AND THEREFORE MATHEMATICALLY  
EQUIVALENT  
31 SUPERVISED LEARNING 235

SVR  
GIVEN TRAINING VECTORS  $x_i \in \mathbb{R}^1$   $i = 1 \dots N$  AND A VECTOR  $y \in \mathbb{R}^1$  SVR SOLVES THE FOLLOWING PRIMAL PROBLEM  
MIN

$$\frac{1}{2} \sum_{i=1}^N (y_i - f(x_i))^2$$

SUBJECT TO  $f(x) = w^T x + b$   
 $w^T x + b \leq y_i$   
 $w^T x + b \geq y_i$

ITS DUAL IS  
MIN

$$\frac{1}{2} \sum_{i=1}^N (y_i - \alpha_i)^2$$

SUBJECT TO  $\alpha_i \geq 0$   
 $\sum \alpha_i = 0$   
 $\alpha_i \leq 1$

WHERE  $\mathbf{1}$  IS THE VECTOR OF ALL ONES  $\alpha$  IS THE UPPER BOUND  $\mathbf{K}$  IS AN  $n$  BY  $n$  POSITIVE SEMIDEFINITE MATRIX  $\mathbf{K}_{ij} = k(x_i, x_j)$  IS THE KERNEL HERE TRAINING VECTORS ARE IMPLICITLY MAPPED INTO A HIGHER MAYBE INFINITE DIMENSIONAL SPACE BY THE FUNCTION  $\phi$   
THE DECISION FUNCTION IS

$$f(x) = \sum_{i=1}^N \alpha_i (y_i - \phi(x_i)^T \phi(x) + b)$$

THESE PARAMETERS CAN BE ACCESSED THROUGH THE MEMBERS DUALCOEF WHICH HOLDS THE DIFFERENCE  $y_i - f(x_i)$   
SUPPORTVECTORS WHICH HOLDS THE SUPPORT VECTORS AND INTERCEPT WHICH HOLDS THE INDEPENDENT TERM  $b$

REFERENCES  
• “A TUTORIAL ON SUPPORT VECTOR REGRESSION” ALEX J SMOLA BERNHARD SCHÖLKOPF STATISTICS AND COMPUTING ARCHIVE VOLUME 14 ISSUE 3 AUGUST 2004 P 199222

IMPLEMENTATION DETAILS  
INTERNALLY WE USE LIBSVM AND LIBLINEAR TO HANDLE ALL COMPUTATIONS THESE LIBRARIES ARE WRAPPED USING C AND CYTHON  
REFERENCES

FOR A DESCRIPTION OF THE IMPLEMENTATION AND DETAILS OF THE ALGORITHMS USED PLEASE REFER TO  
• LIBSVM A LIBRARY FOR SUPPORT VECTOR MACHINES  
• LIBLINEAR - A LIBRARY FOR LARGE LINEAR CLASSIFICATION

315 STOCHASTIC GRADIENT DESCENT  
STOCHASTIC GRADIENT DESCENT SGD IS A SIMPLE YET VERY EFFICIENT APPROACH TO DISCRIMINATIVE LEARNING OF LINEAR CLASSIFIERS UNDER CONVEX LOSS FUNCTIONS SUCH AS LINEAR SUPPORT VECTOR MACHINES AND LOGISTIC REGRESSION EVEN THOUGH  
236 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

SGD HAS BEEN AROUND IN THE MACHINE LEARNING COMMUNITY FOR A LONG TIME IT HAS RECEIVED A CONSIDERABLE AMOUNT OF ATTENTION JUST RECENTLY IN THE CONTEXT OF LARGESCALE LEARNING

SGD HAS BEEN SUCCESSFULLY APPLIED TO LARGESCALE AND SPARSE MACHINE LEARNING PROBLEMS OFTEN ENCOUNTERED IN TEXT CLASSIFICATION AND NATURAL LANGUAGE PROCESSING GIVEN THAT THE DATA IS SPARSE THE CLASSIFIERS IN THIS MODULE EASILY SCALE TO PROBLEMS WITH MORE THAN 105 TRAINING EXAMPLES AND MORE THAN 105 FEATURES

THE ADVANTAGES OF STOCHASTIC GRADIENT DESCENT ARE

- EFFICIENCY
- EASE OF IMPLEMENTATION LOTS OF OPPORTUNITIES FOR CODE TUNING

THE DISADVANTAGES OF STOCHASTIC GRADIENT DESCENT INCLUDE

- SGD REQUIRES A NUMBER OF HYPERPARAMETERS SUCH AS THE REGULARIZATION PARAMETER AND THE NUMBER OF ITERATIONS
- SGD IS SENSITIVE TO FEATURE SCALING

CLASSIFICATION

WARNING MAKE SURE YOU PERMUTE SHUFFLE YOUR TRAINING DATA BEFORE FITTING THE MODEL OR USE SHUFFLETRUE TO SHUFFLE AFTER EACH ITERATION

THE CLASSSGDCLASSIFIER IMPLEMENTS A PLAIN STOCHASTIC GRADIENT DESCENT LEARNING ROUTINE WHICH SUPPORTS DIFFERENT LOSS FUNCTIONS AND PENALTIES FOR CLASSIFICATION

AS OTHER CLASSIFIERS SGD HAS TO BE FITTED WITH TWO ARRAYS AN ARRAY X OF SIZE NSAMPLES NFEATURES HOLDING THE TRAINING SAMPLES AND AN ARRAY Y OF SIZE NSAMPLES HOLDING THE TARGET VALUES CLASS LABELS FOR THE TRAINING SAMPLES

```
FROM SKLEARNLINEARMODEL IMPORT SGDCLASSIFIER
X 0 0 1 1
```

31 SUPERVISED LEARNING 237

SCIKITLEARN USER GUIDE RELEASE 0213

Y 0 1

CLF SGDCLASSIFIERLOSSHINGE PENALTYL2 MAXITER5

CLFFITX Y

SGDCLASSIFIERALPHA00001 AVERAGEFALSE CLASSWEIGHTNONE

EARLYSTOPPINGFALSE EPSILON01 ETA000 FITINTERCEPTTRUE

L1RATIO015 LEARNINGRATEOPTIMAL LOSSHINGE MAXITER5

NITERNOCHANGE5 NJOBSNONE PENALTYL2 POWERT05

RANDOMSTATENONE SHUFFLETRUE TOL0001

VALIDATIONFRACTION01 VERBOSE0 WARMSTARTFALSE

AFTER BEING FITTED THE MODEL CAN THEN BE USED TO PREDICT NEW VALUES

CLFPREDICT2 2

ARRAY1

SGD FITS A LINEAR MODEL TO THE TRAINING DATA THE MEMBER COEF HOLDS THE MODEL PARAMETERS

CLFCOEF

ARRAY99 99

MEMBERINTERCEPT HOLDS THE INTERCEPT AKA OFFSET OR BIAS

CLFINTERCEPT

ARRAY99

WHETHER OR NOT THE MODEL SHOULD USE AN INTERCEPT IE A BIASED HYPERPLANE IS CONTROLLED BY THE PARAMETER

FITINTERCEPT

TO GET THE SIGNED DISTANCE TO THE HYPERPLANE USE SGDCLASSIFIERDECISIONFUNCTION

CLFDECISIONFUNCTION2 2

ARRAY296

THE CONCRETE LOSS FUNCTION CAN BE SET VIA THE LOSS PARAMETER SGDCLASSIFIER SUPPORTS THE FOLLOWING LOSS FUNCTIONS

- LOSSHINGE SOFTMARGIN LINEAR SUPPORT VECTOR MACHINE
- LOSSMODIFIEDHUBER SMOOTHED HINGE LOSS
- LOSSLOG LOGISTIC REGRESSION
- AND ALL REGRESSION LOSSES BELOW

THE FIRST TWO LOSS FUNCTIONS ARE LAZY THEY ONLY UPDATE THE MODEL PARAMETERS IF AN EXAMPLE VIOLATES THE MARGIN CONSTRAINT WHICH MAKES TRAINING VERY EFFICIENT AND MAY RESULT IN SPARSER MODELS EVEN WHEN L2 PENALTY IS USED

USINGLOSSLOG ORLOSSMODIFIEDHUBER ENABLES THE PREDICTPROBA METHOD WHICH GIVES A VECTOR

OF PROBABILITY ESTIMATES [ ]PER SAMPLE [ ]

CLF SGDCLASSIFIERLOSSLOG MAXITER5FITX Y

CLFPREDICTPROBA1 1

ARRAY000 099

THE CONCRETE PENALTY CAN BE SET VIA THE PENALTY PARAMETER SGD SUPPORTS THE FOLLOWING PENALTIES

- PENALTYL2 L2 NORM PENALTY ON COEF
- PENALTYL1 L1 NORM PENALTY ON COEF

238 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

•PENALTYELASTICNET CONVEX COMBINATION OF L2 AND L1 1 L1RATIO L2  
L1RATIO L1

THE DEFAULT SETTING IS PENALTYL2 THE L1 PENALTY LEADS TO SPARSE SOLUTIONS DRIVING MOST COEFFICIENTS TO ZERO  
THE ELASTIC NET SOLVES SOME DEFICIENCIES OF THE L1 PENALTY IN THE PRESENCE OF HIGHLY CORRELATED ATTRIBUTES THE PARAM  
ETERL1RATIO CONTROLS THE CONVEX COMBINATION OF L1 AND L2 PENALTY

SGDCCLASSIFIER SUPPORTS MULTICLASS CLASSIFICATION BY COMBINING MULTIPLE BINARY CLASSIFIERS IN A “ONE VERSUS ALL”  
OV A SCHEME FOR EACH OF THE KCLASSES A BINARY CLASSIFIER IS LEARNED THAT DISCRIMINATES BETWEEN THAT AND ALL OTHER  
K-1CLASSES AT TESTING TIME WE COMPUTE THE CONFIDENCE SCORE IE THE SIGNED DISTANCES TO THE HYPERPLANE FOR EACH  
CLASSIFIER AND CHOOSE THE CLASS WITH THE HIGHEST CONFIDENCE THE FIGURE BELOW ILLUSTRATES THE OV A APPROACH ON THE IRIS  
DATASET THE DASHED LINES REPRESENT THE THREE OV A CLASSIFIERS THE BACKGROUND COLORS SHOW THE DECISION SURFACE INDUCED  
BY THE THREE CLASSIFIERS

IN THE CASE OF MULTICLASS CLASSIFICATION COEF IS A TWODIMENSIONAL ARRAY OF SHAPENCLASSES  
NFEATURES ANDINTERCEPT IS A ONEDIMENSIONAL ARRAY OF SHAPENCLASSES THE ITH ROW OF COEF  
HOLDS THE WEIGHT VECTOR OF THE OV A CLASSIFIER FOR THE ITH CLASS CLASSES ARE INDEXED IN ASCENDING ORDER SEE AT  
TRIBUTECLASSES NOTE THAT IN PRINCIPLE SINCE THEY ALLOW TO CREATE A PROBABILITY MODEL LOSSLOG AND  
LOSSMODIFIEDHUBER ARE MORE SUITABLE FOR ONEVSALL CLASSIFICATION  
SGDCCLASSIFIER SUPPORTS BOTH WEIGHTED CLASSES AND WEIGHTED INSTANCES VIA THE FIT PARAMETERS CLASSWEIGHT  
ANDSAMPLEWEIGHT SEE THE EXAMPLES BELOW AND THE DOCSTRING OF SGDCCLASSIFIERFIT FOR FURTHER INFORMA  
TION

- EXAMPLES
- SGD MAXIMUM MARGIN SEPARATING HYPERPLANE
  - PLOT MULTICLASS SGD ON THE IRIS DATASET
  - SGD WEIGHTED SAMPLES
  - COMPARING VARIOUS ONLINE SOLVERS
- 31 SUPERVISED LEARNING 239

SCIKITLEARN USER GUIDE RELEASE 0213

•SVM SEPARATING HYPERPLANE FOR UNBALANCED CLASSES SEE THENOTE

SGDCLASSIFIER SUPPORTS AVERAGED SGD ASGD AVERAGING CAN BE ENABLED BY SETTING AVERAGETRUE

ASGD WORKS BY AVERAGING THE COEFFICIENTS OF THE PLAIN SGD OVER EACH ITERATION OVER A SAMPLE WHEN USING ASGD

THE LEARNING RATE CAN BE LARGER AND EVEN CONSTANT LEADING ON SOME DATASETS TO A SPEED UP IN TRAINING TIME

FOR CLASSIFICATION WITH A LOGISTIC LOSS ANOTHER VARIANT OF SGD WITH AN AVERAGING STRATEGY IS AVAILABLE WITH STOCHASTIC AVERAGE GRADIENT SAG ALGORITHM AVAILABLE AS A SOLVER IN LOGISTICREGRESSION

THE CLASSSGDREGRESSOR IMPLEMENTS A PLAIN STOCHASTIC GRADIENT DESCENT LEARNING ROUTINE WHICH SUPPORTS DIFFERENT LOSS FUNCTIONS AND PENALTIES TO FIT LINEAR REGRESSION MODELS SGDREGRESSOR IS WELL SUITED FOR REGRESSION PROBLEMS WITH A LARGE NUMBER OF TRAINING SAMPLES 10000 FOR OTHER PROBLEMS WE RECOMMEND RIDGE LASSO OR ELASTICNET

THE CONCRETE LOSS FUNCTION CAN BE SET VIA THE LOSS PARAMETERSGDREGRESSOR SUPPORTS THE FOLLOWING LOSS FUNCTIONS

- LOSSSQUAREDLOSS ORDINARY LEAST SQUARES
- LOSSHUBER HUBER LOSS FOR ROBUST REGRESSION
- LOSSEPSILONINSENSITIVE LINEAR SUPPORT VECTOR REGRESSION

THE HUBER AND EPSILONINSENSITIVE LOSS FUNCTIONS CAN BE USED FOR ROBUST REGRESSION THE WIDTH OF THE INSENSITIVE REGION HAS TO BE SPECIFIED VIA THE PARAMETER EPSILON THIS PARAMETER DEPENDS ON THE SCALE OF THE TARGET VARIABLES

SGDREGRESSOR SUPPORTS AVERAGED SGD AS SGDCLASSIFIER AVERAGING CAN BE ENABLED BY SETTING AVERAGETRUE

FOR REGRESSION WITH A SQUARED LOSS AND A L2 PENALTY ANOTHER VARIANT OF SGD WITH AN AVERAGING STRATEGY IS AVAILABLE WITH STOCHASTIC AVERAGE GRADIENT SAG ALGORITHM AVAILABLE AS A SOLVER IN RIDGE

STOCHASTIC GRADIENT DESCENT FOR SPARSE DATA

NOTE THE SPARSE IMPLEMENTATION PRODUCES SLIGHTLY DIFFERENT RESULTS THAN THE DENSE IMPLEMENTATION DUE TO A SHRUNK LEARNING RATE FOR THE INTERCEPT

THERE IS BUILTIN SUPPORT FOR SPARSE DATA GIVEN IN ANY MATRIX IN A FORMAT SUPPORTED BY SCIPYSPARSE FOR MAXIMUM EFFICIENCY HOWEVER USE THE CSR MATRIX FORMAT AS DEFINED IN SCIPYSPARSECSRMATRIX

EXAMPLES

•CLASSIFICATION OF TEXT DOCUMENTS USING SPARSE FEATURES

COMPLEXITY

THE MAJOR ADVANTAGE OF SGD IS ITS EFFICIENCY WHICH IS BASICALLY LINEAR IN THE NUMBER OF TRAINING EXAMPLES IF X IS A MATRIX OF SIZE N P TRAINING HAS A COST OF  $O(knp)$  WHERE K IS THE NUMBER OF ITERATIONS EPOCHS AND  $p$  IS THE AVERAGE NUMBER OF NONZERO ATTRIBUTES PER SAMPLE

RECENT THEORETICAL RESULTS HOWEVER SHOW THAT THE RUNTIME TO GET SOME DESIRED OPTIMIZATION ACCURACY DOES NOT INCREASE AS THE TRAINING SET SIZE INCREASES

STOPPING CRITERION

THE CLASSES SGDCLASSIFIER ANDSGDREGRESSOR PROVIDE TWO CRITERIA TO STOP THE ALGORITHM WHEN A GIVEN LEVEL OF CONVERGENCE IS REACHED

- WITHEARLYSTOPPINGTRUE THE INPUT DATA IS SPLIT INTO A TRAINING SET AND A VALIDATION SET THE MODEL IS THEN FITTED ON THE TRAINING SET AND THE STOPPING CRITERION IS BASED ON THE PREDICTION SCORE COMPUTED ON THE VALIDATION SET THE SIZE OF THE VALIDATION SET CAN BE CHANGED WITH THE PARAMETER VALIDATIONFRACTION
- WITHEARLYSTOPPINGFALSE THE MODEL IS FITTED ON THE ENTIRE INPUT DATA AND THE STOPPING CRITERION IS BASED ON THE OBJECTIVE FUNCTION COMPUTED ON THE INPUT DATA

IN BOTH CASES THE CRITERION IS EVALUATED ONCE BY EPOCH AND THE ALGORITHM STOPS WHEN THE CRITERION DOES NOT IMPROVE NITERNOCHANGE TIMES IN A ROW THE IMPROVEMENT IS EVALUATED WITH A TOLERANCE TOL AND THE ALGORITHM STOPS IN ANY CASE AFTER A MAXIMUM NUMBER OF ITERATION MAXITER

TIPS ON PRACTICAL USE

- STOCHASTIC GRADIENT DESCENT IS SENSITIVE TO FEATURE SCALING SO IT IS HIGHLY RECOMMENDED TO SCALE YOUR DATA FOR EXAMPLE SCALE EACH ATTRIBUTE ON THE INPUT VECTOR X TO 01 OR 11 OR STANDARDIZE IT TO HAVE MEAN 0 AND VARIANCE 1 NOTE THAT THE SAME SCALING MUST BE APPLIED TO THE TEST VECTOR TO OBTAIN MEANINGFUL RESULTS THIS CAN BE EASILY DONE USING STANDARDSCALER

FROM SKLEARNPREPROCESSING IMPORT STANDARDSCALER

SCALER STANDARDSCALER

SCALERFITXTRAIN DONT CHEAT FIT ONLY ON TRAINING DATA

XTRAIN SCALERTRANSFORMXTRAIN

XTEST SCALERTRANSFORMXTEST APPLY SAME TRANSFORMATION TO TEST DATA

IF YOUR ATTRIBUTES HAVE AN INTRINSIC SCALE EG WORD FREQUENCIES OR INDICATOR FEATURES SCALING IS NOT NEEDED

- FINDING A REASONABLE REGULARIZATION TERM  $\lambda$  IS BEST DONE USING GRIDSEARCHCV USUALLY IN THE RANGE 10 ONPARANGE17

- EMPIRICALLY WE FOUND THAT SGD CONVERGES AFTER OBSERVING APPROX 106 TRAINING SAMPLES THUS A REASONABLE FIRST GUESS FOR THE NUMBER OF ITERATIONS IS MAXITER NPCEIL10 6 N WHERE N IS THE SIZE OF THE TRAINING SET

- IF YOU APPLY SGD TO FEATURES EXTRACTED USING PCA WE FOUND THAT IT IS OFTEN WISE TO SCALE THE FEATURE VALUES BY SOME CONSTANT CSUCH THAT THE AVERAGE L2 NORM OF THE TRAINING DATA EQUALS ONE

- WE FOUND THAT AVERAGED SGD WORKS BEST WITH A LARGER NUMBER OF FEATURES AND A HIGHER ETA0

REFERENCES

- “EFFICIENT BACKPROP” Y LECUN L BOTTOU G ORR K MÜLLER IN NEURAL NETWORKS TRICKS OF THE TRADE 1998

MATHEMATICAL FORMULATION

GIVEN A SET OF TRAINING EXAMPLES  $\{x_i, y_i\}_{i=1}^n$  WHERE  $x_i \in \mathbb{R}^d$  AND  $y_i \in \{-1, 1\}$  OUR GOAL IS TO LEARN A LINEAR SCORING FUNCTION  $f(x) = w^T x + b$  WITH MODEL PARAMETERS  $w \in \mathbb{R}^d$  AND INTERCEPT  $b \in \mathbb{R}$  IN ORDER TO MAKE PREDICTIONS WE SIMPLY LOOK AT THE SIGN OF  $f(x)$  A COMMON CHOICE TO FIND THE MODEL PARAMETERS IS BY MINIMIZING THE REGULARIZED TRAINING ERROR GIVEN BY

$$J(w, b) = \frac{1}{n} \sum_{i=1}^n \max(0, 1 - y_i f(x_i)) + \frac{\lambda}{2} \|w\|^2$$

$$\|w\|^2 = w^T w$$

$$\lambda \geq 0$$

WHERE  $\ell$  IS A LOSS FUNCTION THAT MEASURES MODEL MISFIT AND  $\lambda$  IS A REGULARIZATION TERM AKA PENALTY THAT PENALIZES MODEL COMPLEXITY  $\lambda \geq 0$  IS A NONNEGATIVE HYPERPARAMETER DIFFERENT CHOICES FOR  $\ell$  ENTAIL DIFFERENT CLASSIFIERS SUCH AS

- HINGE SOFTMARGIN SUPPORT VECTOR MACHINES
- LOG LOGISTIC REGRESSION
- LEASTSQUARES RIDGE REGRESSION
- EPSILONINSENSITIVE SOFTMARGIN SUPPORT VECTOR REGRESSION

ALL OF THE ABOVE LOSS FUNCTIONS CAN BE REGARDED AS AN UPPER BOUND ON THE MISCLASSIFICATION ERROR ZEROONE LOSS AS SHOWN IN THE FIGURE BELOW

POPULAR CHOICES FOR THE REGULARIZATION TERM  $\lambda$  INCLUDE

- L2 NORM  $\lambda \sum_{i=1}^n w_i^2$

$\lambda \sum_{i=1}^n |w_i|$

- L1 NORM  $\lambda \sum_{i=1}^n |w_i|$

$\lambda \sum_{i=1}^n |w_i|$  WHICH LEADS TO SPARSE SOLUTIONS

- ELASTIC NET  $\lambda \sum_{i=1}^n (w_i^2 + |w_i|)$

$\lambda \sum_{i=1}^n (w_i^2 + |w_i|)$

$\lambda \sum_{i=1}^n (1 - |w_i|)$

$\lambda \sum_{i=1}^n (1 - |w_i|)$  A CONVEX COMBINATION OF L2 AND L1 WHERE  $\lambda$  IS GIVEN

BY  $\lambda = \lambda_1 \lambda_2$

THE FIGURE BELOW SHOWS THE CONTOURS OF THE DIFFERENT REGULARIZATION TERMS IN THE PARAMETER SPACE WHEN  $\lambda_1 = 1$  SGD

STOCHASTIC GRADIENT DESCENT IS AN OPTIMIZATION METHOD FOR UNCONSTRAINED OPTIMIZATION PROBLEMS IN CONTRAST TO BATCH GRADIENT DESCENT SGD APPROXIMATES THE TRUE GRADIENT OF  $\ell$  BY CONSIDERING A SINGLE TRAINING EXAMPLE AT A TIME THE CLASSSGDClassifier IMPLEMENTS A FIRSTORDER SGD LEARNING ROUTINE THE ALGORITHM ITERATES OVER THE TRAINING





EXAMPLES AND FOR EACH EXAMPLE UPDATES THE MODEL PARAMETERS ACCORDING TO THE UPDATE RULE GIVEN BY

$$\theta \leftarrow \theta - \eta \nabla_{\theta} J(\theta)$$
$$\theta_0 \leftarrow \theta_0 - \eta \nabla_{\theta_0} J(\theta)$$

WHERE  $\eta$  IS THE LEARNING RATE WHICH CONTROLS THE STEPSIZE IN THE PARAMETER SPACE THE INTERCEPT  $\theta_0$  IS UPDATED SIMILARLY BUT WITHOUT REGULARIZATION

THE LEARNING RATE  $\eta$  CAN BE EITHER CONSTANT OR GRADUALLY DECAYING FOR CLASSIFICATION THE DEFAULT LEARNING RATE SCHEDULE LEARNINGRATEOPTIMAL IS GIVEN BY

$$\eta = \frac{1}{\sqrt{t}}$$
$$\eta = \frac{1}{\sqrt{t}}$$
WHERE  $t$  IS THE TIME STEP THERE ARE A TOTAL OF NSAMPLES NITER TIME STEPS  $t$  IS DETERMINED BASED ON A HEURISTIC PROPOSED BY LÉON BOTTOU SUCH THAT THE EXPECTED INITIAL UPDATES ARE COMPARABLE WITH THE EXPECTED SIZE OF THE WEIGHTS THIS ASSUMING THAT THE NORM OF THE TRAINING SAMPLES IS APPROX 1 THE EXACT DEFINITION CAN BE FOUND IN INITT IN BASESGD

FOR REGRESSION THE DEFAULT LEARNING RATE SCHEDULE IS INVERSE SCALING LEARNINGRATEINVSCALING GIVEN BY

$$\eta = \frac{1}{\sqrt{t}}$$
$$\eta = \frac{1}{\sqrt{t}}$$
WHERE  $\eta$  AND  $\eta_0$  ARE HYPERPARAMETERS CHOSEN BY THE USER VIA ETA0 ANDPOWER0 RESP

FOR A CONSTANT LEARNING RATE USE LEARNINGRATECONSTANT AND USEETA0 TO SPECIFY THE LEARNING RATE  
FOR AN ADAPTIVELY DECREASING LEARNING RATE USE LEARNINGRATEADAPTIVE AND USEETA0 TO SPECIFY THE START  
ING LEARNING RATE WHEN THE STOPPING CRITERION IS REACHED THE LEARNING RATE IS DIVIDED BY 5 AND THE ALGORITHM DOES NOT  
STOP THE ALGORITHM STOPS WHEN THE LEARNING RATE GOES BELOW 1E6

THE MODEL PARAMETERS CAN BE ACCESSED THROUGH THE MEMBERS COEF ANDINTERCEPT

- MEMBERCOEF HOLDS THE WEIGHTS  $\theta$
- MEMBERINTERCEPT HOLDS  $\theta_0$

REFERENCES

- “SOLVING LARGE SCALE LINEAR PREDICTION PROBLEMS USING STOCHASTIC GRADIENT DESCENT ALGORITHMS” T ZHANG IN PROCEEDINGS OF ICML ‘04
- “REGULARIZATION AND VARIABLE SELECTION VIA THE ELASTIC NET” H ZOU T HASTIE JOURNAL OF THE ROYAL STATISTICAL SOCIETY SERIES B 67 2 301320
- “TOWARDS OPTIMAL ONE PASS LARGE SCALE LEARNING WITH AVERAGED STOCHASTIC GRADIENT DESCENT” XU WEI

IMPLEMENTATION DETAILS

THE IMPLEMENTATION OF SGD IS INFLUENCED BY THE STOCHASTIC GRADIENT SVM OF LÉON BOTTOU SIMILAR TO SVMMSGD THE WEIGHT VECTOR IS REPRESENTED AS THE PRODUCT OF A SCALAR AND A VECTOR WHICH ALLOWS AN EFFICIENT WEIGHT UPDATE IN THE CASE OF L2 REGULARIZATION IN THE CASE OF SPARSE FEATURE VECTORS THE INTERCEPT IS UPDATED WITH A SMALLER LEARNING RATE MULTIPLIED BY 0.01 TO ACCOUNT FOR THE FACT THAT IT IS UPDATED MORE FREQUENTLY TRAINING EXAMPLES ARE PICKED UP SEQUENTIALLY AND THE LEARNING RATE IS LOWERED AFTER EACH OBSERVED EXAMPLE WE ADOPTED THE LEARNING RATE SCHEDULE FROM SHALEVSHWARTZ ET AL 2007 FOR MULTICLASS CLASSIFICATION A “ONE VERSUS ALL” APPROACH IS USED WE USE THE TRUNCATED GRADIENT ALGORITHM PROPOSED BY TSURUOKA ET AL 2009 FOR L1 REGULARIZATION AND THE ELASTIC NET THE CODE IS WRITTEN IN CYTHON

SCIKITLEARN USER GUIDE RELEASE 0213

REFERENCES

- “STOCHASTIC GRADIENT DESCENT” L BOTTOU WEBSITE 2010
- “THE TRADEOFFS OF LARGE SCALE MACHINE LEARNING” L BOTTOU WEBSITE 2011
- “PEGASOS PRIMAL ESTIMATED SUBGRADIENT SOLVER FOR SVM” S SHALEVSHWARTZ Y SINGER N SREBRO IN PROCEEDINGS OF ICML ’07
- “STOCHASTIC GRADIENT DESCENT TRAINING FOR L1REGULARIZED LOGLINEAR MODELS WITH CUMULATIVE PENALTY” Y TSURUOKA J TSUJII S ANANIADOU IN PROCEEDINGS OF THE AFNLPACL ’09

316 NEAREST NEIGHBORS

SKLEARNNEIGHBORS PROVIDES FUNCTIONALITY FOR UNSUPERVISED AND SUPERVISED NEIGHBORSBASED LEARNING METHODS UNSUPERVISED NEAREST NEIGHBORS IS THE FOUNDATION OF MANY OTHER LEARNING METHODS NOTABLY MANIFOLD LEARNING AND SPECTRAL CLUSTERING SUPERVISED NEIGHBORSBASED LEARNING COMES IN TWO FLAVORS CLASSIFICATION FOR DATA WITH DISCRETE LABELS AND REGRESSION FOR DATA WITH CONTINUOUS LABELS

THE PRINCIPLE BEHIND NEAREST NEIGHBOR METHODS IS TO FIND A PREDEFINED NUMBER OF TRAINING SAMPLES CLOSEST IN DISTANCE TO THE NEW POINT AND PREDICT THE LABEL FROM THESE THE NUMBER OF SAMPLES CAN BE A USERDEFINED CONSTANT KNEAREST NEIGHBOR LEARNING OR VARY BASED ON THE LOCAL DENSITY OF POINTS RADIUSBASED NEIGHBOR LEARNING THE DISTANCE CAN IN GENERAL BE ANY METRIC MEASURE STANDARD EUCLIDEAN DISTANCE IS THE MOST COMMON CHOICE NEIGHBORSBASED METHODS ARE KNOWN AS NONGENERALIZING MACHINE LEARNING METHODS SINCE THEY SIMPLY “REMEMBER” ALL OF ITS TRAINING DATA POSSIBLY TRANSFORMED INTO A FAST INDEXING STRUCTURE SUCH AS A BALL TREE OR KD TREE

DESPITE ITS SIMPLICITY NEAREST NEIGHBORS HAS BEEN SUCCESSFUL IN A LARGE NUMBER OF CLASSIFICATION AND REGRESSION PROBLEMS INCLUDING HANDWRITTEN DIGITS AND SATELLITE IMAGE SCENES BEING A NONPARAMETRIC METHOD IT IS OFTEN SUCCESSFUL IN CLASSIFICATION SITUATIONS WHERE THE DECISION BOUNDARY IS VERY IRREGULAR

THE CLASSES IN SKLEARNNEIGHBORS CAN HANDLE EITHER NUMPY ARRAYS OR SCIPYSPARSE MATRICES AS INPUT FOR DENSE MATRICES A LARGE NUMBER OF POSSIBLE DISTANCE METRICS ARE SUPPORTED FOR SPARSE MATRICES ARBITRARY MINKOWSKI METRICS ARE SUPPORTED FOR SEARCHES

THERE ARE MANY LEARNING ROUTINES WHICH RELY ON NEAREST NEIGHBORS AT THEIR CORE ONE EXAMPLE IS KERNEL DENSITY ESTIMATION DISCUSSED IN THE DENSITY ESTIMATION SECTION

UNSUPERVISED NEAREST NEIGHBORS

NEARESTNEIGHBORS IMPLEMENTS UNSUPERVISED NEAREST NEIGHBORS LEARNING IT ACTS AS A UNIFORM INTERFACE TO THREE DIFFERENT NEAREST NEIGHBORS ALGORITHMS BALLTREE KDTREE AND A BRUTEFORCE ALGORITHM BASED ON ROUTINES IN SKLEARNMETRICSPAIRWISE THE CHOICE OF NEIGHBORS SEARCH ALGORITHM IS CONTROLLED THROUGH THE KEYWORD ALGORITHM WHICH MUST BE ONE OF AUTO BALLTREE KDTREE BRUTE WHEN THE DEFAULT VALUE AUTO IS PASSED THE ALGORITHM ATTEMPTS TO DETERMINE THE BEST APPROACH FROM THE TRAINING DATA FOR A

DISCUSSION OF THE STRENGTHS AND WEAKNESSES OF EACH OPTION SEE NEAREST NEIGHBOR ALGORITHMS WARNING REGARDING THE NEAREST NEIGHBORS ALGORITHMS IF TWO NEIGHBORS  $x_1$  AND  $x_2$  HAVE IDENTICAL DISTANCES BUT DIFFERENT LABELS THE RESULT WILL DEPEND ON THE ORDERING OF THE TRAINING DATA

31 SUPERVISED LEARNING 245

SCIKITLEARN USER GUIDE RELEASE 0213

FINDING THE NEAREST NEIGHBORS

FOR THE SIMPLE TASK OF FINDING THE NEAREST NEIGHBORS BETWEEN TWO SETS OF DATA THE UNSUPERVISED ALGORITHMS WITHIN SKLEARNNEIGHBORS CAN BE USED

```
FROM SKLEARNNEIGHBORS IMPORT NEARESTNEIGHBORS
IMPORT NUMPY AS NP
X NPARRAY1 1 2 1 3 2 1 1 2 1 3 2
NBRS NEARESTNEIGHBORSNNEIGHBORS2 ALGORITHMBALLTREEFITX
DISTANCES INDICES NBRKNEIGHBORSX
INDICES
ARRAY0 1
1 0
2 1
3 4
4 3
5 4
DISTANCES
ARRAY0 1
0 1
0 141421356
0 1
0 1
0 141421356
```

BECAUSE THE QUERY SET MATCHES THE TRAINING SET THE NEAREST NEIGHBOR OF EACH POINT IS THE POINT ITSELF AT A DISTANCE OF ZERO

IT IS ALSO POSSIBLE TO EFFICIENTLY PRODUCE A SPARSE GRAPH SHOWING THE CONNECTIONS BETWEEN NEIGHBORING POINTS

```
NBRKNEIGHBORSGRAPHXTOARRAY
ARRAY1 1 0 0 0 0
1 1 0 0 0 0
0 1 1 0 0 0
0 0 0 1 1 0
0 0 0 1 1 0
0 0 0 0 1 1
```

THE DATASET IS STRUCTURED SUCH THAT POINTS NEARBY IN INDEX ORDER ARE NEARBY IN PARAMETER SPACE LEADING TO AN APPROXIMATELY BLOCKDIAGONAL MATRIX OF KNEAREST NEIGHBORS SUCH A SPARSE GRAPH IS USEFUL IN A VARIETY OF CIRCUMSTANCES WHICH MAKE USE OF SPATIAL RELATIONSHIPS BETWEEN POINTS FOR UNSUPERVISED LEARNING IN PARTICULAR SEESKLEARNMANIFOLDISOMAP SKLEARNMANIFOLDLOCALLYLINEAREMBEDDING ANDSKLEARNCLUSTERSPECTRALCLUSTERING

KDTREE AND BALLTREE CLASSES

ALTERNATIVELY ONE CAN USE THE KDTREE ORBALLTREE CLASSES DIRECTLY TO FIND NEAREST NEIGHBORS THIS IS THE FUNCTIONALITY WRAPPED BY THE NEARESTNEIGHBORS CLASS USED ABOVE THE BALL TREE AND KD TREE HAVE THE SAME INTERFACE

WE’LL SHOW AN EXAMPLE OF USING THE KD TREE HERE

```
FROM SKLEARNNEIGHBORS IMPORT KDTREE
IMPORT NUMPY AS NP
X NPARRAY1 1 2 1 3 2 1 1 2 1 3 2
KDT KDTREEX LEAFSIZE30 METRICEUCLIDEAN
KDTQUERYX K2 RETURNDISTANCE FALSE
ARRAY0 1
```

246 CHAPTER 3 USER GUIDE

1 0  
2 1  
3 4  
4 3  
5 4

REFER TO THE KD TREE AND BALL TREE CLASS DOCUMENTATION FOR MORE INFORMATION ON THE OPTIONS AVAILABLE FOR NEAREST NEIGHBORS SEARCHES INCLUDING SPECIFICATION OF QUERY STRATEGIES DISTANCE METRICS ETC FOR A LIST OF AVAILABLE METRICS SEE THE DOCUMENTATION OF THE DISTANCE METRIC CLASS

NEAREST NEIGHBORS CLASSIFICATION

NEIGHBORS BASED CLASSIFICATION IS A TYPE OF INSTANCE BASED LEARNING OR NON GENERALIZING LEARNING IT DOES NOT ATTEMPT TO CONSTRUCT A GENERAL INTERNAL MODEL BUT SIMPLY STORES INSTANCES OF THE TRAINING DATA CLASSIFICATION IS COMPUTED FROM A SIMPLE MAJORITY VOTE OF THE NEAREST NEIGHBORS OF EACH POINT A QUERY POINT IS ASSIGNED THE DATA CLASS WHICH HAS THE MOST REPRESENTATIVES WITHIN THE NEAREST NEIGHBORS OF THE POINT

SCIKITLEARN IMPLEMENTS TWO DIFFERENT NEAREST NEIGHBORS CLASSIFIERS KNEIGHBORS CLASSIFIER IMPLEMENTS LEARNING BASED ON THE  $k$  NEAREST NEIGHBORS OF EACH QUERY POINT WHERE  $k$  IS AN INTEGER VALUE SPECIFIED BY THE USER

RADIUS NEIGHBORS CLASSIFIER IMPLEMENTS LEARNING BASED ON THE NUMBER OF NEIGHBORS WITHIN A FIXED RADIUS

$r$  OF EACH TRAINING POINT WHERE  $r$  IS A FLOATING POINT VALUE SPECIFIED BY THE USER

THE  $k$  NEIGHBORS CLASSIFICATION IN KNEIGHBORS CLASSIFIER IS THE MOST COMMONLY USED TECHNIQUE THE OPTIMAL CHOICE OF THE VALUE  $k$  IS HIGHLY DATA DEPENDENT IN GENERAL A LARGER  $k$  SUPPRESSES THE EFFECTS OF NOISE BUT MAKES THE CLASSIFICATION BOUNDARIES LESS DISTINCT

IN CASES WHERE THE DATA IS NOT UNIFORMLY SAMPLED RADIUS BASED NEIGHBORS CLASSIFICATION IN

RADIUS NEIGHBORS CLASSIFIER CAN BE A BETTER CHOICE THE USER SPECIFIES A FIXED RADIUS  $r$  SUCH THAT

POINTS IN SPARSE NEIGHBORHOODS USE FEWER NEAREST NEIGHBORS FOR THE CLASSIFICATION FOR HIGH DIMENSIONAL PARAMETER SPACES THIS METHOD BECOMES LESS EFFECTIVE DUE TO THE SO CALLED "CURSE OF DIMENSIONALITY"

THE BASIC NEAREST NEIGHBORS CLASSIFICATION USES UNIFORM WEIGHTS THAT IS THE VALUE ASSIGNED TO A QUERY POINT IS COMPUTED

FROM A SIMPLE MAJORITY VOTE OF THE NEAREST NEIGHBORS UNDER SOME CIRCUMSTANCES IT IS BETTER TO WEIGHT THE NEIGHBORS

SUCH THAT NEARER NEIGHBORS CONTRIBUTE MORE TO THE FIT THIS CAN BE ACCOMPLISHED THROUGH THE WEIGHTS KEYWORD THE

DEFAULT VALUE WEIGHTS UNIFORM ASSIGNS UNIFORM WEIGHTS TO EACH NEIGHBOR WEIGHTS DISTANCE

ASSIGNS WEIGHTS PROPORTIONAL TO THE INVERSE OF THE DISTANCE FROM THE QUERY POINT ALTERNATIVELY A USER DEFINED FUNCTION

OF THE DISTANCE CAN BE SUPPLIED TO COMPUTE THE WEIGHTS

EXAMPLES

•NEAREST NEIGHBORS CLASSIFICATION AN EXAMPLE OF CLASSIFICATION USING NEAREST NEIGHBORS

NEAREST NEIGHBORS REGRESSION

NEIGHBORSBASED REGRESSION CAN BE USED IN CASES WHERE THE DATA LABELS ARE CONTINUOUS RATHER THAN DISCRETE VARIABLES THE LABEL ASSIGNED TO A QUERY POINT IS COMPUTED BASED ON THE MEAN OF THE LABELS OF ITS NEAREST NEIGHBORS

SCIKITLEARN IMPLEMENTS TWO DIFFERENT NEIGHBORS REGRESSORS KNEIGHBORSREGRESSOR IMPLEMENTS LEARNING

BASED ON THE  $k$  NEAREST NEIGHBORS OF EACH QUERY POINT WHERE  $k$  IS AN INTEGER VALUE SPECIFIED BY THE USER

RADIUSNEIGHBORSREGRESSOR IMPLEMENTS LEARNING BASED ON THE NEIGHBORS WITHIN A FIXED RADIUS  $r$  OF THE QUERY POINT WHERE  $r$  IS A FLOATINGPOINT VALUE SPECIFIED BY THE USER

THE BASIC NEAREST NEIGHBORS REGRESSION USES UNIFORM WEIGHTS THAT IS EACH POINT IN THE LOCAL NEIGHBORHOOD CONTRIBUTES UNIFORMLY TO THE CLASSIFICATION OF A QUERY POINT UNDER SOME CIRCUMSTANCES IT CAN BE ADVANTAGEOUS TO WEIGHT POINTS

SUCH THAT NEARBY POINTS CONTRIBUTE MORE TO THE REGRESSION THAN FARAWAY POINTS THIS CAN BE ACCOMPLISHED THROUGH THE WEIGHTS KEYWORD THE DEFAULT VALUE WEIGHTS UNIFORM ASSIGNS EQUAL WEIGHTS TO ALL POINTS WEIGHTS

DISTANCE ASSIGNS WEIGHTS PROPORTIONAL TO THE INVERSE OF THE DISTANCE FROM THE QUERY POINT ALTERNATIVELY A USERDEFINED FUNCTION OF THE DISTANCE CAN BE SUPPLIED WHICH WILL BE USED TO COMPUTE THE WEIGHTS

THE USE OF MULTIOUTPUT NEAREST NEIGHBORS FOR REGRESSION IS DEMONSTRATED IN FACE COMPLETION WITH A MULTIOUTPUT ESTIMATORS IN THIS EXAMPLE THE INPUTS X ARE THE PIXELS OF THE UPPER HALF OF FACES AND THE OUTPUTS Y ARE THE PIXELS OF THE LOWER HALF OF THOSE FACES

EXAMPLES

•NEAREST NEIGHBORS REGRESSION AN EXAMPLE OF REGRESSION USING NEAREST NEIGHBORS



•FACE COMPLETION WITH A MULTIOUTPUT ESTIMATORS AN EXAMPLE OF MULTIOUTPUT REGRESSION USING NEAREST NEIGHBORS

NEAREST NEIGHBOR ALGORITHMS

BRUTE FORCE

FAST COMPUTATION OF NEAREST NEIGHBORS IS AN ACTIVE AREA OF RESEARCH IN MACHINE LEARNING THE MOST NAIVE NEIGHBOR SEARCH IMPLEMENTATION INVOLVES THE BRUTEFORCE COMPUTATION OF DISTANCES BETWEEN ALL PAIRS OF POINTS IN THE DATASET FOR  $n$  SAMPLES IN  $d$  DIMENSIONS THIS APPROACH SCALES AS  $O(n^2d)$  EFFICIENT BRUTEFORCE NEIGHBORS SEARCHES CAN BE VERY COMPETITIVE FOR SMALL DATA SAMPLES HOWEVER AS THE NUMBER OF SAMPLES  $n$  GROWS THE BRUTEFORCE APPROACH QUICKLY BECOMES INFEASIBLE IN THE CLASSES WITHIN SKLEARNNEIGHBORS BRUTEFORCE NEIGHBORS SEARCHES ARE SPECIFIED USING THE KEYWORD ALGORITHM BRUTE AND ARE COMPUTED USING THE ROUTINES AVAILABLE IN SKLEARNMETRICS

PAIRWISE

KD TREE

TO ADDRESS THE COMPUTATIONAL INEFFICIENCIES OF THE BRUTEFORCE APPROACH A VARIETY OF TREEBASED DATA STRUCTURES HAVE BEEN INVENTED IN GENERAL THESE STRUCTURES ATTEMPT TO REDUCE THE REQUIRED NUMBER OF DISTANCE CALCULATIONS BY EFFICIENTLY ENCODING AGGREGATE DISTANCE INFORMATION FOR THE SAMPLE THE BASIC IDEA IS THAT IF POINT  $p_i$  IS VERY DISTANT FROM POINT  $p_j$  AND POINT  $p_j$  IS VERY CLOSE TO POINT  $p_k$  THEN WE KNOW THAT POINTS  $p_i$  AND  $p_k$  ARE VERY DISTANT WITHOUT HAVING TO EXPLICITLY CALCULATE THEIR DISTANCE IN THIS WAY THE COMPUTATIONAL COST OF A NEAREST NEIGHBORS SEARCH CAN BE REDUCED TO  $O(n \log n)$  OR BETTER THIS IS A SIGNIFICANT IMPROVEMENT OVER BRUTEFORCE FOR LARGE  $n$

AN EARLY APPROACH TO TAKING ADVANTAGE OF THIS AGGREGATE INFORMATION WAS THE KD TREE DATA STRUCTURE SHORT FOR K DIMENSIONAL TREE WHICH GENERALIZES TWODIMENSIONAL QUADTREES AND 3DIMENSIONAL OCTTREES TO AN ARBITRARY NUMBER OF DIMENSIONS THE KD TREE IS A BINARY TREE STRUCTURE WHICH RECURSIVELY PARTITIONS THE PARAMETER SPACE ALONG THE DATA AXES DIVIDING IT INTO NESTED ORTHOTROPIC REGIONS INTO WHICH DATA POINTS ARE FILED THE CONSTRUCTION OF A KD TREE IS VERY FAST BECAUSE PARTITIONING IS PERFORMED ONLY ALONG THE DATA AXES NO  $d$  DIMENSIONAL DISTANCES NEED TO BE COMPUTED ONCE CONSTRUCTED THE NEAREST NEIGHBOR OF A QUERY POINT CAN BE DETERMINED WITH ONLY  $O(\log n)$  DISTANCE COMPUTATIONS THOUGH THE KD TREE APPROACH IS VERY FAST FOR LOWDIMENSIONAL  $d \leq 20$  NEIGHBORS SEARCHES IT BECOMES INEFFICIENT AS  $d$  GROWS VERY LARGE THIS IS ONE MANIFESTATION OF THE SOCALLED “CURSE OF DIMENSIONALITY” IN SCIKITLEARN KD TREE NEIGHBORS SEARCHES ARE SPECIFIED USING THE KEYWORD ALGORITHM KDTREE AND ARE COMPUTED USING THE CLASS KDTREE

REFERENCES

- “MULTIDIMENSIONAL BINARY SEARCH TREES USED FOR ASSOCIATIVE SEARCHING” BENTLEY JL COMMUNICATIONS OF THE ACM 1975

BALL TREE

TO ADDRESS THE INEFFICIENCIES OF KD TREES IN HIGHER DIMENSIONS THE BALL TREE DATA STRUCTURE WAS DEVELOPED WHERE KD TREES PARTITION DATA ALONG CARTESIAN AXES BALL TREES PARTITION DATA IN A SERIES OF NESTING HYPERSPHERES THIS MAKES TREE CONSTRUCTION MORE COSTLY THAN THAT OF THE KD TREE BUT RESULTS IN A DATA STRUCTURE WHICH CAN BE VERY EFFICIENT ON HIGHLY STRUCTURED DATA EVEN IN VERY HIGH DIMENSIONS

A BALL TREE RECURSIVELY DIVIDES THE DATA INTO NODES DEFINED BY A CENTROID  $c$  AND RADIUS  $r$  SUCH THAT EACH POINT IN THE NODE LIES WITHIN THE HYPERSPHERE DEFINED BY  $c$  AND  $r$  THE NUMBER OF CANDIDATE POINTS FOR A NEIGHBOR SEARCH IS REDUCED

250 CHAPTER 3 USER GUIDE



SCIKITLEARN USER GUIDE RELEASE 0213  
THROUGH USE OF THE TRIANGLE INEQUALITY  
 $d(p, c) \leq d(p, q)$

WITH THIS SETUP A SINGLE DISTANCE CALCULATION BETWEEN A TEST POINT AND THE CENTROID IS SUFFICIENT TO DETERMINE A LOWER AND UPPER BOUND ON THE DISTANCE TO ALL POINTS WITHIN THE NODE BECAUSE OF THE SPHERICAL GEOMETRY OF THE BALL TREE NODES IT CAN OUTPERFORM A KDTREE IN HIGH DIMENSIONS THOUGH THE ACTUAL PERFORMANCE IS HIGHLY DEPENDENT ON THE STRUCTURE OF THE TRAINING DATA IN SCIKITLEARN BALLTREEBASED NEIGHBORS SEARCHES ARE SPECIFIED USING THE KEYWORD ALGORITHM BALLTREE AND ARE COMPUTED USING THE CLASS SKLEARNNEIGHBORSBALLTREE ALTERNATIVELY THE USER CAN WORK WITH THE BALLTREE CLASS DIRECTLY

REFERENCES

- “FIVE BALLTREE CONSTRUCTION ALGORITHMS” OMOHUNDRO SM INTERNATIONAL COMPUTER SCIENCE INSTITUTE TECHNICAL REPORT 1989

CHOICE OF NEAREST NEIGHBORS ALGORITHM

THE OPTIMAL ALGORITHM FOR A GIVEN DATASET IS A COMPLICATED CHOICE AND DEPENDS ON A NUMBER OF FACTORS

- NUMBER OF SAMPLES  $n$  IENSAMPLES AND DIMENSIONALITY  $d$  IENFEATURES

-BRUTE FORCE QUERY TIME GROWS AS  $n^2$

-BALL TREE QUERY TIME GROWS AS APPROXIMATELY  $n \log n$

-KD TREE QUERY TIME CHANGES WITH  $d$  IN A WAY THAT IS DIFFICULT TO PRECISELY CHARACTERISE FOR SMALL  $d$  LESS THAN 20 OR SO THE COST IS APPROXIMATELY  $n \log d$  AND THE KD TREE QUERY CAN BE VERY EFFICIENT FOR LARGER  $d$  THE COST INCREASES TO NEARLY  $n^2$  AND THE OVERHEAD DUE TO THE TREE STRUCTURE CAN LEAD TO QUERIES WHICH ARE SLOWER THAN BRUTE FORCE

FOR SMALL DATA SETS  $d$  LESS THAN 30 OR SO  $\log d$  IS COMPARABLE TO  $d$  AND BRUTE FORCE ALGORITHMS CAN BE MORE EFFICIENT THAN A TREEBASED APPROACH BOTH KDTREE ANDBALLTREE ADDRESS THIS THROUGH PROVIDING A LEAF SIZE PARAMETER THIS CONTROLS THE NUMBER OF SAMPLES AT WHICH A QUERY SWITCHES TO BRUTEFORCE THIS ALLOWS BOTH ALGORITHMS TO APPROACH THE EFFICIENCY OF A BRUTEFORCE COMPUTATION FOR SMALL  $d$

- DATA STRUCTURE INTRINSIC DIMENSIONALITY OF THE DATA ANDOR SPARSITY OF THE DATA INTRINSIC DIMENSIONALITY REFERS TO THE DIMENSION  $k \leq d$  OF A MANIFOLD ON WHICH THE DATA LIES WHICH CAN BE LINEARLY OR NONLINEARLY EMBEDDED IN THE PARAMETER SPACE SPARSITY REFERS TO THE DEGREE TO WHICH THE DATA FILLS THE PARAMETER SPACE THIS IS TO BE DISTINGUISHED FROM THE CONCEPT AS USED IN “SPARSE” MATRICES THE DATA MATRIX MAY HAVE NO ZERO ENTRIES BUT THE STRUCTURE CAN STILL BE “SPARSE” IN THIS SENSE

-BRUTE FORCE QUERY TIME IS UNCHANGED BY DATA STRUCTURE

-BALL TREE ANDKD TREE QUERY TIMES CAN BE GREATLY INFLUENCED BY DATA STRUCTURE IN GENERAL SPARSER DATA WITH A SMALLER INTRINSIC DIMENSIONALITY LEADS TO FASTER QUERY TIMES BECAUSE THE KD TREE INTERNAL REPRESENTATION IS ALIGNED WITH THE PARAMETER AXES IT WILL NOT GENERALLY SHOW AS MUCH IMPROVEMENT AS BALL TREE FOR ARBITRARILY STRUCTURED DATA

DATASETS USED IN MACHINE LEARNING TEND TO BE VERY STRUCTURED AND ARE VERY WELLSUITED FOR TREEBASED QUERIES

- NUMBER OF NEIGHBORS  $k$  REQUESTED FOR A QUERY POINT

-BRUTE FORCE QUERY TIME IS LARGELY UNAFFECTED BY THE VALUE OF  $k$

-BALL TREE ANDKD TREE QUERY TIME WILL BECOME SLOWER AS  $k$  INCREASES THIS IS DUE TO TWO EFFECTS FIRST A LARGER  $k$  LEADS TO THE NECESSITY TO SEARCH A LARGER PORTION OF THE PARAMETER SPACE SECOND USING  $k = 1$  REQUIRES INTERNAL QUEUEING OF RESULTS AS THE TREE IS TRAVERSED

SCIKITLEARN USER GUIDE RELEASE 0213

AS  $n$  BECOMES LARGE COMPARED TO  $n$  THE ABILITY TO PRUNE BRANCHES IN A TREEBASED QUERY IS REDUCED IN THIS SITUATION BRUTE FORCE QUERIES CAN BE MORE EFFICIENT

- NUMBER OF QUERY POINTS BOTH THE BALL TREE AND THE KD TREE REQUIRE A CONSTRUCTION PHASE THE COST OF THIS CONSTRUCTION BECOMES NEGLIGIBLE WHEN AMORTIZED OVER MANY QUERIES IF ONLY A SMALL NUMBER OF QUERIES WILL BE PERFORMED HOWEVER THE CONSTRUCTION CAN MAKE UP A SIGNIFICANT FRACTION OF THE TOTAL COST IF VERY FEW QUERY POINTS WILL BE REQUIRED BRUTE FORCE IS BETTER THAN A TREEBASED METHOD

CURRENTLY ALGORITHM AUTO SELECTS BRUTE IF  $n \leq 2$  THE INPUT DATA IS SPARSE OR EFFECTIVEMETRIC ISN'T IN THE VALIDMETRICS LIST FOR EITHER KDTREE OR BALLTREE OTHERWISE IT SELECTS THE FIRST OUT OF KDTREE AND BALLTREE THAT HASEFFECTIVEMETRIC IN ITSVALIDMETRICS LIST THIS CHOICE IS BASED ON THE ASSUMPTION THAT THE NUMBER OF QUERY POINTS IS AT LEAST THE SAME ORDER AS THE NUMBER OF TRAINING POINTS AND THAT LEAFSIZE IS CLOSE TO ITS DEFAULT VALUE OF 30

EFFECT OF LEAFSIZE

AS NOTED ABOVE FOR SMALL SAMPLE SIZES A BRUTE FORCE SEARCH CAN BE MORE EFFICIENT THAN A TREEBASED QUERY THIS FACT IS ACCOUNTED FOR IN THE BALL TREE AND KD TREE BY INTERNALLY SWITCHING TO BRUTE FORCE SEARCHES WITHIN LEAF NODES THE LEVEL OF THIS SWITCH CAN BE SPECIFIED WITH THE PARAMETER LEAFSIZE THIS PARAMETER CHOICE HAS MANY EFFECTS

CONSTRUCTION TIME A LARGER LEAFSIZE LEADS TO A FASTER TREE CONSTRUCTION TIME BECAUSE FEWER NODES NEED TO BE CREATED

QUERY TIME BOTH A LARGE OR SMALL LEAFSIZE CAN LEAD TO SUBOPTIMAL QUERY COST FOR LEAFSIZE APPROACHING 1 THE OVERHEAD INVOLVED IN TRAVERSING NODES CAN SIGNIFICANTLY SLOW QUERY TIMES FOR LEAFSIZE APPROACHING THE SIZE OF THE TRAINING SET QUERIES BECOME ESSENTIALLY BRUTE FORCE A GOOD COMPROMISE BETWEEN THESE IS LEAFSIZE 30 THE DEFAULT VALUE OF THE PARAMETER

MEMORY AS LEAFSIZE INCREASES THE MEMORY REQUIRED TO STORE A TREE STRUCTURE DECREASES THIS IS ESPECIALLY IMPORTANT IN THE CASE OF BALL TREE WHICH STORES A  $n$ -DIMENSIONAL CENTROID FOR EACH NODE THE REQUIRED STORAGE SPACE FOR BALLTREE IS APPROXIMATELY  $1 \cdot \text{LEAFSIZE}$  TIMES THE SIZE OF THE TRAINING SET

LEAFSIZE IS NOT REFERENCED FOR BRUTE FORCE QUERIES

NEAREST CENTROID CLASSIFIER

THE NEAREST CENTROID CLASSIFIER IS A SIMPLE ALGORITHM THAT REPRESENTS EACH CLASS BY THE CENTROID OF ITS MEMBERS IN EFFECT THIS MAKES IT SIMILAR TO THE LABEL UPDATING PHASE OF THE SKLEARN KMEANS ALGORITHM IT ALSO HAS NO PARAMETERS TO CHOOSE MAKING IT A GOOD BASELINE CLASSIFIER IT DOES HOWEVER SUFFER ON NONCONVEX CLASSES AS WELL AS WHEN CLASSES HAVE DRASTICALLY DIFFERENT VARIANCES AS EQUAL VARIANCE IN ALL DIMENSIONS IS ASSUMED SEE LINEAR DISCRIMINANT ANALYSIS SKLEARN DISCRIMINANT ANALYSIS LINEAR DISCRIMINANT ANALYSIS AND QUADRATIC DISCRIMINANT ANALYSIS SKLEARN DISCRIMINANT ANALYSIS QUADRATIC DISCRIMINANT ANALYSIS FOR MORE COMPLEX METHODS THAT DO NOT MAKE THIS ASSUMPTION USAGE OF THE DEFAULT NEAREST CENTROID IS SIMPLE

```
FROM SKLEARN NEIGHBORS NEARESTCENTROID IMPORT NEARESTCENTROID
import numpy as np
X = np.array([1, 2, 1, 3, 2, 1, 1, 2, 1, 3, 2])
Y = np.array([1, 1, 1, 2, 2, 2])
clf = NearestCentroid()
clf.fit(X, Y)
nearest_centroid_metric = euclidean
shrink_threshold = None
print(clf.predict([0, 1]))
```

1

252 CHAPTER 3 USER GUIDE

NEAREST SHRUNKEN CENTROID

THE NEAREST CENTROID CLASSIFIER HAS A SHRINKTHRESHOLD PARAMETER WHICH IMPLEMENTS THE NEAREST SHRUNKEN CENTROID CLASSIFIER. IN EFFECT THE VALUE OF EACH FEATURE FOR EACH CENTROID IS DIVIDED BY THE WITHINCLASS VARIANCE OF THAT FEATURE. THE FEATURE VALUES ARE THEN REDUCED BY SHRINKTHRESHOLD. MOST NOTABLY IF A PARTICULAR FEATURE VALUE CROSSES ZERO IT IS SET TO ZERO. IN EFFECT THIS REMOVES THE FEATURE FROM AFFECTING THE CLASSIFICATION. THIS IS USEFUL FOR EXAMPLE FOR REMOVING NOISY FEATURES.

IN THE EXAMPLE BELOW USING A SMALL SHRINK THRESHOLD INCREASES THE ACCURACY OF THE MODEL FROM 0.81 TO 0.82.

EXAMPLES

• NEAREST CENTROID CLASSIFICATION: AN EXAMPLE OF CLASSIFICATION USING NEAREST CENTROID WITH DIFFERENT SHRINK THRESHOLDS.

NEIGHBORHOOD COMPONENTS ANALYSIS

NEIGHBORHOOD COMPONENTS ANALYSIS (NCA) NEIGHBORHOOD COMPONENTS ANALYSIS IS A DISTANCE METRIC LEARNING ALGORITHM WHICH AIMS TO IMPROVE THE ACCURACY OF NEAREST NEIGHBORS CLASSIFICATION COMPARED TO THE STANDARD EUCLIDEAN DISTANCE. THE ALGORITHM DIRECTLY MAXIMIZES A STOCHASTIC VARIANT OF THE LEAVEONEOUT K-NEAREST NEIGHBORS (KNN) SCORE ON THE TRAINING SET. IT CAN ALSO LEARN A LOW-DIMENSIONAL LINEAR PROJECTION OF DATA THAT CAN BE USED FOR DATA VISUALIZATION AND FAST CLASSIFICATION.

SCIKITLEARN USER GUIDE RELEASE 0213

IN THE ABOVE ILLUSTRATING FIGURE WE CONSIDER SOME POINTS FROM A RANDOMLY GENERATED DATASET WE FOCUS ON THE STOCHASTIC KNN CLASSIFICATION OF POINT NO 3 THE THICKNESS OF A LINK BETWEEN SAMPLE 3 AND ANOTHER POINT IS PROPORTIONAL TO THEIR DISTANCE AND CAN BE SEEN AS THE RELATIVE WEIGHT OR PROBABILITY THAT A STOCHASTIC NEAREST NEIGHBOR PREDICTION RULE WOULD ASSIGN TO THIS POINT IN THE ORIGINAL SPACE SAMPLE 3 HAS MANY STOCHASTIC NEIGHBORS FROM VARIOUS CLASSES SO THE RIGHT CLASS IS NOT VERY LIKELY HOWEVER IN THE PROJECTED SPACE LEARNED BY NCA THE ONLY STOCHASTIC NEIGHBORS WITH NON NEGLIGIBLE WEIGHT ARE FROM THE SAME CLASS AS SAMPLE 3 GUARANTEEING THAT THE LATTER WILL BE WELL CLASSIFIED SEE THE MATHEMATICAL FORMULATION FOR MORE DETAILS

CLASSIFICATION

COMBINED WITH A NEAREST NEIGHBORS CLASSIFIER KNEIGHBORSCLASSIFIER NCA IS ATTRACTIVE FOR CLASSIFICATION BE CAUSE IT CAN NATURALLY HANDLE MULTICLASS PROBLEMS WITHOUT ANY INCREASE IN THE MODEL SIZE AND DOES NOT INTRODUCE ADDITIONAL PARAMETERS THAT REQUIRE FINETUNING BY THE USER

NCA CLASSIFICATION HAS BEEN SHOWN TO WORK WELL IN PRACTICE FOR DATA SETS OF VARYING SIZE AND DIFFICULTY IN CONTRAST TO RELATED METHODS SUCH AS LINEAR DISCRIMINANT ANALYSIS NCA DOES NOT MAKE ANY ASSUMPTIONS ABOUT THE CLASS DISTRIBUTIONS THE NEAREST NEIGHBOR CLASSIFICATION CAN NATURALLY PRODUCE HIGHLY IRREGULAR DECISION BOUNDARIES

TO USE THIS MODEL FOR CLASSIFICATION ONE NEEDS TO COMBINE A NEIGHBORHOODCOMPONENTSANALYSIS INSTANCE THAT LEARNS THE OPTIMAL TRANSFORMATION WITH A KNEIGHBORSCLASSIFIER INSTANCE THAT PERFORMS THE CLASSIFICATION IN THE PROJECTED SPACE HERE IS AN EXAMPLE USING THE TWO CLASSES

```
FROM SKLEARNNEIGHBORS IMPORT NEIGHBORHOODCOMPONENTSANALYSIS
KNEIGHBORSCLASSIFIER
FROM SKLEARNDATASETS IMPORT LOADIRIS
FROM SKLEARNMODELSELECTION IMPORT TRAINTESTSPLIT
FROM SKLEARNPIPELINE IMPORT PIPELINE
X Y LOADIRISRETURNXY TRUE
XTRAIN XTEST YTRAIN YTEST TRAINTESTSPLITX Y
STRATIFY TESTSIZE07 RANDOMSTATE42
NCA NEIGHBORHOODCOMPONENTSANALYSISRANDOMSTATE42
KNN KNEIGHBORSCLASSIFIERNNEIGHBORS3
NCAPIPE PIPELINENCA NCA KNN KNN
NCAPIPEFITXTRAIN YTRAIN
PIPELINE
PRINTNCAPIPESCOREXTEST YTEST
096190476
```

THE PLOT SHOWS DECISION BOUNDARIES FOR NEAREST NEIGHBOR CLASSIFICATION AND NEIGHBORHOOD COMPONENTS ANALYSIS CLASSIFICATION ON THE IRIS DATASET WHEN TRAINING AND SCORING ON ONLY TWO FEATURES FOR VISUALISATION PURPOSES

254 CHAPTER 3 USER GUIDE

DIMENSIONALITY REDUCTION

NCA CAN BE USED TO PERFORM SUPERVISED DIMENSIONALITY REDUCTION THE INPUT DATA ARE PROJECTED ONTO A LINEAR SUB SPACE CONSISTING OF THE DIRECTIONS WHICH MINIMIZE THE NCA OBJECTIVE THE DESIRED DIMENSIONALITY CAN BE SET USING THE PARAMETER NCOMPONENTS FOR INSTANCE THE FOLLOWING FIGURE SHOWS A COMPARISON OF DIMENSIONALITY REDUCTION WITH PRINCIPAL COMPONENT ANALYSIS SKLEARNDECOMPOSITIONPCA LINEAR DISCRIMINANT ANALYSIS SKLEARNDISCRIMINANTANALYSISLINEARDISCRIMINANTANALYSIS AND NEIGHBORHOOD COMPONENT ANALYSIS NEIGHBORHOODCOMPONENTSANALYSIS ON THE DIGITS DATASET A DATASET WITH SIZE 1797 AND 64 THE DATA SET IS SPLIT INTO A TRAINING AND A TEST SET OF EQUAL SIZE THEN STANDARDIZED FOR EVALUATION THE 3NEAREST NEIGHBOR CLASSIFICATION ACCURACY IS COMPUTED ON THE 2DIMENSIONAL PROJECTED POINTS FOUND BY EACH METHOD EACH DATA SAMPLE BELONGS TO ONE OF 10 CLASSES

EXAMPLES

- COMPARING NEAREST NEIGHBORS WITH AND WITHOUT NEIGHBORHOOD COMPONENTS ANALYSIS
- DIMENSIONALITY REDUCTION WITH NEIGHBORHOOD COMPONENTS ANALYSIS
- MANIFOLD LEARNING ON HANDWRITTEN DIGITS LOCALLY LINEAR EMBEDDING ISOMAP

MATHEMATICAL FORMULATION

THE GOAL OF NCA IS TO LEARN AN OPTIMAL LINEAR TRANSFORMATION MATRIX OF SIZE NCOMPONENTS NFEATURES WHICH MAXIMISES THE SUM OVER ALL SAMPLES OF THE PROBABILITY THATIS CORRECTLY CLASSIFIED IE

ARG MAX

$$\sum_{i=1}^n \sum_{j=1}^K p_{ij}$$

WITH NSAMPLES AND THE PROBABILITY OF SAMPLE BEING CORRECTLY CLASSIFIED ACCORDING TO A STOCHASTIC NEAREST NEIGHBORS RULE IN THE LEARNED EMBEDDED SPACE

$$p_{ij} = \frac{1}{\sum_{k=1}^K \exp(-\frac{\|x_i - x_{j,k}\|^2}{2\sigma^2})}$$

$$x_i \in \mathbb{R}^D$$

WHERE IS THE SET OF POINTS IN THE SAME CLASS AS SAMPLE AND IS THE SOFTMAX OVER EUCLIDEAN DISTANCES IN THE EMBEDDED SPACE

$$p_{ij} = \frac{\exp(-\frac{\|x_i - x_{j,k}\|^2}{2\sigma^2})}{\sum_{k=1}^K \exp(-\frac{\|x_i - x_{j,k}\|^2}{2\sigma^2})}$$

$$\sigma^2 = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^K \|x_i - x_{j,k}\|^2$$

SCIKITLEARN USER GUIDE RELEASE 0213

MAHALANOBIS DISTANCE

NCA CAN BE SEEN AS LEARNING A SQUARED MAHALANOBIS DISTANCE METRIC

$\mathbf{K} = \mathbf{K}^2 \mathbf{K} - \mathbf{K} \mathbf{K} \mathbf{K} - \mathbf{K}$

WHERE  $\mathbf{K}$  IS A SYMMETRIC POSITIVE SEMIDEFINITE MATRIX OF SIZE  $N \times N$

IMPLEMENTATION

THIS IMPLEMENTATION FOLLOWS WHAT IS EXPLAINED IN THE ORIGINAL PAPER<sup>1</sup> FOR THE OPTIMISATION METHOD IT CURRENTLY USES SCIPY'S LBFGSB WITH A FULL GRADIENT COMPUTATION AT EACH ITERATION TO AVOID TO TUNE THE LEARNING RATE AND PROVIDE STABLE LEARNING

SEE THE EXAMPLES BELOW AND THE DOCSTRING OF NEIGHBORHOODCOMPONENTSANALYSISFIT FOR FURTHER INFORMATION

COMPLEXITY

TRAINING

NCA STORES A MATRIX OF PAIRWISE DISTANCES TAKING  $N \times N$  MEMORY TIME COMPLEXITY DEPENDS ON THE NUMBER OF ITERATIONS DONE BY THE OPTIMISATION ALGORITHM HOWEVER ONE CAN SET THE MAXIMUM NUMBER OF ITERATIONS WITH THE ARGUMENT MAXITER FOR EACH ITERATION TIME COMPLEXITY IS  $O(N \times N \times \text{MAXITER})$

TRANSFORM

HERE THE TRANSFORM OPERATION RETURNS  $N \times N$  THEREFORE ITS TIME COMPLEXITY EQUALS  $O(N \times N)$

REFERENCES

17 GAUSSIAN PROCESSES

GAUSSIAN PROCESSES GP ARE A GENERIC SUPERVISED LEARNING METHOD DESIGNED TO SOLVE REGRESSION ANDPROBABILISTIC CLASSIFICATION PROBLEMS

THE ADVANTAGES OF GAUSSIAN PROCESSES ARE

- THE PREDICTION INTERPOLATES THE OBSERVATIONS AT LEAST FOR REGULAR KERNELS
- THE PREDICTION IS PROBABILISTIC GAUSSIAN SO THAT ONE CAN COMPUTE EMPIRICAL CONFIDENCE INTERVALS AND DECIDE BASED ON THOSE IF ONE SHOULD REFIT ONLINE FITTING ADAPTIVE FITTING THE PREDICTION IN SOME REGION OF INTEREST
- VERSATILE DIFFERENT KERNELS CAN BE SPECIFIED COMMON KERNELS ARE PROVIDED BUT IT IS ALSO POSSIBLE TO SPECIFY CUSTOM KERNELS

1"NEIGHBOURHOOD COMPONENTS ANALYSIS" ADVANCES IN NEURAL INFORMATION" J GOLDBERGER G HINTON S ROWEIS R SALAKHUTDINOV ADV IN NEURAL INFORMATION PROCESSING SYSTEMS V OL 17 MAY 2005 PP 513520

256 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

THE DISADVANTAGES OF GAUSSIAN PROCESSES INCLUDE

- THEY ARE NOT SPARSE IE THEY USE THE WHOLE SAMPLESFEATURES INFORMATION TO PERFORM THE PREDICTION
- THEY LOSE EFFICIENCY IN HIGH DIMENSIONAL SPACES – NAMELY WHEN THE NUMBER OF FEATURES EXCEEDS A FEW DOZENS

GAUSSIAN PROCESS REGRESSION GPR

THEGAUSSIANPROCESSREGRESSOR IMPLEMENTS GAUSSIAN PROCESSES GP FOR REGRESSION PURPOSES FOR THIS THE PRIOR OF THE GP NEEDS TO BE SPECIFIED THE PRIOR MEAN IS ASSUMED TO BE CONSTANT AND ZERO FOR NORMALIZEYFALSE OR THE TRAINING DATA’S MEAN FOR NORMALIZEYTRUE THE PRIOR’S COVARIANCE IS SPECIFIED BY PASSING A KERNEL OBJECT THE HYPERPARAMETERS OF THE KERNEL ARE OPTIMIZED DURING FITTING OF GAUSSIANPROCESSREGRESSOR BY MAXIMIZING THE LOG MARGINALLIKELIHOOD LML BASED ON THE PASSED OPTIMIZER AS THE LML MAY HAVE MULTIPLE LOCAL OPTIMA THE OPTIMIZER CAN BE STARTED REPEATEDLY BY SPECIFYING NRESTARTSOPTIMIZER THE FIRST RUN IS ALWAYS CONDUCTED STARTING FROM THE INITIAL HYPERPARAMETER VALUES OF THE KERNEL SUBSEQUENT RUNS ARE CONDUCTED FROM HYPERPARAMETER VALUES THAT HAVE BEEN CHOSEN RANDOMLY FROM THE RANGE OF ALLOWED VALUES IF THE INITIAL HYPERPARAMETERS SHOULD BE KEPT FIXEDNONE CAN BE PASSED AS OPTIMIZER

THE NOISE LEVEL IN THE TARGETS CAN BE SPECIFIED BY PASSING IT VIA THE PARAMETER ALPHA EITHER GLOBALLY AS A SCALAR OR PER DATAPPOINT NOTE THAT A MODERATE NOISE LEVEL CAN ALSO BE HELPFUL FOR DEALING WITH NUMERIC ISSUES DURING FITTING AS IT IS EFFECTIVELY IMPLEMENTED AS TIKHONOV REGULARIZATION IE BY ADDING IT TO THE DIAGONAL OF THE KERNEL MATRIX AN ALTERNATIVE TO SPECIFYING THE NOISE LEVEL EXPLICITLY IS TO INCLUDE A WHITEKERNEL COMPONENT INTO THE KERNEL WHICH CAN ESTIMATE THE GLOBAL NOISE LEVEL FROM THE DATA SEE EXAMPLE BELOW

THE IMPLEMENTATION IS BASED ON ALGORITHM 21 OF RW2006 IN ADDITION TO THE API OF STANDARD SCIKITLEARN ESTIMATORS GAUSSIANPROCESSREGRESSOR

- ALLOWS PREDICTION WITHOUT PRIOR FITTING BASED ON THE GP PRIOR
- PROVIDES AN ADDITIONAL METHOD SAMPLEYX WHICH EVALUATES SAMPLES DRAWN FROM THE GPR PRIOR OR POSTERIOR AT GIVEN INPUTS
- EXPOSES A METHOD LOGMARGINALLIKELIHOODTHETA WHICH CAN BE USED EXTERNALLY FOR OTHER WAYS OF SELECTING HYPERPARAMETERS EG VIA MARKOV CHAIN MONTE CARLO

GPR EXAMPLES

GPR WITH NOISELEVEL ESTIMATION

THIS EXAMPLE ILLUSTRATES THAT GPR WITH A SUMKERNEL INCLUDING A WHITEKERNEL CAN ESTIMATE THE NOISE LEVEL OF DATA AN ILLUSTRATION OF THE LOGMARGINALLIKELIHOOD LML LANDSCAPE SHOWS THAT THERE EXIST TWO LOCAL MAXIMA OF LML THE FIRST CORRESPONDS TO A MODEL WITH A HIGH NOISE LEVEL AND A LARGE LENGTH SCALE WHICH EXPLAINS ALL VARIATIONS IN THE DATA BY NOISE

THE SECOND ONE HAS A SMALLER NOISE LEVEL AND SHORTER LENGTH SCALE WHICH EXPLAINS MOST OF THE VARIATION BY THE NOISE FREE FUNCTIONAL RELATIONSHIP THE SECOND MODEL HAS A HIGHER LIKELIHOOD HOWEVER DEPENDING ON THE INITIAL VALUE FOR THE HYPERPARAMETERS THE GRADIENTBASED OPTIMIZATION MIGHT ALSO CONVERGE TO THE HIGHNOISE SOLUTION IT IS THUS IMPORTANT TO REPEAT THE OPTIMIZATION SEVERAL TIMES FOR DIFFERENT INITIALIZATIONS

COMPARISON OF GPR AND KERNEL RIDGE REGRESSION

BOTH KERNEL RIDGE REGRESSION KRR AND GPR LEARN A TARGET FUNCTION BY EMPLOYING INTERNALLY THE “KERNEL TRICK” KRR LEARNS A LINEAR FUNCTION IN THE SPACE INDUCED BY THE RESPECTIVE KERNEL WHICH CORRESPONDS TO A NONLINEAR FUNCTION IN THE ORIGINAL SPACE THE LINEAR FUNCTION IN THE KERNEL SPACE IS CHOSEN BASED ON THE MEANSQUARED ERROR LOSS WITH RIDGE REGULARIZATION GPR USES THE KERNEL TO DEFINE THE COVARIANCE OF A PRIOR DISTRIBUTION OVER THE TARGET FUNCTIONS AND USES









SCIKITLEARN USER GUIDE RELEASE 0213

THE OBSERVED TRAINING DATA TO DEFINE A LIKELIHOOD FUNCTION BASED ON BAYES THEOREM A GAUSSIAN POSTERIOR DISTRIBUTION OVER TARGET FUNCTIONS IS DEFINED WHOSE MEAN IS USED FOR PREDICTION

A MAJOR DIFFERENCE IS THAT GPR CAN CHOOSE THE KERNEL'S HYPERPARAMETERS BASED ON GRADIENTASCENT ON THE MARGINAL LIKELIHOOD FUNCTION WHILE KRR NEEDS TO PERFORM A GRID SEARCH ON A CROSSVALIDATED LOSS FUNCTION MEANSQUARED ERROR LOSS A FURTHER DIFFERENCE IS THAT GPR LEARNS A GENERATIVE PROBABILISTIC MODEL OF THE TARGET FUNCTION AND CAN THUS PROVIDE MEANINGFUL CONFIDENCE INTERVALS AND POSTERIOR SAMPLES ALONG WITH THE PREDICTIONS WHILE KRR ONLY PROVIDES PREDICTIONS

THE FOLLOWING FIGURE ILLUSTRATES BOTH METHODS ON AN ARTIFICIAL DATASET WHICH CONSISTS OF A SINUSOIDAL TARGET FUNCTION AND STRONG NOISE THE FIGURE COMPARES THE LEARNED MODEL OF KRR AND GPR BASED ON A EXPSINESQUARED KERNEL WHICH IS SUITED FOR LEARNING PERIODIC FUNCTIONS THE KERNEL'S HYPERPARAMETERS CONTROL THE SMOOTHNESS LENGTHSCALE AND PERIODICITY OF THE KERNEL PERIODICITY MOREOVER THE NOISE LEVEL OF THE DATA IS LEARNED EXPLICITLY BY GPR BY AN ADDITIONAL WHITEKERNEL COMPONENT IN THE KERNEL AND BY THE REGULARIZATION PARAMETER ALPHA OF KRR

THE FIGURE SHOWS THAT BOTH METHODS LEARN REASONABLE MODELS OF THE TARGET FUNCTION GPR CORRECTLY IDENTIFIES THE PERI ODICITY OF THE FUNCTION TO BE ROUGHLY 2π628 WHILE KRR CHOOSES THE DOUBLED PERIODICITY 4π BESIDES THAT GPR PROVIDES REASONABLE CONFIDENCE BOUNDS ON THE PREDICTION WHICH ARE NOT AVAILABLE FOR KRR A MAJOR DIFFERENCE BETWEEN THE TWO METHODS IS THE TIME REQUIRED FOR FITTING AND PREDICTING WHILE FITTING KRR IS FAST IN PRINCIPLE THE GRIDSEARCH FOR HYPERPARAMETER OPTIMIZATION SCALES EXPONENTIALLY WITH THE NUMBER OF HYPERPARAMETERS "CURSE OF DIMENSIONAL ITY" THE GRADIENTBASED OPTIMIZATION OF THE PARAMETERS IN GPR DOES NOT SUFFER FROM THIS EXPONENTIAL SCALING AND IS THUS CONSIDERABLE FASTER ON THIS EXAMPLE WITH 3DIMENSIONAL HYPERPARAMETER SPACE THE TIME FOR PREDICTING IS SIMILAR HOWEVER GENERATING THE VARIANCE OF THE PREDICTIVE DISTRIBUTION OF GPR TAKES CONSIDERABLE LONGER THAN JUST PREDICTING THE MEAN

GPR ON MAUNA LOA CO2 DATA

THIS EXAMPLE IS BASED ON SECTION 543 OF RW2006 IT ILLUSTRATES AN EXAMPLE OF COMPLEX KERNEL ENGINEERING AND HYPERPARAMETER OPTIMIZATION USING GRADIENT ASCENT ON THE LOGMARGINALLIKELIHOOD THE DATA CONSISTS OF THE MONTHLY AVERAGE ATMOSPHERIC CO2 CONCENTRATIONS IN PARTS PER MILLION BY VOLUME PPMV COLLECTED AT THE MAUNA LOA OBSERVATORY IN HAWAII BETWEEN 1958 AND 1997 THE OBJECTIVE IS TO MODEL THE CO2 CONCENTRATION AS A FUNCTION OF THE TIME

T  
THE KERNEL IS COMPOSED OF SEVERAL TERMS THAT ARE RESPONSIBLE FOR EXPLAINING DIFFERENT PROPERTIES OF THE SIGNAL

SCIKITLEARN USER GUIDE RELEASE 0213

- A LONG TERM SMOOTH RISING TREND IS TO BE EXPLAINED BY AN RBF KERNEL THE RBF KERNEL WITH A LARGE LENGTHSCALE ENFORCES THIS COMPONENT TO BE SMOOTH IT IS NOT ENFORCED THAT THE TREND IS RISING WHICH LEAVES THIS CHOICE TO THE GP THE SPECIFIC LENGTHSCALE AND THE AMPLITUDE ARE FREE HYPERPARAMETERS
- A SEASONAL COMPONENT WHICH IS TO BE EXPLAINED BY THE PERIODIC EXPSINESQUARED KERNEL WITH A FIXED PERIODICITY OF 1 YEAR THE LENGTHSCALE OF THIS PERIODIC COMPONENT CONTROLLING ITS SMOOTHNESS IS A FREE PARAMETER IN ORDER TO ALLOW DECAYING AWAY FROM EXACT PERIODICITY THE PRODUCT WITH AN RBF KERNEL IS TAKEN THE LENGTHSCALE OF THIS RBF COMPONENT CONTROLS THE DECAY TIME AND IS A FURTHER FREE PARAMETER
- SMALLER MEDIUM TERM IRREGULARITIES ARE TO BE EXPLAINED BY A RATIONALQUADRATIC KERNEL COMPONENT WHOSE LENGTH SCALE AND ALPHA PARAMETER WHICH DETERMINES THE DIFFUSENESS OF THE LENGTHSCALES ARE TO BE DETERMINED ACCORDING TO RW2006 THESE IRREGULARITIES CAN BETTER BE EXPLAINED BY A RATIONALQUADRATIC THAN AN RBF KERNEL COMPONENT PROBABLY BECAUSE IT CAN ACCOMMODATE SEVERAL LENGTHSCALES
- A “NOISE” TERM CONSISTING OF AN RBF KERNEL CONTRIBUTION WHICH SHALL EXPLAIN THE CORRELATED NOISE COMPONENTS SUCH AS LOCAL WEATHER PHENOMENA AND A WHITEKERNEL CONTRIBUTION FOR THE WHITE NOISE THE RELATIVE AMPLITUDES AND THE RBF’S LENGTH SCALE ARE FURTHER FREE PARAMETERS

MAXIMIZING THE LOGMARGINALLIKELIHOOD AFTER SUBTRACTING THE TARGET’S MEAN YIELDS THE FOLLOWING KERNEL WITH AN LML OF 83214

3442RBFLENGTHSCALE418  
3272RBFLENGTHSCALE180 EXPSINESQUAREDLENGTHSCALE144  
PERIODICITY1  
04462RATIONALQUADRATICALPHA177 LENGTHSCALE0957  
01972RBFLENGTHSCALE0138 WHITEKERNELNOISELEVEL00336

THUS MOST OF THE TARGET SIGNAL 344PPM IS EXPLAINED BY A LONGTERM RISING TREND LENGTHSCALE 418 YEARS THE PERIODIC COMPONENT HAS AN AMPLITUDE OF 327PPM A DECAY TIME OF 180 YEARS AND A LENGTHSCALE OF 144 THE LONG DECAY TIME INDICATES THAT WE HAVE A LOCALLY VERY CLOSE TO PERIODIC SEASONAL COMPONENT THE CORRELATED NOISE HAS AN AMPLITUDE OF 0197PPM WITH A LENGTH SCALE OF 0138 YEARS AND A WHITENOISE CONTRIBUTION OF 0197PPM THUS THE OVERALL NOISE LEVEL IS VERY SMALL INDICATING THAT THE DATA CAN BE VERY WELL EXPLAINED BY THE MODEL THE FIGURE SHOWS ALSO THAT THE MODEL MAKES VERY CONFIDENT PREDICTIONS UNTIL AROUND 2015

GAUSSIAN PROCESS CLASSIFICATION GPC  
THEGAUSSIANPROCESSCLASSIFIER IMPLEMENTS GAUSSIAN PROCESSES GP FOR CLASSIFICATION PURPOSES MORE SPECIFICALLY FOR PROBABILISTIC CLASSIFICATION WHERE TEST PREDICTIONS TAKE THE FORM OF CLASS PROBABILITIES GAUSSIANPRO CESSCLASSIFIER PLACES A GP PRIOR ON A LATENT FUNCTION  $\mathbf{z}$  WHICH IS THEN SQUASHED THROUGH A LINK FUNCTION TO OBTAIN THE PROBABILISTIC CLASSIFICATION THE LATENT FUNCTION  $\mathbf{z}$  IS A SOCALLED NUISANCE FUNCTION WHOSE VALUES ARE NOT OBSERVED AND ARE NOT RELEVANT BY THEMSELVES ITS PURPOSE IS TO ALLOW A CONVENIENT FORMULATION OF THE MODEL AND  $\mathbf{z}$  IS REMOVED INTEGRATED OUT DURING PREDICTION GAUSSIANPROCESSCLASSIFIER IMPLEMENTS THE LOGISTIC LINK FUNCTION FOR WHICH THE INTEGRAL CANNOT BE COMPUTED ANALYTICALLY BUT IS EASILY APPROXIMATED IN THE BINARY CASE

IN CONTRAST TO THE REGRESSION SETTING THE POSTERIOR OF THE LATENT FUNCTION  $\mathbf{z}$  IS NOT GAUSSIAN EVEN FOR A GP PRIOR SINCE A GAUSSIAN LIKELIHOOD IS INAPPROPRIATE FOR DISCRETE CLASS LABELS RATHER A NONGAUSSIAN LIKELIHOOD CORRESPONDING TO THE LOGISTIC LINK FUNCTION LOGIT IS USED GAUSSIANPROCESSCLASSIFIER APPROXIMATES THE NONGAUSSIAN POSTERIOR WITH A GAUSSIAN BASED ON THE LAPLACE APPROXIMATION MORE DETAILS CAN BE FOUND IN CHAPTER 3 OF RW2006

THE GP PRIOR MEAN IS ASSUMED TO BE ZERO THE PRIOR’S COVARIANCE IS SPECIFIED BY PASSING A KERNEL OBJECT THE HYPER PARAMETERS OF THE KERNEL ARE OPTIMIZED DURING FITTING OF GAUSSIANPROCESSREGRESSOR BY MAXIMIZING THE LOGMARGINAL LIKELIHOOD LML BASED ON THE PASSED OPTIMIZER AS THE LML MAY HAVE MULTIPLE LOCAL OPTIMA THE OPTIMIZER CAN BE STARTED REPEATEDLY BY SPECIFYING NRESTARTSOPTIMIZER THE FIRST RUN IS ALWAYS CONDUCTED STARTING FROM THE INITIAL HYPERPARAMETER VALUES OF THE KERNEL SUBSEQUENT RUNS ARE CONDUCTED FROM HYPERPARAMETER VALUES THAT HAVE BEEN CHOSEN RANDOMLY FROM THE RANGE OF ALLOWED VALUES IF THE INITIAL HYPERPARAMETERS SHOULD BE KEPT FIXED NONE CAN BE PASSED AS OPTIMIZER



SCIKITLEARN USER GUIDE RELEASE 0213

GAUSSIANPROCESSCLASSIFIER SUPPORTS MULTICLASS CLASSIFICATION BY PERFORMING EITHER ONEVERSUSREST OR ONE VERSUSONE BASED TRAINING AND PREDICTION IN ONEVERSUSREST ONE BINARY GAUSSIAN PROCESS CLASSIFIER IS FITTED FOR EACH CLASS WHICH IS TRAINED TO SEPARATE THIS CLASS FROM THE REST IN “ONEVSONE” ONE BINARY GAUSSIAN PROCESS CLASSIFIER IS FITTED FOR EACH PAIR OF CLASSES WHICH IS TRAINED TO SEPARATE THESE TWO CLASSES THE PREDICTIONS OF THESE BINARY PREDICTORS ARE COMBINED INTO MULTICLASS PREDICTIONS SEE THE SECTION ON MULTICLASS CLASSIFICATION FOR MORE DETAILS IN THE CASE OF GAUSSIAN PROCESS CLASSIFICATION “ONEVSONE” MIGHT BE COMPUTATIONALLY CHEAPER SINCE IT HAS TO SOLVE MANY PROBLEMS INVOLVING ONLY A SUBSET OF THE WHOLE TRAINING SET RATHER THAN FEWER PROBLEMS ON THE WHOLE DATASET SINCE GAUSSIAN PROCESS CLASSIFICATION SCALES CUBICALLY WITH THE SIZE OF THE DATASET THIS MIGHT BE CONSIDERABLY FASTER HOW EVER NOTE THAT “ONEVSONE” DOES NOT SUPPORT PREDICTING PROBABILITY ESTIMATES BUT ONLY PLAIN PREDICTIONS MOREOVER NOTE THATGAUSSIANPROCESSCLASSIFIER DOES NOT YET IMPLEMENT A TRUE MULTICLASS LAPLACE APPROXIMATION IN TERNALLY BUT AS DISCUSSED ABOVE IS BASED ON SOLVING SEVERAL BINARY CLASSIFICATION TASKS INTERNALLY WHICH ARE COMBINED USING ONEVERSUSREST OR ONEVERSUSONE

GPC EXAMPLES

PROBABILISTIC PREDICTIONS WITH GPC

THIS EXAMPLE ILLUSTRATES THE PREDICTED PROBABILITY OF GPC FOR AN RBF KERNEL WITH DIFFERENT CHOICES OF THE HYPERPARAM ETERS THE FIRST FIGURE SHOWS THE PREDICTED PROBABILITY OF GPC WITH ARBITRARILY CHOSEN HYPERPARAMETERS AND WITH THE HYPERPARAMETERS CORRESPONDING TO THE MAXIMUM LOGMARGINALLIKELIHOOD LML WHILE THE HYPERPARAMETERS CHOSEN BY OPTIMIZING LML HAVE A CONSIDERABLE LARGER LML THEY PERFORM SLIGHTLY WORSE ACCORDING TO THE LOGLOSS ON TEST DATA THE FIGURE SHOWS THAT THIS IS BECAUSE THEY EXHIBIT A STEEP CHANGE OF THE CLASS PROBABILITIES AT THE CLASS BOUNDARIES WHICH IS GOOD BUT HAVE PREDICTED PROBABILITIES CLOSE TO 05 FAR AWAY FROM THE CLASS BOUNDARIES WHICH IS BAD THIS UNDESIRABLE EFFECT IS CAUSED BY THE LAPLACE APPROXIMATION USED INTERNALLY BY GPC

THE SECOND FIGURE SHOWS THE LOGMARGINALLIKELIHOOD FOR DIFFERENT CHOICES OF THE KERNEL’S HYPERPARAMETERS HIGHLIGHTING THE TWO CHOICES OF THE HYPERPARAMETERS USED IN THE FIRST FIGURE BY BLACK DOTS

ILLUSTRATION OF GPC ON THE XOR DATASET

THIS EXAMPLE ILLUSTRATES GPC ON XOR DATA COMPARED ARE A STATIONARY ISOTROPIC KERNEL RBF AND A NONSTATIONARY KERNEL DOTPRODUCT ON THIS PARTICULAR DATASET THE DOTPRODUCT KERNEL OBTAINS CONSIDERABLY BETTER RESULTS BE CAUSE THE CLASSBOUNDARIES ARE LINEAR AND COINCIDE WITH THE COORDINATE AXES IN PRACTICE HOWEVER STATIONARY KERNELS SUCH ASRBF OFTEN OBTAIN BETTER RESULTS

GAUSSIAN PROCESS CLASSIFICATION GPC ON IRIS DATASET

THIS EXAMPLE ILLUSTRATES THE PREDICTED PROBABILITY OF GPC FOR AN ISOTROPIC AND ANISOTROPIC RBF KERNEL ON A TWO DIMENSIONAL VERSION FOR THE IRISDATASET THIS ILLUSTRATES THE APPLICABILITY OF GPC TO NONBINARY CLASSIFICATION THE ANISOTROPIC RBF KERNEL OBTAINS SLIGHTLY HIGHER LOGMARGINALLIKELIHOOD BY ASSIGNING DIFFERENT LENGTHSCALES TO THE TWO FEATURE DIMENSIONS

KERNELS FOR GAUSSIAN PROCESSES

KERNELS ALSO CALLED “COVARIANCE FUNCTIONS” IN THE CONTEXT OF GPS ARE A CRUCIAL INGREDIENT OF GPS WHICH DETERMINE THE SHAPE OF PRIOR AND POSTERIOR OF THE GP THEY ENCODE THE ASSUMPTIONS ON THE FUNCTION BEING LEARNED BY DEFINING THE “SIMILARITY” OF TWO DATAPOINTS COMBINED WITH THE ASSUMPTION THAT SIMILAR DATAPOINTS SHOULD HAVE SIMILAR TARGET VALUES TWO CATEGORIES OF KERNELS CAN BE DISTINGUISHED STATIONARY KERNELS DEPEND ONLY ON THE DISTANCE OF TWO DATAPOINTS AND NOT ON THEIR ABSOLUTE VALUES  $k(x, x') = k(\|x - x'\|)$  AND ARE THUS INVARIANT TO TRANSLATIONS IN THE INPUT SPACE







SCIKITLEARN USER GUIDE RELEASE 0213

WHILE NONSTATIONARY KERNELS DEPEND ALSO ON THE SPECIFIC VALUES OF THE DATAPOINTS STATIONARY KERNELS CAN FURTHER BE SUBDIVIDED INTO ISOTROPIC AND ANISOTROPIC KERNELS WHERE ISOTROPIC KERNELS ARE ALSO INVARIANT TO ROTATIONS IN THE INPUT SPACE FOR MORE DETAILS WE REFER TO CHAPTER 4 OF RW2006

GAUSSIAN PROCESS KERNEL API

THE MAIN USAGE OF A KERNEL IS TO COMPUTE THE GP’S COVARIANCE BETWEEN DATAPOINTS FOR THIS THE METHOD CALL OF THE KERNEL CAN BE CALLED THIS METHOD CAN EITHER BE USED TO COMPUTE THE “AUTOCOVARIANCE” OF ALL PAIRS OF DATAPOINTS IN A 2D ARRAY X OR THE “CROSSCOVARIANCE” OF ALL COMBINATIONS OF DATAPOINTS OF A 2D ARRAY X WITH DATAPOINTS IN A 2D ARRAY Y THE FOLLOWING IDENTITY HOLDS TRUE FOR ALL KERNELS K EXCEPT FOR THE WHITEKERNEL KX KX YX

IF ONLY THE DIAGONAL OF THE AUTOCOVARIANCE IS BEING USED THE METHOD DIAG OF A KERNEL CAN BE CALLED WHICH IS MORE COMPUTATIONALLY EFFICIENT THAN THE EQUIVALENT CALL TO CALL NPDIAKGX X KDIAGX

KERNELS ARE PARAMETERIZED BY A VECTOR  $\theta$  OF HYPERPARAMETERS THESE HYPERPARAMETERS CAN FOR INSTANCE CONTROL LENGTH SCALES OR PERIODICITY OF A KERNEL SEE BELOW ALL KERNELS SUPPORT COMPUTING ANALYTIC GRADIENTS OF THE KERNEL’S AUTO COVARIANCE WITH RESPECT TO  $\theta$  VIA SETTINGEVALGRADIENTTRUE IN THECALL METHOD THIS GRADIENT IS USED BY THE GAUSSIAN PROCESS BOTH REGRESSOR AND CLASSIFIER IN COMPUTING THE GRADIENT OF THE LOGMARGINALLIKELIHOOD WHICH IN TURN IS USED TO DETERMINE THE VALUE OF  $\theta$  WHICH MAXIMIZES THE LOGMARGINALLIKELIHOOD VIA GRADIENT ASCENT FOR EACH HYPERPARAMETER THE INITIAL VALUE AND THE BOUNDS NEED TO BE SPECIFIED WHEN CREATING AN INSTANCE OF THE KERNEL THE CURRENT VALUE OF  $\theta$  CAN BE GET AND SET VIA THE PROPERTY THETA OF THE KERNEL OBJECT MOREOVER THE BOUNDS OF THE HYPERPARAMETERS CAN BE ACCESSED BY THE PROPERTY BOUNDS OF THE KERNEL NOTE THAT BOTH PROPERTIES THETA AND BOUNDS RETURN LOGTRANSFORMED VALUES OF THE INTERNALLY USED VALUES SINCE THOSE ARE TYPICALLY MORE AMENABLE TO GRADIENTBASED OPTIMIZATION THE SPECIFICATION OF EACH HYPERPARAMETER IS STORED IN THE FORM OF AN INSTANCE OF HYPERPARAMETER IN THE RESPECTIVE KERNEL NOTE THAT A KERNEL USING A HYPERPARAMETER WITH NAME “X” MUST HAVE THE ATTRIBUTES SELF<sub>X</sub> AND SELF<sub>X</sub>BOUNDS

THE ABSTRACT BASE CLASS FOR ALL KERNELS IS KERNEL KERNEL IMPLEMENTS A SIMILAR INTERFACE AS ESTIMATOR PROVIDING THE METHODS GETPARAMS SETPARAMS ANDCLONE THIS ALLOWS SETTING KERNEL VALUES ALSO VIA META ESTIMATORS SUCH AS PIPELINE ORGRIDSEARCH NOTE THAT DUE TO THE NESTED STRUCTURE OF KERNELS BY APPLYING KERNEL OPERATORS SEE BELOW THE NAMES OF KERNEL PARAMETERS MIGHT BECOME RELATIVELY COMPLICATED IN GENERAL FOR A BINARY KERNEL OPERATOR PARAMETERS OF THE LEFT OPERAND ARE PREFIXED WITH K1 AND PARAMETERS OF THE RIGHT OPERAND WITH K2 AN ADDITIONAL CONVENIENCE METHOD IS CLONewithTHETATHETA WHICH RETURNS A CLONED VERSION OF THE KERNEL



AND  $\gamma^2$  AND COMBINES THEM VIA  $\frac{1}{\gamma^2} \exp(-\gamma^2 \frac{\|x - x'\|^2}{2})$  THE EXPONENTIATION KERNEL TAKES ONE  
BASE KERNEL AND A SCALAR PARAMETER  $\gamma$  AND COMBINES THEM VIA  $\exp(-\gamma^2 \frac{\|x - x'\|^2}{2})$  EXPONENT

RADIAL BASIS FUNCTION RBF KERNEL

THE RBF KERNEL IS A STATIONARY KERNEL IT IS ALSO KNOWN AS THE “SQUARED EXPONENTIAL” KERNEL IT IS PARAMETERIZED BY A  
LENGTHSCALE PARAMETER  $\ell > 0$  WHICH CAN EITHER BE A SCALAR ISOTROPIC VARIANT OF THE KERNEL OR A VECTOR WITH THE SAME  
NUMBER OF DIMENSIONS AS THE INPUTS  $\ell$  ANISOTROPIC VARIANT OF THE KERNEL THE KERNEL IS GIVEN BY

$$\frac{1}{(2\pi\ell^2)^{D/2}} \exp\left(-\frac{\|x - x'\|^2}{2\ell^2}\right)$$

THIS KERNEL IS INFINITELY DIFFERENTIABLE WHICH IMPLIES THAT GPS WITH THIS KERNEL AS COVARIANCE FUNCTION HAVE MEAN SQUARE  
DERIVATIVES OF ALL ORDERS AND ARE THUS VERY SMOOTH THE PRIOR AND POSTERIOR OF A GP RESULTING FROM AN RBF KERNEL ARE  
SHOWN IN THE FOLLOWING FIGURE

MATÉRN KERNEL

THE MATÉRN KERNEL IS A STATIONARY KERNEL AND A GENERALIZATION OF THE RBF KERNEL IT HAS AN ADDITIONAL PARAMETER  $\nu$   
WHICH CONTROLS THE SMOOTHNESS OF THE RESULTING FUNCTION IT IS PARAMETERIZED BY A LENGTHSCALE PARAMETER  $\ell > 0$  WHICH  
CAN EITHER BE A SCALAR ISOTROPIC VARIANT OF THE KERNEL OR A VECTOR WITH THE SAME NUMBER OF DIMENSIONS AS THE INPUTS  $\ell$   
ANISOTROPIC VARIANT OF THE KERNEL THE KERNEL IS GIVEN BY

$$\frac{2^{\nu-1/2} \Gamma(\nu)}{\sqrt{\pi} \Gamma(\nu)} \left(\frac{\sqrt{2\nu}}{\ell}\right)^{\nu-1/2} K_{\nu-1/2}\left(\frac{\sqrt{2\nu}}{\ell} \|x - x'\|\right)$$

AS  $\nu \rightarrow \infty$  THE MATÉRN KERNEL CONVERGES TO THE RBF KERNEL WHEN  $\nu = 1/2$  THE MATÉRN KERNEL BECOMES IDENTICAL TO  
THE ABSOLUTE EXPONENTIAL KERNEL IE

$$\frac{1}{2\ell} \exp\left(-\frac{\|x - x'\|}{\ell}\right)$$

IN PARTICULAR  $\nu = 3/2$

$$\frac{1}{2\ell^3} \left(1 + \frac{\sqrt{3}}{2\ell} \|x - x'\|\right) \exp\left(-\frac{\sqrt{3}}{2\ell} \|x - x'\|\right)$$

$$\frac{1}{2\ell^5} \left(1 + \frac{5\sqrt{5}}{4\ell} \|x - x'\| + \frac{3\sqrt{5}}{8\ell^2} \|x - x'\|^2\right) \exp\left(-\frac{\sqrt{5}}{2\ell} \|x - x'\|\right)$$

ARE POPULAR CHOICES FOR LEARNING FUNCTIONS THAT ARE NOT INFINITELY DIFFERENTIABLE AS ASSUMED BY THE RBF KERNEL BUT AT  
LEAST ONCE  $\nu = 3/2$  OR TWICE DIFFERENTIABLE  $\nu = 5/2$

THE FLEXIBILITY OF CONTROLLING THE SMOOTHNESS OF THE LEARNED FUNCTION VIA  $\nu$  ALLOWS ADAPTING TO THE PROPERTIES OF THE  
TRUE UNDERLYING FUNCTIONAL RELATION THE PRIOR AND POSTERIOR OF A GP RESULTING FROM A MATÉRN KERNEL ARE SHOWN IN THE  
FOLLOWING FIGURE

SEERW2006 PP84 FOR FURTHER DETAILS REGARDING THE DIFFERENT VARIANTS OF THE MATÉRN KERNEL





SCIKITLEARN USER GUIDE RELEASE 0213

RATIONAL QUADRATIC KERNEL

THE RATIONAL QUADRATIC KERNEL CAN BE SEEN AS A SCALE MIXTURE AN INFINITE SUM OF RBF KERNELS WITH DIFFERENT CHARACTERISTIC LENGTH SCALES IT IS PARAMETERIZED BY A LENGTH SCALE PARAMETER  $\ell^2$  AND A SCALE MIXTURE PARAMETER  $\kappa^2$  ONLY THE ISOTROPIC VARIANT WHERE  $\ell$  IS A SCALAR IS SUPPORTED AT THE MOMENT THE KERNEL IS GIVEN BY

$$k(\mathbf{x}, \mathbf{y}) = \frac{1}{2\pi^2} \int_0^\infty \frac{1 - \exp(-\frac{1}{2}(\mathbf{x} - \mathbf{y})^T \mathbf{L}(\mathbf{x} - \mathbf{y}))}{\ell^2} \exp(-\frac{1}{2}\kappa^2 \ell^2) d\ell$$

THE PRIOR AND POSTERIOR OF A GP RESULTING FROM A RATIONAL QUADRATIC KERNEL ARE SHOWN IN THE FOLLOWING FIGURE

272 CHAPTER 3 USER GUIDE

EXPSINESQUARED KERNEL

THEEXPSINESQUARED KERNEL ALLOWS MODELING PERIODIC FUNCTIONS IT IS PARAMETERIZED BY A LENGTHSCALE PARAMETER  $\ell$  AND A PERIODICITY PARAMETER  $\omega$  ONLY THE ISOTROPIC VARIANT WHERE  $\omega$  IS A SCALAR IS SUPPORTED AT THE MOMENT THE KERNEL IS GIVEN BY

$$k(x, y) = \exp\left(-\frac{1}{2\ell^2} \|x - y\|^2\right) \sin^2\left(\frac{\omega}{2} \|x - y\|\right)$$

THE PRIOR AND POSTERIOR OF A GP RESULTING FROM AN EXPSINESQUARED KERNEL ARE SHOWN IN THE FOLLOWING FIGURE

DOTPRODUCT KERNEL

THE DOTPRODUCT KERNEL IS NONSTATIONARY AND CAN BE OBTAINED FROM LINEAR REGRESSION BY PUTTING  $\gamma_0$  PRIORS ON THE COEFFICIENTS OF  $\gamma_0$  1  $\gamma$  AND A PRIOR OF  $\gamma_0$   $\gamma^2$

ON THE BIAS THE DOTPRODUCT KERNEL IS INVARIANT TO A

ROTATION OF THE COORDINATES ABOUT THE ORIGIN BUT NOT TRANSLATIONS IT IS PARAMETERIZED BY A PARAMETER  $\gamma^2$  0 FOR  $\gamma^2$  0 0

THE KERNEL IS CALLED THE HOMOGENEOUS LINEAR KERNEL OTHERWISE IT IS INHOMOGENEOUS THE KERNEL IS GIVEN BY  $\gamma_0$   $\gamma$   $\gamma^2$  0  $\gamma$   $\gamma$

THE DOTPRODUCT KERNEL IS COMMONLY COMBINED WITH EXPONENTIATION AN EXAMPLE WITH EXPONENT 2 IS SHOWN IN THE FOLLOWING FIGURE



SCIKITLEARN USER GUIDE RELEASE 0213

REFERENCES

318 CROSS DECOMPOSITION

THE CROSS DECOMPOSITION MODULE CONTAINS TWO MAIN FAMILIES OF ALGORITHMS THE PARTIAL LEAST SQUARES PLS AND THE CANONICAL CORRELATION ANALYSIS CCA THESE FAMILIES OF ALGORITHMS ARE USEFUL TO FIND LINEAR RELATIONS BETWEEN TWO MULTIVARIATE DATASETS THE XANDYARGUMENTS OF THE FIT METHOD ARE 2D ARRAYS CROSS DECOMPOSITION ALGORITHMS FIND THE FUNDAMENTAL RELATIONS BETWEEN TWO MATRICES X AND Y THEY ARE LATENT VARIABLE APPROACHES TO MODELING THE COVARIANCE STRUCTURES IN THESE TWO SPACES THEY WILL TRY TO FIND THE MULTIDIRECTIONAL DIRECTION IN THE X SPACE THAT EXPLAINS THE MAXIMUM MULTIDIMENSIONAL VARIANCE DIRECTION IN THE Y SPACE PLSREGRESSION IS PARTICULARLY SUITED WHEN THE MATRIX OF PREDICTORS HAS MORE VARIABLES THAN OBSERVATIONS AND WHEN THERE IS MULTICOLLINEARITY AMONG X VALUES BY CONTRAST STANDARD REGRESSION WILL FAIL IN THESE CASES CLASSES INCLUDED IN THIS MODULE ARE PLSREGRESSIONPLSCANONICAL CCA ANDPLSSVD

REFERENCE

- JA WEGELIN A SURVEY OF PARTIAL LEAST SQUARES PLS METHODS WITH EMPHASIS ON THE TWOBLOCK CASE

EXAMPLES

31 SUPERVISED LEARNING 275

NAIVE BAYES METHODS ARE A SET OF SUPERVISED LEARNING ALGORITHMS BASED ON APPLYING BAYES’ THEOREM WITH THE “NAIVE” ASSUMPTION OF CONDITIONAL INDEPENDENCE BETWEEN EVERY PAIR OF FEATURES GIVEN THE VALUE OF THE CLASS VARIABLE BAYES’ THEOREM STATES THE FOLLOWING RELATIONSHIP GIVEN CLASS VARIABLE  $c$  AND DEPENDENT FEATURE VECTOR  $x_1$  THROUGH  $x_n$

$$p(c|x_1, \dots, x_n) = \frac{p(c) \prod_{i=1}^n p(x_i|c)}$$

USING THE NAIVE CONDITIONAL INDEPENDENCE ASSUMPTION THAT

$$p(x_i|x_{-i}, c) = p(x_i|c) \quad \text{FOR ALL } i$$

THIS RELATIONSHIP IS SIMPLIFIED TO

$$p(c|x_1, \dots, x_n) \propto p(c) \prod_{i=1}^n p(x_i|c)$$

SINCE  $p(c)$  IS CONSTANT GIVEN THE INPUT WE CAN USE THE FOLLOWING CLASSIFICATION RULE

$$p(c|x_1, \dots, x_n) \propto \prod_{i=1}^n p(x_i|c)$$

$$\hat{c} = \underset{c}{\operatorname{arg\,max}} \prod_{i=1}^n p(x_i|c)$$

AND WE CAN USE MAXIMUM A POSTERIORI MAP ESTIMATION TO ESTIMATE  $p(c)$  AND  $p(x_i|c)$  THE FORMER IS THEN THE RELATIVE FREQUENCY OF CLASS  $c$  IN THE TRAINING SET

THE DIFFERENT NAIVE BAYES CLASSIFIERS DIFFER MAINLY BY THE ASSUMPTIONS THEY MAKE REGARDING THE DISTRIBUTION OF  $x_i$

IN SPITE OF THEIR APPARENTLY OVERSIMPLIFIED ASSUMPTIONS NAIVE BAYES CLASSIFIERS HAVE WORKED QUITE WELL IN MANY REAL WORLD SITUATIONS FAMOUSLY DOCUMENT CLASSIFICATION AND SPAM FILTERING THEY REQUIRE A SMALL AMOUNT OF TRAINING DATA TO ESTIMATE THE NECESSARY PARAMETERS FOR THEORETICAL REASONS WHY NAIVE BAYES WORKS WELL AND ON WHICH TYPES OF DATA IT DOES SEE THE REFERENCES BELOW

NAIVE BAYES LEARNERS AND CLASSIFIERS CAN BE EXTREMELY FAST COMPARED TO MORE SOPHISTICATED METHODS THE DECOUPLING OF THE CLASS CONDITIONAL FEATURE DISTRIBUTIONS MEANS THAT EACH DISTRIBUTION CAN BE INDEPENDENTLY ESTIMATED AS A ONE DIMENSIONAL DISTRIBUTION THIS IN TURN HELPS TO ALLEVIATE PROBLEMS STEMMING FROM THE CURSE OF DIMENSIONALITY ON THE FLIP SIDE ALTHOUGH NAIVE BAYES IS KNOWN AS A DECENT CLASSIFIER IT IS KNOWN TO BE A BAD ESTIMATOR SO THE PROBABILITY OUTPUTS FROM PREDICTPROBA ARE NOT TO BE TAKEN TOO SERIOUSLY

GAUSSIAN NAIVE BAYES

GAUSSIANNB IMPLEMENTS THE GAUSSIAN NAIVE BAYES ALGORITHM FOR CLASSIFICATION THE LIKELIHOOD OF THE FEATURES IS ASSUMED TO BE GAUSSIAN

$$\frac{1}{\sqrt{2\pi^2}} \exp\left(-\frac{x^2}{2}\right)$$

THE PARAMETERS  $\mu$  AND  $\sigma$  ARE ESTIMATED USING MAXIMUM LIKELIHOOD

FROM SKLEARN IMPORT DATASETS

IRIS = DATASETSLOADIRIS

FROM SKLEARNNAIVEBAYES IMPORT GAUSSIANNB

GNB = GAUSSIANNB

YPRED = GNBFITIRISDATA IRISTARGETPREDICTIRISDATA

PRINTNUMBER OF MISLABELED POINTS OUT OF A TOTAL DPOINTS = D

IRISDATASHAPE0IRISTARGET\_YPREDSUM

NUMBER OF MISLABELED POINTS OUT OF A TOTAL 150 POINTS = 6

MULTINOMIAL NAIVE BAYES

MULTINOMIALNB IMPLEMENTS THE NAIVE BAYES ALGORITHM FOR MULTINOMIALLY DISTRIBUTED DATA AND IS ONE OF THE TWO CLASSIC NAIVE BAYES VARIANTS USED IN TEXT CLASSIFICATION WHERE THE DATA ARE TYPICALLY REPRESENTED AS WORD VECTOR COUNTS ALTHOUGH TFIDF VECTORS ARE ALSO KNOWN TO WORK WELL IN PRACTICE THE DISTRIBUTION IS PARAMETRIZED BY VECTORS  $\phi$

$\phi_i = \frac{1}{N}$  FOR EACH CLASS  $i$  WHERE  $N$  IS THE NUMBER OF FEATURES IN TEXT CLASSIFICATION THE SIZE OF THE VOCABULARY

AND  $\phi_{ij}$  IS THE PROBABILITY  $\phi_{ij}$  OF FEATURE  $j$  APPEARING IN A SAMPLE BELONGING TO CLASS  $i$

THE PARAMETERS  $\phi_{ij}$  IS ESTIMATED BY A SMOOTHED VERSION OF MAXIMUM LIKELIHOOD IE RELATIVE FREQUENCY COUNTING

$$\phi_{ij} = \frac{N_{ij} + \alpha}{N_i + \alpha V}$$

$$\alpha = \frac{1}{V}$$

WHERE  $N_{ij}$

$N_{ij} \in \mathbb{N}$  IS THE NUMBER OF TIMES FEATURE  $j$  APPEARS IN A SAMPLE OF CLASS  $i$  IN THE TRAINING SET  $\mathcal{D}$  AND

$$N_i = \sum_j N_{ij}$$

$N_i$  IS THE TOTAL COUNT OF ALL FEATURES FOR CLASS  $i$

THE SMOOTHING PRIORS  $\alpha \geq 0$  ACCOUNTS FOR FEATURES NOT PRESENT IN THE LEARNING SAMPLES AND PREVENTS ZERO PROBABILITIES IN FURTHER COMPUTATIONS SETTING  $\alpha = 1$  IS CALLED LAPLACE SMOOTHING WHILE  $\alpha = 1/V$  IS CALLED LIDSTONE SMOOTHING

COMPLEMENT NAIVE BAYES

COMPLEMENTNB IMPLEMENTS THE COMPLEMENT NAIVE BAYES CNB ALGORITHM CNB IS AN ADAPTATION OF THE STANDARD

MULTINOMIAL NAIVE BAYES MNB ALGORITHM THAT IS PARTICULARLY SUITED FOR IMBALANCED DATA SETS SPECIFICALLY CNB USES

STATISTICS FROM THE COMPLEMENT OF EACH CLASS TO COMPUTE THE MODEL’S WEIGHTS THE INVENTORS OF CNB SHOW EMPIRICALLY

THAT THE PARAMETER ESTIMATES FOR CNB ARE MORE STABLE THAN THOSE FOR MNB FURTHER CNB REGULARLY OUTPERFORMS MNB

OFTEN BY A CONSIDERABLE MARGIN ON TEXT CLASSIFICATION TASKS THE PROCEDURE FOR CALCULATING THE WEIGHTS IS AS FOLLOWS

$$\phi_{ij} = \frac{N_{i\bar{j}} + \alpha}{N_i + \alpha V}$$

$$\alpha = \frac{1}{V}$$

$$N_{i\bar{j}} = N_i - N_{ij}$$

$$N_i = \sum_j N_{i\bar{j}}$$

$$\phi_{ij} = \log \phi_{ij}$$

$$\phi_{ij} = \log \phi_{ij}$$

$$\phi_{ij} = \log \phi_{ij}$$

WHERE THE SUMMATIONS ARE OVER ALL DOCUMENTS  $i$  NOT IN CLASS  $i$   $N_{i\bar{j}}$  IS EITHER THE COUNT OR TFIDF VALUE OF TERM  $j$  IN

DOCUMENT  $i$   $\alpha$  IS A SMOOTHING HYPERPARAMETER LIKE THAT FOUND IN MNB AND  $V$

THE SECOND NORMALIZATION

SCIKITLEARN USER GUIDE RELEASE 0213  
ADDRESSES THE TENDENCY FOR LONGER DOCUMENTS TO DOMINATE PARAMETER ESTIMATES IN MNB THE CLASSIFICATION RULE IS  
 $\arg \min_{c \in \mathcal{C}} \sum_{i=1}^n \log p(c | x_i)$

IE A DOCUMENT IS ASSIGNED TO THE CLASS THAT IS THE POOREST COMPLEMENT MATCH  
REFERENCES

- RENNIE J D SHIH L TEEVAN J KARGER D R 2003 TACKLING THE POOR ASSUMPTIONS OF NAIVE BAYES TEXT CLASSIFIERS IN ICML V OL 3 PP 616623

BERNOULLI NAIVE BAYES  
BERNOULLINB IMPLEMENTS THE NAIVE BAYES TRAINING AND CLASSIFICATION ALGORITHMS FOR DATA THAT IS DISTRIBUTED ACCORDING TO MULTIVARIATE BERNOLLI DISTRIBUTIONS IE THERE MAY BE MULTIPLE FEATURES BUT EACH ONE IS ASSUMED TO BE A BINARYVALUED BERNOLLI BOOLEAN VARIABLE THEREFORE THIS CLASS REQUIRES SAMPLES TO BE REPRESENTED AS BINARYVALUED FEATURE VECTORS IF HANDED ANY OTHER KIND OF DATA A BERNOLLINB INSTANCE MAY BINARIZE ITS INPUT DEPENDING ON THE BINARIZE PARAMETER  
THE DECISION RULE FOR BERNOLLI NAIVE BAYES IS BASED ON

$$c = \arg \max_{c \in \mathcal{C}} \prod_{i=1}^n p(x_i | c)$$
 WHICH DIFFERS FROM MULTINOMIAL NB’S RULE IN THAT IT EXPLICITLY PENALIZES THE NONOCCURRENCE OF A FEATURE  $x_i$  THAT IS AN INDICATOR FOR CLASS  $c$  WHERE THE MULTINOMIAL VARIANT WOULD SIMPLY IGNORE A NONOCCURRING FEATURE  
IN THE CASE OF TEXT CLASSIFICATION WORD OCCURRENCE VECTORS RATHER THAN WORD COUNT VECTORS MAY BE USED TO TRAIN AND USE THIS CLASSIFIER BERNOLLINB MIGHT PERFORM BETTER ON SOME DATASETS ESPECIALLY THOSE WITH SHORTER DOCUMENTS  
IT IS ADVISABLE TO EVALUATE BOTH MODELS IF TIME PERMITS

REFERENCES  
• CD MANNING P RAGHAVAN AND H SCHÜTZE 2008 INTRODUCTION TO INFORMATION RETRIEVAL CAMBRIDGE UNIVERSITY PRESS PP 234265  
• A MCCALLUM AND K NIGAM 1998 A COMPARISON OF EVENT MODELS FOR NAIVE BAYES TEXT CLASSIFICATION PROCEEDINGS OF THE AAAI/ICML98 WORKSHOP ON LEARNING FOR TEXT CATEGORIZATION PP 4148  
• V METSIS I ANDROUTSOPOULOS AND G PALIOURAS 2006 SPAM FILTERING WITH NAIVE BAYES - WHICH NAIVE BAYES 3RD CONF ON EMAIL AND ANTISPAM CEAS  
OUTOFCORE NAIVE BAYES MODEL FITTING  
NAIVE BAYES MODELS CAN BE USED TO TACKLE LARGE SCALE CLASSIFICATION PROBLEMS FOR WHICH THE FULL TRAINING SET MIGHT NOT FIT IN MEMORY TO HANDLE THIS CASE MULTINOMIALNB BERNOLLINB ANDGAUSSIANNB EXPOSE A PARTIALFIT METHOD THAT CAN BE USED INCREMENTALLY AS DONE WITH OTHER CLASSIFIERS AS DEMONSTRATED IN OUTOFCORE CLASSIFICATION OF TEXT DOCUMENTS ALL NAIVE BAYES CLASSIFIERS SUPPORT SAMPLE WEIGHTING  
CONTRARY TO THE FIT METHOD THE FIRST CALL TO PARTIALFIT NEEDS TO BE PASSED THE LIST OF ALL THE EXPECTED CLASS LABELS  
FOR AN OVERVIEW OF AVAILABLE STRATEGIES IN SCIKITLEARN SEE ALSO THE OUTOFCORE LEARNING DOCUMENTATION  
278 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

NOTE THEPARTIALFIT METHOD CALL OF NAIVE BAYES MODELS INTRODUCES SOME COMPUTATIONAL OVERHEAD IT IS RECOMMENDED TO USE DATA CHUNK SIZES THAT ARE AS LARGE AS POSSIBLE THAT IS AS THE AVAILABLE RAM ALLOWS

3110 DECISION TREES  
DECISION TREES DTS ARE A NONPARAMETRIC SUPERVISED LEARNING METHOD USED FOR CLASSIFICATION ANDREGRESSION THE GOAL IS TO CREATE A MODEL THAT PREDICTS THE VALUE OF A TARGET VARIABLE BY LEARNING SIMPLE DECISION RULES INFERRED FROM THE DATA FEATURES

FOR INSTANCE IN THE EXAMPLE BELOW DECISION TREES LEARN FROM DATA TO APPROXIMATE A SINE CURVE WITH A SET OF IFTHENELSE DECISION RULES THE DEEPER THE TREE THE MORE COMPLEX THE DECISION RULES AND THE FITTER THE MODEL

SOME ADVANTAGES OF DECISION TREES ARE

- SIMPLE TO UNDERSTAND AND TO INTERPRET TREES CAN BE VISUALISED
- REQUIRES LITTLE DATA PREPARATION OTHER TECHNIQUES OFTEN REQUIRE DATA NORMALISATION DUMMY VARIABLES NEED TO BE CREATED AND BLANK VALUES TO BE REMOVED NOTE HOWEVER THAT THIS MODULE DOES NOT SUPPORT MISSING VALUES
- THE COST OF USING THE TREE IE PREDICTING DATA IS LOGARITHMIC IN THE NUMBER OF DATA POINTS USED TO TRAIN THE TREE
- ABLE TO HANDLE BOTH NUMERICAL AND CATEGORICAL DATA OTHER TECHNIQUES ARE USUALLY SPECIALISED IN ANALYSING DATASETS THAT HAVE ONLY ONE TYPE OF VARIABLE SEE ALGORITHMS FOR MORE INFORMATION
- ABLE TO HANDLE MULTIOUTPUT PROBLEMS
- USES A WHITE BOX MODEL IF A GIVEN SITUATION IS OBSERVABLE IN A MODEL THE EXPLANATION FOR THE CONDITION IS EASILY EXPLAINED BY BOOLEAN LOGIC BY CONTRAST IN A BLACK BOX MODEL EG IN AN ARTIFICIAL NEURAL NETWORK RESULTS MAY BE MORE DIFFICULT TO INTERPRET
- POSSIBLE TO VALIDATE A MODEL USING STATISTICAL TESTS THAT MAKES IT POSSIBLE TO ACCOUNT FOR THE RELIABILITY OF THE MODEL

31 SUPERVISED LEARNING 279

SCIKITLEARN USER GUIDE RELEASE 0213

- PERFORMS WELL EVEN IF ITS ASSUMPTIONS ARE SOMEWHAT VIOLATED BY THE TRUE MODEL FROM WHICH THE DATA WERE GENERATED

THE DISADVANTAGES OF DECISION TREES INCLUDE

- DECISIONTREE LEARNERS CAN CREATE OVERCOMPLEX TREES THAT DO NOT GENERALISE THE DATA WELL THIS IS CALLED OVERFITTING MECHANISMS SUCH AS PRUNING NOT CURRENTLY SUPPORTED SETTING THE MINIMUM NUMBER OF SAMPLES REQUIRED AT A LEAF NODE OR SETTING THE MAXIMUM DEPTH OF THE TREE ARE NECESSARY TO AVOID THIS PROBLEM
- DECISION TREES CAN BE UNSTABLE BECAUSE SMALL VARIATIONS IN THE DATA MIGHT RESULT IN A COMPLETELY DIFFERENT TREE BEING GENERATED THIS PROBLEM IS MITIGATED BY USING DECISION TREES WITHIN AN ENSEMBLE
- THE PROBLEM OF LEARNING AN OPTIMAL DECISION TREE IS KNOWN TO BE NPCOMPLETE UNDER SEVERAL ASPECTS OF OPTIMALITY AND EVEN FOR SIMPLE CONCEPTS CONSEQUENTLY PRACTICAL DECISIONTREE LEARNING ALGORITHMS ARE BASED ON HEURISTIC ALGORITHMS SUCH AS THE GREEDY ALGORITHM WHERE LOCALLY OPTIMAL DECISIONS ARE MADE AT EACH NODE SUCH ALGORITHMS CANNOT GUARANTEE TO RETURN THE GLOBALLY OPTIMAL DECISION TREE THIS CAN BE MITIGATED BY TRAINING MULTIPLE TREES IN AN ENSEMBLE LEARNER WHERE THE FEATURES AND SAMPLES ARE RANDOMLY SAMPLED WITH REPLACEMENT
- THERE ARE CONCEPTS THAT ARE HARD TO LEARN BECAUSE DECISION TREES DO NOT EXPRESS THEM EASILY SUCH AS XOR PARITY OR MULTIPLEXER PROBLEMS
- DECISION TREE LEARNERS CREATE BIASED TREES IF SOME CLASSES DOMINATE IT IS THEREFORE RECOMMENDED TO BALANCE THE DATASET PRIOR TO FITTING WITH THE DECISION TREE

CLASSIFICATION

DECISIONTREECLASSIFIER IS A CLASS CAPABLE OF PERFORMING MULTICLASS CLASSIFICATION ON A DATASET AS WITH OTHER CLASSIFIERS DECISIONTREECLASSIFIER TAKES AS INPUT TWO ARRAYS AN ARRAY X SPARSE OR DENSE OF SIZESAMPLES NFEATURES HOLDING THE TRAINING SAMPLES AND AN ARRAY Y OF INTEGER VALUES SIZE NSAMPLES HOLDING THE CLASS LABELS FOR THE TRAINING SAMPLES

```
FROM SKLEARN IMPORT TREE
```

```
X 0 0 1 1
```

```
Y 0 1
```

```
CLF TREEDECISIONTREECLASSIFIER
```

```
CLF CLFFITX Y
```

AFTER BEING FITTED THE MODEL CAN THEN BE USED TO PREDICT THE CLASS OF SAMPLES

```
CLFPREDICT2 2
```

```
ARRAY1
```

ALTERNATIVELY THE PROBABILITY OF EACH CLASS CAN BE PREDICTED WHICH IS THE FRACTION OF TRAINING SAMPLES OF THE SAME CLASS IN A LEAF

```
CLFPREDICTPROBA2 2
```

```
ARRAY0 1
```

DECISIONTREECLASSIFIER IS CAPABLE OF BOTH BINARY WHERE THE LABELS ARE 1 1 CLASSIFICATION AND MULTICLASS WHERE THE LABELS ARE 0 K1 CLASSIFICATION

USING THE IRIS DATASET WE CAN CONSTRUCT A TREE AS FOLLOWS

```
FROM SKLEARNDATASETS IMPORT LOADIRIS
```

```
FROM SKLEARN IMPORT TREE
```

```
IRIS LOADIRIS
```

```
CLF TREEDECISIONTREECLASSIFIER
```

```
CLF CLFFITIRISDATA IRISTARGET
```

280 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

ONCE TRAINED YOU CAN PLOT THE TREE WITH THE PLOTTREE FUNCTION

```
TREEPLOTTREECLFFITIRISDATA IRISTARGET
```

WE CAN ALSO EXPORT THE TREE IN GRAPHVIZ FORMAT USING THE EXPORTGRAPHVIZ EXPORTER IF YOU USE THE CONDA PACKAGE

MANAGER THE GRAPHVIZ BINARIES

AND THE PYTHON PACKAGE CAN BE INSTALLED WITH

```
CONDA INSTALL PYTHONGRAPHVIZ
```

ALTERNATIVELY BINARIES FOR GRAPHVIZ CAN BE DOWNLOADED FROM THE GRAPHVIZ PROJECT HOMEPAGE AND THE PYTHON WRAPPER

INSTALLED FROM PYPI WITH PIP INSTALL GRAPHVIZ

BELOW IS AN EXAMPLE GRAPHVIZ EXPORT OF THE ABOVE TREE TRAINED ON THE ENTIRE IRIS DATASET THE RESULTS ARE SAVED IN AN

OUTPUT FILEIRISPDF

```
IMPORT GRAPHVIZ
DOTDATA TREEEXPORTGRAPHVIZCLF OUTFILE NONE
GRAPH GRAPHVIZSOURCEDOTDATA
GRAPHRENDERIRIS
THEEXPORTGRAPHVIZ EXPORTER ALSO SUPPORTS A VARIETY OF AESTHETIC OPTIONS INCLUDING COLORING NODES BY THEIR CLASS
OR VALUE FOR REGRESSION AND USING EXPLICIT VARIABLE AND CLASS NAMES IF DESIRED JUPYTER NOTEBOOKS ALSO RENDER THESE
PLOTS INLINE AUTOMATICALLY
DOTDATA TREEEXPORTGRAPHVIZCLF OUTFILE NONE
FEATURENAMESIRISFEATURENAMES
CLASSNAMESIRISTARGETNAMES
FILLEDTRUE ROUNDED TRUE
SPECIALCHARACTERS TRUE
GRAPH GRAPHVIZSOURCEDOTDATA
GRAPH
```

31 SUPERVISED LEARNING 281





```
SCIKITLEARN USER GUIDE RELEASE 0213
ALTERNATIVELY THE TREE CAN ALSO BE EXPORTED IN TEXTUAL FORMAT WITH THE FUNCTION EXPORTTEXT THIS METHOD DOESN'T
REQUIRE THE INSTALLATION OF EXTERNAL LIBRARIES AND IS MORE COMPACT
FROM SKLEARNDATASETS IMPORT LOADIRIS
FROM SKLEARNTREE IMPORT DECISIONTREECLASSIFIER
FROM SKLEARNTREEEXPORT IMPORT EXPORTTEXT
IRIS = LOADIRIS
X = IRISDATA
Y = IRISTARGET
DECISIONTREE = DECISIONTREECLASSIFIERRANDOMSTATE0 MAXDEPTH2
DECISIONTREE = DECISIONTREEFITX Y
R = EXPORTTEXTDECISIONTREE FEATURENAMESIRISFEATURENAMES
PRINTR
PETAL WIDTH CM 080
CLASS 0
PETAL WIDTH CM 080
PETAL WIDTH CM 175
CLASS 1
PETAL WIDTH CM 175
CLASS 2
EXAMPLES
•PLOT THE DECISION SURFACE OF A DECISION TREE ON THE IRIS DATASET
•UNDERSTANDING THE DECISION TREE STRUCTURE
REGRESSION
DECISION TREES CAN ALSO BE APPLIED TO REGRESSION PROBLEMS USING THE DECISIONTREEREgressor CLASS
31 SUPERVISED LEARNING 283
```

SCIKITLEARN USER GUIDE RELEASE 0213

AS IN THE CLASSIFICATION SETTING THE FIT METHOD WILL TAKE AS ARGUMENT ARRAYS X AND Y ONLY THAT IN THIS CASE Y IS EXPECTED TO HAVE FLOATING POINT VALUES INSTEAD OF INTEGER VALUES

```
FROM SKLEARN IMPORT TREE
X 0 0 2 2
Y 05 25
CLF TREEDECISIONTREEREgressor
CLF CLFFITX Y
CLFPREDICT1 1
ARRAY05
EXAMPLES
•DECISION TREE REGRESSION
MULTIOUTPUT PROBLEMS
A MULTIOUTPUT PROBLEM IS A SUPERVISED LEARNING PROBLEM WITH SEVERAL OUTPUTS TO PREDICT THAT IS WHEN Y IS A 2D ARRAY OF SIZESAMPLES NOUTPUTS
WHEN THERE IS NO CORRELATION BETWEEN THE OUTPUTS A VERY SIMPLE WAY TO SOLVE THIS KIND OF PROBLEM IS TO BUILD N INDEPENDENT MODELS IE ONE FOR EACH OUTPUT AND THEN TO USE THOSE MODELS TO INDEPENDENTLY PREDICT EACH ONE OF THE N OUTPUTS HOWEVER BECAUSE IT IS LIKELY THAT THE OUTPUT VALUES RELATED TO THE SAME INPUT ARE THEMSELVES CORRELATED AN OFTEN BETTER WAY IS TO BUILD A SINGLE MODEL CAPABLE OF PREDICTING SIMULTANEOUSLY ALL N OUTPUTS FIRST IT REQUIRES LOWER TRAINING TIME SINCE ONLY A SINGLE ESTIMATOR IS BUILT SECOND THE GENERALIZATION ACCURACY OF THE RESULTING ESTIMATOR MAY OFTEN BE INCREASED
WITH REGARD TO DECISION TREES THIS STRATEGY CAN READILY BE USED TO SUPPORT MULTIOUTPUT PROBLEMS THIS REQUIRES THE FOLLOWING CHANGES
• STORE N OUTPUT VALUES IN LEAVES INSTEAD OF 1
• USE SPLITTING CRITERIA THAT COMPUTE THE AVERAGE REDUCTION ACROSS ALL N OUTPUTS
THIS MODULE OFFERS SUPPORT FOR MULTIOUTPUT PROBLEMS BY IMPLEMENTING THIS STRATEGY IN BOTH DECISIONTREECLASSIFIER ANDDECISIONTREEREgressor IF A DECISION TREE IS FIT ON AN OUTPUT ARRAY Y OF SIZE NSAMPLES NOUTPUTS THEN THE RESULTING ESTIMATOR WILL
• OUTPUT NOUTPUT VALUES UPON PREDICT
• OUTPUT A LIST OF NOUTPUT ARRAYS OF CLASS PROBABILITIES UPON PREDICTPROBA
THE USE OF MULTIOUTPUT TREES FOR REGRESSION IS DEMONSTRATED IN MULTIOUTPUT DECISION TREE REGRESSION IN THIS EXAMPLE THE INPUT X IS A SINGLE REAL VALUE AND THE OUTPUTS Y ARE THE SINE AND COSINE OF X
THE USE OF MULTIOUTPUT TREES FOR CLASSIFICATION IS DEMONSTRATED IN FACE COMPLETION WITH A MULTIOUTPUT ESTIMATORS IN THIS EXAMPLE THE INPUTS X ARE THE PIXELS OF THE UPPER HALF OF FACES AND THE OUTPUTS Y ARE THE PIXELS OF THE LOWER HALF OF THOSE FACES
EXAMPLES
•MULTIOUTPUT DECISION TREE REGRESSION
•FACE COMPLETION WITH A MULTIOUTPUT ESTIMATORS
```

284 CHAPTER 3 USER GUIDE

REFERENCES

- M DUMONT ET AL FAST MULTICLASS IMAGE ANNOTATION WITH RANDOM SUBWINDOWS AND MULTIPLE OUTPUT RANDOMIZED TREES INTERNATIONAL CONFERENCE ON COMPUTER VISION THEORY AND APPLICATIONS 2009

COMPLEXITY

IN GENERAL THE RUN TIME COST TO CONSTRUCT A BALANCED BINARY TREE IS  $O(n \log n)$  AND QUERY TIME  $O(\log n)$  ALTHOUGH THE TREE CONSTRUCTION ALGORITHM ATTEMPTS TO GENERATE BALANCED TREES THEY WILL NOT ALWAYS BE BALANCED ASSUMING THAT THE SUBTREES REMAIN APPROXIMATELY BALANCED THE COST AT EACH NODE CONSISTS OF SEARCHING THROUGH  $n$  TO FIND THE FEATURE THAT OFFERS THE LARGEST REDUCTION IN ENTROPY THIS HAS A COST OF  $O(n \log n)$  AT EACH NODE LEADING TO A TOTAL COST OVER THE ENTIRE TREES BY SUMMING THE COST AT EACH NODE OF  $O(n^2 \log n)$

SCIKITLEARN OFFERS A MORE EFFICIENT IMPLEMENTATION FOR THE CONSTRUCTION OF DECISION TREES A NAIVE IMPLEMENTATION AS ABOVE WOULD RECOMPUTE THE CLASS LABEL HISTOGRAMS FOR CLASSIFICATION OR THE MEANS FOR REGRESSION AT FOR EACH NEW SPLIT POINT ALONG A GIVEN FEATURE PRESORTING THE FEATURE OVER ALL RELEVANT SAMPLES AND RETAINING A RUNNING LABEL COUNT WILL REDUCE THE COMPLEXITY AT EACH NODE TO  $O(n \log n)$  WHICH RESULTS IN A TOTAL COST OF  $O(n^2 \log n)$  THIS IS AN OPTION FOR ALL TREE BASED ALGORITHMS BY DEFAULT IT IS TURNED ON FOR GRADIENT BOOSTING WHERE IN GENERAL IT MAKES TRAINING FASTER BUT TURNED OFF FOR ALL OTHER ALGORITHMS AS IT TENDS TO SLOW DOWN TRAINING WHEN TRAINING DEEP TREES

TIPS ON PRACTICAL USE

- DECISION TREES TEND TO OVERFIT ON DATA WITH A LARGE NUMBER OF FEATURES GETTING THE RIGHT RATIO OF SAMPLES TO NUMBER OF FEATURES IS IMPORTANT SINCE A TREE WITH FEW SAMPLES IN HIGH DIMENSIONAL SPACE IS VERY LIKELY TO OVERFIT



SCIKITLEARN USER GUIDE RELEASE 0213

- CONSIDER PERFORMING DIMENSIONALITY REDUCTION PCA OR FEATURE SELECTION BEFOREHAND TO GIVE YOUR TREE A BETTER CHANCE OF FINDING FEATURES THAT ARE DISCRIMINATIVE
- UNDERSTANDING THE DECISION TREE STRUCTURE WILL HELP IN GAINING MORE INSIGHTS ABOUT HOW THE DECISION TREE MAKES PREDICTIONS WHICH IS IMPORTANT FOR UNDERSTANDING THE IMPORTANT FEATURES IN THE DATA
- VISUALISE YOUR TREE AS YOU ARE TRAINING BY USING THE EXPORT FUNCTION USE MAXDEPTH3 AS AN INITIAL TREE DEPTH TO GET A FEEL FOR HOW THE TREE IS FITTING TO YOUR DATA AND THEN INCREASE THE DEPTH
- REMEMBER THAT THE NUMBER OF SAMPLES REQUIRED TO POPULATE THE TREE DOUBLES FOR EACH ADDITIONAL LEVEL THE TREE GROWS TO USE MAXDEPTH TO CONTROL THE SIZE OF THE TREE TO PREVENT OVERFITTING
- USE MINSAMPLESPLIT OR MINSAMPLESLEAF TO ENSURE THAT MULTIPLE SAMPLES INFORM EVERY DECISION IN THE TREE BY CONTROLLING WHICH SPLITS WILL BE CONSIDERED A VERY SMALL NUMBER WILL USUALLY MEAN THE TREE WILL OVERFIT WHEREAS A LARGE NUMBER WILL PREVENT THE TREE FROM LEARNING THE DATA TRY MINSAMPLESLEAF5 AS AN INITIAL VALUE IF THE SAMPLE SIZE VARIES GREATLY A FLOAT NUMBER CAN BE USED AS PERCENTAGE IN THESE TWO PARAMETERS WHILE MINSAMPLESPLIT CAN CREATE ARBITRARILY SMALL LEAVES MINSAMPLESLEAF GUARANTEES THAT EACH LEAF HAS A MINIMUM SIZE AVOIDING LOW VARIANCE OVERFIT LEAF NODES IN REGRESSION PROBLEMS FOR CLASSIFICATION WITH FEW CLASSES MINSAMPLESLEAF1 IS OFTEN THE BEST CHOICE
- BALANCE YOUR DATASET BEFORE TRAINING TO PREVENT THE TREE FROM BEING BIASED TOWARD THE CLASSES THAT ARE DOMINANT CLASS BALANCING CAN BE DONE BY SAMPLING AN EQUAL NUMBER OF SAMPLES FROM EACH CLASS OR PREFERABLY BY NORMALIZING THE SUM OF THE SAMPLE WEIGHTS SAMPLEWEIGHT FOR EACH CLASS TO THE SAME VALUE ALSO NOTE THAT WEIGHTBASED PREPRUNING CRITERIA SUCH AS MINWEIGHTFRACTIONLEAF WILL THEN BE LESS BIASED TOWARD DOMINANT CLASSES THAN CRITERIA THAT ARE NOT AWARE OF THE SAMPLE WEIGHTS LIKE MINSAMPLESLEAF
- IF THE SAMPLES ARE WEIGHTED IT WILL BE EASIER TO OPTIMIZE THE TREE STRUCTURE USING WEIGHTBASED PREPRUNING CRITERION SUCH AS MINWEIGHTFRACTIONLEAF WHICH ENSURE THAT LEAF NODES CONTAIN AT LEAST A FRACTION OF THE OVERALL SUM OF THE SAMPLE WEIGHTS
- ALL DECISION TREES USE NPFLOAT32 ARRAYS INTERNALLY IF TRAINING DATA IS NOT IN THIS FORMAT A COPY OF THE DATASET WILL BE MADE
- IF THE INPUT MATRIX X IS VERY SPARSE IT IS RECOMMENDED TO CONVERT TO SPARSE CSC MATRIX BEFORE CALLING FIT AND SPARSE CSR MATRIX BEFORE CALLING PREDICT TRAINING TIME CAN BE ORDERS OF MAGNITUDE FASTER FOR A SPARSE MATRIX INPUT COMPARED TO A DENSE MATRIX WHEN FEATURES HAVE ZERO VALUES IN MOST OF THE SAMPLES

TREE ALGORITHMS ID3 C45 C50 AND CART  
WHAT ARE ALL THE VARIOUS DECISION TREE ALGORITHMS AND HOW DO THEY DIFFER FROM EACH OTHER WHICH ONE IS IMPLEMENTED IN SCIKITLEARN

ID3 ITERATIVE DICHOTOMISER 3 WAS DEVELOPED IN 1986 BY ROSS QUINLAN THE ALGORITHM CREATES A MULTIWAY TREE FINDING FOR EACH NODE IE IN A GREEDY MANNER THE CATEGORICAL FEATURE THAT WILL YIELD THE LARGEST INFORMATION GAIN FOR CATEGORICAL TARGETS TREES ARE GROWN TO THEIR MAXIMUM SIZE AND THEN A PRUNING STEP IS USUALLY APPLIED TO IMPROVE THE ABILITY OF THE TREE TO GENERALISE TO UNSEEN DATA  
C45 IS THE SUCCESSOR TO ID3 AND REMOVED THE RESTRICTION THAT FEATURES MUST BE CATEGORICAL BY DYNAMICALLY DEFINING A DISCRETE ATTRIBUTE BASED ON NUMERICAL VARIABLES THAT PARTITIONS THE CONTINUOUS ATTRIBUTE VALUE INTO A DISCRETE SET OF INTERVALS C45 CONVERTS THE TRAINED TREES IE THE OUTPUT OF THE ID3 ALGORITHM INTO SETS OF IF THEN RULES THESE ACCURACY OF EACH RULE IS THEN EVALUATED TO DETERMINE THE ORDER IN WHICH THEY SHOULD BE APPLIED PRUNING IS DONE BY REMOVING A RULE'S PRECONDITION IF THE ACCURACY OF THE RULE IMPROVES WITHOUT IT  
C50 IS QUINLAN'S LATEST VERSION RELEASE UNDER A PROPRIETARY LICENSE IT USES LESS MEMORY AND BUILDS SMALLER RULESETS THAN C45 WHILE BEING MORE ACCURATE  
CART CLASSIFICATION AND REGRESSION TREES IS VERY SIMILAR TO C45 BUT IT DIFFERS IN THAT IT SUPPORTS NUMERICAL TARGET VARIABLES REGRESSION AND DOES NOT COMPUTE RULE SETS CART CONSTRUCTS BINARY TREES USING THE FEATURE AND THRESHOLD THAT YIELD THE LARGEST INFORMATION GAIN AT EACH NODE  
31 SUPERVISED LEARNING 287

SCIKITLEARN USER GUIDE RELEASE 0213

SCIKITLEARN USES AN OPTIMISED VERSION OF THE CART ALGORITHM HOWEVER SCIKITLEARN IMPLEMENTATION DOES NOT SUPPORT CATEGORICAL VARIABLES FOR NOW

MATHEMATICAL FORMULATION

GIVEN TRAINING VECTORS  $\{x_i \in \mathbb{R}^L \mid 1 \leq i \leq n\}$  AND A LABEL VECTOR  $y \in \mathbb{R}$  A DECISION TREE RECURSIVELY PARTITIONS THE SPACE SUCH THAT THE SAMPLES WITH THE SAME LABELS ARE GROUPED TOGETHER

LET THE DATA AT NODE  $t$  BE REPRESENTED BY  $D_t$  FOR EACH CANDIDATE SPLIT  $s = (f, \tau)$  CONSISTING OF A FEATURE  $f$  AND THRESHOLD  $\tau$  PARTITION THE DATA INTO  $D_{t, \text{left}}$  AND  $D_{t, \text{right}}$  SUBSETS

$D_{t, \text{left}} = \{x \in D_t \mid f(x) \leq \tau\}$

$D_{t, \text{right}} = \{x \in D_t \mid f(x) > \tau\}$

THE IMPURITY AT  $t$  IS COMPUTED USING AN IMPURITY FUNCTION  $J$  THE CHOICE OF WHICH DEPENDS ON THE TASK BEING SOLVED CLASSIFICATION OR REGRESSION

$J(D_t) = \sum_{k=1}^K p_k \log p_k$

SELECT THE PARAMETERS THAT MINIMISES THE IMPURITY

$(f, \tau) = \underset{(f, \tau)}{\operatorname{argmin}} J(D_t, (f, \tau))$

RECURSE FOR SUBSETS  $D_{t, \text{left}}$  AND  $D_{t, \text{right}}$  UNTIL THE MAXIMUM ALLOWABLE DEPTH IS REACHED  $\max_{\text{depth}}$  OR  $n \leq 1$

CLASSIFICATION CRITERIA

IF A TARGET IS A CLASSIFICATION OUTCOME TAKING ON VALUES  $0, 1, \dots, K-1$  FOR NODE  $t$  REPRESENTING A REGION  $R_t$  WITH  $n_t$  OBSERVATIONS LET

$p_k = \frac{1}{n_t} \sum_{i \in R_t} \mathbb{1}_{y_i = k}$

BE THE PROPORTION OF CLASS  $k$  OBSERVATIONS IN NODE  $t$

COMMON MEASURES OF IMPURITY ARE GINI

$J(D_t) = \sum_{k=1}^K p_k (1 - p_k)$

ENTROPY

$J(D_t) = -\sum_{k=1}^K p_k \log p_k$

AND MISCLASSIFICATION

$J(D_t) = 1 - \max_{k \in \{0, 1, \dots, K-1\}} p_k$

WHERE  $D_t$  IS THE TRAINING DATA IN NODE  $t$

288 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

REGRESSION CRITERIA

IF THE TARGET IS A CONTINUOUS VALUE THEN FOR NODE  $n$  REPRESENTING A REGION  $R_n$  WITH  $N_n$  OBSERVATIONS COMMON CRITERIA TO MINIMISE AS FOR DETERMINING LOCATIONS FOR FUTURE SPLITS ARE MEAN SQUARED ERROR WHICH MINIMIZES THE L2 ERROR USING MEAN VALUES AT TERMINAL NODES AND MEAN ABSOLUTE ERROR WHICH MINIMIZES THE L1 ERROR USING MEDIAN VALUES AT TERMINAL NODES

MEAN SQUARED ERROR

$$\frac{1}{N_n} \sum_{i \in R_n} (y_i - \hat{y}_n)^2$$
$$\frac{1}{N_n} \sum_{i \in R_n} |y_i - \hat{y}_n|$$

MEAN ABSOLUTE ERROR

$$\frac{1}{N_n} \sum_{i \in R_n} |y_i - \hat{y}_n|$$

WHERE  $R_n$  IS THE TRAINING DATA IN NODE  $n$

REFERENCES

- [HTTPSENWIKIPEDIAORGWIKIDECISIONTREELEARNING](http://en.wikipedia.org/wiki/Decision_tree_learning)
- [HTTPSENWIKIPEDIAORGWIKIPREDICTIVEANALYTICS](http://en.wikipedia.org/wiki/Predictive_analytics)
- L BREIMAN J FRIEDMAN R OLSHEN AND C STONE CLASSIFICATION AND REGRESSION TREES WADSWORTH BELMONT CA 1984
- JR QUINLAN C4 5 PROGRAMS FOR MACHINE LEARNING MORGAN KAUFMANN 1993
- T HASTIE R TIBSHIRANI AND J FRIEDMAN ELEMENTS OF STATISTICAL LEARNING SPRINGER 2009

3111 ENSEMBLE METHODS

THE GOAL OF ENSEMBLE METHODS IS TO COMBINE THE PREDICTIONS OF SEVERAL BASE ESTIMATORS BUILT WITH A GIVEN LEARNING ALGORITHM IN ORDER TO IMPROVE GENERALIZABILITY ROBUSTNESS OVER A SINGLE ESTIMATOR

TWO FAMILIES OF ENSEMBLE METHODS ARE USUALLY DISTINGUISHED

- IN AVERAGING METHODS THE DRIVING PRINCIPLE IS TO BUILD SEVERAL ESTIMATORS INDEPENDENTLY AND THEN TO AVERAGE THEIR PREDICTIONS ON AVERAGE THE COMBINED ESTIMATOR IS USUALLY BETTER THAN ANY OF THE SINGLE BASE ESTIMATOR BECAUSE ITS VARIANCE IS REDUCED

EXAMPLES BAGGING METHODS FORESTS OF RANDOMIZED TREES

- BY CONTRAST IN BOOSTING METHODS BASE ESTIMATORS ARE BUILT SEQUENTIALLY AND ONE TRIES TO REDUCE THE BIAS OF THE COMBINED ESTIMATOR THE MOTIVATION IS TO COMBINE SEVERAL WEAK MODELS TO PRODUCE A POWERFUL ENSEMBLE

EXAMPLES ADABOOST GRADIENT TREE BOOSTING

SCIKITLEARN USER GUIDE RELEASE 0213

BAGGING METAESTIMATOR

IN ENSEMBLE ALGORITHMS BAGGING METHODS FORM A CLASS OF ALGORITHMS WHICH BUILD SEVERAL INSTANCES OF A BLACKBOX ESTIMATOR ON RANDOM SUBSETS OF THE ORIGINAL TRAINING SET AND THEN AGGREGATE THEIR INDIVIDUAL PREDICTIONS TO FORM A FINAL PREDICTION THESE METHODS ARE USED AS A WAY TO REDUCE THE VARIANCE OF A BASE ESTIMATOR EG A DECISION TREE BY INTRODUCING RANDOMIZATION INTO ITS CONSTRUCTION PROCEDURE AND THEN MAKING AN ENSEMBLE OUT OF IT IN MANY CASES BAGGING METHODS CONSTITUTE A VERY SIMPLE WAY TO IMPROVE WITH RESPECT TO A SINGLE MODEL WITHOUT MAKING IT NECESSARY TO ADAPT THE UNDERLYING BASE ALGORITHM AS THEY PROVIDE A WAY TO REDUCE OVERFITTING BAGGING METHODS WORK BEST WITH STRONG AND COMPLEX MODELS EG FULLY DEVELOPED DECISION TREES IN CONTRAST WITH BOOSTING METHODS WHICH USUALLY WORK BEST WITH WEAK MODELS EG SHALLOW DECISION TREES

BAGGING METHODS COME IN MANY FLAVOURS BUT MOSTLY DIFFER FROM EACH OTHER BY THE WAY THEY DRAW RANDOM SUBSETS OF THE TRAINING SET

- WHEN RANDOM SUBSETS OF THE DATASET ARE DRAWN AS RANDOM SUBSETS OF THE SAMPLES THEN THIS ALGORITHM IS KNOWN AS PASTING B1999
- WHEN SAMPLES ARE DRAWN WITH REPLACEMENT THEN THE METHOD IS KNOWN AS BAGGING B1996
- WHEN RANDOM SUBSETS OF THE DATASET ARE DRAWN AS RANDOM SUBSETS OF THE FEATURES THEN THE METHOD IS KNOWN AS RANDOM SUBSPACES H1998
- FINALLY WHEN BASE ESTIMATORS ARE BUILT ON SUBSETS OF BOTH SAMPLES AND FEATURES THEN THE METHOD IS KNOWN AS RANDOM PATCHES LG2012

IN SCIKITLEARN BAGGING METHODS ARE OFFERED AS A UNIFIED BAGGINGCLASSIFIER METAESTIMATOR RESP BAGGINGREGRESSOR TAKING AS INPUT A USERSPECIFIED BASE ESTIMATOR ALONG WITH PARAMETERS SPECIFYING THE STRATEGY TO DRAW RANDOM SUBSETS IN PARTICULAR MAXSAMPLES ANDMAXFEATURES CONTROL THE SIZE OF THE SUBSETS IN TERMS OF SAMPLES AND FEATURES WHILE BOOTSTRAP ANDBOOTSTRAPFEATURES CONTROL WHETHER SAMPLES AND FEATURES ARE DRAWN WITH OR WITHOUT REPLACEMENT WHEN USING A SUBSET OF THE AVAILABLE SAMPLES THE GENERALIZATION ACCURACY CAN BE ESTIMATED WITH THE OUTFBAG SAMPLES BY SETTING OOBSCORETRUE AS AN EXAMPLE THE SNIPPET BELOW ILLUSTRATES HOW TO INSTANTIATE A BAGGING ENSEMBLE OF KNEIGHBORSCLASSIFIER BASE ESTIMATORS EACH BUILT ON RANDOM SUBSETS OF 50 OF THE SAMPLES AND 50 OF THE FEATURES

```
FROM SKLEARNENSEMBLE IMPORT BAGGINGCLASSIFIER
FROM SKLEARNNEIGHBORS IMPORT KNEIGHBORSCLASSIFIER
BAGGING BAGGINGCLASSIFIERKNEIGHBORSCLASSIFIER
MAXSAMPLES05 MAXFEATURES05
```

EXAMPLES

- SINGLE ESTIMATOR VERSUS BAGGING BIASVARIANCE DECOMPOSITION

REFERENCES

FORESTS OF RANDOMIZED TREES

THESKLEARNENSEMBLE MODULE INCLUDES TWO AVERAGING ALGORITHMS BASED ON RANDOMIZED DECISION TREES THE RAN DOMFOREST ALGORITHM AND THE EXTRATREES METHOD BOTH ALGORITHMS ARE PERTURBANDCOMBINE TECHNIQUES B1998 SPECIFICALLY DESIGNED FOR TREES THIS MEANS A DIVERSE SET OF CLASSIFIERS IS CREATED BY INTRODUCING RANDOMNESS IN THE CLASSIFIER CONSTRUCTION THE PREDICTION OF THE ENSEMBLE IS GIVEN AS THE AVERAGED PREDICTION OF THE INDIVIDUAL CLASSIFIERS



SCIKITLEARN USER GUIDE RELEASE 0213

AS OTHER CLASSIFIERS FOREST CLASSIFIERS HAVE TO BE FITTED WITH TWO ARRAYS A SPARSE OR DENSE ARRAY X OF SIZE NSAMPLES NFEATURES HOLDING THE TRAINING SAMPLES AND AN ARRAY Y OF SIZE NSAMPLES HOLDING THE TARGET VALUES CLASS LABELS FOR THE TRAINING SAMPLES

```
FROM SKLEARNENSEMBLE IMPORT RANDOMFORESTCLASSIFIER
X 0 0 1 1
Y 0 1
CLF RANDOMFORESTCLASSIFIERNESTIMATORS10
CLF CLFFITX Y
```

LIKE DECISION TREES FORESTS OF TREES ALSO EXTEND TO MULTIOUTPUT PROBLEMS IF Y IS AN ARRAY OF SIZE NSAMPLES

NOOUTPUTS

RANDOM FORESTS

IN RANDOM FORESTS SEE RANDOMFORESTCLASSIFIER ANDRANDOMFORESTREGRESSOR CLASSES EACH TREE IN THE ENSEMBLE IS BUILT FROM A SAMPLE DRAWN WITH REPLACEMENT IE A BOOTSTRAP SAMPLE FROM THE TRAINING SET

FURTHERMORE WHEN SPLITTING EACH NODE DURING THE CONSTRUCTION OF A TREE THE BEST SPLIT IS FOUND EITHER FROM ALL INPUT FEATURES OR A RANDOM SUBSET OF SIZE MAXFEATURES SEE THE PARAMETER TUNING GUIDELINES FOR MORE DETAILS

THE PURPOSE OF THESE TWO SOURCES OF RANDOMNESS IS TO DECREASE THE VARIANCE OF THE FOREST ESTIMATOR INDEED INDIVIDUAL DECISION TREES TYPICALLY EXHIBIT HIGH VARIANCE AND TEND TO OVERFIT THE INJECTED RANDOMNESS IN FORESTS YIELD DECISION TREES WITH SOMEWHAT DECOUPLED PREDICTION ERRORS BY TAKING AN AVERAGE OF THOSE PREDICTIONS SOME ERRORS CAN CANCEL OUT

RANDOM FORESTS ACHIEVE A REDUCED VARIANCE BY COMBINING DIVERSE TREES SOMETIMES AT THE COST OF A SLIGHT INCREASE IN BIAS

IN PRACTICE THE VARIANCE REDUCTION IS OFTEN SIGNIFICANT HENCE YIELDING AN OVERALL BETTER MODEL

IN CONTRAST TO THE ORIGINAL PUBLICATION B2001 THE SCIKITLEARN IMPLEMENTATION COMBINES CLASSIFIERS BY AVERAGING THEIR PROBABILISTIC PREDICTION INSTEAD OF LETTING EACH CLASSIFIER VOTE FOR A SINGLE CLASS

EXTREMELY RANDOMIZED TREES

IN EXTREMELY RANDOMIZED TREES SEE EXTRATREESCLASSIFIER ANDEXTRATREESREGRESSOR CLASSES

RANDOMNESS GOES ONE STEP FURTHER IN THE WAY SPLITS ARE COMPUTED AS IN RANDOM FORESTS A RANDOM SUBSET OF CANDIDATE FEATURES IS USED BUT INSTEAD OF LOOKING FOR THE MOST DISCRIMINATIVE THRESHOLDS THRESHOLDS ARE DRAWN AT RANDOM FOR EACH CANDIDATE FEATURE AND THE BEST OF THESE RANDOMLYGENERATED THRESHOLDS IS PICKED AS THE SPLITTING RULE THIS USUALLY ALLOWS TO REDUCE THE VARIANCE OF THE MODEL A BIT MORE AT THE EXPENSE OF A SLIGHTLY GREATER INCREASE IN BIAS

```
FROM SKLEARNMODELSELECTION IMPORT CROSSVALSCORE
FROM SKLEARNDATASETS IMPORT MAKEBLOBS
FROM SKLEARNENSEMBLE IMPORT RANDOMFORESTCLASSIFIER
FROM SKLEARNENSEMBLE IMPORT EXTRATREESCLASSIFIER
FROM SKLEARNTREE IMPORT DECISIONTREECLASSIFIER
X Y MAKEBLOBSNSAMPLES10000 NFEATURES10 CENTERS100
RANDOMSTATE0
CLF DECISIONTREECLASSIFIERMAXDEPTH NONE MINSAMPLESSPLIT2
RANDOMSTATE0
SCORES CROSSVALSCORECLF X Y CV5
SCORESMEAN
098
CLF RANDOMFORESTCLASSIFIERNESTIMATORS10 MAXDEPTH NONE
MINSAMPLESSPLIT2 RANDOMSTATE0
SCORES CROSSVALSCORECLF X Y CV5
31 SUPERVISED LEARNING 291
```

SCORESMEAN

0999

CLF EXTRATREESCLASSIFIERNESTIMATORS10 MAXDEPTH NONE

MINSAMPLESSPLIT2 RANDOMSTATE0

SCORES CROSSVALSCORECLF X Y CV5

SCORESMEAN 0999

TRUE

PARAMETERS

THE MAIN PARAMETERS TO ADJUST WHEN USING THESE METHODS IS NESTIMATORS ANDMAXFEATURES THE FORMER IS THE NUMBER OF TREES IN THE FOREST THE LARGER THE BETTER BUT ALSO THE LONGER IT WILL TAKE TO COMPUTE IN ADDITION NOTE THAT RESULTS WILL STOP GETTING SIGNIFICANTLY BETTER BEYOND A CRITICAL NUMBER OF TREES THE LATTER IS THE SIZE OF THE RANDOM SUBSETS OF FEATURES TO CONSIDER WHEN SPLITTING A NODE THE LOWER THE GREATER THE REDUCTION OF VARIANCE BUT ALSO THE GREATER THE INCREASE IN BIAS EMPIRICAL GOOD DEFAULT VALUES ARE MAXFEATURESNONE ALWAYS CONSIDERING ALL FEATURES INSTEAD OF A RANDOM SUBSET FOR REGRESSION PROBLEMS AND MAXFEATURESSQRT USING A RANDOM SUBSET OF SIZESQRTNFEATURES FOR CLASSIFICATION TASKS WHERE NFEATURES IS THE NUMBER OF FEATURES IN THE DATA GOOD RESULTS ARE OFTEN ACHIEVED WHEN SETTING MAXDEPTHNONE IN COMBINATION WITH MINSAMPLESSPLIT2 IE WHEN FULLY DEVELOPING THE TREES BEAR IN MIND THOUGH THAT THESE VALUES ARE USUALLY NOT OPTIMAL AND MIGHT RESULT IN MODELS THAT CONSUME A LOT OF RAM THE BEST PARAMETER VALUES SHOULD ALWAYS BE CROSSVALIDATED IN ADDITION NOTE THAT IN RANDOM FORESTS BOOTSTRAP SAMPLES ARE USED BY DEFAULT BOOTSTRAPTRUE WHILE THE DEFAULT STRATEGY FOR EXTRATREES IS TO USE THE WHOLE DATASET BOOTSTRAPFALSE WHEN USING BOOTSTRAP SAMPLING THE GENERALIZATION ACCURACY CAN BE ESTIMATED ON THE LEFT OUT OR OUTOFBAG SAMPLES THIS CAN BE ENABLED BY SETTING OOBSCORETRUE NOTE THE SIZE OF THE MODEL WITH THE DEFAULT PARAMETERS IS  $\frac{N}{M} \times \frac{N}{M}$  WHERE  $N$  IS THE NUMBER OF TREES AND  $M$  IS THE NUMBER OF SAMPLES IN ORDER TO REDUCE THE SIZE OF THE MODEL YOU CAN CHANGE THESE PARAMETERS MINSAMPLESSPLIT MAXLEAFNODES MAXDEPTH ANDMINSAMPLESLEAF

SCIKITLEARN USER GUIDE RELEASE 0213

PARALLELIZATION

FINALLY THIS MODULE ALSO FEATURES THE PARALLEL CONSTRUCTION OF THE TREES AND THE PARALLEL COMPUTATION OF THE PREDICTIONS THROUGH THE NJOBS PARAMETER IF NJOBSK THEN COMPUTATIONS ARE PARTITIONED INTO KJOBS AND RUN ON KCORES OF THE MACHINE IF NJOBS1 THEN ALL CORES AVAILABLE ON THE MACHINE ARE USED NOTE THAT BECAUSE OF INTERPROCESS COMMUNICATION OVERHEAD THE SPEEDUP MIGHT NOT BE LINEAR IE USING KJOBS WILL UNFORTUNATELY NOT BE KTIMES AS FAST SIGNIFICANT SPEEDUP CAN STILL BE ACHIEVED THOUGH WHEN BUILDING A LARGE NUMBER OF TREES OR WHEN BUILDING A SINGLE TREE REQUIRES A FAIR AMOUNT OF TIME EG ON LARGE DATASETS

EXAMPLES

- PLOT THE DECISION SURFACES OF ENSEMBLES OF TREES ON THE IRIS DATASET
- PIXEL IMPORTANCES WITH A PARALLEL FOREST OF TREES
- FACE COMPLETION WITH A MULTIOUTPUT ESTIMATORS

REFERENCES

- P GEURTS D ERNST AND L WEHENKEL “EXTREMELY RANDOMIZED TREES” MACHINE LEARNING 631 342 2006

FEATURE IMPORTANCE EVALUATION

THE RELATIVE RANK IE DEPTH OF A FEATURE USED AS A DECISION NODE IN A TREE CAN BE USED TO ASSESS THE RELATIVE IMPORTANCE OF THAT FEATURE WITH RESPECT TO THE PREDICTABILITY OF THE TARGET VARIABLE FEATURES USED AT THE TOP OF THE TREE CONTRIBUTE TO THE FINAL PREDICTION DECISION OF A LARGER FRACTION OF THE INPUT SAMPLES THE EXPECTED FRACTION OF THE SAMPLES THEY CONTRIBUTE TO CAN THUS BE USED AS AN ESTIMATE OF THE RELATIVE IMPORTANCE OF THE FEATURES IN SCIKITLEARN THE FRACTION OF SAMPLES A FEATURE CONTRIBUTES TO IS COMBINED WITH THE DECREASE IN IMPURITY FROM SPLITTING THEM TO CREATE A NORMALIZED ESTIMATE OF THE PREDICTIVE POWER OF THAT FEATURE

BYAVERAGING THE ESTIMATES OF PREDICTIVE ABILITY OVER SEVERAL RANDOMIZED TREES ONE CAN REDUCE THE VARIANCE OF SUCH AN ESTIMATE AND USE IT FOR FEATURE SELECTION THIS IS KNOWN AS THE MEAN DECREASE IN IMPURITY OR MDI REFER TO L2014 FOR MORE INFORMATION ON MDI AND FEATURE IMPORTANCE EVALUATION WITH RANDOM FORESTS

THE FOLLOWING EXAMPLE SHOWS A COLORCODED REPRESENTATION OF THE RELATIVE IMPORTANCES OF EACH INDIVIDUAL PIXEL FOR A FACE RECOGNITION TASK USING A EXTRATREESCLASSIFIER MODEL

IN PRACTICE THOSE ESTIMATES ARE STORED AS AN ATTRIBUTE NAMED FEATUREIMPORTANCES ON THE FITTED MODEL THIS IS AN ARRAY WITH SHAPE NFEATURES WHOSE VALUES ARE POSITIVE AND SUM TO 10 THE HIGHER THE VALUE THE MORE IMPORTANT IS THE CONTRIBUTION OF THE MATCHING FEATURE TO THE PREDICTION FUNCTION

EXAMPLES

- PIXEL IMPORTANCES WITH A PARALLEL FOREST OF TREES
- FEATURE IMPORTANCES WITH FORESTS OF TREES

REFERENCES

SCIKITLEARN USER GUIDE RELEASE 0213

TOTALLY RANDOM TREES EMBEDDING

RANDOMTREESEMBEDDING IMPLEMENTS AN UNSUPERVISED TRANSFORMATION OF THE DATA USING A FOREST OF COMPLETELY RANDOM TREES RANDOMTREESEMBEDDING ENCODES THE DATA BY THE INDICES OF THE LEAVES A DATA POINT ENDS UP IN THIS INDEX IS THEN ENCODED IN A ONEOFK MANNER LEADING TO A HIGH DIMENSIONAL SPARSE BINARY CODING THIS CODING CAN BE COMPUTED VERY EFFICIENTLY AND CAN THEN BE USED AS A BASIS FOR OTHER LEARNING TASKS THE SIZE AND SPARSITY OF THE CODE CAN BE INFLUENCED BY CHOOSING THE NUMBER OF TREES AND THE MAXIMUM DEPTH PER TREE FOR EACH TREE IN THE ENSEMBLE THE CODING CONTAINS ONE ENTRY OF ONE THE SIZE OF THE CODING IS AT MOST  $2^{MAXDEPTH}$  THE MAXIMUM NUMBER OF LEAVES IN THE FOREST

AS NEIGHBORING DATA POINTS ARE MORE LIKELY TO LIE WITHIN THE SAME LEAF OF A TREE THE TRANSFORMATION PERFORMS AN IMPLICIT NONPARAMETRIC DENSITY ESTIMATION

EXAMPLES

- HASHING FEATURE TRANSFORMATION USING TOTALLY RANDOM TREES
- MANIFOLD LEARNING ON HANDWRITTEN DIGITS LOCALLY LINEAR EMBEDDING ISOMAP COMPARES NONLINEAR DIMENSIONALITY REDUCTION TECHNIQUES ON HANDWRITTEN DIGITS
- FEATURE TRANSFORMATIONS WITH ENSEMBLES OF TREES COMPARES SUPERVISED AND UNSUPERVISED TREE BASED FEATURE TRANSFORMATIONS

SEE ALSO

MANIFOLD LEARNING TECHNIQUES CAN ALSO BE USEFUL TO DERIVE NONLINEAR REPRESENTATIONS OF FEATURE SPACE ALSO THESE APPROACHES FOCUS ALSO ON DIMENSIONALITY REDUCTION

294 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

ADABOOST

THE MODULE SKLEARNENSEMBLE INCLUDES THE POPULAR BOOSTING ALGORITHM ADABOOST INTRODUCED IN 1995 BY FREUND AND SCHAPIRE FS1995

THE CORE PRINCIPLE OF ADABOOST IS TO FIT A SEQUENCE OF WEAK LEARNERS IE MODELS THAT ARE ONLY SLIGHTLY BETTER THAN RANDOM GUESSING SUCH AS SMALL DECISION TREES ON REPEATEDLY MODIFIED VERSIONS OF THE DATA THE PREDICTIONS FROM ALL OF THEM ARE THEN COMBINED THROUGH A WEIGHTED MAJORITY VOTE OR SUM TO PRODUCE THE FINAL PREDICTION THE DATA MODIFICATIONS AT EACH SOCALLED BOOSTING ITERATION CONSIST OF APPLYING WEIGHTS  $\alpha_1, \alpha_2, \dots, \alpha_n$  TO EACH OF THE TRAINING SAMPLES INITIALLY THOSE WEIGHTS ARE ALL SET TO  $\frac{1}{n}$  SO THAT THE FIRST STEP SIMPLY TRAINS A WEAK LEARNER ON THE ORIGINAL DATA FOR EACH SUCCESSIVE ITERATION THE SAMPLE WEIGHTS ARE INDIVIDUALLY MODIFIED AND THE LEARNING ALGORITHM IS REAPPLIED TO THE REWEIGHTED DATA AT A GIVEN STEP THOSE TRAINING EXAMPLES THAT WERE INCORRECTLY PREDICTED BY THE BOOSTED MODEL INDUCED AT THE PREVIOUS STEP HAVE THEIR WEIGHTS INCREASED WHEREAS THE WEIGHTS ARE DECREASED FOR THOSE THAT WERE PREDICTED CORRECTLY AS ITERATIONS PROCEED EXAMPLES THAT ARE DIFFICULT TO PREDICT RECEIVE EVERINCREASING INFLUENCE EACH SUBSEQUENT WEAK LEARNER IS THEREBY FORCED TO CONCENTRATE ON THE EXAMPLES THAT ARE MISSED BY THE PREVIOUS ONES IN THE SEQUENCE HTF

ADABOOST CAN BE USED BOTH FOR CLASSIFICATION AND REGRESSION PROBLEMS

- FOR MULTICLASS CLASSIFICATION ADABOOSTCLASSIFIER IMPLEMENTS ADABOOSTSAMME AND ADABOOSTSAMMER ZZRH2009

- FOR REGRESSION ADABOOSTREGRESSOR IMPLEMENTS ADABOOSTR2 D1997

USAGE

THE FOLLOWING EXAMPLE SHOWS HOW TO FIT AN ADABOOST CLASSIFIER WITH 100 WEAK LEARNERS

```
FROM SKLEARNMODELSELECTION IMPORT CROSSVALSCORE
```

```
FROM SKLEARNDATASETS IMPORT LOADIRIS
```

```
FROM SKLEARNENSEMBLE IMPORT ADABOOSTCLASSIFIER
```

```
31 SUPERVISED LEARNING 295
```

SCIKITLEARN USER GUIDE RELEASE 0213

IRIS LOADIRIS

CLF ADABOOSTCLASSIFIERNESTIMATORS100

SCORES CROSSVALSCORECLF IRISDATA IRISTARGET CV5

SCORESMEAN

09

THE NUMBER OF WEAK LEARNERS IS CONTROLLED BY THE PARAMETER NESTIMATORS THELEARNINGRATE PARAMETER CONTROLS THE CONTRIBUTION OF THE WEAK LEARNERS IN THE FINAL COMBINATION BY DEFAULT WEAK LEARNERS ARE DECISION STUMPS DIFFERENT WEAK LEARNERS CAN BE SPECIFIED THROUGH THE BASEESTIMATOR PARAMETER THE MAIN PARAMETERS TO TUNE TO OBTAIN GOOD RESULTS ARE NESTIMATORS AND THE COMPLEXITY OF THE BASE ESTIMATORS EG ITS DEPTH MAXDEPTH OR MINIMUM REQUIRED NUMBER OF SAMPLES TO CONSIDER A SPLIT MINSAMPLESSPLIT

EXAMPLES

- DISCRETE VERSUS REAL ADABOOST COMPARES THE CLASSIFICATION ERROR OF A DECISION STUMP DECISION TREE AND A BOOSTED DECISION STUMP USING ADABOOSTSAMME AND ADABOOSTSAMMER
- MULTICLASS ADABOOSTED DECISION TREES SHOWS THE PERFORMANCE OF ADABOOSTSAMME AND ADABOOST SAMMER ON A MULTICLASS PROBLEM
- TWOCLASS ADABOOST SHOWS THE DECISION BOUNDARY AND DECISION FUNCTION VALUES FOR A NONLINEARLY SEPARABLE TWOCLASS PROBLEM USING ADABOOSTSAMME
- DECISION TREE REGRESSION WITH ADABOOST DEMONSTRATES REGRESSION WITH THE ADABOOSTR2 ALGORITHM

REFERENCES

GRADIENT TREE BOOSTING

GRADIENT TREE BOOSTING OR GRADIENT BOOSTED REGRESSION TREES GBRT IS A GENERALIZATION OF BOOSTING TO ARBITRARY DIFFERENTIABLE LOSS FUNCTIONS GBRT IS AN ACCURATE AND EFFECTIVE OFFTHESHELF PROCEDURE THAT CAN BE USED FOR BOTH REGRESSION AND CLASSIFICATION PROBLEMS GRADIENT TREE BOOSTING MODELS ARE USED IN A VARIETY OF AREAS INCLUDING WEB SEARCH RANKING AND ECOLOGY

THE ADVANTAGES OF GBRT ARE

- NATURAL HANDLING OF DATA OF MIXED TYPE HETEROGENEOUS FEATURES
- PREDICTIVE POWER
- ROBUSTNESS TO OUTLIERS IN OUTPUT SPACE VIA ROBUST LOSS FUNCTIONS

THE DISADVANTAGES OF GBRT ARE

- SCALABILITY DUE TO THE SEQUENTIAL NATURE OF BOOSTING IT CAN HARDLY BE PARALLELIZED

THE MODULE SKLEARNENSEMBLE PROVIDES METHODS FOR BOTH CLASSIFICATION AND REGRESSION VIA GRADIENT BOOSTED REGRESSION TREES

NOTE SCIKITLEARN 021 INTRODUCES TWO NEW EXPERIMENTAL IMPLEMENTATION OF GRADIENT BOOSTING TREES NAMELY HISTGRADIENTBOOSTINGCLASSIFIER ANDHISTGRADIENTBOOSTINGREGRESSOR INSPIRED BY LIGHT GBM THESE FAST ESTIMATORS FIRST BIN THE INPUT SAMPLES XINTO INTEGERVALUED BINS TYPICALLY 256 BINS WHICH TREMENDOUSLY REDUCES THE NUMBER OF SPLITTING POINTS TO CONSIDER AND ALLOW THE ALGORITHM TO LEVERAGE INTEGERBASED DATA STRUCTURES HISTOGRAMS INSTEAD OF RELYING ON SORTED CONTINUOUS VALUES

296 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

THE NEW HISTOGRAMBASED ESTIMATORS CAN BE ORDERS OF MAGNITUDE FASTER THAN THEIR CONTINUOUS COUNTERPARTS WHEN THE NUMBER OF SAMPLES IS LARGER THAN TENS OF THOUSANDS OF SAMPLES THE API OF THESE NEW ESTIMATORS IS SLIGHTLY DIFFERENT AND SOME OF THE FEATURES FROM GRADIENTBOOSTINGCLASSIFIER ANDGRADIENTBOOSTINGREGRESSOR ARE NOT YET SUPPORTED

THESE NEW ESTIMATORS ARE STILL EXPERIMENTAL FOR NOW THEIR PREDICTIONS AND THEIR API MIGHT CHANGE WITHOUT ANY DEPRECATION CYCLE TO USE THEM YOU NEED TO EXPLICITLY IMPORT ENABLEHISTGRADIENTBOOSTING

EXPLICITLY REQUIRE THIS EXPERIMENTAL FEATURE

FROM SKLEARNEXPERIMENTAL IMPORT ENABLEHISTGRADIENTBOOSTING NOQA

NOW YOU CAN IMPORT NORMALLY FROM ENSEMBLE

FROM SKLEARNENSEMBLE IMPORT HISTGRADIENTBOOSTINGCLASSIFIER

THE FOLLOWING GUIDE FOCUSES ON GRADIENTBOOSTINGCLASSIFIER ANDGRADIENTBOOSTINGREGRESSOR ONLY WHICH MIGHT BE PREFERRED FOR SMALL SAMPLE SIZES SINCE BINNING MAY LEAD TO SPLIT POINTS THAT ARE TOO APPROXIMATE IN THIS SETTING

CLASSIFICATION

GRADIENTBOOSTINGCLASSIFIER SUPPORTS BOTH BINARY AND MULTICLASS CLASSIFICATION THE FOLLOWING EXAMPLE SHOWS HOW TO FIT A GRADIENT BOOSTING CLASSIFIER WITH 100 DECISION STUMPS AS WEAK LEARNERS

```
FROM SKLEARNDATASETS IMPORT MAKEHASTIE102
FROM SKLEARNENSEMBLE IMPORT GRADIENTBOOSTINGCLASSIFIER
X Y MAKEHASTIE102RANDOMSTATE0
XTRAIN XTEST X2000 X2000
YTRAIN YTEST Y2000 Y2000
CLF GRADIENTBOOSTINGCLASSIFIERNESTIMATORS100 LEARNINGRATE10
MAXDEPTH1 RANDOMSTATEOFITXTRAIN YTRAIN
CLFSCOREXTEST YTEST
```

0913

THE NUMBER OF WEAK LEARNERS IE REGRESSION TREES IS CONTROLLED BY THE PARAMETER NESTIMATORS THE SIZE OF EACH TREECAN BE CONTROLLED EITHER BY SETTING THE TREE DEPTH VIA MAXDEPTH OR BY SETTING THE NUMBER OF LEAF NODES VIA MAXLEAFNODES THELEARNINGRATE IS A HYPERPARAMETER IN THE RANGE 00 10 THAT CONTROLS OVERFITTING VIA SHRINKAGE

NOTE CLASSIFICATION WITH MORE THAN 2 CLASSES REQUIRES THE INDUCTION OF NCLASSES REGRESSION TREES AT EACH ITERATION THUS THE TOTAL NUMBER OF INDUCED TREES EQUALS NCLASSES NESTIMATORS FOR DATASETS WITH A LARGE NUMBER OF CLASSES WE STRONGLY RECOMMEND TO USE RANDOMFORESTCLASSIFIER AS AN ALTERNATIVE TO GRADIENTBOOSTINGCLASSIFIER

REGRESSION

GRADIENTBOOSTINGREGRESSOR SUPPORTS A NUMBER OF DIFFERENT LOSS FUNCTIONS FOR REGRESSION WHICH CAN BE SPECIFIED VIA THE ARGUMENT LOSS THE DEFAULT LOSS FUNCTION FOR REGRESSION IS LEAST SQUARES LS

```
IMPORT NUMPY AS NP
FROM SKLEARNMETRICS IMPORT MEANSQUAREDERROR
FROM SKLEARNDATASETS IMPORT MAKEFRIEDMAN1
```

31 SUPERVISED LEARNING 297

SCIKITLEARN USER GUIDE RELEASE 0213

```
FROM SKLEARNENSEMBLE IMPORT GRADIENTBOOSTINGREGRESSOR
X Y MAKEFRIEDMAN1NSAMPLES1200 RANDOMSTATE0 NOISE10
XTRAIN XTEST X200 X200
YTRAIN YTEST Y200 Y200
EST GRADIENTBOOSTINGREGRESSORNESTIMATORS100 LEARNINGRATE01
MAXDEPTH1 RANDOMSTATE0 LOSSLSFITXTRAIN YTRAIN
MEANSQUAREDERRORYTEST ESTPREDICTXTEST
500
```

THE FIGURE BELOW SHOWS THE RESULTS OF APPLYING GRADIENTBOOSTINGREGRESSOR WITH LEAST SQUARES LOSS AND 500 BASE LEARNERS TO THE BOSTON HOUSE PRICE DATASET SKLEARNDATASETSLOADBOSTON THE PLOT ON THE LEFT SHOWS THE TRAIN AND TEST ERROR AT EACH ITERATION THE TRAIN ERROR AT EACH ITERATION IS STORED IN THE TRAINSCORE ATTRIBUTE OF THE GRADIENT BOOSTING MODEL THE TEST ERROR AT EACH ITERATIONS CAN BE OBTAINED VIA THE STAGEDPREDICT METHOD WHICH RETURNS A GENERATOR THAT YIELDS THE PREDICTIONS AT EACH STAGE PLOTS LIKE THESE CAN BE USED TO DETERMINE THE OPTIMAL NUMBER OF TREES IE NESTIMATORS BY EARLY STOPPING THE PLOT ON THE RIGHT SHOWS THE FEATURE IMPORTANCES WHICH CAN BE OBTAINED VIA THE FEATUREIMPORTANCES PROPERTY

EXAMPLES

- GRADIENT BOOSTING REGRESSION
  - GRADIENT BOOSTING OUTFBAG ESTIMATES
- FITTING ADDITIONAL WEAKLEARNERS  
BOTHGRADIENTBOOSTINGREGRESSOR ANDGRADIENTBOOSTINGCLASSIFIER SUPPORT  
WARMSTARTTRUE WHICH ALLOWS YOU TO ADD MORE ESTIMATORS TO AN ALREADY FITTED MODEL



SCIKITLEARN USER GUIDE RELEASE 0213  
ESTSETPARAMSNESTIMATORS200 WARMSTART TRUE SET WARMSTART AND NEW  
<→NR OF TREES  
ESTFITXTRAIN YTRAIN FIT ADDITIONAL 100 TREES TO EST  
MEANSQUAREDERRORYTEST ESTPREDICTXTEST  
384

CONTROLLING THE TREE SIZE  
THE SIZE OF THE REGRESSION TREE BASE LEARNERS DEFINES THE LEVEL OF VARIABLE INTERACTIONS THAT CAN BE CAPTURED BY THE GRADIENT BOOSTING MODEL IN GENERAL A TREE OF DEPTH H CAN CAPTURE INTERACTIONS OF ORDER H THERE ARE TWO WAYS IN WHICH THE SIZE OF THE INDIVIDUAL REGRESSION TREES CAN BE CONTROLLED  
IF YOU SPECIFY MAXDEPTHH THEN COMPLETE BINARY TREES OF DEPTH H WILL BE GROWN SUCH TREES WILL HAVE AT MOST  $2^H$  LEAF NODES AND  $2^H - 1$  SPLIT NODES  
ALTERNATIVELY YOU CAN CONTROL THE TREE SIZE BY SPECIFYING THE NUMBER OF LEAF NODES VIA THE PARAMETER MAXLEAFNODES IN THIS CASE TREES WILL BE GROWN USING BESTFIRST SEARCH WHERE NODES WITH THE HIGHEST IMPROVEMENT IN IMPURITY WILL BE EXPANDED FIRST A TREE WITH MAXLEAFNODES K HAS  $2^{\lceil \log_2 K \rceil} - 1$  SPLIT NODES AND THUS CAN MODEL INTERACTIONS OF UP TO ORDER MAXLEAFNODES  
WE FOUND THAT MAXLEAFNODES K GIVES COMPARABLE RESULTS TO MAXDEPTH K BUT IS SIGNIFICANTLY FASTER TO TRAIN AT THE EXPENSE OF A SLIGHTLY HIGHER TRAINING ERROR THE PARAMETER MAXLEAFNODES CORRESPONDS TO THE VARIABLE J IN THE CHAPTER ON GRADIENT BOOSTING IN F2001 AND IS RELATED TO THE PARAMETER INTERACTIONDEPTH IN R'S GBM PACKAGE WHERE MAXLEAFNODES INTERACTIONDEPTH 1  
MATHEMATICAL FORMULATION  
GBRT CONSIDERS ADDITIVE MODELS OF THE FOLLOWING FORM

$$\hat{y} = \sum_{j=1}^M h_j(x)$$
WHERE  $h_j$  ARE THE BASIS FUNCTIONS WHICH ARE USUALLY CALLED WEAK LEARNERS IN THE CONTEXT OF BOOSTING GRADIENT TREE BOOSTING USES DECISION TREES OF FIXED SIZE AS WEAK LEARNERS DECISION TREES HAVE A NUMBER OF ABILITIES THAT MAKE THEM VALUABLE FOR BOOSTING NAMELY THE ABILITY TO HANDLE DATA OF MIXED TYPE AND THE ABILITY TO MODEL COMPLEX FUNCTIONS SIMILAR TO OTHER BOOSTING ALGORITHMS GBRT BUILDS THE ADDITIVE MODEL IN A GREEDY FASHION

$$\hat{y} = \hat{y}_{j-1} + h_j(x)$$
WHERE THE NEWLY ADDED TREE  $h_j$  TRIES TO MINIMIZE THE LOSS  $L$  GIVEN THE PREVIOUS ENSEMBLE  $\hat{y}_{j-1}$   
$$h_j = \underset{h}{\text{ARG MIN}} L(\hat{y}_{j-1} + h)$$
THE INITIAL MODEL  $\hat{y}_0$  IS PROBLEM SPECIFIC FOR LEAST SQUARES REGRESSION ONE USUALLY CHOOSES THE MEAN OF THE TARGET VALUES  
31 SUPERVISED LEARNING 299

NOTE THE INITIAL MODEL CAN ALSO BE SPECIFIED VIA THE INIT ARGUMENT THE PASSED OBJECT HAS TO IMPLEMENT FIT AND PREDICT

GRADIENT BOOSTING ATTEMPTS TO SOLVE THIS MINIMIZATION PROBLEM NUMERICALLY VIA STEEPEST DESCENT THE STEEPEST DESCENT DIRECTION IS THE NEGATIVE GRADIENT OF THE LOSS FUNCTION EVALUATED AT THE CURRENT MODEL  $\theta_{t-1}$  WHICH CAN BE CALCULATED FOR ANY DIFFERENTIABLE LOSS FUNCTION

$$\eta_t = -\frac{1}{n} \sum_{i=1}^n \nabla_{\theta} \ell(y_i, \theta_{t-1})$$

$$\theta_t = \arg \min_{\theta} \sum_{i=1}^n \ell(y_i, \theta)$$

WHERE THE STEP LENGTH  $\eta_t$  IS CHOSEN USING LINE SEARCH

$$\theta_t = \arg \min_{\theta} \sum_{i=1}^n \ell(y_i, \theta)$$

$$\theta_t = \arg \min_{\theta} \sum_{i=1}^n \ell(y_i, \theta)$$

$$\theta_t = \arg \min_{\theta} \sum_{i=1}^n \ell(y_i, \theta)$$

$$\theta_t = \arg \min_{\theta} \sum_{i=1}^n \ell(y_i, \theta)$$

THE ALGORITHMS FOR REGRESSION AND CLASSIFICATION ONLY DIFFER IN THE CONCRETE LOSS FUNCTION USED

LOSS FUNCTIONS

THE FOLLOWING LOSS FUNCTIONS ARE SUPPORTED AND CAN BE SPECIFIED USING THE PARAMETER LOSS

• REGRESSION

-LEAST SQUARES LS THE NATURAL CHOICE FOR REGRESSION DUE TO ITS SUPERIOR COMPUTATIONAL PROPERTIES THE INITIAL MODEL IS GIVEN BY THE MEAN OF THE TARGET VALUES

-LEAST ABSOLUTE DEVIATION LAD A ROBUST LOSS FUNCTION FOR REGRESSION THE INITIAL MODEL IS GIVEN BY THE MEDIAN OF THE TARGET VALUES

-HUBER HUBER ANOTHER ROBUST LOSS FUNCTION THAT COMBINES LEAST SQUARES AND LEAST ABSOLUTE DEVIATION USE ALPHA TO CONTROL THE SENSITIVITY WITH REGARDS TO OUTLIERS SEE F2001 FOR MORE DETAILS

-QUANTILE QUANTILE A LOSS FUNCTION FOR QUANTILE REGRESSION USE 0 ALPHA 1 TO SPECIFY THE QUANTILE THIS LOSS FUNCTION CAN BE USED TO CREATE PREDICTION INTERVALS SEE PREDICTION INTERVALS FOR GRADIENT BOOSTING REGRESSION

• CLASSIFICATION

-BINOMIAL DEVIANCE DEVIANCE THE NEGATIVE BINOMIAL LOGLIKELIHOOD LOSS FUNCTION FOR BINARY CLASSIFICATION PROVIDES PROBABILITY ESTIMATES THE INITIAL MODEL IS GIVEN BY THE LOG ODDS RATIO

-MULTINOMIAL DEVIANCE DEVIANCE THE NEGATIVE MULTINOMIAL LOGLIKELIHOOD LOSS FUNCTION FOR MULTI CLASS CLASSIFICATION WITH NCLASSES MUTUALLY EXCLUSIVE CLASSES IT PROVIDES PROBABILITY ESTIMATES THE INITIAL MODEL IS GIVEN BY THE PRIOR PROBABILITY OF EACH CLASS AT EACH ITERATION NCLASSES REGRESSION TREES HAVE TO BE CONSTRUCTED WHICH MAKES GBRT RATHER INEFFICIENT FOR DATA SETS WITH A LARGE NUMBER OF CLASSES

-EXPONENTIAL LOSS EXPONENTIAL THE SAME LOSS FUNCTION AS ADABOOSTCLASSIFIER LESS ROBUST TO MISLABELED EXAMPLES THAN DEVIANCE CAN ONLY BE USED FOR BINARY CLASSIFICATION

REGULARIZATION

SHRINKAGE

F2001 PROPOSED A SIMPLE REGULARIZATION STRATEGY THAT SCALES THE CONTRIBUTION OF EACH WEAK LEARNER BY A FACTOR  $\eta$

$\eta = 1 - \eta_{\text{shrinkage}}$

THE PARAMETER  $\eta$  IS ALSO CALLED THE LEARNING RATE BECAUSE IT SCALES THE STEP LENGTH THE GRADIENT DESCENT PROCEDURE IT CAN BE SET VIA THE LEARNINGRATE PARAMETER

THE PARAMETER LEARNINGRATE STRONGLY INTERACTS WITH THE PARAMETER NESTIMATORS THE NUMBER OF WEAK LEARNERS TO FIT SMALLER VALUES OF LEARNINGRATE REQUIRE LARGER NUMBERS OF WEAK LEARNERS TO MAINTAIN A CONSTANT TRAINING ERROR EMPIRICAL EVIDENCE SUGGESTS THAT SMALL VALUES OF LEARNINGRATE FAVOR BETTER TEST ERROR HTF2009 RECOMMEND TO SET THE LEARNING RATE TO A SMALL CONSTANT EG LEARNINGRATE 01 AND CHOOSE NESTIMATORS BY EARLY STOPPING FOR A MORE DETAILED DISCUSSION OF THE INTERACTION BETWEEN LEARNINGRATE ANDNESTIMATORS SEE R2007

SUBSAMPLING

F1999 PROPOSED STOCHASTIC GRADIENT BOOSTING WHICH COMBINES GRADIENT BOOSTING WITH BOOTSTRAP AVERAGING BAGGING AT EACH ITERATION THE BASE CLASSIFIER IS TRAINED ON A FRACTION SUBSAMPLE OF THE AVAILABLE TRAINING DATA THE SUBSAMPLE IS DRAWN WITHOUT REPLACEMENT A TYPICAL VALUE OF SUBSAMPLE IS 05

THE FIGURE BELOW ILLUSTRATES THE EFFECT OF SHRINKAGE AND SUBSAMPLING ON THE GOODNESSOFFIT OF THE MODEL WE CAN CLEARLY SEE THAT SHRINKAGE OUTPERFORMS NOSHRINKAGE SUBSAMPLING WITH SHRINKAGE CAN FURTHER INCREASE THE ACCURACY OF THE MODEL SUBSAMPLING WITHOUT SHRINKAGE ON THE OTHER HAND DOES POORLY

ANOTHER STRATEGY TO REDUCE THE VARIANCE IS BY SUBSAMPLING THE FEATURES ANALOGOUS TO THE RANDOM SPLITS IN RANDOMFORESTCLASSIFIER THE NUMBER OF SUBSAMPLED FEATURES CAN BE CONTROLLED VIA THE MAXFEATURES

PARAMETER

SCIKITLEARN USER GUIDE RELEASE 0213

NOTE USING A SMALL MAXFEATURES VALUE CAN SIGNIFICANTLY DECREASE THE RUNTIME  
STOCHASTIC GRADIENT BOOSTING ALLOWS TO COMPUTE OUTFBAG ESTIMATES OF THE TEST DEVIANC BY COMPUTING THE IMPROVE  
MENT IN DEVIANC ON THE EXAMPLES THAT ARE NOT INCLUDED IN THE BOOTSTRAP SAMPLE IE THE OUTFBAG EXAMPLES THE  
IMPROVEMENTS ARE STORED IN THE ATTRIBUTE OOBIMPROVEMENT OOBIMPROVEMENTI HOLDS THE IMPROVEMENT  
IN TERMS OF THE LOSS ON THE OOB SAMPLES IF YOU ADD THE ITH STAGE TO THE CURRENT PREDICTIONS OUTFBAG ESTIMATES CAN  
BE USED FOR MODEL SELECTION FOR EXAMPLE TO DETERMINE THE OPTIMAL NUMBER OF ITERATIONS OOB ESTIMATES ARE USUALLY  
VERY PESSIMISTIC THUS WE RECOMMEND TO USE CROSSVALIDATION INSTEAD AND ONLY USE OOB IF CROSSVALIDATION IS TOO TIME  
CONSUMING

- EXAMPLES
- GRADIENT BOOSTING REGULARIZATION
  - GRADIENT BOOSTING OUTFBAG ESTIMATES
  - OOB ERRORS FOR RANDOM FORESTS

INTERPRETATION  
INDIVIDUAL DECISION TREES CAN BE INTERPRETED EASILY BY SIMPLY VISUALIZING THE TREE STRUCTURE GRADIENT BOOSTING MODELS  
HOWEVER COMPRISE HUNDREDS OF REGRESSION TREES THUS THEY CANNOT BE EASILY INTERPRETED BY VISUAL INSPECTION OF THE  
INDIVIDUAL TREES FORTUNATELY A NUMBER OF TECHNIQUES HAVE BEEN PROPOSED TO SUMMARIZE AND INTERPRET GRADIENT BOOSTING  
MODELS

FEATURE IMPORTANCE  
OFTEN FEATURES DO NOT CONTRIBUTE EQUALLY TO PREDICT THE TARGET RESPONSE IN MANY SITUATIONS THE MAJORITY OF THE FEATURES  
ARE IN FACT IRRELEVANT WHEN INTERPRETING A MODEL THE FIRST QUESTION USUALLY IS WHAT ARE THOSE IMPORTANT FEATURES AND  
HOW DO THEY CONTRIBUTING IN PREDICTING THE TARGET RESPONSE  
INDIVIDUAL DECISION TREES INTRINSICALLY PERFORM FEATURE SELECTION BY SELECTING APPROPRIATE SPLIT POINTS THIS INFORMATION  
CAN BE USED TO MEASURE THE IMPORTANCE OF EACH FEATURE THE BASIC IDEA IS THE MORE OFTEN A FEATURE IS USED IN THE SPLIT  
POINTS OF A TREE THE MORE IMPORTANT THAT FEATURE IS THIS NOTION OF IMPORTANCE CAN BE EXTENDED TO DECISION TREE ENSEMBLES  
BY SIMPLY AVERAGING THE FEATURE IMPORTANCE OF EACH TREE SEE FEATURE IMPORTANCE EVALUATION FOR MORE DETAILS  
THE FEATURE IMPORTANCE SCORES OF A FIT GRADIENT BOOSTING MODEL CAN BE ACCESSED VIA THE FEATUREIMPORTANCES  
PROPERTY

```
FROM SKLEARNDATASETS IMPORT MAKEHASTIE102
FROM SKLEARNENSEMBLE IMPORT GRADIENTBOOSTINGCLASSIFIER
X Y MAKEHASTIE102RANDOMSTATE0
CLF GRADIENTBOOSTINGCLASSIFIERNESTIMATORS100 LEARNINGRATE10
MAXDEPTH1 RANDOMSTATEOFITX Y
CLFFEATUREIMPORTANCES
ARRAY010 010 011
EXAMPLES
302 CHAPTER 3 USER GUIDE
```

SCIKITLEARN USER GUIDE RELEASE 0213

•GRADIENT BOOSTING REGRESSION

VOTING CLASSIFIER

THE IDEA BEHIND THE VOTINGCLASSIFIER IS TO COMBINE CONCEPTUALLY DIFFERENT MACHINE LEARNING CLASSIFIERS AND USE A MAJORITY VOTE OR THE AVERAGE PREDICTED PROBABILITIES SOFT VOTE TO PREDICT THE CLASS LABELS SUCH A CLASSIFIER CAN BE USEFUL FOR A SET OF EQUALLY WELL PERFORMING MODEL IN ORDER TO BALANCE OUT THEIR INDIVIDUAL WEAKNESSES

MAJORITY CLASS LABELS MAJORITYHARD VOTING

IN MAJORITY VOTING THE PREDICTED CLASS LABEL FOR A PARTICULAR SAMPLE IS THE CLASS LABEL THAT REPRESENTS THE MAJORITY MODE OF THE CLASS LABELS PREDICTED BY EACH INDIVIDUAL CLASSIFIER

EG IF THE PREDICTION FOR A GIVEN SAMPLE IS

- CLASSIFIER 1 CLASS 1
- CLASSIFIER 2 CLASS 1
- CLASSIFIER 3 CLASS 2

THE V OTINGCLASSIFIER WITH VOTINGHARD WOULD CLASSIFY THE SAMPLE AS “CLASS 1” BASED ON THE MAJORITY CLASS LABEL

IN THE CASES OF A TIE THE VOTINGCLASSIFIER WILL SELECT THE CLASS BASED ON THE ASCENDING SORT ORDER EG IN THE FOLLOWING SCENARIO

- CLASSIFIER 1 CLASS 2
- CLASSIFIER 2 CLASS 1

THE CLASS LABEL 1 WILL BE ASSIGNED TO THE SAMPLE

USAGE

THE FOLLOWING EXAMPLE SHOWS HOW TO FIT THE MAJORITY RULE CLASSIFIER

FROM SKLEARN IMPORT DATASETS

FROM SKLEARNMODELSELECTION IMPORT CROSSVALSCORE

FROM SKLEARNLINEARMODEL IMPORT LOGISTICREGRESSION

FROM SKLEARNNAIVEBAYES IMPORT GAUSSIANNB

FROM SKLEARNENSEMBLE IMPORT RANDOMFORESTCLASSIFIER

FROM SKLEARNENSEMBLE IMPORT VOTINGCLASSIFIER

IRIS DATASETSLOADIRIS

X Y IRISDATA 13 IRISTARGET

CLF1 LOGISTICREGRESSIONSOLVERLBFGS MULTICLASSMULTINOMIAL  
RANDOMSTATE1

CLF2 RANDOMFORESTCLASSIFIERNESTIMATORS50 RANDOMSTATE1

CLF3 GAUSSIANNB

ECLF VOTINGCLASSIFIERESTIMATORSCLR CLF1 RF CLF2 GNB CLF3

↪VOTINGHARD

FOR CLF LABEL INZIPCLF1 CLF2 CLF3 ECLF LOGISTIC REGRESSION RANDOM

↪FOREST NAIVE BAYES ENSEMBLE

31 SUPERVISED LEARNING 303

SCIKITLEARN USER GUIDE RELEASE 0213

SCORES CROSSVALSCORECLF X Y CV5 SCORINGACCURACY  
PRINTACCURACY 02F02F S SCORESMEAN SCORESSTD  
'→LABEL

ACCURACY 095 004 LOGISTIC REGRESSION  
ACCURACY 094 004 RANDOM FOREST  
ACCURACY 091 004 NAIVE BAYES  
ACCURACY 095 004 ENSEMBLE

WEIGHTED AVERAGE PROBABILITIES SOFT VOTING  
IN CONTRAST TO MAJORITY VOTING HARD VOTING SOFT VOTING RETURNS THE CLASS LABEL AS ARGMAX OF THE SUM OF PREDICTED PROBABILITIES

SPECIFIC WEIGHTS CAN BE ASSIGNED TO EACH CLASSIFIER VIA THE WEIGHTS PARAMETER WHEN WEIGHTS ARE PROVIDED THE PREDICTED CLASS PROBABILITIES FOR EACH CLASSIFIER ARE COLLECTED MULTIPLIED BY THE CLASSIFIER WEIGHT AND AVERAGED THE FINAL CLASS LABEL IS THEN DERIVED FROM THE CLASS LABEL WITH THE HIGHEST AVERAGE PROBABILITY

TO ILLUSTRATE THIS WITH A SIMPLE EXAMPLE LET'S ASSUME WE HAVE 3 CLASSIFIERS AND A 3CLASS CLASSIFICATION PROBLEMS WHERE WE ASSIGN EQUAL WEIGHTS TO ALL CLASSIFIERS W11 W21 W31

THE WEIGHTED AVERAGE PROBABILITIES FOR A SAMPLE WOULD THEN BE CALCULATED AS FOLLOWS

CLASSIFIER CLASS 1 CLASS 2 CLASS 3  
CLASSIFIER 1 W1 02 W1 05 W1 03  
CLASSIFIER 2 W2 06 W2 03 W2 01  
CLASSIFIER 3 W3 03 W3 04 W3 03  
WEIGHTED AVERAGE 037 04 023

HERE THE PREDICTED CLASS LABEL IS 2 SINCE IT HAS THE HIGHEST AVERAGE PROBABILITY

THE FOLLOWING EXAMPLE ILLUSTRATES HOW THE DECISION REGIONS MAY CHANGE WHEN A SOFT VOTINGCLASSIFIER IS USED BASED ON AN LINEAR SUPPORT VECTOR MACHINE A DECISION TREE AND A KNEAREST NEIGHBOR CLASSIFIER

FROM SKLEARN IMPORT DATASETS  
FROM SKLEARNTREE IMPORT DECISIONTREECLASSIFIER  
FROM SKLEARNNEIGHBORS IMPORT KNEIGHBORSCLASSIFIER  
FROM SKLEARN SVM IMPORT SVC  
FROM ITERTOOLS IMPORT PRODUCT  
FROM SKLEARNENSEMBLE IMPORT VOTINGCLASSIFIER

LOADING SOME EXAMPLE DATA  
IRIS DATASETSLOADIRIS  
X IRISDATA 0 2  
Y IRISTARGET

TRAINING CLASSIFIERS  
CLF1 DECISIONTREECLASSIFIERMAXDEPTH4  
CLF2 KNEIGHBORSCLASSIFIERNNEIGHBORS7  
CLF3 SVCGAMMASCALE KERNELRBF PROBABILITY TRUE  
ECLF VOTINGCLASSIFIERESTIMATORS7 CLF1 KNN CLF2 SVC CLF3  
VOTINGSOFT WEIGHTS2 1 2

CLF1 CLF1FITX Y  
CLF2 CLF2FITX Y  
CLF3 CLF3FITX Y  
ECLF ECLFFITX Y

304 CHAPTER 3 USER GUIDE



SCIKITLEARN USER GUIDE RELEASE 0213

USING THE VOTINGCLASSIFIER WITHGRIDSEARCHCV

THEVOTINGCLASSIFIER CAN ALSO BE USED TOGETHER WITH GRIDSEARCHCV IN ORDER TO TUNE THE HYPERPARAMETERS OF THE INDIVIDUAL ESTIMATORS

```
FROM SKLEARNMODELSELECTION IMPORT GRIDSEARCHCV
CLF1 LOGISTICREGRESSIONSOLVERLBFGS MULTICLASSMULTINOMIAL
RANDOMSTATE1
CLF2 RANDOMFORESTCLASSIFIERRANDOMSTATE1
CLF3 GAUSSIANNB
ECLF VOTINGCLASSIFIERESTIMATORSCLR CLF1 RF CLF2 GNB CLF3
VOTINGSOFT
PARAMS LRC 10 1000 RFNESTIMATORS 20 200
GRID GRIDSEARCHCVESTIMATORECLF PARAMGRIDPARAMS CV5
GRID GRIDFITIRISDATA IRISTARGET
USAGE
IN ORDER TO PREDICT THE CLASS LABELS BASED ON THE PREDICTED CLASSPROBABILITIES SCIKITLEARN ESTIMATORS IN THE V OTINGCLAS
SIFIER MUST SUPPORT PREDICTPROBA METHOD
ECLF VOTINGCLASSIFIERESTIMATORSCLR CLF1 RF CLF2 GNB CLF3
VOTINGSOFT
OPTIONALLY WEIGHTS CAN BE PROVIDED FOR THE INDIVIDUAL CLASSIFIERS
ECLF VOTINGCLASSIFIERESTIMATORSCLR CLF1 RF CLF2 GNB CLF3
VOTINGSOFT WEIGHTS2 5 1
VOTING REGRESSOR
THE IDEA BEHIND THE VOTINGREGRESSOR IS TO COMBINE CONCEPTUALLY DIFFERENT MACHINE LEARNING REGRESSORS AND RETURN
THE AVERAGE PREDICTED VALUES SUCH A REGRESSOR CAN BE USEFUL FOR A SET OF EQUALLY WELL PERFORMING MODELS IN ORDER TO
BALANCE OUT THEIR INDIVIDUAL WEAKNESSES
THE FOLLOWING EXAMPLE SHOWS HOW TO FIT THE V OTINGREGRESSOR
FROM SKLEARN IMPORT DATASETS
FROM SKLEARNENSEMBLE IMPORT GRADIENTBOOSTINGREGRESSOR
FROM SKLEARNENSEMBLE IMPORT RANDOMFORESTREGRESSOR
FROM SKLEARNLINEARMODEL IMPORT LINEARREGRESSION
FROM SKLEARNENSEMBLE IMPORT VOTINGREGRESSOR
LOADING SOME EXAMPLE DATA
BOSTON DATASETSLOADBOSTON
X BOSTONDATA
Y BOSTONTARGET
TRAINING CLASSIFIERS
REG1 GRADIENTBOOSTINGREGRESSORRANDOMSTATE1 NESTIMATORS10
REG2 RANDOMFORESTREGRESSORRANDOMSTATE1 NESTIMATORS10
REG3 LINEARREGRESSION
306 CHAPTER 3 USER GUIDE
```



SCIKITLEARN USER GUIDE RELEASE 0213  
EREG VOTINGREGRESSORESTIMATORSGB REG1 RF REG2 LR REG3  
EREG EREGFITX Y  
EXAMPLES

•PLOT INDIVIDUAL AND VOTING REGRESSION PREDICTIONS

3112 MULTICLASS AND MULTILABEL ALGORITHMS

WARNING ALL CLASSIFIERS IN SCIKITLEARN DO MULTICLASS CLASSIFICATION OUTOFTHEBOX YOU DON'T NEED TO USE THE  
SKLEARNMULTICLASS MODULE UNLESS YOU WANT TO EXPERIMENT WITH DIFFERENT MULTICLASS STRATEGIES

THESKLEARNMULTICLASS MODULE IMPLEMENTS METAESTIMATORS TO SOLVEMULTICLASS ANDMULTILABEL CLAS  
SIFICATION PROBLEMS BY DECOMPOSING SUCH PROBLEMS INTO BINARY CLASSIFICATION PROBLEMS MULTITARGET REGRESSION IS ALSO  
SUPPORTED

•MULTICLASS CLASSIFICATION MEANS A CLASSIFICATION TASK WITH MORE THAN TWO CLASSES EG CLASSIFY A SET OF IMAGES OF  
FRUITS WHICH MAY BE ORANGES APPLES OR PEARS MULTICLASS CLASSIFICATION MAKES THE ASSUMPTION THAT EACH SAMPLE  
IS ASSIGNED TO ONE AND ONLY ONE LABEL A FRUIT CAN BE EITHER AN APPLE OR A PEAR BUT NOT BOTH AT THE SAME TIME

•MULTILABEL CLASSIFICATION ASSIGNS TO EACH SAMPLE A SET OF TARGET LABELS THIS CAN BE THOUGHT AS PREDICTING PROPER  
TIES OF A DATAPOINT THAT ARE NOT MUTUALLY EXCLUSIVE SUCH AS TOPICS THAT ARE RELEVANT FOR A DOCUMENT A TEXT MIGHT  
BE ABOUT ANY OF RELIGION POLITICS FINANCE OR EDUCATION AT THE SAME TIME OR NONE OF THESE

•MULTIOUTPUT REGRESSION ASSIGNS EACH SAMPLE A SET OF TARGET VALUES THIS CAN BE THOUGHT OF AS PREDICTING SEVERAL  
PROPERTIES FOR EACH DATAPOINT SUCH AS WIND DIRECTION AND MAGNITUDE AT A CERTAIN LOCATION

31 SUPERVISED LEARNING 307

SCIKITLEARN USER GUIDE RELEASE 0213

•MULTIOUTPUTMULTICLASS CLASSIFICATION ANDMULTITASK CLASSIFICATION MEANS THAT A SINGLE ESTIMATOR HAS TO HANDLE SEVERAL JOINT CLASSIFICATION TASKS THIS IS BOTH A GENERALIZATION OF THE MULTILABEL CLASSIFICATION TASK WHICH ONLY CONSIDERS BINARY CLASSIFICATION AS WELL AS A GENERALIZATION OF THE MULTICLASS CLASSIFICATION TASK THE OUTPUT FORMAT IS A 2D NUMPY ARRAY OR SPARSE MATRIX  
THE SET OF LABELS CAN BE DIFFERENT FOR EACH OUTPUT VARIABLE FOR INSTANCE A SAMPLE COULD BE ASSIGNED “PEAR” FOR AN OUTPUT VARIABLE THAT TAKES POSSIBLE VALUES IN A FINITE SET OF SPECIES SUCH AS “PEAR” “APPLE” AND “BLUE” OR “GREEN” FOR A SECOND OUTPUT VARIABLE THAT TAKES POSSIBLE VALUES IN A FINITE SET OF COLORS SUCH AS “GREEN” “RED” “BLUE” “YELLOW”  
THIS MEANS THAT ANY CLASSIFIERS HANDLING MULTIOUTPUT MULTICLASS OR MULTITASK CLASSIFICATION TASKS SUPPORT THE MULTILABEL CLASSIFICATION TASK AS A SPECIAL CASE MULTITASK CLASSIFICATION IS SIMILAR TO THE MULTIOUTPUT CLASSIFICATION TASK WITH DIFFERENT MODEL FORMULATIONS FOR MORE INFORMATION SEE THE RELEVANT ESTIMATOR DOCUMENTATION  
ALL SCIKITLEARN CLASSIFIERS ARE CAPABLE OF MULTICLASS CLASSIFICATION BUT THE METAESTIMATORS OFFERED BY SKLEARN MULTICLASS PERMIT CHANGING THE WAY THEY HANDLE MORE THAN TWO CLASSES BECAUSE THIS MAY HAVE AN EFFECT ON CLASSIFIER PERFORMANCE EITHER IN TERMS OF GENERALIZATION ERROR OR REQUIRED COMPUTATIONAL RESOURCES  
BELOW IS A SUMMARY OF THE CLASSIFIERS SUPPORTED BY SCIKITLEARN GROUPED BY STRATEGY YOU DON’T NEED THE METAESTIMATORS IN THIS CLASS IF YOU’RE USING ONE OF THESE UNLESS YOU WANT CUSTOM MULTICLASS BEHAVIOR

- INHERENTLY MULTICLASS
    - SKLEARNNAIVEBAYESBERNOULLINB
    - SKLEARNTREEDECISIONTREECLASSIFIER
    - SKLEARNTREEEXTRATREECLASSIFIER
    - SKLEARNENSEMBLEEXTRATREESCLASSIFIER
    - SKLEARNNAIVEBAYESGAUSSIANNB
    - SKLEARNNEIGHBORSKNEIGHBORSCLASSIFIER
    - SKLEARNSEMISUPERVISEDLABELPROPAGATION
    - SKLEARNSEMISUPERVISEDLABELSPREADING
    - SKLEARNDISCRIMINANTANALYSISLINEARDISCRIMINANTANALYSIS
    - SKLEARNSVMLINEARSVC SETTING MULTICLASS”CRAMMERSINGER”
    - SKLEARNLINEARMODELLOGISTICREGRESSION SETTING MULTICLASS”MULTINOMIAL”
    - SKLEARNLINEARMODELLOGISTICREGRESSIONCV SETTING MULTICLASS”MULTINOMIAL”
    - SKLEARNNEURALNETWORKMLPCLASSIFIER
    - SKLEARNNEIGHBORSNEARESTCENTROID
    - SKLEARNDISCRIMINANTANALYSISQUADRATICDISCRIMINANTANALYSIS
    - SKLEARNNEIGHBORSRADIUSNEIGHBORSCLASSIFIER
    - SKLEARNENSEMBLERANDOMFORESTCLASSIFIER
    - SKLEARNLINEARMODELRIDGECLASSIFIER
    - SKLEARNLINEARMODELRIDGECLASSIFIERCV
  - MULTICLASS AS ONEVSONE
    - SKLEARNSVMNUSVC
    - SKLEARNSVMSVC
- 308 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

-SKLEARNGAUSSIANPROCESSGAUSSIANPROCESSCLASSIFIER SETTING MULTICLASS  
"ONEVSONE"

- MULTICLASS AS ONEVSALL

-SKLEARNENSEMBLEGRADIENTBOOSTINGCLASSIFIER

-SKLEARNGAUSSIANPROCESSGAUSSIANPROCESSCLASSIFIER SETTING MULTICLASS  
"ONEVSREST"

-SKLEARNVMLINEARSVC SETTING MULTICLASS"OVR"

-SKLEARNLINEARMODELLOGISTICREGRESSION SETTING MULTICLASS"OVR"

-SKLEARNLINEARMODELLOGISTICREGRESSIONCV SETTING MULTICLASS"OVR"

-SKLEARNLINEARMODELSGDCLASSIFIER

-SKLEARNLINEARMODELPERCEPTRON

-SKLEARNLINEARMODELPASSIVEAGGRESSIVECLASSIFIER

- SUPPORT MULTILABEL

-SKLEARNTREEDECISIONTREECLASSIFIER

-SKLEARNTREEEXTRATREECLASSIFIER

-SKLEARNENSEMBLEEXTRATREESCLASSIFIER

-SKLEARNNEIGHBORSKNEIGHBORSCLASSIFIER

-SKLEARNNEURALNETWORKMLPCLASSIFIER

-SKLEARNNEIGHBORSRADIUSNEIGHBORSCLASSIFIER

-SKLEARNENSEMBLERANDOMFORESTCLASSIFIER

-SKLEARNLINEARMODELRIDGECLASSIFIERCV

- SUPPORT MULTICLASSMULTIOUTPUT

-SKLEARNTREEDECISIONTREECLASSIFIER

-SKLEARNTREEEXTRATREECLASSIFIER

-SKLEARNENSEMBLEEXTRATREESCLASSIFIER

-SKLEARNNEIGHBORSKNEIGHBORSCLASSIFIER

-SKLEARNNEIGHBORSRADIUSNEIGHBORSCLASSIFIER

-SKLEARNENSEMBLERANDOMFORESTCLASSIFIER

WARNING AT PRESENT NO METRIC IN SKLEARNMETRICS SUPPORTS THE MULTIOUTPUTMULTICLASS CLASSIFICATION TASK

MULTILABEL CLASSIFICATION FORMAT

IN MULTILABEL LEARNING THE JOINT SET OF BINARY CLASSIFICATION TASKS IS EXPRESSED WITH LABEL BINARY INDICATOR ARRAY EACH  
SAMPLE IS ONE ROW OF A 2D ARRAY OF SHAPE NSAMPLES NCLASSES WITH BINARY VALUES THE ONE IE THE NON ZERO ELEMENTS  
CORRESPONDS TO THE SUBSET OF LABELS AN ARRAY SUCH AS NPARRAY1 0 0 0 1 1 0 0 0  
REPRESENTS LABEL 0 IN THE FIRST SAMPLE LABELS 1 AND 2 IN THE SECOND SAMPLE AND NO LABELS IN THE THIRD SAMPLE  
PRODUCING MULTILABEL DATA AS A LIST OF SETS OF LABELS MAY BE MORE INTUITIVE THE MULTILABELBINARIZER TRANSFORMER  
CAN BE USED TO CONVERT BETWEEN A COLLECTION OF COLLECTIONS OF LABELS AND THE INDICATOR FORMAT

31 SUPERVISED LEARNING 309







SCIKITLEARN USER GUIDE RELEASE 0213

- “SOLVING MULTICLASS LEARNING PROBLEMS VIA ERRORCORRECTING OUTPUT CODES” DIETTERICH T BAKIRI G JOURNAL OF ARTIFICIAL INTELLIGENCE RESEARCH 2 1995
- “THE ELEMENTS OF STATISTICAL LEARNING” HASTIE T TIBSHIRANI R FRIEDMAN J PAGE 606 SECONDEDITION 2008

MULTIOUTPUT REGRESSION

MULTIOUTPUT REGRESSION SUPPORT CAN BE ADDED TO ANY REGRESSOR WITH MULTIOUTPUTREGRESSOR THIS STRATEGY CONSISTS OF FITTING ONE REGRESSOR PER TARGET SINCE EACH TARGET IS REPRESENTED BY EXACTLY ONE REGRESSOR IT IS POSSIBLE TO GAIN KNOWLEDGE ABOUT THE TARGET BY INSPECTING ITS CORRESPONDING REGRESSOR AS MULTIOUTPUTREGRESSOR FITS ONE REGRESSOR PER TARGET IT CAN NOT TAKE ADVANTAGE OF CORRELATIONS BETWEEN TARGETS

BELOW IS AN EXAMPLE OF MULTIOUTPUT REGRESSION

```
FROM SKLEARNDATASETS IMPORT MAKEREGRESSION
FROM SKLEARNMULTIOUTPUT IMPORT MULTIOUTPUTREGRESSOR
FROM SKLEARNENSEMBLE IMPORT GRADIENTBOOSTINGREGRESSOR
X Y MAKEREGRESSIONNSAMPLES10 NTARGETS3 RANDOMSTATE1
MULTIOUTPUTREGRESSORGRADIENTBOOSTINGREGRESSORRANDOMSTATE0FITX Y
↪PREDICTX
ARRAY15475474165 14703498585 5003812219
712165031 512914884 8146081961
1878948621 10044373091 1388978285
14162745778 9502891072 19148204257
9703260883 16534867495 13952003279
12392529176 2125719016 784253
12225193977 8516443186 10712274212
30170388 9480956739 1216979946
14072667194 17650941682 1750447799
14937967282 8115699552 572850319
```

MULTIOUTPUT CLASSIFICATION

MULTIOUTPUT CLASSIFICATION SUPPORT CAN BE ADDED TO ANY CLASSIFIER WITH MULTIOUTPUTCLASSIFIER THIS STRATEGY CONSISTS OF FITTING ONE CLASSIFIER PER TARGET THIS ALLOWS MULTIPLE TARGET VARIABLE CLASSIFICATIONS THE PURPOSE OF THIS CLASS IS TO EXTEND ESTIMATORS TO BE ABLE TO ESTIMATE A SERIES OF TARGET FUNCTIONS F1F2F3 FN THAT ARE TRAINED ON A SINGLE X PREDICTOR MATRIX TO PREDICT A SERIES OF RESPONSES Y1Y2Y3 YN

BELOW IS AN EXAMPLE OF MULTIOUTPUT CLASSIFICATION

```
FROM SKLEARNDATASETS IMPORT MAKECLASSIFICATION
FROM SKLEARNMULTIOUTPUT IMPORT MULTIOUTPUTCLASSIFIER
FROM SKLEARNENSEMBLE IMPORT RANDOMFORESTCLASSIFIER
FROM SKLEARNUTILS IMPORT SHUFFLE
IMPORT NUMPY AS NP
X Y1 MAKECLASSIFICATIONNSAMPLES10 NFEATURES100 NINFORMATIVE30 N
↪CLASSES3 RANDOMSTATE1
Y2 SHUFFLEY1 RANDOMSTATE1
Y3 SHUFFLEY1 RANDOMSTATE2
Y NPVSTACKY1 Y2 Y3T
NSAMPLES NFEATURES XSHAPE 10100
NOUTPUTS YSHAPE1 3
NCLASSES 3
FOREST RANDOMFORESTCLASSIFIERNESTIMATORS100 RANDOMSTATE1
MULTITARGETFOREST MULTIOUTPUTCLASSIFIERFOREST NJOBS1
31 SUPERVISED LEARNING 313
```

SCIKITLEARN USER GUIDE RELEASE 0213  
MULTITARGETFORESTFITX YPREDICTX  
ARRAY2 2 0

1 2 1  
2 1 0  
0 0 2  
0 2 1  
0 0 2  
1 1 0  
1 1 1  
0 0 2  
2 0 0

CLASSIFIER CHAIN

CLASSIFIER CHAINS SEE CLASSIFIERCHAIN ARE A WAY OF COMBINING A NUMBER OF BINARY CLASSIFIERS INTO A SINGLE MULTILABEL MODEL THAT IS CAPABLE OF EXPLOITING CORRELATIONS AMONG TARGETS FOR A MULTILABEL CLASSIFICATION PROBLEM WITH N CLASSES N BINARY CLASSIFIERS ARE ASSIGNED AN INTEGER BETWEEN 0 AND N1 THESE INTEGERS DEFINE THE ORDER OF MODELS IN THE CHAIN EACH CLASSIFIER IS THEN FIT ON THE AVAILABLE TRAINING DATA PLUS THE TRUE LABELS OF THE CLASSES WHOSE MODELS WERE ASSIGNED A LOWER NUMBER WHEN PREDICTING THE TRUE LABELS WILL NOT BE AVAILABLE INSTEAD THE PREDICTIONS OF EACH MODEL ARE PASSED ON TO THE SUBSEQUENT MODELS IN THE CHAIN TO BE USED AS FEATURES CLEARLY THE ORDER OF THE CHAIN IS IMPORTANT THE FIRST MODEL IN THE CHAIN HAS NO INFORMATION ABOUT THE OTHER LABELS WHILE THE LAST MODEL IN THE CHAIN HAS FEATURES INDICATING THE PRESENCE OF ALL OF THE OTHER LABELS IN GENERAL ONE DOES NOT KNOW THE OPTIMAL ORDERING OF THE MODELS IN THE CHAIN SO TYPICALLY MANY RANDOMLY ORDERED CHAINS ARE FIT AND THEIR PREDICTIONS ARE AVERAGED TOGETHER

REFERENCES

JESSE READ BERNHARD PFAHRINGER GEOFF HOLMES EIBE FRANK “CLASSIFIER CHAINS FOR MULTILABEL CLASSIFICATION” 2009

REGRESSOR CHAIN

REGRESSOR CHAINS SEE REGRESSORCHAIN IS ANALOGOUS TO CLASSIFIERCHAIN AS A WAY OF COMBINING A NUMBER OF REGRESSIONS INTO A SINGLE MULTITARGET MODEL THAT IS CAPABLE OF EXPLOITING CORRELATIONS AMONG TARGETS

3113 FEATURE SELECTION

THE CLASSES IN THE SKLEARNFEATURESELECTION MODULE CAN BE USED FOR FEATURE SELECTION DIMENSIONALITY REDUCTION ON SAMPLE SETS EITHER TO IMPROVE ESTIMATORS’ ACCURACY SCORES OR TO BOOST THEIR PERFORMANCE ON VERY HIGH DIMENSIONAL DATASETS

REMOVING FEATURES WITH LOW VARIANCE

VARIANCETHRESHOLD IS A SIMPLE BASELINE APPROACH TO FEATURE SELECTION IT REMOVES ALL FEATURES WHOSE VARIANCE DOESN’T MEET SOME THRESHOLD BY DEFAULT IT REMOVES ALL ZEROVARIANCE FEATURES IE FEATURES THAT HAVE THE SAME VALUE IN ALL SAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213

AS AN EXAMPLE SUPPOSE THAT WE HAVE A DATASET WITH BOOLEAN FEATURES AND WE WANT TO REMOVE ALL FEATURES THAT ARE EITHER ONE OR ZERO ON OR OFF IN MORE THAN 80 OF THE SAMPLES BOOLEAN FEATURES ARE BERNOULLI RANDOM VARIABLES AND THE VARIANCE OF SUCH VARIABLES IS GIVEN BY

```
VAR = 1 -
SO WE CAN SELECT USING THE THRESHOLD 81 8
FROM SKLEARNFEATURESELECTION IMPORT VARIANCETHRESHOLD
X 0 0 1 0 1 0 1 0 0 0 1 1 0 1 0 0 1 1
SEL VARIANCETHRESHOLDTHRESHOLD8 1 8
SELFITTRANSFORMX
ARRAY0 1
1 0
0 0
1 1
1 0
1 1
```

AS EXPECTED VARIANCETHRESHOLD HAS REMOVED THE FIRST COLUMN WHICH HAS A PROBABILITY 56 80F CONTAINING A ZERO

UNIVARIATE FEATURE SELECTION

UNIVARIATE FEATURE SELECTION WORKS BY SELECTING THE BEST FEATURES BASED ON UNIVARIATE STATISTICAL TESTS IT CAN BE SEEN AS A PREPROCESSING STEP TO AN ESTIMATOR SCIKITLEARN EXPOSES FEATURE SELECTION ROUTINES AS OBJECTS THAT IMPLEMENT THE TRANSFORM METHOD

- SELECTKBEST REMOVES ALL BUT THE HIGHEST SCORING FEATURES
- SELECTPERCENTILE REMOVES ALL BUT A USERSPECIFIED HIGHEST SCORING PERCENTAGE OF FEATURES
- USING COMMON UNIVARIATE STATISTICAL TESTS FOR EACH FEATURE FALSE POSITIVE RATE SELECTFPR FALSE DISCOVERY RATE SELECTFDR OR FAMILY WISE ERROR SELECTFWE
- GENERICUNIVARIATESELECT ALLOWS TO PERFORM UNIVARIATE FEATURE SELECTION WITH A CONFIGURABLE STRATEGY THIS ALLOWS TO SELECT THE BEST UNIVARIATE SELECTION STRATEGY WITH HYPERPARAMETER SEARCH ESTIMATOR FOR INSTANCE WE CAN PERFORM A K2TEST TO THE SAMPLES TO RETRIEVE ONLY THE TWO BEST FEATURES AS FOLLOWS

```
FROM SKLEARNDATASETS IMPORT LOADIRIS
FROM SKLEARNFEATURESELECTION IMPORT SELECTKBEST
FROM SKLEARNFEATURESELECTION IMPORT CHI2
IRIS LOADIRIS
X Y IRISDATA IRISTARGET
XSHAPE
150 4
XNEW SELECTKBESTCHI2 K2FITTRANSFORMX Y
XNEWSHAPE
150 2
```

THESE OBJECTS TAKE AS INPUT A SCORING FUNCTION THAT RETURNS UNIVARIATE SCORES AND PVALUES OR ONLY SCORES FOR SELECTKBEST ANDSELECTPERCENTILE

- FOR REGRESSION FREGRESSION MUTUALINFOREGRESSION
- FOR CLASSIFICATION CHI2 FCLASSIF MUTUALINFOCLASSIF

SCIKITLEARN USER GUIDE RELEASE 0213

THE METHODS BASED ON FTEST ESTIMATE THE DEGREE OF LINEAR DEPENDENCY BETWEEN TWO RANDOM VARIABLES ON THE OTHER HAND MUTUAL INFORMATION METHODS CAN CAPTURE ANY KIND OF STATISTICAL DEPENDENCY BUT BEING NONPARAMETRIC THEY REQUIRE MORE SAMPLES FOR ACCURATE ESTIMATION

FEATURE SELECTION WITH SPARSE DATA

IF YOU USE SPARSE DATA IE DATA REPRESENTED AS SPARSE MATRICES CHI2 MUTUALINFOREGRESSION

MUTUALINFOCLASSIF WILL DEAL WITH THE DATA WITHOUT MAKING IT DENSE

WARNING BEWARE NOT TO USE A REGRESSION SCORING FUNCTION WITH A CLASSIFICATION PROBLEM YOU WILL GET USELESS RESULTS

EXAMPLES

- UNIVARIATE FEATURE SELECTION
  - COMPARISON OF FTEST AND MUTUAL INFORMATION
- RECURSIVE FEATURE ELIMINATION

GIVEN AN EXTERNAL ESTIMATOR THAT ASSIGNS WEIGHTS TO FEATURES EG THE COEFFICIENTS OF A LINEAR MODEL RECURSIVE FEATURE ELIMINATION RFE IS TO SELECT FEATURES BY RECURSIVELY CONSIDERING SMALLER AND SMALLER SETS OF FEATURES FIRST THE ESTIMATOR IS TRAINED ON THE INITIAL SET OF FEATURES AND THE IMPORTANCE OF EACH FEATURE IS OBTAINED EITHER THROUGH A COEF ATTRIBUTE OR THROUGH A FEATUREIMPORTANCES ATTRIBUTE THEN THE LEAST IMPORTANT FEATURES ARE PRUNED FROM CURRENT SET OF FEATURESTHAT PROCEDURE IS RECURSIVELY REPEATED ON THE PRUNED SET UNTIL THE DESIRED NUMBER OF FEATURES TO SELECT IS EVENTUALLY REACHED

RFECV PERFORMS RFE IN A CROSSVALIDATION LOOP TO FIND THE OPTIMAL NUMBER OF FEATURES

EXAMPLES

- RECURSIVE FEATURE ELIMINATION A RECURSIVE FEATURE ELIMINATION EXAMPLE SHOWING THE RELEVANCE OF PIXELS IN A DIGIT CLASSIFICATION TASK
- RECURSIVE FEATURE ELIMINATION WITH CROSSVALIDATION A RECURSIVE FEATURE ELIMINATION EXAMPLE WITH AUTOMATIC TUNING OF THE NUMBER OF FEATURES SELECTED WITH CROSSVALIDATION

FEATURE SELECTION USING SELECTFROMMODEL

SELECTFROMMODEL IS A METATransformer THAT CAN BE USED ALONG WITH ANY ESTIMATOR THAT HAS A COEF OR FEATUREIMPORTANCES ATTRIBUTE AFTER FITTING THE FEATURES ARE CONSIDERED UNIMPORTANT AND REMOVED IF THE CORRESPONDING COEF ORFEATUREIMPORTANCES VALUES ARE BELOW THE PROVIDED THRESHOLD PARAMETER APART FROM SPECIFYING THE THRESHOLD NUMERICALLY THERE ARE BUILTIN HEURISTICS FOR FINDING A THRESHOLD USING A STRING ARGUMENT AVAILABLE HEURISTICS ARE “MEAN” “MEDIAN” AND FLOAT MULTIPLES OF THESE LIKE “01MEAN” FOR EXAMPLES ON HOW IT IS TO BE USED REFER TO THE SECTIONS BELOW

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES

•FEATURE SELECTION USING SELECTFROMMODEL AND LASSOCV SELECTING THE TWO MOST IMPORTANT FEATURES FROM THE BOSTON DATASET WITHOUT KNOWING THE THRESHOLD BEFOREHAND

L1BASED FEATURE SELECTION

LINEAR MODELS PENALIZED WITH THE L1 NORM HAVE SPARSE SOLUTIONS MANY OF THEIR ESTIMATED COEFFICIENTS ARE ZERO WHEN THE GOAL IS TO REDUCE THE DIMENSIONALITY OF THE DATA TO USE WITH ANOTHER CLASSIFIER THEY CAN BE USED ALONG WITHFEATURESELECTIONSELECTFROMMODEL TO SELECT THE NONZERO COEFFICIENTS IN PARTICULAR SPARSE ESTIMATORS USEFUL FOR THIS PURPOSE ARE THE LINEARMODELLASSO FOR REGRESSION AND OF LINEARMODEL LOGISTICREGRESSION ANDSVMLINEARSVC FOR CLASSIFICATION

```
FROM SKLEARN SVM IMPORT LINEARSVC
FROM SKLEARN DATASETS IMPORT LOADIRIS
FROM SKLEARN FEATURESELECTION IMPORT SELECTFROMMODEL
```

```
IRIS LOADIRIS
X Y IRIS DATA IRIS TARGET
X SHAPE
```

```
150 4
LSVC LINEARSVC C001 PENALTY L1 DUAL FALSE FIT X Y
MODEL SELECTFROMMODEL SVC PREPIT TRUE
X NEW MODEL TRANSFORM X
X NEW SHAPE
```

150 3  
WITH SVMs AND LOGISTIC REGRESSION THE PARAMETER C CONTROLS THE SPARSITY THE SMALLER C THE FEWER FEATURES SELECTED WITH LASSO THE HIGHER THE ALPHA PARAMETER THE FEWER FEATURES SELECTED

EXAMPLES

•CLASSIFICATION OF TEXT DOCUMENTS USING SPARSE FEATURES COMPARISON OF DIFFERENT ALGORITHMS FOR DOCUMENT CLASSIFICATION INCLUDING L1BASED FEATURE SELECTION

L1RECOVERY AND COMPRESSIVE SENSING

FOR A GOOD CHOICE OF ALPHA THE LASSO CAN FULLY RECOVER THE EXACT SET OF NONZERO VARIABLES USING ONLY FEW OBSERVATIONS PROVIDED CERTAIN SPECIFIC CONDITIONS ARE MET IN PARTICULAR THE NUMBER OF SAMPLES SHOULD BE “SUFFICIENTLY LARGE” OR L1 MODELS WILL PERFORM AT RANDOM WHERE “SUFFICIENTLY LARGE” DEPENDS ON THE NUMBER OF NONZERO COEFFICIENTS THE LOGARITHM OF THE NUMBER OF FEATURES THE AMOUNT OF NOISE THE SMALLEST ABSOLUTE VALUE OF NONZERO COEFFICIENTS AND THE STRUCTURE OF THE DESIGN MATRIX X IN ADDITION THE DESIGN MATRIX MUST DISPLAY CERTAIN SPECIFIC PROPERTIES SUCH AS NOT BEING TOO CORRELATED

THERE IS NO GENERAL RULE TO SELECT AN ALPHA PARAMETER FOR RECOVERY OF NONZERO COEFFICIENTS IT CAN BE SET BY CROSS VALIDATION LASSOCV OR LASSO LARS CV THOUGH THIS MAY LEAD TO UNDERPENALIZED MODELS INCLUDING A SMALL NUMBER OF NONRELEVANT VARIABLES IS NOT DETRIMENTAL TO PREDICTION SCORE BIC LASSO LARS IC TENDS ON THE OPPOSITE TO SET HIGH VALUES OF ALPHA

REFERENCE RICHARD G BARANIUK “COMPRESSIVE SENSING” IEEE SIGNAL PROCESSING MAGAZINE 120 JULY 2007 HTTP  
USERS/ISRI/STUTLPTAGUI/ARCS/NOTES/PDF  
31 SUPERVISED LEARNING 317

SCIKITLEARN USER GUIDE RELEASE 0213

TREEBASED FEATURE SELECTION

TREEBASED ESTIMATORS SEE THE SKLEARNNTREE MODULE AND FOREST OF TREES IN THE SKLEARNENSEMBLE MODULE

CAN BE USED TO COMPUTE FEATURE IMPORTANCES WHICH IN TURN CAN BE USED TO DISCARD IRRELEVANT FEATURES WHEN COUPLED

WITH THESKLEARNFEATURESELECTIONSELECTFROMMODEL METATransformer

```
FROM SKLEARNENSEMBLE IMPORT EXTRATREESCLASSIFIER
FROM SKLEARNDATASETS IMPORT LOADIRIS
FROM SKLEARNFEATURESELECTION IMPORT SELECTFROMMODEL
IRIS = LOADIRIS
X, Y = IRIS.data, IRIS.target
X.shape
150, 4
clf = EXTRATREESCLASSIFIER(n_estimators=50)
clf.fit(X, Y)
clf.feature_importances_
array([0.04, 0.05, 0.04, 0.04])
model = SELECTFROMMODEL(clf, prefit=True)
X_new = model.transform(X)
X_new.shape
150, 2
```

EXAMPLES

- FEATURE IMPORTANCES WITH FORESTS OF TREES EXAMPLE ON SYNTHETIC DATA SHOWING THE RECOVERY OF THE ACTUALLY MEANINGFUL FEATURES
- PIXEL IMPORTANCES WITH A PARALLEL FOREST OF TREES EXAMPLE ON FACE RECOGNITION DATA

FEATURE SELECTION AS PART OF A PIPELINE

FEATURE SELECTION IS USUALLY USED AS A PREPROCESSING STEP BEFORE DOING THE ACTUAL LEARNING THE RECOMMENDED WAY TO DO THIS IN SCIKITLEARN IS TO USE A SKLEARNPIPELINEPIPELINE

```
clf = PIPELINE(
    FEATURESELECTION=SELECTFROMMODEL(LINEARSVCPENALTY=L1),
    CLASSIFICATION=RANDOMFORESTCLASSIFIER)

clf.fit(X, Y)
```

IN THIS SNIPPET WE MAKE USE OF A SKLEARNsvmlLINEARSVC COUPLED WITH SKLEARNFEATURESELECTION

SELECTFROMMODEL TO EVALUATE FEATURE IMPORTANCES AND SELECT THE MOST RELEVANT FEATURES THEN A SKLEARN

ENSEMBLERANDOMFORESTCLASSIFIER IS TRAINED ON THE TRANSFORMED OUTPUT IE USING ONLY RELEVANT FEATURES

YOU CAN PERFORM SIMILAR OPERATIONS WITH THE OTHER FEATURE SELECTION METHODS AND ALSO CLASSIFIERS THAT PROVIDE A WAY TO EVALUATE FEATURE IMPORTANCES OF COURSE SEE THE SKLEARNPIPELINEPIPELINE EXAMPLES FOR MORE DETAILS

3114 SEMISUPERVISED

SEMISUPERVISED LEARNING IS A SITUATION IN WHICH IN YOUR TRAINING DATA SOME OF THE SAMPLES ARE NOT LABELED THE SEMI

SUPERVISED ESTIMATORS IN SKLEARNSEMISUPERVISED ARE ABLE TO MAKE USE OF THIS ADDITIONAL UNLABELED DATA TO

BETTER CAPTURE THE SHAPE OF THE UNDERLYING DATA DISTRIBUTION AND GENERALIZE BETTER TO NEW SAMPLES THESE ALGORITHMS

CAN PERFORM WELL WHEN WE HAVE A VERY SMALL AMOUNT OF LABELED POINTS AND A LARGE AMOUNT OF UNLABELED POINTS

318 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

UNLABELED ENTRIES IN Y

IT IS IMPORTANT TO ASSIGN AN IDENTIFIER TO UNLABELED POINTS ALONG WITH THE LABELED DATA WHEN TRAINING THE MODEL WITH THEFIT METHOD THE IDENTIFIER THAT THIS IMPLEMENTATION USES IS THE INTEGER VALUE -1

LABEL PROPAGATION

LABEL PROPAGATION DENOTES A FEW VARIATIONS OF SEMISUPERVISED GRAPH INFERENCE ALGORITHMS

A FEW FEATURES AVAILABLE IN THIS MODEL

- CAN BE USED FOR CLASSIFICATION AND REGRESSION TASKS
- KERNEL METHODS TO PROJECT DATA INTO ALTERNATE DIMENSIONAL SPACES

SCIKITLEARN PROVIDES TWO LABEL PROPAGATION MODELS LABELPROPAGATION ANDLABELSPREADING BOTH WORK BY CONSTRUCTING A SIMILARITY GRAPH OVER ALL ITEMS IN THE INPUT DATASET

FIG 31 AN ILLUSTRATION OF LABELPROPAGATION THE STRUCTURE OF UNLABELED OBSERVATIONS IS CONSISTENT WITH THE CLASS STRUCTURE AND THUS THE CLASS LABEL CAN BE PROPAGATED TO THE UNLABELED OBSERVATIONS OF THE TRAINING SET LABELPROPAGATION ANDLABELSPREADING DIFFER IN MODIFICATIONS TO THE SIMILARITY MATRIX THAT GRAPH AND THE CLAMPING EFFECT ON THE LABEL DISTRIBUTIONS CLAMPING ALLOWS THE ALGORITHM TO CHANGE THE WEIGHT OF THE TRUE GROUND LABELED DATA TO SOME DEGREE THE LABELPROPAGATION ALGORITHM PERFORMS HARD CLAMPING OF INPUT LABELS WHICH MEANS 0 THIS CLAMPING FACTOR CAN BE RELAXED TO SAY 0.2 WHICH MEANS THAT WE WILL ALWAYS RETAIN 80 PERCENT OF OUR ORIGINAL LABEL DISTRIBUTION BUT THE ALGORITHM GETS TO CHANGE ITS CONFIDENCE OF THE DISTRIBUTION WITHIN 20 PERCENT LABELPROPAGATION USES THE RAW SIMILARITY MATRIX CONSTRUCTED FROM THE DATA WITH NO MODIFICATIONS IN CONTRAST LABELSPREADING MINIMIZES A LOSS FUNCTION THAT HAS REGULARIZATION PROPERTIES AS SUCH IT IS OFTEN MORE ROBUST TO NOISE THE ALGORITHM ITERATES ON A MODIFIED VERSION OF THE ORIGINAL GRAPH AND NORMALIZES THE EDGE WEIGHTS BY COMPUTING THE NORMALIZED GRAPH LAPLACIAN MATRIX THIS PROCEDURE IS ALSO USED IN SPECTRAL CLUSTERING LABEL PROPAGATION MODELS HAVE TWO BUILTIN KERNEL METHODS CHOICE OF KERNEL EFFECTS BOTH SCALABILITY AND PERFORMANCE OF THE ALGORITHMS THE FOLLOWING ARE AVAILABLE

- RBF EXP- $\gamma$ - $\gamma^2$   $\gamma$  IS SPECIFIED BY KEYWORD GAMMA
- KNN  $1/\sqrt{k}$   $k$  IS SPECIFIED BY KEYWORD NNEIGHBORS

THE RBF KERNEL WILL PRODUCE A FULLY CONNECTED GRAPH WHICH IS REPRESENTED IN MEMORY BY A DENSE MATRIX THIS MATRIX MAY BE VERY LARGE AND COMBINED WITH THE COST OF PERFORMING A FULL MATRIX MULTIPLICATION CALCULATION FOR EACH ITERATION OF THE ALGORITHM CAN LEAD TO PROHIBITIVELY LONG RUNNING TIMES ON THE OTHER HAND THE KNN KERNEL WILL PRODUCE A MUCH MORE MEMORYFRIENDLY SPARSE MATRIX WHICH CAN DRASTICALLY REDUCE RUNNING TIMES

EXAMPLES

- DECISION BOUNDARY OF LABEL PROPAGATION VERSUS SVM ON THE IRIS DATASET
- LABEL PROPAGATION LEARNING A COMPLEX STRUCTURE
- LABEL PROPAGATION DIGITS DEMONSTRATING PERFORMANCE
- LABEL PROPAGATION DIGITS ACTIVE LEARNING

REFERENCES

1 YOSHUA BENGIO OLIVIER DELALLEAU NICOLAS LE ROUX IN SEMISUPERVISED LEARNING 2006 PP 193216  
2 OLIVIER DELALLEAU YOSHUA BENGIO NICOLAS LE ROUX EFFICIENT NONPARAMETRIC FUNCTION INDUCTION IN SEMI SUPERVISED LEARNING AISTAT 2005 [HTTPSRESEARCHMICROSOFTCOMENUSPEOPLENICOLASLEFFICIENTSSLPDF](https://research.microsoft.com/en-us/people/nicolas/efficientssl.pdf)

3115 ISOTONIC REGRESSION

THE CLASSISOTONICREGRESSION FITS A NONDECREASING FUNCTION TO DATA IT SOLVES THE FOLLOWING PROBLEM

MINIMIZE  $\sum_{i=1}^n (y_i - f(x_i))^2$

SUBJECT TO  $f(x_1) \leq f(x_2) \leq \dots \leq f(x_n)$

WHERE EACH  $y_i$  IS STRICTLY POSITIVE AND EACH  $x_i$  IS AN ARBITRARY REAL NUMBER IT YIELDS THE VECTOR WHICH IS COMPOSED OF NONDECREASING ELEMENTS THE CLOSEST IN TERMS OF MEAN SQUARED ERROR IN PRACTICE THIS LIST OF ELEMENTS FORMS A FUNCTION THAT IS PIECEWISE LINEAR

3116 PROBABILITY CALIBRATION

WHEN PERFORMING CLASSIFICATION YOU OFTEN WANT NOT ONLY TO PREDICT THE CLASS LABEL BUT ALSO OBTAIN A PROBABILITY OF THE RESPECTIVE LABEL THIS PROBABILITY GIVES YOU SOME KIND OF CONFIDENCE ON THE PREDICTION SOME MODELS CAN GIVE YOU POOR ESTIMATES OF THE CLASS PROBABILITIES AND SOME EVEN DO NOT SUPPORT PROBABILITY PREDICTION THE CALIBRATION MODULE ALLOWS YOU TO BETTER CALIBRATE THE PROBABILITIES OF A GIVEN MODEL OR TO ADD SUPPORT FOR PROBABILITY PREDICTION WELL CALIBRATED CLASSIFIERS ARE PROBABILISTIC CLASSIFIERS FOR WHICH THE OUTPUT OF THE PREDICTPROBA METHOD CAN BE DIRECTLY INTERPRETED AS A CONFIDENCE LEVEL FOR INSTANCE A WELL CALIBRATED BINARY CLASSIFIER SHOULD CLASSIFY THE SAMPLES SUCH THAT AMONG THE SAMPLES TO WHICH IT GAVE A PREDICTPROBA VALUE CLOSE TO 0.8 APPROXIMATELY 80% ACTUALLY BELONG TO THE POSITIVE CLASS THE FOLLOWING PLOT COMPARES HOW WELL THE PROBABILISTIC PREDICTIONS OF DIFFERENT CLASSIFIERS ARE CALIBRATED LOGISTICREGRESSION RETURNS WELL CALIBRATED PREDICTIONS BY DEFAULT AS IT DIRECTLY OPTIMIZES LOGLOSS IN CONTRAST THE OTHER METHODS RETURN BIASED PROBABILITIES WITH DIFFERENT BIASES PER METHOD

- GAUSSIANNB TENDS TO PUSH PROBABILITIES TO 0 OR 1 NOTE THE COUNTS IN THE HISTOGRAMS THIS IS MAINLY BECAUSE IT MAKES THE ASSUMPTION THAT FEATURES ARE CONDITIONALLY INDEPENDENT GIVEN THE CLASS WHICH IS NOT THE CASE IN THIS DATASET WHICH CONTAINS 2 REDUNDANT FEATURES







SCIKITLEARN USER GUIDE RELEASE 0213

•RANDOMFORESTCLASSIFIER SHOWS THE OPPOSITE BEHAVIOR THE HISTOGRAMS SHOW PEAKS AT APPROXIMATELY 0.2 AND 0.9 PROBABILITY WHILE PROBABILITIES CLOSE TO 0 OR 1 ARE VERY RARE AN EXPLANATION FOR THIS IS GIVEN BY NICULESCUMIZIL AND CARUANA4 “METHODS SUCH AS BAGGING AND RANDOM FORESTS THAT AVERAGE PREDICTIONS FROM A BASE SET OF MODELS CAN HAVE DIFFICULTY MAKING PREDICTIONS NEAR 0 AND 1 BECAUSE VARIANCE IN THE UNDERLYING BASE MODELS WILL BIAS PREDICTIONS THAT SHOULD BE NEAR ZERO OR ONE AWAY FROM THESE VALUES BECAUSE PREDICTIONS ARE RESTRICTED TO THE INTERVAL [0,1] ERRORS CAUSED BY VARIANCE TEND TO BE ONESIDED NEAR ZERO AND ONE FOR EXAMPLE IF A MODEL SHOULD PREDICT 0 FOR A CASE THE ONLY WAY BAGGING CAN ACHIEVE THIS IS IF ALL BAGGED TREES PREDICT ZERO IF WE ADD NOISE TO THE TREES THAT BAGGING IS AVERAGING OVER THIS NOISE WILL CAUSE SOME TREES TO PREDICT VALUES LARGER THAN 0 FOR THIS CASE THUS MOVING THE AVERAGE PREDICTION OF THE BAGGED ENSEMBLE AWAY FROM 0 WE OBSERVE THIS EFFECT MOST STRONGLY WITH RANDOM FORESTS BECAUSE THE BASELEVEL TREES TRAINED WITH RANDOM FORESTS HAVE RELATIVELY HIGH VARIANCE DUE TO FEATURE SUBSETTING” AS A RESULT THE CALIBRATION CURVE ALSO REFERRED TO AS THE RELIABILITY DIAGRAM WILKS 19955 SHOWS A CHARACTERISTIC SIGMOID SHAPE INDICATING THAT THE CLASSIFIER COULD TRUST ITS “INTUITION” MORE AND RETURN PROBABILITIES CLOSER TO 0 OR 1 TYPICALLY

• LINEAR SUPPORT VECTOR CLASSIFICATION LINEARSVC SHOWS AN EVEN MORE SIGMOID CURVE AS THE RANDOMFOREST CLASSIFIER WHICH IS TYPICAL FOR MAXIMUMMARGIN METHODS COMPARE NICULESCUMIZIL AND CARUANA4 WHICH FOCUS ON HARD SAMPLES THAT ARE CLOSE TO THE DECISION BOUNDARY THE SUPPORT VECTORS

TWO APPROACHES FOR PERFORMING CALIBRATION OF PROBABILISTIC PREDICTIONS ARE PROVIDED A PARAMETRIC APPROACH BASED ON PLATT’S SIGMOID MODEL AND A NONPARAMETRIC APPROACH BASED ON ISOTONIC REGRESSION SKLEARNISOTONIC PROBABILITY CALIBRATION SHOULD BE DONE ON NEW DATA NOT USED FOR MODEL FITTING THE CLASS CALIBRATEDCLASSIFIERCV USES A CROSSVALIDATION GENERATOR AND ESTIMATES FOR EACH SPLIT THE MODEL PARAMETER ON THE TRAIN SAMPLES AND THE CALIBRATION OF THE TEST SAMPLES THE PROBABILITIES PREDICTED FOR THE FOLDS ARE THEN AVERAGED ALREADY FITTED CLASSIFIERS CAN BE CALIBRATED BY CALIBRATEDCLASSIFIERCV VIA THE PARAMETER CV”PREFIT” IN THIS CASE THE USER HAS TO TAKE CARE MANUALLY THAT DATA FOR MODEL FITTING AND CALIBRATION ARE DISJOINT

THE FOLLOWING IMAGES DEMONSTRATE THE BENEFIT OF PROBABILITY CALIBRATION THE FIRST IMAGE PRESENT A DATASET WITH 2 CLASSES AND 3 BLOBS OF DATA THE BLOB IN THE MIDDLE CONTAINS RANDOM SAMPLES OF EACH CLASS THE PROBABILITY FOR THE SAMPLES IN THIS BLOB SHOULD BE 0.5

THE FOLLOWING IMAGE SHOWS ON THE DATA ABOVE THE ESTIMATED PROBABILITY USING A GAUSSIAN NAIVE BAYES CLASSIFIER WITHOUT CALIBRATION WITH A SIGMOID CALIBRATION AND WITH A NONPARAMETRIC ISOTONIC CALIBRATION ONE CAN OBSERVE THAT THE NON PARAMETRIC MODEL PROVIDES THE MOST ACCURATE PROBABILITY ESTIMATES FOR SAMPLES IN THE MIDDLE IE 0.5

THE FOLLOWING EXPERIMENT IS PERFORMED ON AN ARTIFICIAL DATASET FOR BINARY CLASSIFICATION WITH 100000 SAMPLES 1000 OF THEM ARE USED FOR MODEL FITTING WITH 20 FEATURES OF THE 20 FEATURES ONLY 2 ARE INFORMATIVE AND 10 ARE REDUNDANT THE FIGURE SHOWS THE ESTIMATED PROBABILITIES OBTAINED WITH LOGISTIC REGRESSION A LINEAR SUPPORTVECTOR CLASSIFIER SVC AND LINEAR SVC WITH BOTH ISOTONIC CALIBRATION AND SIGMOID CALIBRATION THE BRIER SCORE IS A METRIC WHICH IS A COMBINATION OF CALIBRATION LOSS AND REFINEMENT LOSS BRIERSCORELOSS REPORTED IN THE LEGEND THE SMALLER THE BETTER CALIBRATION LOSS IS DEFINED AS THE MEAN SQUARED DEVIATION FROM EMPIRICAL PROBABILITIES DERIVED FROM THE SLOPE OF ROC SEGMENTS REFINEMENT LOSS CAN BE DEFINED AS THE EXPECTED OPTIMAL LOSS AS MEASURED BY THE AREA UNDER THE OPTIMAL COST CURVE

ONE CAN OBSERVE HERE THAT LOGISTIC REGRESSION IS WELL CALIBRATED AS ITS CURVE IS NEARLY DIAGONAL LINEAR SVC’S CALIBRATION CURVE OR RELIABILITY DIAGRAM HAS A SIGMOID CURVE WHICH IS TYPICAL FOR AN UNDERCONFIDENT CLASSIFIER IN THE CASE OF LINEARSVC THIS IS CAUSED BY THE MARGIN PROPERTY OF THE HINGE LOSS WHICH LETS THE MODEL FOCUS ON HARD SAMPLES THAT ARE CLOSE TO THE DECISION BOUNDARY THE SUPPORT VECTORS BOTH KINDS OF CALIBRATION CAN FIX THIS ISSUE AND YIELD NEARLY IDENTICAL RESULTS THE NEXT FIGURE SHOWS THE CALIBRATION CURVE OF GAUSSIAN NAIVE BAYES ON THE SAME DATA WITH BOTH KINDS OF CALIBRATION AND ALSO WITHOUT CALIBRATION

ONE CAN SEE THAT GAUSSIAN NAIVE BAYES PERFORMS VERY BADLY BUT DOES SO IN AN OTHER WAY THAN LINEAR SVC WHILE LINEAR SVC EXHIBITED A SIGMOID CALIBRATION CURVE GAUSSIAN NAIVE BAYES’ CALIBRATION CURVE HAS A TRANSPOSEDSIGMOID SHAPE THIS IS TYPICAL FOR AN OVERCONFIDENT CLASSIFIER IN THIS CASE THE CLASSIFIER’S OVERCONFIDENCE IS CAUSED BY THE REDUNDANT FEATURES WHICH VIOLATE THE NAIVE BAYES ASSUMPTION OF FEATUREINDEPENDENCE

4PREDICTING GOOD PROBABILITIES WITH SUPERVISED LEARNING A NICULESCUMIZIL R CARUANA ICML 2005

5ON THE COMBINATION OF FORECAST PROBABILITIES FOR CONSECUTIVE PRECIPITATION PERIODS WEA FORECASTING 5 640-650 WILKS D S 1990.

31 SUPERVISED LEARNING 323









SCIKITLEARN USER GUIDE RELEASE 0213

CALIBRATION OF THE PROBABILITIES OF GAUSSIAN NAIVE BAYES WITH ISOTONIC REGRESSION CAN FIX THIS ISSUE AS CAN BE SEEN FROM THE NEARLY DIAGONAL CALIBRATION CURVE SIGMOID CALIBRATION ALSO IMPROVES THE BRIER SCORE SLIGHTLY ALBEIT NOT AS STRONGLY AS THE NONPARAMETRIC ISOTONIC CALIBRATION THIS IS AN INTRINSIC LIMITATION OF SIGMOID CALIBRATION WHOSE PARAMETRIC FORM ASSUMES A SIGMOID RATHER THAN A TRANSPOSEDSIGMOID CURVE THE NONPARAMETRIC ISOTONIC CALIBRATION MODEL HOWEVER MAKES NO SUCH STRONG ASSUMPTIONS AND CAN DEAL WITH EITHER SHAPE PROVIDED THAT THERE IS SUFFICIENT CALIBRATION DATA IN GENERAL SIGMOID CALIBRATION IS PREFERABLE IN CASES WHERE THE CALIBRATION CURVE IS SIGMOID AND WHERE THERE IS LIMITED CALIBRATION DATA WHILE ISOTONIC CALIBRATION IS PREFERABLE FOR NONSIGMOID CALIBRATION CURVES AND IN SITUATIONS WHERE LARGE AMOUNTS OF DATA ARE AVAILABLE FOR CALIBRATION

CALIBRATEDCLASSIFIERCV CAN ALSO DEAL WITH CLASSIFICATION TASKS THAT INVOLVE MORE THAN TWO CLASSES IF THE BASE ESTIMATOR CAN DO SO IN THIS CASE THE CLASSIFIER IS CALIBRATED FIRST FOR EACH CLASS SEPARATELY IN AN ONEVSREST FASHION WHEN PREDICTING PROBABILITIES FOR UNSEEN DATA THE CALIBRATED PROBABILITIES FOR EACH CLASS ARE PREDICTED SEPARATELY AS THOSE PROBABILITIES DO NOT NECESSARILY SUM TO ONE A POSTPROCESSING IS PERFORMED TO NORMALIZE THEM THE NEXT IMAGE ILLUSTRATES HOW SIGMOID CALIBRATION CHANGES PREDICTED PROBABILITIES FOR A 3CLASS CLASSIFICATION PROBLEM ILLUSTRATED IS THE STANDARD 2SIMPLEX WHERE THE THREE CORNERS CORRESPOND TO THE THREE CLASSES ARROWS POINT FROM THE PROBABILITY VECTORS PREDICTED BY AN UNCALIBRATED CLASSIFIER TO THE PROBABILITY VECTORS PREDICTED BY THE SAME CLASSIFIER AFTER SIGMOID CALIBRATION ON A HOLDOUT VALIDATION SET COLORS INDICATE THE TRUE CLASS OF AN INSTANCE RED CLASS 1 GREEN CLASS 2 BLUE CLASS 3

THE BASE CLASSIFIER IS A RANDOM FOREST CLASSIFIER WITH 25 BASE ESTIMATORS TREES IF THIS CLASSIFIER IS TRAINED ON ALL 800 TRAINING DATAPPOINTS IT IS OVERLY CONFIDENT IN ITS PREDICTIONS AND THUS INCURS A LARGE LOGLOSS CALIBRATING AN IDENTICAL CLASSIFIER WHICH WAS TRAINED ON 600 DATAPPOINTS WITH METHOD'SIGMOID' ON THE REMAINING 200 DATAPPOINTS REDUCES THE CONFIDENCE OF THE PREDICTIONS IE MOVES THE PROBABILITY VECTORS FROM THE EDGES OF THE SIMPLEX TOWARDS THE CENTER

328 CHAPTER 3 USER GUIDE



SCIKITLEARN USER GUIDE RELEASE 0213

THIS CALIBRATION RESULTS IN A LOWER LOGLOSS NOTE THAT AN ALTERNATIVE WOULD HAVE BEEN TO INCREASE THE NUMBER OF BASE ESTIMATORS WHICH WOULD HAVE RESULTED IN A SIMILAR DECREASE IN LOGLOSS

REFERENCES

- OBTAINING CALIBRATED PROBABILITY ESTIMATES FROM DECISION TREES AND NAIVE BAYESIAN CLASSIFIERS B ZADROZNY C ELKAN ICML 2001
- TRANSFORMING CLASSIFIER SCORES INTO ACCURATE MULTICLASS PROBABILITY ESTIMATES B ZADROZNY C ELKAN KDD 2002
- PROBABILISTIC OUTPUTS FOR SUPPORT VECTOR MACHINES AND COMPARISONS TO REGULARIZED LIKELIHOOD METHODS J PLATT 1999

3117 NEURAL NETWORK MODELS SUPERVISED

WARNING THIS IMPLEMENTATION IS NOT INTENDED FOR LARGESCALE APPLICATIONS IN PARTICULAR SCIKITLEARN OFFERS NO GPU SUPPORT FOR MUCH FASTER GPUBASED IMPLEMENTATIONS AS WELL AS FRAMEWORKS OFFERING MUCH MORE FLEXIBILITY TO BUILD DEEP LEARNING ARCHITECTURES SEE RELATED PROJECTS

MULTILAYER PERCEPTRON

MULTILAYER PERCEPTRON MLP IS A SUPERVISED LEARNING ALGORITHM THAT LEARNS A FUNCTION  $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$  BY TRAINING ON A DATASET WHERE  $n$  IS THE NUMBER OF DIMENSIONS FOR INPUT AND  $m$  IS THE NUMBER OF DIMENSIONS FOR OUTPUT GIVEN A SET OF FEATURES  $x_1, x_2, \dots, x_n$  AND A TARGET  $y$  IT CAN LEARN A NONLINEAR FUNCTION APPROXIMATOR FOR EITHER CLASSIFICATION OR REGRESSION IT IS DIFFERENT FROM LOGISTIC REGRESSION IN THAT BETWEEN THE INPUT AND THE OUTPUT LAYER THERE CAN BE ONE OR MORE NONLINEAR LAYERS CALLED HIDDEN LAYERS FIGURE 1 SHOWS A ONE HIDDEN LAYER MLP WITH SCALAR OUTPUT

FIG 32 FIGURE 1 ONE HIDDEN LAYER MLP

THE LEFTMOST LAYER KNOWN AS THE INPUT LAYER CONSISTS OF A SET OF NEURONS  $x_1, x_2, \dots, x_n$  REPRESENTING THE INPUT FEATURES EACH NEURON IN THE HIDDEN LAYER TRANSFORMS THE VALUES FROM THE PREVIOUS LAYER WITH A WEIGHTED LINEAR SUMMATION  $z_1, z_2, \dots, z_m$  FOLLOWED BY A NONLINEAR ACTIVATION FUNCTION  $f: \mathbb{R} \rightarrow \mathbb{R}$  LIKE THE HYPERBOLIC TAN FUNCTION THE OUTPUT LAYER RECEIVES THE VALUES FROM THE LAST HIDDEN LAYER AND TRANSFORMS THEM INTO OUTPUT VALUES



SCIKITLEARN USER GUIDE RELEASE 0213

THE MODULE CONTAINS THE PUBLIC ATTRIBUTES COEFS ANDINTERCEPTS COEFS IS A LIST OF WEIGHT MATRICES WHERE WEIGHT MATRIX AT INDEX  $i$  REPRESENTS THE WEIGHTS BETWEEN LAYER  $i$  AND LAYER  $i+1$  INTERCEPTS IS A LIST OF BIAS VECTORS WHERE THE VECTOR AT INDEX  $i$  REPRESENTS THE BIAS VALUES ADDED TO LAYER  $i$

THE ADVANTAGES OF MULTILAYER PERCEPTRON ARE

- CAPABILITY TO LEARN NONLINEAR MODELS
- CAPABILITY TO LEARN MODELS IN REALTIME ONLINE LEARNING USING PARTIALFIT

THE DISADVANTAGES OF MULTILAYER PERCEPTRON MLP INCLUDE

- MLP WITH HIDDEN LAYERS HAVE A NONCONVEX LOSS FUNCTION WHERE THERE EXISTS MORE THAN ONE LOCAL MINIMUM THEREFORE DIFFERENT RANDOM WEIGHT INITIALIZATIONS CAN LEAD TO DIFFERENT VALIDATION ACCURACY
- MLP REQUIRES TUNING A NUMBER OF HYPERPARAMETERS SUCH AS THE NUMBER OF HIDDEN NEURONS LAYERS AND ITERATIONS
- MLP IS SENSITIVE TO FEATURE SCALING

PLEASE SEE TIPS ON PRACTICAL USE SECTION THAT ADDRESSES SOME OF THESE DISADVANTAGES

CLASSIFICATION

CLASSMLPCLASSIFIER IMPLEMENTS A MULTILAYER PERCEPTRON MLP ALGORITHM THAT TRAINS USING BACKPROPAGATION MLP TRAINS ON TWO ARRAYS ARRAY X OF SIZE NSAMPLES NFEATURES WHICH HOLDS THE TRAINING SAMPLES REPRESENTED AS FLOATING POINT FEATURE VECTORS AND ARRAY Y OF SIZE NSAMPLES WHICH HOLDS THE TARGET VALUES CLASS LABELS FOR THE TRAINING SAMPLES

```
FROM SKLEARNNEURALNETWORK IMPORT MLPCLASSIFIER
```

```
X 0 0 1 1
```

```
Y 0 1
```

```
CLF MLPCLASSIFIERSOLVERLBFGS ALPHA1E5
```

```
HIDDENLAYERSIZES5 2 RANDOMSTATE1
```

```
CLFFITX Y
```

```
MLPCLASSIFIERACTIVATIONRELU ALPHA1E05 BATCHSIZEAUTO
```

```
BETA109 BETA20999 EARLYSTOPPINGFALSE
```

```
EPSILON1E08 HIDDENLAYERSIZES5 2
```

```
LEARNINGRATECONSTANT LEARNINGRATEINIT0001
```

```
MAXITER200 MOMENTUM09 NITERNOCHANGE10
```

```
NESTEROVSMOMENTUMTRUE POWERT05 RANDOMSTATE1
```

```
SHUFFLETRUE SOLVERLBFGS TOL00001
```

```
VALIDATIONFRACTION01 VERBOSEFALSE WARMSTARTFALSE
```

AFTER FITTING TRAINING THE MODEL CAN PREDICT LABELS FOR NEW SAMPLES

```
CLFPREDICT2 2 1 2
```

```
ARRAY1 0
```

MLP CAN FIT A NONLINEAR MODEL TO THE TRAINING DATA CLFCOEFS CONTAINS THE WEIGHT MATRICES THAT CONSTITUTE THE MODEL PARAMETERS

```
COEFSHAPE FORCOEFINCLFCOEFS
```

```
2 5 5 2 2 1
```

CURRENTLYMLPCLASSIFIER SUPPORTS ONLY THE CROSSENTROPY LOSS FUNCTION WHICH ALLOWS PROBABILITY ESTIMATES BY RUNNING THE PREDICTPROBA METHOD

31 SUPERVISED LEARNING 331

SCIKITLEARN USER GUIDE RELEASE 0213

MLP TRAINS USING BACKPROPAGATION MORE PRECISELY IT TRAINS USING SOME FORM OF GRADIENT DESCENT AND THE GRADIENTS ARE CALCULATED USING BACKPROPAGATION FOR CLASSIFICATION IT MINIMIZES THE CROSSENTROPY LOSS FUNCTION GIVING A VECTOR OF PROBABILITY ESTIMATES [ ]PER SAMPLE [ ]

CLFPREDICTPROBA2 2 1 2

ARRAY1967E04 999801

1967E04 999801

MLPCLASSIFIER SUPPORTS MULTICLASS CLASSIFICATION BY APPLYING SOFTMAX AS THE OUTPUT FUNCTION

FURTHER THE MODEL SUPPORTS MULTILABEL CLASSIFICATION IN WHICH A SAMPLE CAN BELONG TO MORE THAN ONE CLASS FOR EACH CLASS THE RAW OUTPUT PASSES THROUGH THE LOGISTIC FUNCTION VALUES LARGER OR EQUAL TO 0.5 ARE ROUNDED TO 1 OTHERWISE TO 0 FOR A PREDICTED OUTPUT OF A SAMPLE THE INDICES WHERE THE VALUE IS 1 REPRESENTS THE ASSIGNED CLASSES OF THAT SAMPLE

X 0 0 1 1

Y 0 1 1 1

CLF MLPCLASSIFIERSOLVERLBFGS ALPHA1E5

HIDDENLAYERSIZES15 RANDOMSTATE1

CLFFITX Y

MLPCLASSIFIERACTIVATIONRELU ALPHA1E05 BATCHSIZEAUTO

BETA109 BETA20999 EARLYSTOPPINGFALSE

EPSILON1E08 HIDDENLAYERSIZES15

LEARNINGRATECONSTANT LEARNINGRATEINIT0001

MAXITER200 MOMENTUM09 NITERNOCHANGE10

NESTEROVMOMENTUMTRUE POWERT05 RANDOMSTATE1

SHUFFLETRUE SOLVERLBFGS TOL00001

VALIDATIONFRACTION01 VERBOSEFALSE WARMSTARTFALSE

CLFPREDICT1 2

ARRAY1 1

CLFPREDICT0 0

ARRAY0 1

SEE THE EXAMPLES BELOW AND THE DOCSTRING OF MLPCLASSIFIERFIT FOR FURTHER INFORMATION

EXAMPLES

- COMPARE STOCHASTIC LEARNING STRATEGIES FOR MLPCLASSIFIER
- VISUALIZATION OF MLP WEIGHTS ON MNIST

REGRESSION

CLASSMLPREGRESSOR IMPLEMENTS A MULTILAYER PERCEPTRON MLP THAT TRAINS USING BACKPROPAGATION WITH NO ACTIVATION FUNCTION IN THE OUTPUT LAYER WHICH CAN ALSO BE SEEN AS USING THE IDENTITY FUNCTION AS ACTIVATION FUNCTION THEREFORE IT USES THE SQUARE ERROR AS THE LOSS FUNCTION AND THE OUTPUT IS A SET OF CONTINUOUS VALUES

MLPREGRESSOR ALSO SUPPORTS MULTIOUTPUT REGRESSION IN WHICH A SAMPLE CAN HAVE MORE THAN ONE TARGET

REGULARIZATION

BOTHMLPREGRESSOR ANDMLPCLASSIFIER USE PARAMETER ALPHA FOR REGULARIZATION L2 REGULARIZATION TERM WHICH HELPS IN AVOIDING OVERFITTING BY PENALIZING WEIGHTS WITH LARGE MAGNITUDES FOLLOWING PLOT DISPLAYS VARYING DECISION FUNCTION WITH VALUE OF ALPHA

332 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213  
SEE THE EXAMPLES BELOW FOR FURTHER INFORMATION  
EXAMPLES

•VARYING REGULARIZATION IN MULTILAYER PERCEPTRON  
ALGORITHMS

MLP TRAINS USING STOCHASTIC GRADIENT DESCENT ADAM OR LBFGS STOCHASTIC GRADIENT DESCENT SGD UPDATES PA  
RAMETERS USING THE GRADIENT OF THE LOSS FUNCTION WITH RESPECT TO A PARAMETER THAT NEEDS ADAPTATION IE

$\theta \leftarrow \theta - \eta \nabla_{\theta} L(\theta)$   
 $\eta$  is the learning rate

WHERE  $\eta$  IS THE LEARNING RATE WHICH CONTROLS THE STEPSIZE IN THE PARAMETER SPACE SEARCH  $L(\theta)$  IS THE LOSS FUNCTION USED  
FOR THE NETWORK

MORE DETAILS CAN BE FOUND IN THE DOCUMENTATION OF SGD

ADAM IS SIMILAR TO SGD IN A SENSE THAT IT IS A STOCHASTIC OPTIMIZER BUT IT CAN AUTOMATICALLY ADJUST THE AMOUNT TO UPDATE  
PARAMETERS BASED ON ADAPTIVE ESTIMATES OF LOWERORDER MOMENTS  
WITH SGD OR ADAM TRAINING SUPPORTS ONLINE AND MINIBATCH LEARNING

LBFGS IS A SOLVER THAT APPROXIMATES THE HESSIAN MATRIX WHICH REPRESENTS THE SECONDDORDER PARTIAL DERIVATIVE OF A  
FUNCTION FURTHER IT APPROXIMATES THE INVERSE OF THE HESSIAN MATRIX TO PERFORM PARAMETER UPDATES THE IMPLEMENTATION  
USES THE SCIPY VERSION OF LBFGS

IF THE SELECTED SOLVER IS 'LBFGS' TRAINING DOES NOT SUPPORT ONLINE NOR MINIBATCH LEARNING

COMPLEXITY

SUPPOSE THERE ARE  $n$  TRAINING SAMPLES  $m$  FEATURES  $h$  HIDDEN LAYERS EACH CONTAINING  $h$  NEURONS FOR SIMPLICITY AND  $1$  OUTPUT NEURONS THE TIME COMPLEXITY OF BACKPROPAGATION IS  $O(n \cdot m \cdot h \cdot n \cdot n)$  WHERE  $n$  IS THE NUMBER OF ITERATIONS SINCE BACKPROPAGATION HAS A HIGH TIME COMPLEXITY IT IS ADVISABLE TO START WITH SMALLER NUMBER OF HIDDEN NEURONS AND FEW HIDDEN LAYERS FOR TRAINING

MATHEMATICAL FORMULATION

GIVEN A SET OF TRAINING EXAMPLES  $x_1, x_2, \dots, x_n$  WHERE  $x_i \in \mathbb{R}^m$  AND  $y_i \in \{0, 1\}$  A ONE HIDDEN LAYER

ONE HIDDEN NEURON MLP LEARNS THE FUNCTION  $f(x) = w_1 x + b_1$

WHERE  $w_1 \in \mathbb{R}^m$  AND  $b_1 \in \mathbb{R}$  ARE

MODEL PARAMETERS  $w_1, b_1$  REPRESENT THE WEIGHTS OF THE INPUT LAYER AND HIDDEN LAYER RESPECTIVELY AND  $w_2, b_2$  REPRESENT THE BIAS ADDED TO THE HIDDEN LAYER AND THE OUTPUT LAYER RESPECTIVELY  $\sigma(\cdot)$  IS THE ACTIVATION FUNCTION SET BY DEFAULT AS THE HYPERBOLIC TAN IT IS GIVEN AS

$$\sigma(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

FOR BINARY CLASSIFICATION  $\sigma$  PASSES THROUGH THE LOGISTIC FUNCTION  $\sigma(x) = \frac{1}{1 + e^{-x}}$  TO OBTAIN OUTPUT VALUES BETWEEN ZERO AND ONE A THRESHOLD SET TO 0.5 WOULD ASSIGN SAMPLES OF OUTPUTS LARGER OR EQUAL 0.5 TO THE POSITIVE CLASS AND THE REST TO THE NEGATIVE CLASS

IF THERE ARE MORE THAN TWO CLASSES  $\sigma$  ITSELF WOULD BE A VECTOR OF SIZE  $N_{CLASSES}$  INSTEAD OF PASSING THROUGH LOGISTIC FUNCTION IT PASSES THROUGH THE SOFTMAX FUNCTION WHICH IS WRITTEN AS

$$\text{SOFTMAX}(x) = \frac{\exp(x_i)}{\sum_j \exp(x_j)}$$

$$\sum_j$$

$$\exp(x_j)$$

WHERE  $x_i$  REPRESENTS THE  $i$ TH ELEMENT OF THE INPUT TO SOFTMAX WHICH CORRESPONDS TO CLASS  $i$  AND  $N$  IS THE NUMBER OF CLASSES THE RESULT IS A VECTOR CONTAINING THE PROBABILITIES THAT SAMPLE  $x$  BELONG TO EACH CLASS THE OUTPUT IS THE CLASS WITH THE HIGHEST PROBABILITY

IN REGRESSION THE OUTPUT REMAINS AS  $\sigma$  THEREFORE OUTPUT ACTIVATION FUNCTION IS JUST THE IDENTITY FUNCTION

MLP USES DIFFERENT LOSS FUNCTIONS DEPENDING ON THE PROBLEM TYPE THE LOSS FUNCTION FOR CLASSIFICATION IS CROSS ENTROPY WHICH IN BINARY CASE IS GIVEN AS

$$L(y, \hat{y}) = -\frac{1}{n} \sum_i [y_i \ln \hat{y}_i + (1 - y_i) \ln (1 - \hat{y}_i)]$$

$$2$$

WHERE  $\hat{y}_i$

IS AN L2 REGULARIZATION TERM AKA PENALTY THAT PENALIZES COMPLEX MODELS AND  $\lambda$  IS A NONNEGATIVE

HYPERPARAMETER THAT CONTROLS THE MAGNITUDE OF THE PENALTY

FOR REGRESSION MLP USES THE SQUARE ERROR LOSS FUNCTION WRITTEN AS

$$L(y, \hat{y}) = \frac{1}{n} \sum_i (y_i - \hat{y}_i)^2$$

$$2$$

$$2$$

$$2$$

$$2$$

STARTING FROM INITIAL RANDOM WEIGHTS MULTILAYER PERCEPTRON MLP MINIMIZES THE LOSS FUNCTION BY REPEATEDLY UPDATING THESE WEIGHTS AFTER COMPUTING THE LOSS A BACKWARD PASS PROPAGATES IT FROM THE OUTPUT LAYER TO THE PREVIOUS LAYERS PROVIDING EACH WEIGHT PARAMETER WITH AN UPDATE VALUE MEANT TO DECREASE THE LOSS

IN GRADIENT DESCENT THE GRADIENT  $\nabla_{w_i} L$  OF THE LOSS WITH RESPECT TO THE WEIGHTS IS COMPUTED AND DEDUCTED FROM  $w_i$  MORE FORMALLY THIS IS EXPRESSED AS

$$w_i^{t+1} = w_i^t - \eta \nabla_{w_i} L$$

$$t$$

WHERE  $t$  IS THE ITERATION STEP AND  $\eta$  IS THE LEARNING RATE WITH A VALUE LARGER THAN 0

THE ALGORITHM STOPS WHEN IT REACHES A PRESET MAXIMUM NUMBER OF ITERATIONS OR WHEN THE IMPROVEMENT IN LOSS IS BELOW A CERTAIN SMALL NUMBER

SCIKITLEARN USER GUIDE RELEASE 0213

TIPS ON PRACTICAL USE

• MULTILAYER PERCEPTRON IS SENSITIVE TO FEATURE SCALING SO IT IS HIGHLY RECOMMENDED TO SCALE YOUR DATA FOR EXAMPLE SCALE EACH ATTRIBUTE ON THE INPUT VECTOR X TO 0 1 OR 1 1 OR STANDARDIZE IT TO HAVE MEAN 0 AND VARIANCE 1 NOTE THAT YOU MUST APPLY THE SAME SCALING TO THE TEST SET FOR MEANINGFUL RESULTS YOU CAN USE STANDARDSCALER FOR STANDARDIZATION

FROM SKLEARNPREPROCESSING IMPORT STANDARDSCALER

SCALER STANDARDSCALER

DONT CHEAT FIT ONLY ON TRAINING DATA

SCALERFITXTRAIN

XTRAIN SCALERTRANSFORMXTRAIN

APPLY SAME TRANSFORMATION TO TEST DATA

XTEST SCALERTRANSFORMXTEST

AN ALTERNATIVE AND RECOMMENDED APPROACH IS TO USE STANDARDSCALER IN APIPELINE

• FINDING A REASONABLE REGULARIZATION PARAMETER  $\alpha$  IS BEST DONE USING GRIDSEARCHCV USUALLY IN THE RANGE 100 NPARANGE1 7

• EMPIRICALLY WE OBSERVED THAT LBFGS CONVERGES FASTER AND WITH BETTER SOLUTIONS ON SMALL DATASETS FOR RELATIVELY LARGE DATASETS HOWEVER ADAM IS VERY ROBUST IT USUALLY CONVERGES QUICKLY AND GIVES PRETTY GOOD PERFORMANCE STOCHASTIC GRADIENT DESCENT WITH MOMENTUM OR NESTEROV’S MOMENTUM ON THE OTHER HAND CAN PERFORM BETTER THAN THOSE TWO ALGORITHMS IF LEARNING RATE IS CORRECTLY TUNED

MORE CONTROL WITH WARMSTART

IF YOU WANT MORE CONTROL OVER STOPPING CRITERIA OR LEARNING RATE IN SGD OR WANT TO DO ADDITIONAL MONITORING USING WARMSTARTTRUE ANDMAXITER1 AND ITERATING YOURSELF CAN BE HELPFUL

X 0 0 1 1

Y 0 1

CLF MLPCLASSIFIERHIDDENLAYERSIZES15 RANDOMSTATE1 MAXITER1 WARM

↪STARTTRUE

FOR IINRANGE10

CLFFITX Y

ADDITIONAL MONITORING INSPECTION

MLPCLASSIFIER

REFERENCES

- “LEARNING REPRESENTATIONS BY BACKPROPAGATING ERRORS” RUMELHART DAVID E GEOFFREY E HINTON AND RONALD J WILLIAMS
- “STOCHASTIC GRADIENT DESCENT” L BOTTOU WEBSITE 2010
- “BACKPROPAGATION” ANDREW NG JIQUAN NGIAM CHUAN YU FOO YIFAN MAI CAROLINE SUEN WEBSITE 2011
- “EFFICIENT BACKPROP” Y LECUN L BOTTOU G ORR K MÜLLER IN NEURAL NETWORKS TRICKS OF THE TRADE 1998
- “ADAM A METHOD FOR STOCHASTIC OPTIMIZATION” KINGMA DIEDERIK AND JIMMY BA ARXIV PREPRINT ARXIV14126980 2014

SCIKITLEARN USER GUIDE RELEASE 0213

32 UNSUPERVISED LEARNING

321 GAUSSIAN MIXTURE MODELS

SKLEARNMIXTURE IS A PACKAGE WHICH ENABLES ONE TO LEARN GAUSSIAN MIXTURE MODELS DIAGONAL SPHERICAL TIED AND FULL COVARIANCE MATRICES SUPPORTED SAMPLE THEM AND ESTIMATE THEM FROM DATA FACILITIES TO HELP DETERMINE THE APPROPRIATE NUMBER OF COMPONENTS ARE ALSO PROVIDED

FIG 33 TWOCOMPONENT GAUSSIAN MIXTURE MODEL DATA POINTS AND EQUIPROBABILITY SURFACES OF THE MODEL

A GAUSSIAN MIXTURE MODEL IS A PROBABILISTIC MODEL THAT ASSUMES ALL THE DATA POINTS ARE GENERATED FROM A MIXTURE OF A FINITE NUMBER OF GAUSSIAN DISTRIBUTIONS WITH UNKNOWN PARAMETERS ONE CAN THINK OF MIXTURE MODELS AS GENERALIZING KMEANS CLUSTERING TO INCORPORATE INFORMATION ABOUT THE COVARIANCE STRUCTURE OF THE DATA AS WELL AS THE CENTERS OF THE LATENT GAUSSIANS

SCIKITLEARN IMPLEMENTS DIFFERENT CLASSES TO ESTIMATE GAUSSIAN MIXTURE MODELS THAT CORRESPOND TO DIFFERENT ESTIMATION STRATEGIES DETAILED BELOW

GAUSSIAN MIXTURE

THEGAUSSIANMIXTURE OBJECT IMPLEMENTS THE EXPECTATIONMAXIMIZATION EM ALGORITHM FOR FITTING MIXTUREOF GAUSSIAN MODELS IT CAN ALSO DRAW CONFIDENCE ELLIPSOIDS FOR MULTIVARIATE MODELS AND COMPUTE THE BAYESIAN INFORMATION CRITERION TO ASSESS THE NUMBER OF CLUSTERS IN THE DATA A GAUSSIANMIXTUREFIT METHOD IS PROVIDED THAT LEARNS A GAUSSIAN MIXTURE MODEL FROM TRAIN DATA GIVEN TEST DATA IT CAN ASSIGN TO EACH SAMPLE THE GAUSSIAN IT MOSTLY PROBABLY BELONG TO USING THE GAUSSIANMIXTUREPREDICT METHOD

THEGAUSSIANMIXTURE COMES WITH DIFFERENT OPTIONS TO CONSTRAIN THE COVARIANCE OF THE DIFFERENCE CLASSES ESTIMATED SPHERICAL DIAGONAL TIED OR FULL COVARIANCE

EXAMPLES

- SEE GMM COVARIANCES FOR AN EXAMPLE OF USING THE GAUSSIAN MIXTURE AS CLUSTERING ON THE IRIS DATASET
- SEE DENSITY ESTIMATION FOR A GAUSSIAN MIXTURE FOR AN EXAMPLE ON PLOTTING THE DENSITY ESTIMATION

336 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

PROS AND CONS OF CLASS GAUSSIANMIXTURE

PROS

SPEED IT IS THE FASTEST ALGORITHM FOR LEARNING MIXTURE MODELS

AGNOSTIC AS THIS ALGORITHM MAXIMIZES ONLY THE LIKELIHOOD IT WILL NOT BIAS THE MEANS TOWARDS ZERO OR BIAS THE CLUSTER SIZES TO HAVE SPECIFIC STRUCTURES THAT MIGHT OR MIGHT NOT APPLY

CONS

SINGULARITIES WHEN ONE HAS INSUFFICIENTLY MANY POINTS PER MIXTURE ESTIMATING THE COVARIANCE MATRICES BECOMES DIFFICULT AND THE ALGORITHM IS KNOWN TO DIVERGE AND FIND SOLUTIONS WITH INFINITE LIKELIHOOD UNLESS ONE REGULARIZES THE COVARIANCES ARTIFICIALLY

NUMBER OF COMPONENTS THIS ALGORITHM WILL ALWAYS USE ALL THE COMPONENTS IT HAS ACCESS TO NEEDING HELDOUT DATA OR INFORMATION THEORETICAL CRITERIA TO DECIDE HOW MANY COMPONENTS TO USE IN THE ABSENCE OF EXTERNAL CUES

SELECTING THE NUMBER OF COMPONENTS IN A CLASSICAL GAUSSIAN MIXTURE MODEL

THE BIC CRITERION CAN BE USED TO SELECT THE NUMBER OF COMPONENTS IN A GAUSSIAN MIXTURE IN AN EFFICIENT WAY IN THEORY IT RECOVERS THE TRUE NUMBER OF COMPONENTS ONLY IN THE ASYMPTOTIC REGIME IE IF MUCH DATA IS AVAILABLE AND ASSUMING

32 UNSUPERVISED LEARNING 337

SCIKITLEARN USER GUIDE RELEASE 0213

THAT THE DATA WAS ACTUALLY GENERATED IID FROM A MIXTURE OF GAUSSIAN DISTRIBUTION NOTE THAT USING A VARIATIONAL BAYESIAN GAUSSIAN MIXTURE AVOIDS THE SPECIFICATION OF THE NUMBER OF COMPONENTS FOR A GAUSSIAN MIXTURE MODEL EXAMPLES

- SEE GAUSSIAN MIXTURE MODEL SELECTION FOR AN EXAMPLE OF MODEL SELECTION PERFORMED WITH CLASSICAL GAUSSIAN MIXTURE

ESTIMATION ALGORITHM EXPECTATIONMAXIMIZATION

THE MAIN DIFFICULTY IN LEARNING GAUSSIAN MIXTURE MODELS FROM UNLABELED DATA IS THAT IT IS ONE USUALLY DOESN'T KNOW WHICH POINTS CAME FROM WHICH LATENT COMPONENT IF ONE HAS ACCESS TO THIS INFORMATION IT GETS VERY EASY TO FIT A SEPARATE GAUSSIAN DISTRIBUTION TO EACH SET OF POINTS EXPECTATIONMAXIMIZATION IS A WELLFOUNDED STATISTICAL ALGORITHM TO GET AROUND THIS PROBLEM BY AN ITERATIVE PROCESS FIRST ONE ASSUMES RANDOM COMPONENTS RANDOMLY CENTERED ON DATA POINTS LEARNED FROM KMEANS OR EVEN JUST NORMALLY DISTRIBUTED AROUND THE ORIGIN AND COMPUTES FOR EACH POINT A PROBABILITY OF BEING GENERATED BY EACH COMPONENT OF THE MODEL THEN ONE TWEAKS THE PARAMETERS TO MAXIMIZE THE LIKELIHOOD OF THE DATA GIVEN THOSE ASSIGNMENTS REPEATING THIS PROCESS IS GUARANTEED TO ALWAYS CONVERGE TO A LOCAL OPTIMUM

VARIATIONAL BAYESIAN GAUSSIAN MIXTURE

THEBAYESIANGAUSSIANMIXTURE OBJECT IMPLEMENTS A VARIANT OF THE GAUSSIAN MIXTURE MODEL WITH VARIATIONAL INFERENCE ALGORITHMS THE API IS SIMILAR AS THE ONE DEFINED BY GAUSSIANMIXTURE

ESTIMATION ALGORITHM VARIATIONAL INFERENCE

VARIATIONAL INFERENCE IS AN EXTENSION OF EXPECTATIONMAXIMIZATION THAT MAXIMIZES A LOWER BOUND ON MODEL EVIDENCE INCLUDING PRIORS INSTEAD OF DATA LIKELIHOOD THE PRINCIPLE BEHIND VARIATIONAL METHODS IS THE SAME AS EXPECTATION MAXIMIZATION THAT IS BOTH ARE ITERATIVE ALGORITHMS THAT ALTERNATE BETWEEN FINDING THE PROBABILITIES FOR EACH POINT TO

338 CHAPTER 3 USER GUIDE



BE GENERATED BY EACH MIXTURE AND FITTING THE MIXTURE TO THESE ASSIGNED POINTS BUT VARIATIONAL METHODS ADD REGULARIZATION BY INTEGRATING INFORMATION FROM PRIOR DISTRIBUTIONS THIS AVOIDS THE SINGULARITIES OFTEN FOUND IN EXPECTATION MAXIMIZATION SOLUTIONS BUT INTRODUCES SOME SUBTLE BIASES TO THE MODEL INFERENCE IS OFTEN NOTABLY SLOWER BUT NOT USUALLY AS MUCH SO AS TO RENDER USAGE UNPRACTICAL

DUE TO ITS BAYESIAN NATURE THE VARIATIONAL ALGORITHM NEEDS MORE HYPER PARAMETERS THAN EXPECTATIONMAXIMIZATION THE MOST IMPORTANT OF THESE BEING THE CONCENTRATION PARAMETER WEIGHTCONCENTRATIONPRIOR SPECIFYING A LOW VALUE FOR THE CONCENTRATION PRIOR WILL MAKE THE MODEL PUT MOST OF THE WEIGHT ON FEW COMPONENTS SET THE REMAINING COMPONENTS WEIGHTS VERY CLOSE TO ZERO HIGH VALUES OF THE CONCENTRATION PRIOR WILL ALLOW A LARGER NUMBER OF COMPONENTS TO BE ACTIVE IN THE MIXTURE

THE PARAMETERS IMPLEMENTATION OF THE BAYESIANGAUSSIANMIXTURE CLASS PROPOSES TWO TYPES OF PRIOR FOR THE WEIGHTS DISTRIBUTION A FINITE MIXTURE MODEL WITH DIRICHLET DISTRIBUTION AND AN INFINITE MIXTURE MODEL WITH THE DIRICHLET PROCESS IN PRACTICE DIRICHLET PROCESS INFERENCE ALGORITHM IS APPROXIMATED AND USES A TRUNCATED DISTRIBUTION WITH A FIXED MAXIMUM NUMBER OF COMPONENTS CALLED THE STICKBREAKING REPRESENTATION THE NUMBER OF COMPONENTS ACTUALLY USED ALMOST ALWAYS DEPENDS ON THE DATA

THE NEXT FIGURE COMPARES THE RESULTS OBTAINED FOR THE DIFFERENT TYPE OF THE WEIGHT CONCENTRATION PRIOR PARAMETER WEIGHTCONCENTRATIONPRIORTYPE FOR DIFFERENT VALUES OF WEIGHTCONCENTRATIONPRIOR HERE WE CAN SEE THE VALUE OF THE WEIGHTCONCENTRATIONPRIOR PARAMETER HAS A STRONG IMPACT ON THE EFFECTIVE NUMBER OF ACTIVE COMPONENTS OBTAINED WE CAN ALSO NOTICE THAT LARGE VALUES FOR THE CONCENTRATION WEIGHT PRIOR LEAD TO MORE UNIFORM WEIGHTS WHEN THE TYPE OF PRIOR IS 'DIRICHLETDISTRIBUTION' WHILE THIS IS NOT NECESSARILY THE CASE FOR THE 'DIRICHLETPROCESS' TYPE USED BY DEFAULT

32 UNSUPERVISED LEARNING 339

SCIKITLEARN USER GUIDE RELEASE 0213

THE EXAMPLES BELOW COMPARE GAUSSIAN MIXTURE MODELS WITH A FIXED NUMBER OF COMPONENTS TO THE VARIATIONAL GAUSSIAN MIXTURE MODELS WITH A DIRICHLET PROCESS PRIOR HERE A CLASSICAL GAUSSIAN MIXTURE IS FITTED WITH 5 COMPONENTS ON A DATASET COMPOSED OF 2 CLUSTERS WE CAN SEE THAT THE VARIATIONAL GAUSSIAN MIXTURE WITH A DIRICHLET PROCESS PRIOR IS ABLE TO LIMIT ITSELF TO ONLY 2 COMPONENTS WHEREAS THE GAUSSIAN MIXTURE FITS THE DATA WITH A FIXED NUMBER OF COMPONENTS THAT HAS TO BE SET A PRIORI BY THE USER IN THIS CASE THE USER HAS SELECTED NCOMPONENTS5 WHICH DOES NOT MATCH THE TRUE GENERATIVE DISTRIBUTION OF THIS TOY DATASET NOTE THAT WITH VERY LITTLE OBSERVATIONS THE VARIATIONAL GAUSSIAN MIXTURE MODELS WITH A DIRICHLET PROCESS PRIOR CAN TAKE A CONSERVATIVE STAND AND FIT ONLY ONE COMPONENT ON THE FOLLOWING FIGURE WE ARE FITTING A DATASET NOT WELLDEPICTED BY A GAUSSIAN MIXTURE ADJUSTING THE WEIGHTCONCENTRATIONPRIOR PARAMETER OF THE BAYESIANGAUSSIANMIXTURE CONTROLS THE NUMBER OF

340 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

COMPONENTS USED TO FIT THIS DATA WE ALSO PRESENT ON THE LAST TWO PLOTS A RANDOM SAMPLING GENERATED FROM THE TWO RESULTING MIXTURES

EXAMPLES

- SEE GAUSSIAN MIXTURE MODEL ELLIPSOIDS FOR AN EXAMPLE ON PLOTTING THE CONFIDENCE ELLIPSOIDS FOR BOTH GAUSSIANMIXTURE ANDBAYESIANGAUSSIANMIXTURE
- GAUSSIAN MIXTURE MODEL SINE CURVE SHOWS USING GAUSSIANMIXTURE AND BAYESIANGAUSSIANMIXTURE TO FIT A SINE WAVE
- SEE CONCENTRATION PRIOR TYPE ANALYSIS OF VARIATION BAYESIAN GAUSSIAN MIXTURE FOR AN EX

AMPLE PLOTTING THE CONFIDENCE ELLIPSOIDS FOR THE BAYESIANGAUSSIANMIXTURE WITH DIF

32 UNSUPERVISED LEARNING 341

SCIKITLEARN USER GUIDE RELEASE 0213

FERENTWEIGHTCONCENTRATIONPRIORTYPE FOR DIFFERENT VALUES OF THE PARAMETER  
WEIGHTCONCENTRATIONPRIOR

PROS AND CONS OF VARIATIONAL INFERENCE WITH BAYESIANGAUSSIANMIXTURE

PROS

AUTOMATIC SELECTION WHENWEIGHTCONCENTRATIONPRIOR IS SMALL ENOUGH AND  
NCOMPONENTS IS LARGER THAN WHAT IS FOUND NECESSARY BY THE MODEL THE VARIATIONAL BAYESIAN  
MIXTURE MODEL HAS A NATURAL TENDENCY TO SET SOME MIXTURE WEIGHTS VALUES CLOSE TO ZERO THIS MAKES  
IT POSSIBLE TO LET THE MODEL CHOOSE A SUITABLE NUMBER OF EFFECTIVE COMPONENTS AUTOMATICALLY ONLY AN  
UPPER BOUND OF THIS NUMBER NEEDS TO BE PROVIDED NOTE HOWEVER THAT THE “IDEAL” NUMBER OF ACTIVE  
COMPONENTS IS VERY APPLICATION SPECIFIC AND IS TYPICALLY ILLDEFINED IN A DATA EXPLORATION SETTING  
LESS SENSITIVITY TO THE NUMBER OF PARAMETERS UNLIKE FINITE MODELS WHICH WILL ALMOST ALWAYS USE  
ALL COMPONENTS AS MUCH AS THEY CAN AND HENCE WILL PRODUCE WILDLY DIFFERENT SOLUTIONS FOR  
DIFFERENT NUMBERS OF COMPONENTS THE VARIATIONAL INFERENCE WITH A DIRICHLET PROCESS PRIOR  
WEIGHTCONCENTRATIONPRIORTYPE DIRICHLETPROCESS WON’T CHANGE MUCH  
WITH CHANGES TO THE PARAMETERS LEADING TO MORE STABILITY AND LESS TUNING  
REGULARIZATION DUE TO THE INCORPORATION OF PRIOR INFORMATION VARIATIONAL SOLUTIONS HAVE LESS PATHOLOGICAL  
SPECIAL CASES THAN EXPECTATIONMAXIMIZATION SOLUTIONS

CONS

SPEED THE EXTRA PARAMETRIZATION NECESSARY FOR VARIATIONAL INFERENCE MAKE INFERENCE SLOWER ALTHOUGH NOT  
BY MUCH

HYPERPARAMETERS THIS ALGORITHM NEEDS AN EXTRA HYPERPARAMETER THAT MIGHT NEED EXPERIMENTAL TUNING VIA  
CROSSVALIDATION

BIAS THERE ARE MANY IMPLICIT BIASES IN THE INFERENCE ALGORITHMS AND ALSO IN THE DIRICHLET PROCESS IF USED  
AND WHENEVER THERE IS A MISMATCH BETWEEN THESE BIASES AND THE DATA IT MIGHT BE POSSIBLE TO FIT BETTER  
MODELS USING A FINITE MIXTURE

THE DIRICHLET PROCESS

HERE WE DESCRIBE VARIATIONAL INFERENCE ALGORITHMS ON DIRICHLET PROCESS MIXTURE THE DIRICHLET PROCESS IS A PRIOR  
PROBABILITY DISTRIBUTION ON CLUSTERINGS WITH AN INFINITE UNBOUNDED NUMBER OF PARTITIONS VARIATIONAL TECHNIQUES LET US  
INCORPORATE THIS PRIOR STRUCTURE ON GAUSSIAN MIXTURE MODELS AT ALMOST NO PENALTY IN INFERENCE TIME COMPARING WITH A  
FINITE GAUSSIAN MIXTURE MODEL

AN IMPORTANT QUESTION IS HOW CAN THE DIRICHLET PROCESS USE AN INFINITE UNBOUNDED NUMBER OF CLUSTERS AND STILL BE  
CONSISTENT WHILE A FULL EXPLANATION DOESN’T FIT THIS MANUAL ONE CAN THINK OF ITS STICK BREAKING PROCESS ANALOGY TO HELP  
UNDERSTANDING IT THE STICK BREAKING PROCESS IS A GENERATIVE STORY FOR THE DIRICHLET PROCESS WE START WITH A UNITLENGTH  
STICK AND IN EACH STEP WE BREAK OFF A PORTION OF THE REMAINING STICK EACH TIME WE ASSOCIATE THE LENGTH OF THE PIECE OF  
THE STICK TO THE PROPORTION OF POINTS THAT FALLS INTO A GROUP OF THE MIXTURE AT THE END TO REPRESENT THE INFINITE MIXTURE  
WE ASSOCIATE THE LAST REMAINING PIECE OF THE STICK TO THE PROPORTION OF POINTS THAT DON’T FALL INTO ALL THE OTHER GROUPS THE  
LENGTH OF EACH PIECE IS A RANDOM VARIABLE WITH PROBABILITY PROPORTIONAL TO THE CONCENTRATION PARAMETER SMALLER VALUE  
OF THE CONCENTRATION WILL DIVIDE THE UNITLENGTH INTO LARGER PIECES OF THE STICK DEFINING MORE CONCENTRATED DISTRIBUTION  
LARGER CONCENTRATION VALUES WILL CREATE SMALLER PIECES OF THE STICK INCREASING THE NUMBER OF COMPONENTS WITH NON  
ZERO WEIGHTS

342 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

VARIATIONAL INFERENCE TECHNIQUES FOR THE DIRICHLET PROCESS STILL WORK WITH A FINITE APPROXIMATION TO THIS INFINITE MIXTURE MODEL BUT INSTEAD OF HAVING TO SPECIFY A PRIORI HOW MANY COMPONENTS ONE WANTS TO USE ONE JUST SPECIFIES THE CONCEN TRATION PARAMETER AND AN UPPER BOUND ON THE NUMBER OF MIXTURE COMPONENTS THIS UPPER BOUND ASSUMING IT IS HIGHER THAN THE “TRUE” NUMBER OF COMPONENTS AFFECTS ONLY ALGORITHMIC COMPLEXITY NOT THE ACTUAL NUMBER OF COMPONENTS USED

322 MANIFOLD LEARNING

LOOK FOR THE BARE NECESSITIES

THE SIMPLE BARE NECESSITIES

FORGET ABOUT YOUR WORRIES AND YOUR STRIFE

I MEAN THE BARE NECESSITIES

OLD MOTHER NATURE’S RECIPES

THAT BRING THE BARE NECESSITIES OF LIFE

– BALOO’S SONG THE JUNGLE BOOK

MANIFOLD LEARNING IS AN APPROACH TO NONLINEAR DIMENSIONALITY REDUCTION ALGORITHMS FOR THIS TASK ARE BASED ON THE IDEA THAT THE DIMENSIONALITY OF MANY DATA SETS IS ONLY ARTIFICIALLY HIGH

INTRODUCTION

HIGHDIMENSIONAL DATASETS CAN BE VERY DIFFICULT TO VISUALIZE WHILE DATA IN TWO OR THREE DIMENSIONS CAN BE PLOTTED TO SHOW THE INHERENT STRUCTURE OF THE DATA EQUIVALENT HIGHDIMENSIONAL PLOTS ARE MUCH LESS INTUITIVE TO AID VISUALIZATION OF THE STRUCTURE OF A DATASET THE DIMENSION MUST BE REDUCED IN SOME WAY

THE SIMPLEST WAY TO ACCOMPLISH THIS DIMENSIONALITY REDUCTION IS BY TAKING A RANDOM PROJECTION OF THE DATA THOUGH THIS ALLOWS SOME DEGREE OF VISUALIZATION OF THE DATA STRUCTURE THE RANDOMNESS OF THE CHOICE LEAVES MUCH TO BE DESIRED IN A RANDOM PROJECTION IT IS LIKELY THAT THE MORE INTERESTING STRUCTURE WITHIN THE DATA WILL BE LOST

32 UNSUPERVISED LEARNING 343

SCIKITLEARN USER GUIDE RELEASE 0213

TO ADDRESS THIS CONCERN A NUMBER OF SUPERVISED AND UNSUPERVISED LINEAR DIMENSIONALITY REDUCTION FRAMEWORKS HAVE BEEN DESIGNED SUCH AS PRINCIPAL COMPONENT ANALYSIS PCA INDEPENDENT COMPONENT ANALYSIS LINEAR DISCRIMINANT ANALYSIS AND OTHERS THESE ALGORITHMS DEFINE SPECIFIC RUBRICS TO CHOOSE AN “INTERESTING” LINEAR PROJECTION OF THE DATA THESE METHODS CAN BE POWERFUL BUT OFTEN MISS IMPORTANT NONLINEAR STRUCTURE IN THE DATA MANIFOLD LEARNING CAN BE THOUGHT OF AS AN ATTEMPT TO GENERALIZE LINEAR FRAMEWORKS LIKE PCA TO BE SENSITIVE TO NON LINEAR STRUCTURE IN DATA THOUGH SUPERVISED VARIANTS EXIST THE TYPICAL MANIFOLD LEARNING PROBLEM IS UNSUPERVISED IT LEARNS THE HIGHDIMENSIONAL STRUCTURE OF THE DATA FROM THE DATA ITSELF WITHOUT THE USE OF PREDETERMINED CLASSIFICATIONS

- EXAMPLES
- SEE MANIFOLD LEARNING ON HANDWRITTEN DIGITS LOCALLY LINEAR EMBEDDING ISOMAP FOR AN EXAMPLE OF DIMENSIONALITY REDUCTION ON HANDWRITTEN DIGITS
- SEE COMPARISON OF MANIFOLD LEARNING METHODS FOR AN EXAMPLE OF DIMENSIONALITY REDUCTION ON A TOY “S CURVE” DATASET

THE MANIFOLD LEARNING IMPLEMENTATIONS AVAILABLE IN SCIKITLEARN ARE SUMMARIZED BELOW

ISOMAP

ONE OF THE EARLIEST APPROACHES TO MANIFOLD LEARNING IS THE ISOMAP ALGORITHM SHORT FOR ISOMETRIC MAPPING ISOMAP CAN BE VIEWED AS AN EXTENSION OF MULTIDIMENSIONAL SCALING MDS OR KERNEL PCA ISOMAP SEEKS A LOWERDIMENSIONAL

SCIKITLEARN USER GUIDE RELEASE 0213

EMBEDDING WHICH MAINTAINS GEODESIC DISTANCES BETWEEN ALL POINTS ISOMAP CAN BE PERFORMED WITH THE OBJECT ISOMAP

COMPLEXITY

THE ISOMAP ALGORITHM COMPRISES THREE STAGES

1NEAREST NEIGHBOR SEARCH ISOMAP USES SKLEARNNEIGHBORSBALLTREE FOR EFFICIENT NEIGHBOR SEARCH

THE COST IS APPROXIMATELY  $\mathcal{O}(\log n \log d)$  FOR  $k$  NEAREST NEIGHBORS OF  $n$  POINTS IN  $d$  DIMENSIONS

2SHORTESTPATH GRAPH SEARCH THE MOST EFFICIENT KNOWN ALGORITHMS FOR THIS ARE DIJKSTRA’S ALGORITHM WHICH IS APPROXIMATELY  $\mathcal{O}(n^2 \log n)$  OR THE FLOYDWARSHALL ALGORITHM WHICH IS  $\mathcal{O}(n^3)$  THE ALGORITHM CAN BE SELECTED BY THE USER WITH THE PATHMETHOD KEYWORD OF ISOMAP IF UNSPECIFIED THE CODE ATTEMPTS TO CHOOSE THE BEST ALGORITHM FOR THE INPUT DATA

3PARTIAL EIGENVALUE DECOMPOSITION THE EMBEDDING IS ENCODED IN THE EIGENVECTORS CORRESPONDING TO THE  $k$  LARGEST EIGENVALUES OF THE  $k \times k$  ISOMAP KERNEL FOR A DENSE SOLVER THE COST IS APPROXIMATELY  $\mathcal{O}(n^2)$  THIS COST CAN OFTEN BE IMPROVED USING THE ARPACK SOLVER THE EIGENSOLVER CAN BE SPECIFIED BY THE USER WITH THE PATHMETHOD KEYWORD OF ISOMAP IF UNSPECIFIED THE CODE ATTEMPTS TO CHOOSE THE BEST ALGORITHM FOR THE INPUT DATA

THE OVERALL COMPLEXITY OF ISOMAP IS  $\mathcal{O}(\log n \log d + n^2 \log n + k^3)$

- $n$  NUMBER OF TRAINING DATA POINTS
- $d$  INPUT DIMENSION
- $k$  NUMBER OF NEAREST NEIGHBORS
- $m$  OUTPUT DIMENSION

REFERENCES

- “A GLOBAL GEOMETRIC FRAMEWORK FOR NONLINEAR DIMENSIONALITY REDUCTION” TENENBAUM JB DE SILVA V LANGFORD JC SCIENCE 290 5500

32 UNSUPERVISED LEARNING 345

SCIKITLEARN USER GUIDE RELEASE 0213

LOCALLY LINEAR EMBEDDING

LOCALLY LINEAR EMBEDDING LLE SEEKS A LOWERDIMENSIONAL PROJECTION OF THE DATA WHICH PRESERVES DISTANCES WITHIN LOCAL NEIGHBORHOODS IT CAN BE THOUGHT OF AS A SERIES OF LOCAL PRINCIPAL COMPONENT ANALYSES WHICH ARE GLOBALLY COMPARED TO FIND THE BEST NONLINEAR EMBEDDING

LOCALLY LINEAR EMBEDDING CAN BE PERFORMED WITH FUNCTION LOCALLYLINEAREMBEDDING OR ITS OBJECTORIENTED COUNTERPART LOCALLYLINEAREMBEDDING COMPLEXITY

THE STANDARD LLE ALGORITHM COMPRISES THREE STAGES

1NEAREST NEIGHBORS SEARCH SEE DISCUSSION UNDER ISOMAP ABOVE

2WEIGHT MATRIX CONSTRUCTION  $\mathcal{O}(n^3)$  THE CONSTRUCTION OF THE LLE WEIGHT MATRIX INVOLVES THE SOLUTION OF  $A \times W$  LINEAR EQUATION FOR EACH OF THE  $k$  LOCAL NEIGHBORHOODS

3PARTIAL EIGENVALUE DECOMPOSITION SEE DISCUSSION UNDER ISOMAP ABOVE

THE OVERALL COMPLEXITY OF STANDARD LLE IS  $\mathcal{O}(n \log n \log k)$   $\mathcal{O}(n^3)$   $\mathcal{O}(n^2)$

- $n$  NUMBER OF TRAINING DATA POINTS
- $d$  INPUT DIMENSION
- $k$  NUMBER OF NEAREST NEIGHBORS
- $m$  OUTPUT DIMENSION

REFERENCES

- “NONLINEAR DIMENSIONALITY REDUCTION BY LOCALLY LINEAR EMBEDDING” ROWEIS S SAUL L SCIENCE 2902323 2000

MODIFIED LOCALLY LINEAR EMBEDDING

ONE WELLKNOWN ISSUE WITH LLE IS THE REGULARIZATION PROBLEM WHEN THE NUMBER OF NEIGHBORS IS GREATER THAN THE NUMBER OF INPUT DIMENSIONS THE MATRIX DEFINING EACH LOCAL NEIGHBORHOOD IS RANKDEFICIENT TO ADDRESS THIS STANDARD 346 CHAPTER 3 USER GUIDE



SCIKITLEARN USER GUIDE RELEASE 0213

LLE APPLIES AN ARBITRARY REGULARIZATION PARAMETER  $\lambda$  WHICH IS CHOSEN RELATIVE TO THE TRACE OF THE LOCAL WEIGHT MATRIX  
THOUGH IT CAN BE SHOWN FORMALLY THAT AS  $\lambda \rightarrow 0$  THE SOLUTION CONVERGES TO THE DESIRED EMBEDDING THERE IS NO GUARANTEE  
THAT THE OPTIMAL SOLUTION WILL BE FOUND FOR  $\lambda = 0$  THIS PROBLEM MANIFESTS ITSELF IN EMBEDDINGS WHICH DISTORT THE  
UNDERLYING GEOMETRY OF THE MANIFOLD  
ONE METHOD TO ADDRESS THE REGULARIZATION PROBLEM IS TO USE MULTIPLE WEIGHT VECTORS IN EACH NEIGHBORHOOD  
THIS IS THE ESSENCE OF MODIFIED LOCALLY LINEAR EMBEDDING MLE MLE CAN BE PERFORMED WITH FUNCTION  
LOCALLYLINEAREMBEDDING OR ITS OBJECTORIENTED COUNTERPART LOCALLYLINEAREMBEDDING WITH THE KEY  
WORDMETHOD MODIFIED IT REQUIRES NNEIGHBORS NCOMPONENTS  
COMPLEXITY

THE MLE ALGORITHM COMPRISES THREE STAGES

1NEAREST NEIGHBORS SEARCH SAME AS STANDARD LLE

2WEIGHT MATRIX CONSTRUCTION APPROXIMATELY  $O(n^3k) - O(n^2k)$  THE FIRST TERM IS EXACTLY EQUIVALENT

TO THAT OF STANDARD LLE THE SECOND TERM HAS TO DO WITH CONSTRUCTING THE WEIGHT MATRIX FROM MULTIPLE WEIGHTS  
IN PRACTICE THE ADDED COST OF CONSTRUCTING THE MLE WEIGHT MATRIX IS RELATIVELY SMALL COMPARED TO THE COST OF  
STEPS 1 AND 3

3PARTIAL EIGENVALUE DECOMPOSITION SAME AS STANDARD LLE

THE OVERALL COMPLEXITY OF MLE IS  $O(n \log n \log k + n^3k - n^2k)$

- $n$  NUMBER OF TRAINING DATA POINTS
- $d$  INPUT DIMENSION
- $k$  NUMBER OF NEAREST NEIGHBORS
- $m$  OUTPUT DIMENSION

REFERENCES

- “MLE MODIFIED LOCALLY LINEAR EMBEDDING USING MULTIPLE WEIGHTS” ZHANG Z WANG J

SCIKITLEARN USER GUIDE RELEASE 0213

HESSIAN EIGENMAPPING

HESSIAN EIGENMAPPING ALSO KNOWN AS HESSIANBASED LLE HLLE IS ANOTHER METHOD OF SOLVING THE REGULARIZATION PROBLEM OF LLE IT REVOLVES AROUND A HESSIANBASED QUADRATIC FORM AT EACH NEIGHBORHOOD WHICH IS USED TO RECOVER THE LOCALLY LINEAR STRUCTURE THOUGH OTHER IMPLEMENTATIONS NOTE ITS POOR SCALING WITH DATA SIZE SKLEARN IMPLEMENTS SOME ALGORITHMIC IMPROVEMENTS WHICH MAKE ITS COST COMPARABLE TO THAT OF OTHER LLE VARIANTS FOR SMALL OUTPUT DIMENSION HLLE CAN BE PERFORMED WITH FUNCTION LOCALLYLINEAREMBEDDING OR ITS OBJECTORIENTED COUNTER PARTLOCALLYLINEAREMBEDDING WITH THE KEYWORD METHOD HESSIAN IT REQUIRES NNEIGHBORS

NCOMPONENTS NCOMPONENTS 3 2

COMPLEXITY

THE HLLE ALGORITHM COMPRISES THREE STAGES

1NEAREST NEIGHBORS SEARCH SAME AS STANDARD LLE

2WEIGHT MATRIX CONSTRUCTION APPROXIMATELY  $O(n^3)$   $O(n^6)$  THE FIRST TERM REFLECTS A SIMILAR COST TO THAT OF STANDARD LLE THE SECOND TERM COMES FROM A QR DECOMPOSITION OF THE LOCAL HESSIAN ESTIMATOR

3PARTIAL EIGENVALUE DECOMPOSITION SAME AS STANDARD LLE

THE OVERALL COMPLEXITY OF STANDARD HLLE IS  $O(n \log n \log n)$   $O(n^3)$   $O(n^6)$   $O(n^2)$

- $n$  NUMBER OF TRAINING DATA POINTS
- $n$  INPUT DIMENSION
- $k$  NUMBER OF NEAREST NEIGHBORS
- $m$  OUTPUT DIMENSION

REFERENCES

- “HESSIAN EIGENMAPS LOCALLY LINEAR EMBEDDING TECHNIQUES FOR HIGHDIMENSIONAL DATA” DONOHO D GRIMES C PROC NATL ACAD SCI USA 1005591 2003

SCIKITLEARN USER GUIDE RELEASE 0213

SPECTRAL EMBEDDING

SPECTRAL EMBEDDING IS AN APPROACH TO CALCULATING A NONLINEAR EMBEDDING SCIKITLEARN IMPLEMENTS LAPLACIAN EIGENMAPS WHICH FINDS A LOW DIMENSIONAL REPRESENTATION OF THE DATA USING A SPECTRAL DECOMPOSITION OF THE GRAPH LAPLACIAN THE GRAPH GENERATED CAN BE CONSIDERED AS A DISCRETE APPROXIMATION OF THE LOW DIMENSIONAL MANIFOLD IN THE HIGH DIMENSIONAL SPACE MINIMIZATION OF A COST FUNCTION BASED ON THE GRAPH ENSURES THAT POINTS CLOSE TO EACH OTHER ON THE MANIFOLD ARE MAPPED CLOSE TO EACH OTHER IN THE LOW DIMENSIONAL SPACE PRESERVING LOCAL DISTANCES SPECTRAL EMBEDDING CAN BE PERFORMED WITH THE FUNCTION SPECTRALEMBEDDING OR ITS OBJECTORIENTED COUNTERPART SPECTRALEMBEDDING

COMPLEXITY THE SPECTRAL EMBEDDING LAPLACIAN EIGENMAPS ALGORITHM COMPRISES THREE STAGES

1WEIGHTED GRAPH CONSTRUCTION TRANSFORM THE RAW INPUT DATA INTO GRAPH REPRESENTATION USING AFFINITY ADJACENCY MATRIX REPRESENTATION

2GRAPH LAPLACIAN CONSTRUCTION UNNORMALIZED GRAPH LAPLACIAN IS CONSTRUCTED AS  $A - \frac{1}{n} \mathbf{1}\mathbf{1}^T$  FOR AND NORMALIZED ONE AS  $\frac{1}{\sqrt{D}}(A - \frac{1}{n} \mathbf{1}\mathbf{1}^T)\frac{1}{\sqrt{D}}$

2

3PARTIAL EIGENVALUE DECOMPOSITION EIGENVALUE DECOMPOSITION IS DONE ON GRAPH LAPLACIAN

THE OVERALL COMPLEXITY OF SPECTRAL EMBEDDING IS  $O(n^3 \log n)$

- $n$  NUMBER OF TRAINING DATA POINTS
- $d$  INPUT DIMENSION
- $k$  NUMBER OF NEAREST NEIGHBORS
- $m$  OUTPUT DIMENSION

REFERENCES

- “LAPLACIAN EIGENMAPS FOR DIMENSIONALITY REDUCTION AND DATA REPRESENTATION” M BELKIN P NIYOGI NEURAL COMPUTATION JUNE 2003 15 613731396

LOCAL TANGENT SPACE ALIGNMENT

THOUGH NOT TECHNICALLY A VARIANT OF LLE LOCAL TANGENT SPACE ALIGNMENT LTSA IS ALGORITHMICALLY SIMILAR ENOUGH TO LLE THAT IT CAN BE PUT IN THIS CATEGORY RATHER THAN FOCUSING ON PRESERVING NEIGHBORHOOD DISTANCES AS IN LLE LTSA SEEKS TO CHARACTERIZE THE LOCAL GEOMETRY AT EACH NEIGHBORHOOD VIA ITS TANGENT SPACE AND PERFORMS A GLOBAL OPTIMIZATION TO ALIGN THESE LOCAL TANGENT SPACES TO LEARN THE EMBEDDING LTSA CAN BE PERFORMED WITH FUNCTION LOCALLYLINEAREMBEDDING OR ITS OBJECTORIENTED COUNTERPART LOCALLYLINEAREMBEDDING WITH THE KEY

WORDMETHOD LTSA

COMPLEXITY

THE LTSA ALGORITHM COMPRISES THREE STAGES

1NEAREST NEIGHBORS SEARCH SAME AS STANDARD LLE

2WEIGHT MATRIX CONSTRUCTION APPROXIMATELY  $O(n^3)$  THE FIRST TERM REFLECTS A SIMILAR COST TO THAT OF STANDARD LLE

3PARTIAL EIGENVALUE DECOMPOSITION SAME AS STANDARD LLE

32 UNSUPERVISED LEARNING 349

SCIKITLEARN USER GUIDE RELEASE 0213

THE OVERALL COMPLEXITY OF STANDARD LTSA IS  $O(n \log n \log n)$  WHERE  $n$  IS THE NUMBER OF TRAINING DATA POINTS

- $n$  NUMBER OF TRAINING DATA POINTS
- $n$  INPUT DIMENSION
- $n$  NUMBER OF NEAREST NEIGHBORS
- $n$  OUTPUT DIMENSION

REFERENCES

- “PRINCIPAL MANIFOLDS AND NONLINEAR DIMENSIONALITY REDUCTION VIA TANGENT SPACE ALIGNMENT” ZHANG Z ZHA

H JOURNAL OF SHANGHAI UNIV 8406 2004

MULTIDIMENSIONAL SCALING MDS

MULTIDIMENSIONAL SCALING MDS SEEKS A LOWDIMENSIONAL REPRESENTATION OF THE DATA IN WHICH THE DISTANCES RESPECT WELL THE DISTANCES IN THE ORIGINAL HIGHDIMENSIONAL SPACE

IN GENERAL IS A TECHNIQUE USED FOR ANALYZING SIMILARITY OR DISSIMILARITY DATA MDS ATTEMPTS TO MODEL SIMILARITY OR DISSIMILARITY DATA AS DISTANCES IN A GEOMETRIC SPACES THE DATA CAN BE RATINGS OF SIMILARITY BETWEEN OBJECTS INTERACTION FREQUENCIES OF MOLECULES OR TRADE INDICES BETWEEN COUNTRIES

THERE EXISTS TWO TYPES OF MDS ALGORITHM METRIC AND NON METRIC IN THE SCIKITLEARN THE CLASS MDS IMPLEMENTS BOTH IN METRIC MDS THE INPUT SIMILARITY MATRIX ARISES FROM A METRIC AND THUS RESPECTS THE TRIANGULAR INEQUALITY THE DISTANCES BETWEEN OUTPUT TWO POINTS ARE THEN SET TO BE AS CLOSE AS POSSIBLE TO THE SIMILARITY OR DISSIMILARITY DATA IN THE NONMETRIC VERSION THE ALGORITHMS WILL TRY TO PRESERVE THE ORDER OF THE DISTANCES AND HENCE SEEK FOR A MONOTONIC RELATIONSHIP BETWEEN THE DISTANCES IN THE EMBEDDED SPACE AND THE SIMILARITIESDISSIMILARITIES

LET  $S$  BE THE SIMILARITY MATRIX AND  $X$  THE COORDINATES OF THE  $n$  INPUT POINTS DISPARITIES  $d_{ij}$  ARE TRANSFORMATION OF THE SIMILARITIES CHOSEN IN SOME OPTIMAL WAYS THE OBJECTIVE CALLED THE STRESS IS THEN DEFINED BY 
$$\text{STRESS} = \sqrt{\frac{1}{n(n-1)} \sum_{i < j} (d_{ij} - \delta_{ij})^2}$$

METRIC MDS THE SIMPLEST METRIC MDS MODEL CALLED ABSOLUTE MDS DISPARITIES ARE DEFINED BY  $d_{ij} = |x_i - x_j|$  WITH ABSOLUTE MDS THE VALUE  $d_{ij}$  SHOULD THEN CORRESPOND EXACTLY TO THE DISTANCE BETWEEN POINT  $i$  AND  $j$  IN THE EMBEDDING POINT

SCIKITLEARN USER GUIDE RELEASE 0213

MOST COMMONLY DISPARITIES ARE SET TO 1000000

NONMETRIC MDS

NON METRIC MDS FOCUSES ON THE ORDINATION OF THE DATA IF 1000 1000 THEN THE EMBEDDING SHOULD ENFORCE 1000 1000

A SIMPLE ALGORITHM TO ENFORCE THAT IS TO USE A MONOTONIC REGRESSION OF 1000ON1000 YIELDING DISPARITIES 1000IN THE SAME ORDER AS 1000

A TRIVIAL SOLUTION TO THIS PROBLEM IS TO SET ALL THE POINTS ON THE ORIGIN IN ORDER TO AVOID THAT THE DISPARITIES 1000ARE NORMALIZED

REFERENCES

• “MODERN MULTIDIMENSIONAL SCALING THEORY AND APPLICATIONS” BORG I GROENEN P SPRINGER SERIES IN STATIS TICS 1997

• “NONMETRIC MULTIDIMENSIONAL SCALING A NUMERICAL METHOD” KRUSKAL J PSYCHOMETRIKA 29 1964

32 UNSUPERVISED LEARNING 351

SCIKITLEARN USER GUIDE RELEASE 0213

- “MULTIDIMENSIONAL SCALING BY OPTIMIZING GOODNESS OF FIT TO A NONMETRIC HYPOTHESIS” KRUSKAL J PSYCHOMETRIKA 29 1964

TDISTRIBUTED STOCHASTIC NEIGHBOR EMBEDDING TSNE

TSNE TSNE CONVERTS AFFINITIES OF DATA POINTS TO PROBABILITIES THE AFFINITIES IN THE ORIGINAL SPACE ARE REPRESENTED BY GAUSSIAN JOINT PROBABILITIES AND THE AFFINITIES IN THE EMBEDDED SPACE ARE REPRESENTED BY STUDENT’S TDISTRIBUTIONS THIS ALLOWS TSNE TO BE PARTICULARLY SENSITIVE TO LOCAL STRUCTURE AND HAS A FEW OTHER ADVANTAGES OVER EXISTING TECHNIQUES

- REVEALING THE STRUCTURE AT MANY SCALES ON A SINGLE MAP
- REVEALING DATA THAT LIE IN MULTIPLE DIFFERENT MANIFOLDS OR CLUSTERS
- REDUCING THE TENDENCY TO CROWD POINTS TOGETHER AT THE CENTER

WHILE ISOMAP LLE AND VARIANTS ARE BEST SUITED TO UNFOLD A SINGLE CONTINUOUS LOW DIMENSIONAL MANIFOLD TSNE WILL FOCUS ON THE LOCAL STRUCTURE OF THE DATA AND WILL TEND TO EXTRACT CLUSTERED LOCAL GROUPS OF SAMPLES AS HIGHLIGHTED ON THE SCURVE EXAMPLE THIS ABILITY TO GROUP SAMPLES BASED ON THE LOCAL STRUCTURE MIGHT BE BENEFICIAL TO VISUALLY DISENTANGLE A DATASET THAT COMPRISES SEVERAL MANIFOLDS AT ONCE AS IS THE CASE IN THE DIGITS DATASET

THE KULLBACKLEIBLER KL DIVERGENCE OF THE JOINT PROBABILITIES IN THE ORIGINAL SPACE AND THE EMBEDDED SPACE WILL BE MINIMIZED BY GRADIENT DESCENT NOTE THAT THE KL DIVERGENCE IS NOT CONVEX IE MULTIPLE RESTARTS WITH DIFFERENT INITIALIZATIONS WILL END UP IN LOCAL MINIMA OF THE KL DIVERGENCE HENCE IT IS SOMETIMES USEFUL TO TRY DIFFERENT SEEDS AND SELECT THE EMBEDDING WITH THE LOWEST KL DIVERGENCE

THE DISADVANTAGES TO USING TSNE ARE ROUGHLY

- TSNE IS COMPUTATIONALLY EXPENSIVE AND CAN TAKE SEVERAL HOURS ON MILLIONSAMPLE DATASETS WHERE PCA WILL FINISH IN SECONDS OR MINUTES
- THE BARNESHUT TSNE METHOD IS LIMITED TO TWO OR THREE DIMENSIONAL EMBEDDINGS
- THE ALGORITHM IS STOCHASTIC AND MULTIPLE RESTARTS WITH DIFFERENT SEEDS CAN YIELD DIFFERENT EMBEDDINGS HOWEVER IT IS PERFECTLY LEGITIMATE TO PICK THE EMBEDDING WITH THE LEAST ERROR
- GLOBAL STRUCTURE IS NOT EXPLICITLY PRESERVED THIS PROBLEM IS MITIGATED BY INITIALIZING POINTS WITH PCA USING INITPCA

OPTIMIZING TSNE

THE MAIN PURPOSE OF TSNE IS VISUALIZATION OF HIGHDIMENSIONAL DATA HENCE IT WORKS BEST WHEN THE DATA WILL BE EMBEDDED ON TWO OR THREE DIMENSIONS

OPTIMIZING THE KL DIVERGENCE CAN BE A LITTLE BIT TRICKY SOMETIMES THERE ARE FIVE PARAMETERS THAT CONTROL THE OPTIMIZATION OF TSNE AND THEREFORE POSSIBLY THE QUALITY OF THE RESULTING EMBEDDING

- PERPLEXITY
- EARLY EXAGGERATION FACTOR
- LEARNING RATE
- MAXIMUM NUMBER OF ITERATIONS
- ANGLE NOT USED IN THE EXACT METHOD

THE PERPLEXITY IS DEFINED AS  $\frac{1}{\sum_{j=1}^k p_{ij}}$  WHERE  $p_i$  IS THE SHANNON ENTROPY OF THE CONDITIONAL PROBABILITY DISTRIBUTION THE PERPLEXITY OF A  $\eta$ -SIDED DIE IS  $\eta$  SO THAT  $p_i$  IS EFFECTIVELY THE NUMBER OF NEAREST NEIGHBORS TSNE CONSIDERS WHEN GENERATING THE CONDITIONAL PROBABILITIES LARGER PERPLEXITIES LEAD TO MORE NEAREST NEIGHBORS AND LESS SENSITIVE TO SMALL STRUCTURE CONVERSELY A LOWER PERPLEXITY CONSIDERS A SMALLER NUMBER OF NEIGHBORS AND THUS IGNORES MORE GLOBAL INFORMATION IN FAVOUR OF THE LOCAL NEIGHBORHOOD AS DATASET SIZES GET LARGER MORE POINTS WILL BE REQUIRED TO GET A REASONABLE SAMPLE OF THE LOCAL NEIGHBORHOOD AND HENCE LARGER PERPLEXITIES MAY BE REQUIRED SIMILARLY NOISIER DATASETS WILL REQUIRE LARGER PERPLEXITY VALUES TO ENCOMPASS ENOUGH LOCAL NEIGHBORS TO SEE BEYOND THE BACKGROUND NOISE THE MAXIMUM NUMBER OF ITERATIONS IS USUALLY HIGH ENOUGH AND DOES NOT NEED ANY TUNING THE OPTIMIZATION CONSISTS OF TWO PHASES THE EARLY EXAGGERATION PHASE AND THE FINAL OPTIMIZATION DURING EARLY EXAGGERATION THE JOINT PROBABILITIES IN THE ORIGINAL SPACE WILL BE ARTIFICIALLY INCREASED BY MULTIPLICATION WITH A GIVEN FACTOR LARGER FACTORS RESULT IN LARGER GAPS BETWEEN NATURAL CLUSTERS IN THE DATA IF THE FACTOR IS TOO HIGH THE KL DIVERGENCE COULD INCREASE DURING THIS PHASE USUALLY IT DOES NOT HAVE TO BE TUNED A CRITICAL PARAMETER IS THE LEARNING RATE IF IT IS TOO LOW GRADIENT DESCENT WILL GET STUCK IN A BAD LOCAL MINIMUM IF IT IS TOO HIGH THE KL DIVERGENCE WILL INCREASE DURING OPTIMIZATION MORE TIPS CAN BE FOUND IN LAURENS VAN DER MAATEN'S FAQ SEE REFERENCES THE LAST PARAMETER ANGLE IS A TRADEOFF BETWEEN PERFORMANCE AND ACCURACY LARGER ANGLES IMPLY THAT WE CAN APPROXIMATE LARGER REGIONS BY A SINGLE POINT LEADING TO BETTER SPEED BUT LESS ACCURATE RESULTS

"HOW TO USE TSNE EFFECTIVELY" PROVIDES A GOOD DISCUSSION OF THE EFFECTS OF THE VARIOUS PARAMETERS AS WELL AS INTERACTIVE PLOTS TO EXPLORE THE EFFECTS OF DIFFERENT PARAMETERS

BARNESHUT TSNE

THE BARNESHUT TSNE THAT HAS BEEN IMPLEMENTED HERE IS USUALLY MUCH SLOWER THAN OTHER MANIFOLD LEARNING ALGORITHMS THE OPTIMIZATION IS QUITE DIFFICULT AND THE COMPUTATION OF THE GRADIENT IS  $O(n^2 \log n)$  WHERE  $n$  IS THE NUMBER OF OUTPUT DIMENSIONS AND  $n$  IS THE NUMBER OF SAMPLES THE BARNESHUT METHOD IMPROVES ON THE EXACT METHOD WHERE TSNE COMPLEXITY IS  $O(n^2)$  BUT HAS SEVERAL OTHER NOTABLE DIFFERENCES

- THE BARNESHUT IMPLEMENTATION ONLY WORKS WHEN THE TARGET DIMENSIONALITY IS 3 OR LESS THE 2D CASE IS TYPICAL WHEN BUILDING VISUALIZATIONS
- BARNESHUT ONLY WORKS WITH DENSE INPUT DATA SPARSE DATA MATRICES CAN ONLY BE EMBEDDED WITH THE EXACT METHOD OR CAN BE APPROXIMATED BY A DENSE LOW RANK PROJECTION FOR INSTANCE USING SKLEARNDECOMPOSITION TRUNCATEDSVD
- BARNESHUT IS AN APPROXIMATION OF THE EXACT METHOD THE APPROXIMATION IS PARAMETERIZED WITH THE ANGLE PARAMETER THEREFORE THE ANGLE PARAMETER IS UNUSED WHEN METHOD="EXACT"
- BARNESHUT IS SIGNIFICANTLY MORE SCALABLE BARNESHUT CAN BE USED TO EMBED HUNDRED OF THOUSANDS OF DATA POINTS WHILE THE EXACT METHOD CAN HANDLE THOUSANDS OF SAMPLES BEFORE BECOMING COMPUTATIONALLY INTRACTABLE

SCIKITLEARN USER GUIDE RELEASE 0213

FOR VISUALIZATION PURPOSE WHICH IS THE MAIN USE CASE OF TSNE USING THE BARNESHUT METHOD IS STRONGLY RECOMMENDED THE EXACT TSNE METHOD IS USEFUL FOR CHECKING THE THEORETICALLY PROPERTIES OF THE EMBEDDING POSSIBLY IN HIGHER DIMENSIONAL SPACE BUT LIMIT TO SMALL DATASETS DUE TO COMPUTATIONAL CONSTRAINTS ALSO NOTE THAT THE DIGITS LABELS ROUGHLY MATCH THE NATURAL GROUPING FOUND BY TSNE WHILE THE LINEAR 2D PROJECTION OF THE PCA MODEL YIELDS A REPRESENTATION WHERE LABEL REGIONS LARGELY OVERLAP THIS IS A STRONG CLUE THAT THIS DATA CAN BE WELL SEPARATED BY NON LINEAR METHODS THAT FOCUS ON THE LOCAL STRUCTURE EG AN SVM WITH A GAUSSIAN RBF KERNEL HOWEVER FAILING TO VISUALIZE WELL SEPARATED HOMOGENEOUSLY LABELED GROUPS WITH TSNE IN 2D DOES NOT NECESSARILY IMPLY THAT THE DATA CANNOT BE CORRECTLY CLASSIFIED BY A SUPERVISED MODEL IT MIGHT BE THE CASE THAT 2 DIMENSIONS ARE NOT LOW ENOUGH TO ACCURATELY REPRESENTS THE INTERNAL STRUCTURE OF THE DATA

REFERENCES

- “VISUALIZING HIGHDIMENSIONAL DATA USING TSNE” VAN DER MAATEN LJP HINTON G JOURNAL OF MACHINE LEARNING RESEARCH 2008
- “TDISTRIBUTED STOCHASTIC NEIGHBOR EMBEDDING” VAN DER MAATEN LJP
- “ACCELERATING TSNE USING TREEBASED ALGORITHMS” LJP VAN DER MAATEN JOURNAL OF MACHINE LEARNING RESEARCH 15OCT32213245 2014

TIPS ON PRACTICAL USE

- MAKE SURE THE SAME SCALE IS USED OVER ALL FEATURES BECAUSE MANIFOLD LEARNING METHODS ARE BASED ON A NEAREST NEIGHBOR SEARCH THE ALGORITHM MAY PERFORM POORLY OTHERWISE SEE STANDARDSCALER FOR CONVENIENT WAYS OF SCALING HETEROGENEOUS DATA
- THE RECONSTRUCTION ERROR COMPUTED BY EACH ROUTINE CAN BE USED TO CHOOSE THE OPTIMAL OUTPUT DIMENSION FOR A  $n$ -DIMENSIONAL MANIFOLD EMBEDDED IN A  $d$ -DIMENSIONAL PARAMETER SPACE THE RECONSTRUCTION ERROR WILL DECREASE AS  $n$  COMPONENTS IS INCREASED UNTIL  $n$  COMPONENTS  $= d$
- NOTE THAT NOISY DATA CAN “SHORTCIRCUIT” THE MANIFOLD IN ESSENCE ACTING AS A BRIDGE BETWEEN PARTS OF THE MANIFOLD THAT WOULD OTHERWISE BE WELLSEPARATED MANIFOLD LEARNING ON NOISY ANDOR INCOMPLETE DATA IS AN ACTIVE AREA OF RESEARCH
- CERTAIN INPUT CONFIGURATIONS CAN LEAD TO SINGULAR WEIGHT MATRICES FOR EXAMPLE WHEN MORE THAN TWO POINTS IN THE DATASET ARE IDENTICAL OR WHEN THE DATA IS SPLIT INTO DISJOINTED GROUPS IN THIS CASE SOLVERARPACK WILL FAIL TO FIND THE NULL SPACE THE EASIEST WAY TO ADDRESS THIS IS TO USE SOLVERDENSE WHICH WILL WORK ON A SINGULAR MATRIX THOUGH IT MAY BE VERY SLOW DEPENDING ON THE NUMBER OF INPUT POINTS ALTERNATIVELY ONE CAN ATTEMPT TO UNDERSTAND THE SOURCE OF THE SINGULARITY IF IT IS DUE TO DISJOINT SETS INCREASING  $n$ NEIGHBORS MAY HELP IF IT IS DUE TO IDENTICAL POINTS IN THE DATASET REMOVING THESE POINTS MAY HELP

SEE ALSO  
TOTALLY RANDOM TREES EMBEDDING CAN ALSO BE USEFUL TO DERIVE NONLINEAR REPRESENTATIONS OF FEATURE SPACE ALSO IT DOES NOT PERFORM DIMENSIONALITY REDUCTION

323 CLUSTERING

CLUSTERING OF UNLABELED DATA CAN BE PERFORMED WITH THE MODULE SKLEARNCLUSTER EACH CLUSTERING ALGORITHM COMES IN TWO VARIANTS A CLASS THAT IMPLEMENTS THE FIT METHOD TO LEARN THE CLUSTERS ON TRAIN DATA AND A FUNCTION THAT GIVEN TRAIN DATA RETURNS AN ARRAY OF INTEGER LABELS CORRESPONDING TO THE DIFFERENT CLUSTERS FOR THE CLASS THE LABELS OVER THE TRAINING DATA CAN BE FOUND IN THE LABELS ATTRIBUTE

354 CHAPTER 3 USER GUIDE



INPUT DATA

ONE IMPORTANT THING TO NOTE IS THAT THE ALGORITHMS IMPLEMENTED IN THIS MODULE CAN TAKE DIFFERENT KINDS OF MATRIX AS INPUT ALL THE METHODS ACCEPT STANDARD DATA MATRICES OF SHAPE NSAMPLES NFEATURES THESE CAN BE OBTAINED FROM THE CLASSES IN THE SKLEARNFEATUREEXTRACTION MODULE FOR AFFINITYPROPAGATION SPECTRALCLUSTERING ANDDBSCAN ONE CAN ALSO INPUT SIMILARITY MATRICES OF SHAPE NSAMPLES NSAMPLES THESE CAN BE OBTAINED FROM THE FUNCTIONS IN THE SKLEARNMETRICSPAIRWISE MODULE

OVERVIEW OF CLUSTERING METHODS

FIG 34 A COMPARISON OF THE CLUSTERING ALGORITHMS IN SCIKITLEARN



POINTS FROM  $\mu$  ALTHOUGH THEY LIVE IN THE SAME SPACE

THE KMEANS ALGORITHM AIMS TO CHOOSE CENTROIDS THAT MINIMISE THE INERTIA OR WITHIN CLUSTER SUM OF SQUARES CRITERION

$$\sum$$

$$\sum 0 \text{ MIN}$$

$$\mu \in \mu_1 - \mu_2$$

INERTIA CAN BE RECOGNIZED AS A MEASURE OF HOW INTERNALLY COHERENT CLUSTERS ARE IT SUFFERS FROM VARIOUS DRAWBACKS

- INERTIA MAKES THE ASSUMPTION THAT CLUSTERS ARE CONVEX AND ISOTROPIC WHICH IS NOT ALWAYS THE CASE IT RESPONDS POORLY TO ELONGATED CLUSTERS OR MANIFOLDS WITH IRREGULAR SHAPES
- INERTIA IS NOT A NORMALIZED METRIC WE JUST KNOW THAT LOWER VALUES ARE BETTER AND ZERO IS OPTIMAL BUT IN VERY HIGH DIMENSIONAL SPACES EUCLIDEAN DISTANCES TEND TO BECOME INFLATED THIS IS AN INSTANCE OF THE SO CALLED “CURSE OF DIMENSIONALITY” RUNNING A DIMENSIONALITY REDUCTION ALGORITHM SUCH AS PRINCIPAL COMPONENT ANALYSIS PCA PRIOR TO KMEANS CLUSTERING CAN ALLEVIATE THIS PROBLEM AND SPEED UP THE COMPUTATIONS

K

MEANS IS OFTEN REFERRED TO AS LLOYD’S ALGORITHM IN BASIC TERMS THE ALGORITHM HAS THREE STEPS THE FIRST STEP CHOOSES THE INITIAL CENTROIDS WITH THE MOST BASIC METHOD BEING TO CHOOSE  $\mu$  SAMPLES FROM THE DATASET  $\mu$  AFTER INITIALIZATION

32 UNSUPERVISED LEARNING 357

SCIKITLEARN USER GUIDE RELEASE 0213

KMEANS CONSISTS OF LOOPING BETWEEN THE TWO OTHER STEPS THE FIRST STEP ASSIGNS EACH SAMPLE TO ITS NEAREST CENTROID THE SECOND STEP CREATES NEW CENTROIDS BY TAKING THE MEAN VALUE OF ALL OF THE SAMPLES ASSIGNED TO EACH PREVIOUS CENTROID THE DIFFERENCE BETWEEN THE OLD AND THE NEW CENTROIDS ARE COMPUTED AND THE ALGORITHM REPEATS THESE LAST TWO STEPS UNTIL THIS VALUE IS LESS THAN A THRESHOLD IN OTHER WORDS IT REPEATS UNTIL THE CENTROIDS DO NOT MOVE SIGNIFICANTLY KMEANS IS EQUIVALENT TO THE EXPECTATIONMAXIMIZATION ALGORITHM WITH A SMALL ALLEQUAL DIAGONAL COVARIANCE MATRIX

THE ALGORITHM CAN ALSO BE UNDERSTOOD THROUGH THE CONCEPT OF V ORONOI DIAGRAMS FIRST THE V ORONOI DIAGRAM OF THE POINTS IS CALCULATED USING THE CURRENT CENTROIDS EACH SEGMENT IN THE V ORONOI DIAGRAM BECOMES A SEPARATE CLUSTER SECONDLY THE CENTROIDS ARE UPDATED TO THE MEAN OF EACH SEGMENT THE ALGORITHM THEN REPEATS THIS UNTIL A STOPPING CRITERION IS FULFILLED USUALLY THE ALGORITHM STOPS WHEN THE RELATIVE DECREASE IN THE OBJECTIVE FUNCTION BETWEEN ITERATIONS IS LESS THAN THE GIVEN TOLERANCE VALUE THIS IS NOT THE CASE IN THIS IMPLEMENTATION ITERATION STOPS WHEN CENTROIDS MOVE LESS THAN THE TOLERANCE

GIVEN ENOUGH TIME KMEANS WILL ALWAYS CONVERGE HOWEVER THIS MAY BE TO A LOCAL MINIMUM THIS IS HIGHLY DEPENDENT ON THE INITIALIZATION OF THE CENTROIDS AS A RESULT THE COMPUTATION IS OFTEN DONE SEVERAL TIMES WITH DIFFERENT INITIALIZATIONS OF THE CENTROIDS ONE METHOD TO HELP ADDRESS THIS ISSUE IS THE KMEANS INITIALIZATION SCHEME WHICH HAS BEEN IMPLEMENTED IN SCIKITLEARN USE THE INITKMEANS PARAMETER THIS INITIALIZES THE CENTROIDS TO BE GENERALLY DISTANT FROM EACH OTHER LEADING TO PROVABLY BETTER RESULTS THAN RANDOM INITIALIZATION AS SHOWN IN THE REFERENCE

THE ALGORITHM SUPPORTS SAMPLE WEIGHTS WHICH CAN BE GIVEN BY A PARAMETER SAMPLEWEIGHT THIS ALLOWS TO ASSIGN MORE WEIGHT TO SOME SAMPLES WHEN COMPUTING CLUSTER CENTERS AND VALUES OF INERTIA FOR EXAMPLE ASSIGNING A WEIGHT OF 2 TO A SAMPLE IS EQUIVALENT TO ADDING A DUPLICATE OF THAT SAMPLE TO THE DATASET

A PARAMETER CAN BE GIVEN TO ALLOW KMEANS TO BE RUN IN PARALLEL CALLED NJOBS GIVING THIS PARAMETER A POSITIVE VALUE USES THAT MANY PROCESSORS DEFAULT 1 A VALUE OF 1 USES ALL AVAILABLE PROCESSORS WITH 2 USING ONE LESS AND SO ON PARALLELIZATION GENERALLY SPEEDS UP COMPUTATION AT THE COST OF MEMORY IN THIS CASE MULTIPLE COPIES OF CENTROIDS NEED TO BE STORED ONE FOR EACH JOB

WARNING THE PARALLEL VERSION OF KMEANS IS BROKEN ON OS X WHEN NUMPY USES THEACCELERATE FRAMEWORK THIS IS EXPECTED BEHAVIOR ACCELERATE CAN BE CALLED AFTER A FORK BUT YOU NEED TO EXECV THE SUBPROCESS WITH THE PYTHON BINARY WHICH MULTIPROCESSING DOES NOT DO UNDER POSIX

KMEANS CAN BE USED FOR VECTOR QUANTIZATION THIS IS ACHIEVED USING THE TRANSFORM METHOD OF A TRAINED MODEL OF KMEANS

EXAMPLES

- DEMONSTRATION OF KMEANS ASSUMPTIONS DEMONSTRATING WHEN KMEANS PERFORMS INTUITIVELY AND WHEN IT DOES NOT
- A DEMO OF KMEANS CLUSTERING ON THE HANDWRITTEN DIGITS DATA CLUSTERING HANDWRITTEN DIGITS

358 CHAPTER 3 USER GUIDE

REFERENCES

- “KMEANS THE ADVANTAGES OF CAREFUL SEEDING” ARTHUR DAVID AND SERGEI VASSILVITSKII PROCEEDINGS OF THE EIGHTEENTH ANNUAL ACMSIAM SYMPOSIUM ON DISCRETE ALGORITHMS SOCIETY FOR INDUSTRIAL AND APPLIED MATHEMATICS 2007

MINI BATCH KMEANS

THEMINIBATCHKMEANS IS A VARIANT OF THE KMEANS ALGORITHM WHICH USES MINIBATCHES TO REDUCE THE COMPUTATION TIME WHILE STILL ATTEMPTING TO OPTIMISE THE SAME OBJECTIVE FUNCTION MINIBATCHES ARE SUBSETS OF THE INPUT DATA RANDOMLY SAMPLED IN EACH TRAINING ITERATION THESE MINIBATCHES DRASTICALLY REDUCE THE AMOUNT OF COMPUTATION REQUIRED TO CONVERGE TO A LOCAL SOLUTION IN CONTRAST TO OTHER ALGORITHMS THAT REDUCE THE CONVERGENCE TIME OF KMEANS MINIBATCH KMEANS PRODUCES RESULTS THAT ARE GENERALLY ONLY SLIGHTLY WORSE THAN THE STANDARD ALGORITHM

THE ALGORITHM ITERATES BETWEEN TWO MAJOR STEPS SIMILAR TO VANILLA KMEANS IN THE FIRST STEP □SAMPLES ARE DRAWN RANDOMLY FROM THE DATASET TO FORM A MINIBATCH THESE ARE THEN ASSIGNED TO THE NEAREST CENTROID IN THE SECOND STEP THE CENTROIDS ARE UPDATED IN CONTRAST TO KMEANS THIS IS DONE ON A PERSAMPLE BASIS FOR EACH SAMPLE IN THE MINIBATCH THE ASSIGNED CENTROID IS UPDATED BY TAKING THE STREAMING AVERAGE OF THE SAMPLE AND ALL PREVIOUS SAMPLES ASSIGNED TO THAT CENTROID THIS HAS THE EFFECT OF DECREASING THE RATE OF CHANGE FOR A CENTROID OVER TIME THESE STEPS ARE PERFORMED UNTIL CONVERGENCE OR A PREDETERMINED NUMBER OF ITERATIONS IS REACHED

MINIBATCHKMEANS CONVERGES FASTER THAN KMEANS BUT THE QUALITY OF THE RESULTS IS REDUCED IN PRACTICE THIS DIFFERENCE IN QUALITY CAN BE QUITE SMALL AS SHOWN IN THE EXAMPLE AND CITED REFERENCE

EXAMPLES

- COMPARISON OF THE KMEANS AND MINIBATCHKMEANS CLUSTERING ALGORITHMS COMPARISON OF KMEANS AND MINIBATCHKMEANS

- CLUSTERING TEXT DOCUMENTS USING KMEANS DOCUMENT CLUSTERING USING SPARSE MINIBATCHKMEANS

- ONLINE LEARNING OF A DICTIONARY OF PARTS OF FACES

REFERENCES

- “WEB SCALE KMEANS CLUSTERING” D SCULLEY PROCEEDINGS OF THE 19TH INTERNATIONAL CONFERENCE ON WORLD WIDE WEB2010

AFFINITY PROPAGATION

AFFINITYPROPAGATION CREATES CLUSTERS BY SENDING MESSAGES BETWEEN PAIRS OF SAMPLES UNTIL CONVERGENCE A DATASET IS THEN DESCRIBED USING A SMALL NUMBER OF EXEMPLARS WHICH ARE IDENTIFIED AS THOSE MOST REPRESENTATIVE OF OTHER SAMPLES THE MESSAGES SENT BETWEEN PAIRS REPRESENT THE SUITABILITY FOR ONE SAMPLE TO BE THE EXEMPLAR OF THE OTHER WHICH IS UPDATED IN RESPONSE TO THE VALUES FROM OTHER PAIRS THIS UPDATING HAPPENS ITERATIVELY UNTIL CONVERGENCE AT WHICH POINT THE FINAL EXEMPLARS ARE CHOSEN AND HENCE THE FINAL CLUSTERING IS GIVEN AFFINITY PROPAGATION CAN BE INTERESTING AS IT CHOOSES THE NUMBER OF CLUSTERS BASED ON THE DATA PROVIDED FOR THIS PURPOSE THE TWO IMPORTANT PARAMETERS ARE THE PREFERENCE WHICH CONTROLS HOW MANY EXEMPLARS ARE USED AND THE DAMPING FACTOR WHICH DAMPS THE RESPONSIBILITY AND AVAILABILITY MESSAGES TO AVOID NUMERICAL OSCILLATIONS WHEN UPDATING THESE MESSAGES

THE MAIN DRAWBACK OF AFFINITY PROPAGATION IS ITS COMPLEXITY THE ALGORITHM HAS A TIME COMPLEXITY OF THE ORDER  $O(n^2)$  WHERE  $n$  IS THE NUMBER OF SAMPLES AND  $k$  IS THE NUMBER OF ITERATIONS UNTIL CONVERGENCE FURTHER THE MEMORY COMPLEXITY IS OF THE ORDER  $O(n^2)$ IF A DENSE SIMILARITY MATRIX IS USED BUT REDUCIBLE IF A SPARSE SIMILARITY MATRIX IS USED THIS MAKES AFFINITY PROPAGATION MOST APPROPRIATE FOR SMALL TO MEDIUM SIZED DATASETS

EXAMPLES

- DEMO OF AFFINITY PROPAGATION CLUSTERING ALGORITHM AFFINITY PROPAGATION ON A SYNTHETIC 2D DATASETS WITH 3 CLASSES
- VISUALIZING THE STOCK MARKET STRUCTURE AFFINITY PROPAGATION ON FINANCIAL TIME SERIES TO FIND GROUPS OF COMPANIES

ALGORITHM DESCRIPTION THE MESSAGES SENT BETWEEN POINTS BELONG TO ONE OF TWO CATEGORIES THE FIRST IS THE RESPONSIBILITY  $r_{ij}$  WHICH IS THE ACCUMULATED EVIDENCE THAT SAMPLE  $i$  SHOULD BE THE EXEMPLAR FOR SAMPLE  $j$  THE SECOND IS THE AVAILABILITY  $a_{ij}$  WHICH IS THE ACCUMULATED EVIDENCE THAT SAMPLE  $j$  SHOULD CHOOSE SAMPLE  $i$  TO BE ITS EXEMPLAR AND CONSIDERS THE VALUES FOR ALL OTHER SAMPLES THAT  $j$  SHOULD BE AN EXEMPLAR IN THIS WAY EXEMPLARS ARE CHOSEN BY SAMPLES IF THEY ARE 1 SIMILAR ENOUGH 2 CHOSEN BY MANY SAMPLES TO BE REPRESENTATIVE OF THEMSELVES

360 CHAPTER 3 USER GUIDE

MORE FORMALLY THE RESPONSIBILITY OF A SAMPLE  $i$  TO BE THE EXEMPLAR OF SAMPLE  $j$  IS GIVEN BY

$$r_{ij} \leftarrow \frac{1 - \sum_{k \neq j} r_{ik}}{\sum_{k \neq j} r_{ik}}$$

WHERE  $r_{ij}$  IS THE SIMILARITY BETWEEN SAMPLES  $i$  AND  $j$  THE AVAILABILITY OF SAMPLE  $i$  TO BE THE EXEMPLAR OF SAMPLE  $j$  IS GIVEN BY

$$r_{ij} \leftarrow \frac{1}{\sum_{k \neq j} r_{ik}}$$

$$r_{ij} \in [0, 1]$$

TO BEGIN WITH ALL VALUES FOR  $r_{ij}$  AND  $r_{ji}$  ARE SET TO ZERO AND THE CALCULATION OF EACH ITERATES UNTIL CONVERGENCE AS DISCUSSED ABOVE IN ORDER TO AVOID NUMERICAL OSCILLATIONS WHEN UPDATING THE MESSAGES THE DAMPING FACTOR  $\alpha$  IS INTRODUCED TO ITERATION PROCESS

$$r_{ij}^{t+1} \leftarrow \alpha \cdot r_{ij}^t + (1 - \alpha) \cdot \tilde{r}_{ij}^t$$

$$r_{ji}^{t+1} \leftarrow \alpha \cdot r_{ji}^t + (1 - \alpha) \cdot \tilde{r}_{ji}^t$$

WHERE  $t$  INDICATES THE ITERATION TIMES

MEAN SHIFT

MEANSHIFT CLUSTERING AIMS TO DISCOVER BLOBS IN A SMOOTH DENSITY OF SAMPLES IT IS A CENTROID BASED ALGORITHM WHICH WORKS BY UPDATING CANDIDATES FOR CENTROIDS TO BE THE MEAN OF THE POINTS WITHIN A GIVEN REGION THESE CANDIDATES ARE THEN FILTERED IN A POSTPROCESSING STAGE TO ELIMINATE NEARDUPLICATES TO FORM THE FINAL SET OF CENTROIDS GIVEN A CANDIDATE CENTROID  $\mu_i$  FOR ITERATION  $t$  THE CANDIDATE IS UPDATED ACCORDING TO THE FOLLOWING EQUATION

$$\mu_i^{t+1} = \frac{1}{n_i} \sum_{x \in N(\mu_i, h)} x$$

$$n_i = |N(\mu_i, h)|$$

$$N(\mu_i, h)$$

WHERE  $N(\mu_i, h)$  IS THE NEIGHBORHOOD OF SAMPLES WITHIN A GIVEN DISTANCE AROUND  $\mu_i$  AND  $\mu_i$  IS THE MEAN SHIFT VECTOR THAT IS COMPUTED FOR EACH CENTROID THAT POINTS TOWARDS A REGION OF THE MAXIMUM INCREASE IN THE DENSITY OF POINTS THIS IS COMPUTED USING THE FOLLOWING EQUATION EFFECTIVELY UPDATING A CENTROID TO BE THE MEAN OF THE SAMPLES WITHIN ITS NEIGHBORHOOD

$$\mu_i^{t+1} = \mu_i^t + \frac{1}{n_i} \sum_{x \in N(\mu_i, h)} (x - \mu_i^t)$$

$$\mu_i^{t+1} \leftarrow \mu_i^t + \frac{1}{n_i} \sum_{x \in N(\mu_i, h)} (x - \mu_i^t)$$

THE ALGORITHM AUTOMATICALLY SETS THE NUMBER OF CLUSTERS INSTEAD OF RELYING ON A PARAMETER BANDWIDTH WHICH DICTATES THE SIZE OF THE REGION TO SEARCH THROUGH THIS PARAMETER CAN BE SET MANUALLY BUT CAN BE ESTIMATED USING THE PROVIDED ESTIMATEBANDWIDTH FUNCTION WHICH IS CALLED IF THE BANDWIDTH IS NOT SET

THE ALGORITHM IS NOT HIGHLY SCALABLE AS IT REQUIRES MULTIPLE NEAREST NEIGHBOR SEARCHES DURING THE EXECUTION OF THE ALGORITHM THE ALGORITHM IS GUARANTEED TO CONVERGE HOWEVER THE ALGORITHM WILL STOP ITERATING WHEN THE CHANGE IN CENTROIDS IS SMALL

LABELLING A NEW SAMPLE IS PERFORMED BY FINDING THE NEAREST CENTROID FOR A GIVEN SAMPLE

EXAMPLES

- A DEMO OF THE MEANSHIFT CLUSTERING ALGORITHM MEAN SHIFT CLUSTERING ON A SYNTHETIC 2D DATASETS WITH 3 CLASSES

REFERENCES

- “MEAN SHIFT A ROBUST APPROACH TOWARD FEATURE SPACE ANALYSIS” D COMANICIU AND P MEER IEEE TRANSACTIONS ON PATTERN ANALYSIS AND MACHINE INTELLIGENCE 2002

SPECTRAL CLUSTERING

SPECTRALCLUSTERING DOES A LOWDIMENSION EMBEDDING OF THE AFFINITY MATRIX BETWEEN SAMPLES FOLLOWED BY A KMEANS IN THE LOW DIMENSIONAL SPACE IT IS ESPECIALLY EFFICIENT IF THE AFFINITY MATRIX IS SPARSE AND THE PYAMG MODULE IS INSTALLED SPECTRALCLUSTERING REQUIRES THE NUMBER OF CLUSTERS TO BE SPECIFIED IT WORKS WELL FOR A SMALL NUMBER OF CLUSTERS BUT IS NOT ADVISED WHEN USING MANY CLUSTERS

FOR TWO CLUSTERS IT SOLVES A CONVEX RELAXATION OF THE NORMALISED CUTS PROBLEM ON THE SIMILARITY GRAPH CUTTING THE GRAPH IN TWO SO THAT THE WEIGHT OF THE EDGES CUT IS SMALL COMPARED TO THE WEIGHTS OF THE EDGES INSIDE EACH CLUSTER THIS CRITERIA IS ESPECIALLY INTERESTING WHEN WORKING ON IMAGES GRAPH VERTICES ARE PIXELS AND EDGES OF THE SIMILARITY GRAPH ARE A FUNCTION OF THE GRADIENT OF THE IMAGE

WARNING TRANSFORMING DISTANCE TO WELLBEHAVED SIMILARITIES

NOTE THAT IF THE VALUES OF YOUR SIMILARITY MATRIX ARE NOT WELL DISTRIBUTED EG WITH NEGATIVE VALUES OR WITH A DISTANCE 362 CHAPTER 3 USER GUIDE



SCIKITLEARN USER GUIDE RELEASE 0213  
MATRIX RATHER THAN A SIMILARITY THE SPECTRAL PROBLEM WILL BE SINGULAR AND THE PROBLEM NOT SOLVABLE IN WHICH CASE IT IS ADVISED TO APPLY A TRANSFORMATION TO THE ENTRIES OF THE MATRIX FOR INSTANCE IN THE CASE OF A SIGNED DISTANCE MATRIX IS COMMON TO APPLY A HEAT KERNEL  
SIMILARITY NPEXPBETA DISTANCE DISTANCESTD  
SEE THE EXAMPLES FOR SUCH AN APPLICATION

EXAMPLES

- SPECTRAL CLUSTERING FOR IMAGE SEGMENTATION SEGMENTING OBJECTS FROM A NOISY BACKGROUND USING SPECTRAL CLUSTERING
- SEGMENTING THE PICTURE OF GREEK COINS IN REGIONS SPECTRAL CLUSTERING TO SPLIT THE IMAGE OF COINS IN REGIONS

DIFFERENT LABEL ASSIGNMENT STRATEGIES  
DIFFERENT LABEL ASSIGNMENT STRATEGIES CAN BE USED CORRESPONDING TO THE ASSIGNLABELS PARAMETER OF SPECTRALCLUSTERING THEKMEANS STRATEGY CAN MATCH FINER DETAILS OF THE DATA BUT IT CAN BE MORE UNSTABLE IN PARTICULAR UNLESS YOU CONTROL THE RANDOMSTATE IT MAY NOT BE REPRODUCIBLE FROM RUNTORUN AS IT DEPENDS ON A RANDOM INITIALIZATION ON THE OTHER HAND THE DISCRETIZE STRATEGY IS 100 REPRODUCIBLE BUT IT TENDS TO CREATE PARCELS OF FAIRLY EVEN AND GEOMETRICAL SHAPE  
ASSIGNLABELSKMEANS ASSIGNLABELSDISCRETIZE  
SPECTRAL CLUSTERING GRAPHS  
SPECTRAL CLUSTERING CAN ALSO BE USED TO CLUSTER GRAPHS BY THEIR SPECTRAL EMBEDDINGS IN THIS CASE THE AFFINITY MATRIX IS THE ADJACENCY MATRIX OF THE GRAPH AND SPECTRALCLUSTERING IS INITIALIZED WITH AFFINITYPRECOMPUTED  
32 UNSUPERVISED LEARNING 363

SCIKITLEARN USER GUIDE RELEASE 0213  
FROM SKLEARNCLUSTER IMPORT SPECTRALCLUSTERING  
SC SPECTRALCLUSTERING3 AFFINITYPRECOMPUTED NINIT100  
ASSIGNLABELSDISCRETIZE  
SCFITPREDICTADJACENCYMATRIX  
REFERENCES

- “A TUTORIAL ON SPECTRAL CLUSTERING” ULRIKE VON LUXBURG 2007
- “NORMALIZED CUTS AND IMAGE SEGMENTATION” JIANBO SHI JITENDRA MALIK 2000
- “A RANDOM WALKS VIEW OF SPECTRAL SEGMENTATION” MARINA MEILA JIANBO SHI 2001
- “ON SPECTRAL CLUSTERING ANALYSIS AND AN ALGORITHM” ANDREW Y NG MICHAEL I JORDAN YAIR WEISS 2001

HIERARCHICAL CLUSTERING  
HIERARCHICAL CLUSTERING IS A GENERAL FAMILY OF CLUSTERING ALGORITHMS THAT BUILD NESTED CLUSTERS BY MERGING OR SPLITTING THEM SUCCESSIVELY THIS HIERARCHY OF CLUSTERS IS REPRESENTED AS A TREE OR DENDROGRAM THE ROOT OF THE TREE IS THE UNIQUE CLUSTER THAT GATHERS ALL THE SAMPLES THE LEAVES BEING THE CLUSTERS WITH ONLY ONE SAMPLE SEE THE WIKIPEDIA PAGE FOR MORE DETAILS

THEAGGLOMERATIVECLUSTERING OBJECT PERFORMS A HIERARCHICAL CLUSTERING USING A BOTTOM UP APPROACH EACH OBSERVATION STARTS IN ITS OWN CLUSTER AND CLUSTERS ARE SUCCESSIVELY MERGED TOGETHER THE LINKAGE CRITERIA DETERMINES THE METRIC USED FOR THE MERGE STRATEGY

- WARD MINIMIZES THE SUM OF SQUARED DIFFERENCES WITHIN ALL CLUSTERS IT IS A VARIANCEMINIMIZING APPROACH AND IN THIS SENSE IS SIMILAR TO THE KMEANS OBJECTIVE FUNCTION BUT TACKLED WITH AN AGGLOMERATIVE HIERARCHICAL APPROACH
- MAXIMUM ORCOMPLETE LINKAGE MINIMIZES THE MAXIMUM DISTANCE BETWEEN OBSERVATIONS OF PAIRS OF CLUSTERS
- AVERAGE LINKAGE MINIMIZES THE AVERAGE OF THE DISTANCES BETWEEN ALL OBSERVATIONS OF PAIRS OF CLUSTERS
- SINGLE LINKAGE MINIMIZES THE DISTANCE BETWEEN THE CLOSEST OBSERVATIONS OF PAIRS OF CLUSTERS

AGGLOMERATIVECLUSTERING CAN ALSO SCALE TO LARGE NUMBER OF SAMPLES WHEN IT IS USED JOINTLY WITH A CONNECTIVITY MATRIX BUT IS COMPUTATIONALLY EXPENSIVE WHEN NO CONNECTIVITY CONSTRAINTS ARE ADDED BETWEEN SAMPLES IT CONSIDERS AT EACH STEP ALL THE POSSIBLE MERGES

FEATUREAGGLOMERATION  
THEFEATUREAGGLOMERATION USES AGGLOMERATIVE CLUSTERING TO GROUP TOGETHER FEATURES THAT LOOK VERY SIMILAR THUS DECREASING THE NUMBER OF FEATURES IT IS A DIMENSIONALITY REDUCTION TOOL SEE UNSUPERVISED DIMENSIONALITY REDUCTION

DIFFERENT LINKAGE TYPE WARD COMPLETE AVERAGE AND SINGLE LINKAGE  
AGGLOMERATIVECLUSTERING SUPPORTS WARD SINGLE AVERAGE AND COMPLETE LINKAGE STRATEGIES  
364 CHAPTER 3 USER GUIDE

GLOMERATIVE CLUSTER HAS A “RICH GET RICHER” BEHAVIOR THAT LEADS TO UNEVEN CLUSTER SIZES IN THIS REGARD SINGLE LINKAGE IS THE WORST STRATEGY AND WARD GIVES THE MOST REGULAR SIZES HOWEVER THE AFFINITY OR DISTANCE USED IN CLUSTERING CANNOT BE VARIED WITH WARD THUS FOR NON EUCLIDEAN METRICS AVERAGE LINKAGE IS A GOOD ALTERNATIVE SINGLE LINKAGE WHILE NOT ROBUST TO NOISY DATA CAN BE COMPUTED VERY EFFICIENTLY AND CAN THEREFORE BE USEFUL TO PROVIDE HIERARCHICAL CLUSTERING OF LARGER DATASETS SINGLE LINKAGE CAN ALSO PERFORM WELL ON NONGLOBULAR DATA

EXAMPLES

- VARIOUS AGGLOMERATIVE CLUSTERING ON A 2D EMBEDDING OF DIGITS EXPLORATION OF THE DIFFERENT LINKAGE STRATEGIES IN A REAL DATASET

SCIKITLEARN USER GUIDE RELEASE 0213

ADDING CONNECTIVITY CONSTRAINTS

AN INTERESTING ASPECT OF AGGLOMERATIVE CLUSTERING IS THAT CONNECTIVITY CONSTRAINTS CAN BE ADDED TO THIS ALGORITHM ONLY ADJACENT CLUSTERS CAN BE MERGED TOGETHER THROUGH A CONNECTIVITY MATRIX THAT DEFINES FOR EACH SAMPLE THE NEIGHBORING SAMPLES FOLLOWING A GIVEN STRUCTURE OF THE DATA FOR INSTANCE IN THE SWISSROLL EXAMPLE BELOW THE CONNECTIVITY CONSTRAINTS FORBID THE MERGING OF POINTS THAT ARE NOT ADJACENT ON THE SWISS ROLL AND THUS AVOID FORMING CLUSTERS THAT EXTEND ACROSS OVERLAPPING FOLDS OF THE ROLL THESE CONSTRAINT ARE USEFUL TO IMPOSE A CERTAIN LOCAL STRUCTURE BUT THEY ALSO MAKE THE ALGORITHM FASTER ESPECIALLY WHEN THE NUMBER OF THE SAMPLES IS HIGH

THE CONNECTIVITY CONSTRAINTS ARE IMPOSED VIA AN CONNECTIVITY MATRIX A SCIPY SPARSE MATRIX THAT HAS ELEMENTS ONLY AT THE INTERSECTION OF A ROW AND A COLUMN WITH INDICES OF THE DATASET THAT SHOULD BE CONNECTED THIS MATRIX CAN BE CONSTRUCTED FROM APRIORI INFORMATION FOR INSTANCE YOU MAY WISH TO CLUSTER WEB PAGES BY ONLY MERGING PAGES WITH A LINK POINTING FROM ONE TO ANOTHER IT CAN ALSO BE LEARNED FROM THE DATA FOR INSTANCE USING SKLEARN NEIGHBORSKNEIGHBORSGRAPH TO RESTRICT MERGING TO NEAREST NEIGHBORS AS IN THIS EXAMPLE OR USINGSKLEARN FEATUREEXTRACTIONIMAGEGRIDTOGRAPH TO ENABLE ONLY MERGING OF NEIGHBORING PIXELS ON AN IMAGE AS IN THE COIN EXAMPLE

EXAMPLES

- A DEMO OF STRUCTURED WARD HIERARCHICAL CLUSTERING ON AN IMAGE OF COINS WARD CLUSTERING TO SPLIT THE IMAGE OF COINS IN REGIONS
- HIERARCHICAL CLUSTERING STRUCTURED VS UNSTRUCTURED WARD EXAMPLE OF WARD ALGORITHM ON A SWISSROLL COMPARISON OF STRUCTURED APPROACHES VERSUS UNSTRUCTURED APPROACHES
- FEATURE AGGLOMERATION VS UNIVARIATE SELECTION EXAMPLE OF DIMENSIONALITY REDUCTION WITH FEATURE AGGLOMERATION BASED ON WARD HIERARCHICAL CLUSTERING
- AGGLOMERATIVE CLUSTERING WITH AND WITHOUT STRUCTURE

WARNING CONNECTIVITY CONSTRAINTS WITH SINGLE AVERAGE AND COMPLETE LINKAGE  
CONNECTIVITY CONSTRAINTS AND SINGLE COMPLETE OR AVERAGE LINKAGE CAN ENHANCE THE 'RICH GETTING RICHER' ASPECT OF AGGLOMERATIVE CLUSTERING PARTICULARLY SO IF THEY ARE BUILT WITH SKLEARNNEIGHBORSKNEIGHBORSGRAPH IN THE LIMIT OF A SMALL NUMBER OF CLUSTERS THEY TEND TO GIVE A FEW MACROSCOPICALLY OCCUPIED CLUSTERS AND ALMOST EMPTY ONES SEE THE DISCUSSION IN AGGLOMERATIVE CLUSTERING WITH AND WITHOUT STRUCTURE SINGLE LINKAGE IS THE MOST BRITTLE LINKAGE OPTION WITH REGARD TO THIS ISSUE

366 CHAPTER 3 USER GUIDE

VARYING THE METRIC

SINGLE AVERAGE AND COMPLETE LINKAGE CAN BE USED WITH A VARIETY OF DISTANCES OR AFFINITIES IN PARTICULAR EUCLIDEAN DISTANCE L2 MANHATTAN DISTANCE OR CITYBLOCK OR L1 COSINE DISTANCE OR ANY PRECOMPUTED AFFINITY MATRIX

•L1DISTANCE IS OFTEN GOOD FOR SPARSE FEATURES OR SPARSE NOISE IE MANY OF THE FEATURES ARE ZERO AS IN TEXT MINING USING OCCURRENCES OF RARE WORDS

•COSINE DISTANCE IS INTERESTING BECAUSE IT IS INVARIANT TO GLOBAL SCALINGS OF THE SIGNAL

THE GUIDELINES FOR CHOOSING A METRIC IS TO USE ONE THAT MAXIMIZES THE DISTANCE BETWEEN SAMPLES IN DIFFERENT CLASSES AND MINIMIZES THAT WITHIN EACH CLASS

EXAMPLES

•AGGLOMERATIVE CLUSTERING WITH DIFFERENT METRICS

DBSCAN

THEDBSCAN ALGORITHM VIEWS CLUSTERS AS AREAS OF HIGH DENSITY SEPARATED BY AREAS OF LOW DENSITY DUE TO THIS RATHER GENERIC VIEW CLUSTERS FOUND BY DBSCAN CAN BE ANY SHAPE AS OPPOSED TO KMEANS WHICH ASSUMES THAT CLUSTERS ARE CONVEX SHAPED THE CENTRAL COMPONENT TO THE DBSCAN IS THE CONCEPT OF CORE SAMPLES WHICH ARE SAMPLES THAT ARE IN AREAS OF HIGH DENSITY A CLUSTER IS THEREFORE A SET OF CORE SAMPLES EACH CLOSE TO EACH OTHER MEASURED BY SOME DISTANCE MEASURE AND A SET OF NONCORE SAMPLES THAT ARE CLOSE TO A CORE SAMPLE BUT ARE NOT THEMSELVES CORE SAMPLES THERE ARE TWO PARAMETERS TO THE ALGORITHM MINSAMPLES ANDEPS WHICH DEFINE FORMALLY WHAT WE MEAN WHEN WE SAY DENSE HIGHERMINSAMPLES OR LOWEREPS INDICATE HIGHER DENSITY NECESSARY TO FORM A CLUSTER MORE FORMALLY WE DEFINE A CORE SAMPLE AS BEING A SAMPLE IN THE DATASET SUCH THAT THERE EXIST MINSAMPLES OTHER SAMPLES WITHIN A DISTANCE OF EPS WHICH ARE DEFINED AS NEIGHBORS OF THE CORE SAMPLE THIS TELLS US THAT THE CORE SAMPLE IS IN A DENSE AREA OF THE VECTOR SPACE A CLUSTER IS A SET OF CORE SAMPLES THAT CAN BE BUILT BY RECURSIVELY TAKING A CORE SAMPLE FINDING ALL OF ITS NEIGHBORS THAT ARE CORE SAMPLES FINDING ALL OF THEIR NEIGHBORS THAT ARE CORE SAMPLES AND SO ON A CLUSTER ALSO HAS A SET OF NONCORE SAMPLES WHICH ARE SAMPLES THAT ARE NEIGHBORS OF A CORE SAMPLE IN THE CLUSTER BUT ARE NOT THEMSELVES CORE SAMPLES INTUITIVELY THESE SAMPLES ARE ON THE FRINGES OF A CLUSTER

SCIKITLEARN USER GUIDE RELEASE 0213

ANY CORE SAMPLE IS PART OF A CLUSTER BY DEFINITION ANY SAMPLE THAT IS NOT A CORE SAMPLE AND IS AT LEAST EPS IN DISTANCE FROM ANY CORE SAMPLE IS CONSIDERED AN OUTLIER BY THE ALGORITHM

WHILE THE PARAMETER MINSAMPLES PRIMARILY CONTROLS HOW TOLERANT THE ALGORITHM IS TOWARDS NOISE ON NOISY AND LARGE DATA SETS IT MAY BE DESIABLE TO INCREASE THIS PARAMETER THE PARAMETER EPS IS CRUCIAL TO CHOOSE APPROPRIATELY FOR THE DATA SET AND DISTANCE FUNCTION AND USUALLY CANNOT BE LEFT AT THE DEFAULT VALUE IT CONTROLS THE LOCAL NEIGHBORHOOD OF THE POINTS WHEN CHOSEN TOO SMALL MOST DATA WILL NOT BE CLUSTERED AT ALL AND LABELED AS 1 FOR “NOISE” WHEN CHOSEN TOO LARGE IT CAUSES CLOSE CLUSTERS TO BE MERGED INTO ONE CLUSTER AND EVENTUALLY THE ENTIRE DATA SET TO BE RETURNED AS A SINGLE CLUSTER SOME HEURISTICS FOR CHOOSING THIS PARAMETER HAVE BEEN DISCUSSED IN LITERATURE FOR EXAMPLE BASED ON A KNEE IN THE NEAREST NEIGHBOR DISTANCES PLOT AS DISCUSSED IN THE REFERENCES BELOW

IN THE FIGURE BELOW THE COLOR INDICATES CLUSTER MEMBERSHIP WITH LARGE CIRCLES INDICATING CORE SAMPLES FOUND BY THE ALGORITHM SMALLER CIRCLES ARE NONCORE SAMPLES THAT ARE STILL PART OF A CLUSTER MOREOVER THE OUTLIERS ARE INDICATED BY BLACK POINTS BELOW

EXAMPLES

- DEMO OF DBSCAN CLUSTERING ALGORITHM IMPLEMENTATION

THE DBSCAN ALGORITHM IS DETERMINISTIC ALWAYS GENERATING THE SAME CLUSTERS WHEN GIVEN THE SAME DATA IN THE SAME ORDER HOWEVER THE RESULTS CAN DIFFER WHEN DATA IS PROVIDED IN A DIFFERENT ORDER FIRST EVEN THOUGH THE CORE SAMPLES WILL ALWAYS BE ASSIGNED TO THE SAME CLUSTERS THE LABELS OF THOSE CLUSTERS WILL DEPEND ON THE ORDER IN WHICH THOSE SAMPLES ARE ENCOUNTERED IN THE DATA SECOND AND MORE IMPORTANTLY THE CLUSTERS TO WHICH NONCORE SAMPLES ARE ASSIGNED CAN DIFFER DEPENDING ON THE DATA ORDER THIS WOULD HAPPEN WHEN A NONCORE SAMPLE HAS A DISTANCE LOWER THAN EPS TO TWO CORE SAMPLES IN DIFFERENT CLUSTERS BY THE TRIANGULAR INEQUALITY THOSE TWO CORE SAMPLES MUST BE MORE DISTANT THAN EPS FROM EACH OTHER OR THEY WOULD BE IN THE SAME CLUSTER THE NONCORE SAMPLE IS ASSIGNED TO WHICHEVER CLUSTER IS GENERATED FIRST IN A PASS THROUGH THE DATA AND SO THE RESULTS WILL DEPEND ON THE DATA ORDERING

THE CURRENT IMPLEMENTATION USES BALL TREES AND KD TREES TO DETERMINE THE NEIGHBORHOOD OF POINTS WHICH AVOIDS CALCULATING THE FULL DISTANCE MATRIX AS WAS DONE IN SCIKITLEARN VERSIONS BEFORE 014 THE POSSIBILITY TO USE CUSTOM METRICS IS RETAINED FOR DETAILS SEE NEAREST NEIGHBORS

MEMORY CONSUMPTION FOR LARGE SAMPLE SIZES

368 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

THIS IMPLEMENTATION IS BY DEFAULT NOT MEMORY EFFICIENT BECAUSE IT CONSTRUCTS A FULL PAIRWISE SIMILARITY MATRIX IN THE CASE WHERE KD TREES OR BALL TREES CANNOT BE USED EG WITH SPARSE MATRICES THIS MATRIX WILL CONSUME  $N^2$  FLOATS A COUPLE OF MECHANISMS FOR GETTING AROUND THIS ARE

- USE OPTICS CLUSTERING IN CONJUNCTION WITH THE EXTRACT DBSCAN METHOD OPTICS CLUSTERING ALSO CALCULATES THE FULL PAIRWISE MATRIX BUT ONLY KEEPS ONE ROW IN MEMORY AT A TIME MEMORY COMPLEXITY  $N$
- A SPARSE RADIUS NEIGHBORHOOD GRAPH WHERE MISSING ENTRIES ARE PRESUMED TO BE OUT OF EPS CAN BE PRECOMPUTED IN A MEMORY EFFICIENT WAY AND DBSCAN CAN BE RUN OVER THIS WITH METRIC PRECOMPUTED SEE SKLEARN NEIGHBORS NEAREST NEIGHBORS RADIUS NEIGHBORS GRAPH

- THE DATASET CAN BE COMPRESSED EITHER BY REMOVING EXACT DUPLICATES IF THESE OCCUR IN YOUR DATA OR BY USING BIRCH THEN YOU ONLY HAVE A RELATIVELY SMALL NUMBER OF REPRESENTATIVES FOR A LARGE NUMBER OF POINTS YOU CAN THEN PROVIDE A SAMPLE WEIGHT WHEN FITTING DBSCAN

REFERENCES

- “A DENSITY BASED ALGORITHM FOR DISCOVERING CLUSTERS IN LARGE SPATIAL DATABASES WITH NOISE” ESTER M H P KRIEDEL J SANDER AND X XU IN PROCEEDINGS OF THE 2ND INTERNATIONAL CONFERENCE ON KNOWLEDGE DISCOVERY AND DATA MINING PORTLAND OR AAAI PRESS PP 226-231 1996

- “DBSCAN REVISITED REVISITED WHY AND HOW YOU SHOULD STILL USE DBSCAN SCHUBERT E SANDER J ESTER M KRIEDEL H P XU X 2017 IN ACM TRANSACTIONS ON DATABASE SYSTEMS TODS 423 19

OPTICS

THE OPTICS ALGORITHM SHARES MANY SIMILARITIES WITH THE DBSCAN ALGORITHM AND CAN BE CONSIDERED A GENERALIZATION OF DBSCAN THAT RELAXES THE EPS REQUIREMENT FROM A SINGLE VALUE TO A VALUE RANGE THE KEY DIFFERENCE BETWEEN DBSCAN AND OPTICS IS THAT THE OPTICS ALGORITHM BUILDS A REACHABILITY GRAPH WHICH ASSIGNS EACH SAMPLE BOTH A REACHABILITY DISTANCE AND A SPOT WITHIN THE CLUSTER ORDERING ATTRIBUTE THESE TWO ATTRIBUTES ARE ASSIGNED WHEN THE MODEL IS FITTED AND ARE USED TO DETERMINE CLUSTER MEMBERSHIP IF OPTICS IS RUN WITH THE DEFAULT VALUE OF INF SET FOR MAX EPS THEN DBSCAN STYLE CLUSTER EXTRACTION CAN BE PERFORMED REPEATEDLY IN LINEAR TIME FOR ANY GIVEN EPS VALUE USING THE CLUSTER OPTICS DBSCAN METHOD SETTING MAX EPS TO A LOWER VALUE WILL RESULT IN SHORTER RUN TIMES AND CAN BE THOUGHT OF AS THE MAXIMUM NEIGHBORHOOD RADIUS FROM EACH POINT TO FIND OTHER POTENTIAL REACHABLE POINTS

SCIKITLEARN USER GUIDE RELEASE 0213

THE REACHABILITY DISTANCES GENERATED BY OPTICS ALLOW FOR VARIABLE DENSITY EXTRACTION OF CLUSTERS WITHIN A SINGLE DATA SET AS SHOWN IN THE ABOVE PLOT. COMBINING REACHABILITY DISTANCES AND DATA SET ORDERING PRODUCES A REACHABILITY PLOT WHERE POINT DENSITY IS REPRESENTED ON THE Y-Axis AND POINTS ARE ORDERED SUCH THAT NEARBY POINTS ARE ADJACENT. 'CUTTING' THE REACHABILITY PLOT AT A SINGLE VALUE PRODUCES DBSCAN LIKE RESULTS. ALL POINTS ABOVE THE 'CUT' ARE CLASSIFIED AS NOISE AND EACH TIME THAT THERE IS A BREAK WHEN READING FROM LEFT TO RIGHT SIGNIFIES A NEW CLUSTER. THE DEFAULT CLUSTER EXTRACTION WITH OPTICS LOOKS AT THE STEEP SLOPES WITHIN THE GRAPH TO FIND CLUSTERS AND THE USER CAN DEFINE WHAT COUNTS AS A STEEP SLOPE USING THE PARAMETER  $\xi$ . THERE ARE ALSO OTHER POSSIBILITIES FOR ANALYSIS ON THE GRAPH ITSELF SUCH AS GENERATING HIERARCHICAL REPRESENTATIONS OF THE DATA THROUGH REACHABILITY PLOT DENDROGRAMS AND THE HIERARCHY OF CLUSTERS DETECTED BY THE ALGORITHM CAN BE ACCESSED THROUGH THE CLUSTER HIERARCHY PARAMETER. THE PLOT ABOVE HAS BEEN COLOR CODED SO THAT CLUSTER COLORS IN PLANAR SPACE MATCH THE LINEAR SEGMENT CLUSTERS OF THE REACHABILITY PLOT. NOTE THAT THE BLUE AND RED CLUSTERS ARE ADJACENT IN THE REACHABILITY PLOT AND CAN BE HIERARCHICALLY REPRESENTED AS CHILDREN OF A LARGER PARENT CLUSTER.

EXAMPLES

- DEMO OF OPTICS CLUSTERING ALGORITHM
- COMPARISON WITH DBSCAN

THE RESULTS FROM OPTICS CLUSTER OPTICS DBSCAN METHOD AND DBSCAN ARE VERY SIMILAR BUT NOT ALWAYS IDENTICAL. SPECIFICALLY LABELING OF PERIPHERY AND NOISE POINTS. THIS IS IN PART BECAUSE THE FIRST SAMPLES OF EACH DENSE AREA PROCESSED BY OPTICS HAVE A LARGE REACHABILITY VALUE WHILE BEING CLOSE TO OTHER POINTS IN THEIR AREA AND WILL THUS SOMETIMES BE MARKED AS NOISE RATHER THAN PERIPHERY. THIS AFFECTS ADJACENT POINTS WHEN THEY ARE CONSIDERED AS CANDIDATES FOR BEING MARKED AS EITHER PERIPHERY OR NOISE. NOTE THAT FOR ANY SINGLE VALUE OF EPS DBSCAN WILL TEND TO HAVE A SHORTER RUN TIME THAN OPTICS HOWEVER FOR REPEATED RUNS AT VARYING EPS VALUES A SINGLE RUN OF OPTICS MAY REQUIRE LESS CUMULATIVE RUNTIME THAN DBSCAN. IT IS ALSO IMPORTANT TO NOTE THAT OPTICS' OUTPUT IS CLOSE TO DBSCAN'S ONLY IF EPS AND MAXEPS ARE CLOSE.

370 CHAPTER 3 USER GUIDE



SCIKITLEARN USER GUIDE RELEASE 0213

COMPUTATIONAL COMPLEXITY

SPATIAL INDEXING TREES ARE USED TO AVOID CALCULATING THE FULL DISTANCE MATRIX AND ALLOW FOR EFFICIENT MEMORY USAGE ON LARGE SETS OF SAMPLES DIFFERENT DISTANCE METRICS CAN BE SUPPLIED VIA THE METRIC KEYWORD FOR LARGE DATASETS SIMILAR BUT NOT IDENTICAL RESULTS CAN BE OBTAINED VIA HDBSCAN THE HDBSCAN IMPLEMENTATION IS MULTITHREADED AND HAS BETTER ALGORITHMIC RUNTIME COMPLEXITY THAN OPTICS AT THE COST OF WORSE MEMORY SCALING FOR EXTREMELY LARGE DATASETS THAT EXHAUST SYSTEM MEMORY USING HDBSCAN OPTICS WILL MAINTAIN NAS OPPOSED TO N2 MEMORY SCALING HOWEVER TUNING OF THE MAXEPS PARAMETER WILL LIKELY NEED TO BE USED TO GIVE A SOLUTION IN A REASONABLE AMOUNT OF WALL TIME

REFERENCES

- “OPTICS ORDERING POINTS TO IDENTIFY THE CLUSTERING STRUCTURE” ANKERST MIHAEL MARKUS M BREUNIG HANS PETER KRIEGLER AND JÖRG SANDER IN ACM SIGMOD RECORD VOL 28 NO 2 PP 4960 ACM 1999

BIRCH

THE BIRCH BUILDS A TREE CALLED THE CHARACTERISTIC FEATURE TREE CFT FOR THE GIVEN DATA THE DATA IS ESSENTIALLY LOSSY COMPRESSED TO A SET OF CHARACTERISTIC FEATURE NODES CF NODES THE CF NODES HAVE A NUMBER OF SUBCLUSTERS CALLED CHARACTERISTIC FEATURE SUBCLUSTERS CF SUBCLUSTERS AND THESE CF SUBCLUSTERS LOCATED IN THE NONTERMINAL CF NODES CAN HAVE CF NODES AS CHILDREN

THE CF SUBCLUSTERS HOLD THE NECESSARY INFORMATION FOR CLUSTERING WHICH PREVENTS THE NEED TO HOLD THE ENTIRE INPUT DATA IN MEMORY THIS INFORMATION INCLUDES

- NUMBER OF SAMPLES IN A SUBCLUSTER
- LINEAR SUM A NDIMENSIONAL VECTOR HOLDING THE SUM OF ALL SAMPLES
- SQUARED SUM SUM OF THE SQUARED L2 NORM OF ALL SAMPLES
- CENTROIDS TO AVOID RECALCULATION LINEAR SUM NSAMPLES
- SQUARED NORM OF THE CENTROIDS

THE BIRCH ALGORITHM HAS TWO PARAMETERS THE THRESHOLD AND THE BRANCHING FACTOR THE BRANCHING FACTOR LIMITS THE NUMBER OF SUBCLUSTERS IN A NODE AND THE THRESHOLD LIMITS THE DISTANCE BETWEEN THE ENTERING SAMPLE AND THE EXISTING SUBCLUSTERS

THIS ALGORITHM CAN BE VIEWED AS AN INSTANCE OR DATA REDUCTION METHOD SINCE IT REDUCES THE INPUT DATA TO A SET OF SUBCLUSTERS WHICH ARE OBTAINED DIRECTLY FROM THE LEAVES OF THE CFT THIS REDUCED DATA CAN BE FURTHER PROCESSED BY FEEDING IT INTO A GLOBAL CLUSTERER THIS GLOBAL CLUSTERER CAN BE SET BY NCLUSTERS IF NCLUSTERS IS SET TO NONE THE SUBCLUSTERS FROM THE LEAVES ARE DIRECTLY READ OFF OTHERWISE A GLOBAL CLUSTERING STEP LABELS THESE SUBCLUSTERS INTO GLOBAL CLUSTERS LABELS AND THE SAMPLES ARE MAPPED TO THE GLOBAL LABEL OF THE NEAREST SUBCLUSTER

ALGORITHM DESCRIPTION

- A NEW SAMPLE IS INSERTED INTO THE ROOT OF THE CF TREE WHICH IS A CF NODE IT IS THEN MERGED WITH THE SUBCLUSTER OF THE ROOT THAT HAS THE SMALLEST RADIUS AFTER MERGING CONSTRAINED BY THE THRESHOLD AND BRANCHING FACTOR CONDITIONS IF THE SUBCLUSTER HAS ANY CHILD NODE THEN THIS IS DONE REPEATEDLY TILL IT REACHES A LEAF AFTER FINDING THE NEAREST SUBCLUSTER IN THE LEAF THE PROPERTIES OF THIS SUBCLUSTER AND THE PARENT SUBCLUSTERS ARE RECURSIVELY UPDATED
- IF THE RADIUS OF THE SUBCLUSTER OBTAINED BY MERGING THE NEW SAMPLE AND THE NEAREST SUBCLUSTER IS GREATER THAN THE SQUARE OF THE THRESHOLD AND IF THE NUMBER OF SUBCLUSTERS IS GREATER THAN THE BRANCHING FACTOR THEN A SPACE IS TEMPORARILY ALLOCATED TO THIS NEW SAMPLE THE TWO FARTHEST SUBCLUSTERS ARE TAKEN AND THE SUBCLUSTERS ARE DIVIDED INTO TWO GROUPS ON THE BASIS OF THE DISTANCE BETWEEN THESE SUBCLUSTERS

SCIKITLEARN USER GUIDE RELEASE 0213

- IF THIS SPLIT NODE HAS A PARENT SUBCLUSTER AND THERE IS ROOM FOR A NEW SUBCLUSTER THEN THE PARENT IS SPLIT INTO TWO IF THERE IS NO ROOM THEN THIS NODE IS AGAIN SPLIT INTO TWO AND THE PROCESS IS CONTINUED RECURSIVELY TILL IT REACHES THE ROOT

BIRCH OR MINIBATCHKMEANS

- BIRCH DOES NOT SCALE VERY WELL TO HIGH DIMENSIONAL DATA AS A RULE OF THUMB IF NFEATURES IS GREATER THAN TWENTY IT IS GENERALLY BETTER TO USE MINIBATCHKMEANS
- IF THE NUMBER OF INSTANCES OF DATA NEEDS TO BE REDUCED OR IF ONE WANTS A LARGE NUMBER OF SUBCLUSTERS EITHER AS A PREPROCESSING STEP OR OTHERWISE BIRCH IS MORE USEFUL THAN MINIBATCHKMEANS

HOW TO USE PARTIALFIT

TO AVOID THE COMPUTATION OF GLOBAL CLUSTERING FOR EVERY CALL OF PARTIALFIT THE USER IS ADVISED

1 TO SETNCLUSTERSNONE INITIALLY

2 TRAIN ALL DATA BY MULTIPLE CALLS TO PARTIALFIT

3 SETNCLUSTERS TO A REQUIRED VALUE USING BRCSETPARAMSNCLUSTERSNCLUSTERS

4 CALLPARTIALFIT FINALLY WITH NO ARGUMENTS IE BRCPARTIALFIT WHICH PERFORMS THE GLOBAL CLUS

TERING

REFERENCES

- TIAN ZHANG RAGHU RAMAKRISHNAN MARON LIVNY BIRCH AN EFFICIENT DATA CLUSTERING METHOD FOR LARGE DATABASES [HTTPSWWWCSFUCACOURSECENTRAL459HANPAPERSZHANG96PDF](https://www.cs.fucacoursecentral459.hanpapers.org/hang96.pdf)

- ROBERTO PERDISCI JBIRCH JAVA IMPLEMENTATION OF BIRCH CLUSTERING ALGORITHM [HTTPSCODEGOOGLECOM ARCHIVEPJBIRCH](http://code.google.com/archive/p/jbirch/)

CLUSTERING PERFORMANCE EVALUATION

EVALUATING THE PERFORMANCE OF A CLUSTERING ALGORITHM IS NOT AS TRIVIAL AS COUNTING THE NUMBER OF ERRORS OR THE PRECISION AND RECALL OF A SUPERVISED CLASSIFICATION ALGORITHM IN PARTICULAR ANY EVALUATION METRIC SHOULD NOT TAKE THE ABSOLUTE VALUES OF THE CLUSTER LABELS INTO ACCOUNT BUT RATHER IF THIS CLUSTERING DEFINE SEPARATIONS OF THE DATA SIMILAR TO SOME GROUND TRUTH SET OF CLASSES OR SATISFYING SOME ASSUMPTION SUCH THAT MEMBERS BELONG TO THE SAME CLASS ARE MORE SIMILAR THAT MEMBERS OF DIFFERENT CLASSES ACCORDING TO SOME SIMILARITY METRIC

SCIKITLEARN USER GUIDE RELEASE 0213

ADJUSTED RAND INDEX

GIVEN THE KNOWLEDGE OF THE GROUND TRUTH CLASS ASSIGNMENTS LABELTRUE AND OUR CLUSTERING ALGORITHM ASSIGNMENTS OF THE SAME SAMPLES LABELSPRED THE ADJUSTED RAND INDEX IS A FUNCTION THAT MEASURES THE SIMILARITY OF THE TWO ASSIGNMENTS IGNORING PERMUTATIONS AND WITH CHANCE NORMALIZATION

FROM SKLEARN IMPORT METRICS

LABELTRUE 0 0 0 1 1 1

LABELSPRED 0 0 1 1 2 2

METRICSADJUSTEDRANDSCORELABELTRUE LABELSPRED

024

ONE CAN PERMUTE 0 AND 1 IN THE PREDICTED LABELS RENAME 2 TO 3 AND GET THE SAME SCORE

LABELSPRED 1 1 0 0 3 3

METRICSADJUSTEDRANDSCORELABELTRUE LABELSPRED

024

FURTHERMORE ADJUSTEDRANDSCORE ISSYMMETRIC SWAPPING THE ARGUMENT DOES NOT CHANGE THE SCORE IT CAN THUS BE USED AS A CONSENSUS MEASURE

METRICSADJUSTEDRANDSCORELABELSPRED LABELTRUE

024

PERFECT LABELING IS SCORED 10

LABELSPRED LABELTRUE

METRICSADJUSTEDRANDSCORELABELTRUE LABELSPRED

10

BAD EG INDEPENDENT LABELINGS HAVE NEGATIVE OR CLOSE TO 00 SCORES

LABELTRUE 0 1 2 0 3 4 5 1

LABELSPRED 1 1 0 0 2 2 2 2

METRICSADJUSTEDRANDSCORELABELTRUE LABELSPRED

012

ADVANTAGES

- RANDOM UNIFORM LABEL ASSIGNMENTS HAVE A ARI SCORE CLOSE TO 00 FOR ANY VALUE OF NCLUSTERS AND NSAMPLES WHICH IS NOT THE CASE FOR RAW RAND INDEX OR THE VMEASURE FOR INSTANCE
- BOUNDED RANGE 1 1 NEGATIVE VALUES ARE BAD INDEPENDENT LABELINGS SIMILAR CLUSTERINGS HAVE A POSITIVE ARI 10 IS THE PERFECT MATCH SCORE
- NO ASSUMPTION IS MADE ON THE CLUSTER STRUCTURE CAN BE USED TO COMPARE CLUSTERING ALGORITHMS SUCH AS K MEANS WHICH ASSUMES ISOTROPIC BLOB SHAPES WITH RESULTS OF SPECTRAL CLUSTERING ALGORITHMS WHICH CAN FIND CLUSTER WITH “FOLDED” SHAPES

DRAWBACKS

- CONTRARY TO INERTIA ARI REQUIRES KNOWLEDGE OF THE GROUND TRUTH CLASSES WHILE IS ALMOST NEVER AVAILABLE IN PRACTICE OR REQUIRES MANUAL ASSIGNMENT BY HUMAN ANNOTATORS AS IN THE SUPERVISED LEARNING SETTING

32 UNSUPERVISED LEARNING 373

HOWEVER ARI CAN ALSO BE USEFUL IN A PURELY UNSUPERVISED SETTING AS A BUILDING BLOCK FOR A CONSENSUS INDEX THAT CAN BE USED FOR CLUSTERING MODEL SELECTION TODO

EXAMPLES

- ADJUSTMENT FOR CHANCE IN CLUSTERING PERFORMANCE EVALUATION ANALYSIS OF THE IMPACT OF THE DATASET SIZE ON THE VALUE OF CLUSTERING MEASURES FOR RANDOM ASSIGNMENTS

MATHEMATICAL FORMULATION

IF C IS A GROUND TRUTH CLASS ASSIGNMENT AND K THE CLUSTERING LET US DEFINE  $\alpha$  AND  $\beta$  AS

- $\alpha$  THE NUMBER OF PAIRS OF ELEMENTS THAT ARE IN THE SAME SET IN C AND IN THE SAME SET IN K
- $\beta$  THE NUMBER OF PAIRS OF ELEMENTS THAT ARE IN DIFFERENT SETS IN C AND IN DIFFERENT SETS IN K

THE RAW UNADJUSTED RAND INDEX IS THEN GIVEN BY

$$RI = \frac{\alpha}{\beta}$$

WHERE  $\beta$

2 IS THE TOTAL NUMBER OF POSSIBLE PAIRS IN THE DATASET WITHOUT ORDERING

HOWEVER THE RI SCORE DOES NOT GUARANTEE THAT RANDOM LABEL ASSIGNMENTS WILL GET A VALUE CLOSE TO ZERO ESP IF THE NUMBER OF CLUSTERS IS IN THE SAME ORDER OF MAGNITUDE AS THE NUMBER OF SAMPLES

TO COUNTER THIS EFFECT WE CAN DISCOUNT THE EXPECTED RI  $\alpha$  OF RANDOM LABELINGS BY DEFINING THE ADJUSTED RAND INDEX AS FOLLOWS

$$ARI = \frac{RI - \alpha}{1 - \alpha}$$

$$MAX RI = \alpha$$

REFERENCES

- COMPARING PARTITIONS L HUBERT AND P ARABIE JOURNAL OF CLASSIFICATION 1985
- WIKIPEDIA ENTRY FOR THE ADJUSTED RAND INDEX

MUTUAL INFORMATION BASED SCORES

GIVEN THE KNOWLEDGE OF THE GROUND TRUTH CLASS ASSIGNMENTS LABELSTRUE AND OUR CLUSTERING ALGORITHM ASSIGNMENTS OF THE SAME SAMPLES LABELSPRED THE MUTUAL INFORMATION IS A FUNCTION THAT MEASURES THE AGREEMENT OF THE TWO ASSIGNMENTS IGNORING PERMUTATIONS TWO DIFFERENT NORMALIZED VERSIONS OF THIS MEASURE ARE AVAILABLE NORMALIZED MUTUAL INFORMATION NMI ANDADJUSTED MUTUAL INFORMATION AMI NMI IS OFTEN USED IN THE LITERATURE WHILE AMI WAS PROPOSED MORE RECENTLY AND IS NORMALIZED AGAINST CHANCE

FROM SKLEARN IMPORT METRICS

LABELSTRUE 0 0 0 1 1 1

LABELSPRED 0 0 1 1 2 2

METRICSADJUSTEDMUTUALINFOSCORELABELSTRUE LABELSPRED

022504

ONE CAN PERMUTE 0 AND 1 IN THE PREDICTED LABELS RENAME 2 TO 3 AND GET THE SAME SCORE

SCIKITLEARN USER GUIDE RELEASE 0213

LABELSPRED 1 1 0 0 3 3

METRICSADJUSTEDMUTUALINFOSCORELABELSTRUE LABELSPRED  
022504

ALLMUTUALINFOSCORE ADJUSTEDMUTUALINFOSCORE ANDNORMALIZEDMUTUALINFOSCORE  
ARE SYMMETRIC SWAPPING THE ARGUMENT DOES NOT CHANGE THE SCORE THUS THEY CAN BE USED AS A CONSENSUS MEASURE  
METRICSADJUSTEDMUTUALINFOSCORELABELSPRED LABELSTRUE  
022504

PERFECT LABELING IS SCORED 10

LABELSPRED LABELSTRUE

METRICSADJUSTEDMUTUALINFOSCORELABELSTRUE LABELSPRED  
10

METRICSNORMALIZEDMUTUALINFOSCORELABELSTRUE LABELSPRED  
10

THIS IS NOT TRUE FOR MUTUALINFOSCORE WHICH IS THEREFORE HARDER TO JUDGE

METRICSMUTUALINFOSCORELABELSTRUE LABELSPRED  
069

BAD EG INDEPENDENT LABELINGS HAVE NONPOSITIVE SCORES

LABELSTRUE 0 1 2 0 3 4 5 1

LABELSPRED 1 1 0 0 2 2 2 2

METRICSADJUSTEDMUTUALINFOSCORELABELSTRUE LABELSPRED  
010526

ADVANTAGES

- RANDOM UNIFORM LABEL ASSIGNMENTS HAVE A AMI SCORE CLOSE TO 00 FOR ANY VALUE OF NCLUSTERS AND NSAMPLES WHICH IS NOT THE CASE FOR RAW MUTUAL INFORMATION OR THE VMEASURE FOR INSTANCE
- UPPER BOUND OF 1 VALUES CLOSE TO ZERO INDICATE TWO LABEL ASSIGNMENTS THAT ARE LARGELY INDEPENDENT WHILE VALUES CLOSE TO ONE INDICATE SIGNIFICANT AGREEMENT FURTHER AN AMI OF EXACTLY 1 INDICATES THAT THE TWO LABEL ASSIGNMENTS ARE EQUAL WITH OR WITHOUT PERMUTATION

DRAWBACKS

- CONTRARY TO INERTIA MIBASED MEASURES REQUIRE THE KNOWLEDGE OF THE GROUND TRUTH CLASSES WHILE ALMOST NEVER AVAILABLE IN PRACTICE OR REQUIRES MANUAL ASSIGNMENT BY HUMAN ANNOTATORS AS IN THE SUPERVISED LEARNING SETTING

HOWEVER MIBASED MEASURES CAN ALSO BE USEFUL IN PURELY UNSUPERVISED SETTING AS A BUILDING BLOCK FOR A CONSEN  
SUS INDEX THAT CAN BE USED FOR CLUSTERING MODEL SELECTION

- NMI AND MI ARE NOT ADJUSTED AGAINST CHANCE

EXAMPLES

32 UNSUPERVISED LEARNING 375

•ADJUSTMENT FOR CHANCE IN CLUSTERING PERFORMANCE EVALUATION ANALYSIS OF THE IMPACT OF THE DATASET SIZE ON THE VALUE OF CLUSTERING MEASURES FOR RANDOM ASSIGNMENTS THIS EXAMPLE ALSO INCLUDES THE ADJUSTED RAND INDEX MATHEMATICAL FORMULATION

ASSUME TWO LABEL ASSIGNMENTS OF THE SAME N OBJECTS [AND] THEIR ENTROPY IS THE AMOUNT OF UNCERTAINTY FOR A PARTITION SET DEFINED BY

$$H(C) = -\sum_{i=1}^K p_i \log p_i$$

WHERE  $p_i$  IS THE PROBABILITY THAT AN OBJECT PICKED AT RANDOM FROM  $C$  FALLS INTO CLASS  $i$  LIKEWISE FOR  $C'$

$$H(C, C') = -\sum_{i=1}^K p_i \log p_i$$

WITH  $p_i$  THE MUTUAL INFORMATION MI BETWEEN  $C$  AND  $C'$  IS CALCULATED BY

$$MI(C, C') = -\sum_{i=1}^K p_i \log p_i$$

WHERE  $p_i$  IS THE PROBABILITY THAT AN OBJECT PICKED AT RANDOM FALLS INTO BOTH CLASSES  $C$  AND  $C'$

IT ALSO CAN BE EXPRESSED IN SET CARDINALITY FORMULATION

$$MI(C, C') = -\sum_{i=1}^K p_i \log p_i$$

THE NORMALIZED MUTUAL INFORMATION IS DEFINED AS

$$NMI(C, C') = MI(C, C') / H(C, C')$$

THIS VALUE OF THE MUTUAL INFORMATION AND ALSO THE NORMALIZED VARIANT IS NOT ADJUSTED FOR CHANCE AND WILL TEND TO INCREASE AS THE NUMBER OF DIFFERENT LABELS CLUSTERS INCREASES REGARDLESS OF THE ACTUAL AMOUNT OF “MUTUAL INFORMATION” BETWEEN THE LABEL ASSIGNMENTS

THE EXPECTED VALUE FOR THE MUTUAL INFORMATION CAN BE CALCULATED USING THE FOLLOWING EQUATION VEB2009 IN THIS EQUATION  $n_1$  THE NUMBER OF ELEMENTS IN  $C$  AND  $n_2$  THE NUMBER OF ELEMENTS IN  $C'$

$$E[MI(C, C')] = \frac{1}{n_1 n_2} \sum_{i=1}^{n_1} \sum_{j=1}^{n_2} \log \frac{n_{ij}}{n_i n_j}$$

USING THE EXPECTED VALUE THE ADJUSTED MUTUAL INFORMATION CAN THEN BE CALCULATED USING A SIMILAR FORM TO THAT OF THE ADJUSTED RAND INDEX

$$AMI(C, C') = \frac{MI(C, C') - E[MI(C, C')]}{H(C, C') - E[MI(C, C')]}$$

FOR NORMALIZED MUTUAL INFORMATION AND ADJUSTED MUTUAL INFORMATION THE NORMALIZING VALUE IS TYPICALLY SOME GENERALIZED MEAN OF THE ENTROPIES OF EACH CLUSTERING VARIOUS GENERALIZED MEANS EXIST AND NO FIRM RULES EXIST FOR PREFERRING ONE OVER THE OTHERS THE DECISION IS LARGELY A FIELD-BY-FIELD BASIS FOR INSTANCE IN COMMUNITY DETECTION THE ARITHMETIC MEAN IS MOST COMMON EACH NORMALIZING METHOD PROVIDES “QUALITATIVELY SIMILAR BEHAVIOURS” YAT2016 IN OUR IMPLEMENTATION THIS IS CONTROLLED BY THE AVERAGEMETHOD PARAMETER

VINH ET AL 2010 NAMED VARIANTS OF NMI AND AMI BY THEIR AVERAGING METHOD VEB2010 THEIR ‘SQRT’ AND ‘SUM’ AVERAGES ARE THE GEOMETRIC AND ARITHMETIC MEANS WE USE THESE MORE BROADLY COMMON NAMES

SCIKITLEARN USER GUIDE RELEASE 0213

REFERENCES

- STREHL ALEXANDER AND JOYDEEP GHOSH 2002 “CLUSTER ENSEMBLES – A KNOWLEDGE REUSE FRAME  
WORK FOR COMBINING MULTIPLE PARTITIONS” JOURNAL OF MACHINE LEARNING RESEARCH 3 583–617  
DOI101162153244303321897735

- WIKIPEDIA ENTRY FOR THE NORMALIZED MUTUAL INFORMATION
  - WIKIPEDIA ENTRY FOR THE ADJUSTED MUTUAL INFORMATION
- HOMOGENEITY COMPLETENESS AND VMEASURE

GIVEN THE KNOWLEDGE OF THE GROUND TRUTH CLASS ASSIGNMENTS OF THE SAMPLES IT IS POSSIBLE TO DEFINE SOME INTUITIVE METRIC  
USING CONDITIONAL ENTROPY ANALYSIS

IN PARTICULAR ROSENBERG AND HIRSCHBERG 2007 DEFINE THE FOLLOWING TWO DESIRABLE OBJECTIVES FOR ANY CLUSTER ASSIGN  
MENT

- HOMOGENEITY EACH CLUSTER CONTAINS ONLY MEMBERS OF A SINGLE CLASS
- COMPLETENESS ALL MEMBERS OF A GIVEN CLASS ARE ASSIGNED TO THE SAME CLUSTER

WE CAN TURN THOSE CONCEPT AS SCORES HOMOGENEITYSCORE ANDCOMPLETENESSSCORE BOTH ARE BOUNDED  
BELOW BY 00 AND ABOVE BY 10 HIGHER IS BETTER

FROM SKLEARN IMPORT METRICS

LABELSTRUE 0 0 0 1 1 1

LABELSPRED 0 0 1 1 2 2

METRICSHOMOGENEITYSCORELABELSTRUE LABELSPRED

066

METRICSCOMPLETENESSSCORELABELSTRUE LABELSPRED

042

THEIR HARMONIC MEAN CALLED VMEASURE IS COMPUTED BY VMEASURESCORE

METRICSVMEASURESCORELABELSTRUE LABELSPRED

051

THIS FUNCTION’S FORMULA IS AS FOLLOWS

MATH  $V = \frac{1}{\beta \times \text{TEXTHOMOGENEITY} + (1 - \beta) \times \text{TEXTCOMPLETENESS}}$

↪ $\beta \times \text{TEXTHOMOGENEITY} + (1 - \beta) \times \text{TEXTCOMPLETENESS}$

BETA DEFAULTS TO A VALUE OF 10 BUT FOR USING A VALUE LESS THAN 1 FOR BETA

METRICSVMEASURESCORELABELSTRUE LABELSPRED BETA06

054

MORE WEIGHT WILL BE ATTRIBUTED TO HOMOGENEITY AND USING A VALUE GREATER THAN 1

METRICSVMEASURESCORELABELSTRUE LABELSPRED BETA18

048

MORE WEIGHT WILL BE ATTRIBUTED TO COMPLETENESS

32 UNSUPERVISED LEARNING 377

SCIKITLEARN USER GUIDE RELEASE 0213  
THE VMEASURE IS ACTUALLY EQUIVALENT TO THE MUTUAL INFORMATION NMI DISCUSSED ABOVE WITH THE AGGREGATION FUNCTION BEING THE ARITHMETIC MEAN B2011  
HOMOGENEITY COMPLETENESS AND VMEASURE CAN BE COMPUTED AT ONCE USING  
HOMOGENEITYCOMPLETENESSVMEASURE AS FOLLOWS  
METRICSHOMOGENEITYCOMPLETENESSVMEASURELABELSTRUE LABELSPRED

066 042 051  
THE FOLLOWING CLUSTERING ASSIGNMENT IS SLIGHTLY BETTER SINCE IT IS HOMOGENEOUS BUT NOT COMPLETE  
LABELSPRED 0 0 0 1 2 2  
METRICSHOMOGENEITYCOMPLETENESSVMEASURELABELSTRUE LABELSPRED

10 068 081  
NOTEVMEASURESCORE ISSYMMETRIC IT CAN BE USED TO EVALUATE THE AGREEMENT OF TWO INDEPENDENT ASSIGNMENTS ON THE SAME DATASET  
THIS IS NOT THE CASE FOR COMPLETENESSSCORE ANDHOMOGENEITYSCORE BOTH ARE BOUND BY THE RELATIONSHIP  
HOMOGENEITYSCOREA B COMPLETENESSSCOREB A  
ADVANTAGES  
•BOUNDED SCORES 00 IS AS BAD AS IT CAN BE 10 IS A PERFECT SCORE  
• INTUITIVE INTERPRETATION CLUSTERING WITH BAD VMEASURE CAN BE QUALITATIVELY ANALYZED IN TERMS OF HOMOGENEITY AND COMPLETENESS TO BETTER FEEL WHAT ‘KIND’ OF MISTAKES IS DONE BY THE ASSIGNMENT  
•NO ASSUMPTION IS MADE ON THE CLUSTER STRUCTURE CAN BE USED TO COMPARE CLUSTERING ALGORITHMS SUCH AS K MEANS WHICH ASSUMES ISOTROPIC BLOB SHAPES WITH RESULTS OF SPECTRAL CLUSTERING ALGORITHMS WHICH CAN FIND CLUSTER WITH “FOLDED” SHAPES  
DRAWBACKS  
• THE PREVIOUSLY INTRODUCED METRICS ARE NOT NORMALIZED WITH REGARDS TO RANDOM LABELING THIS MEANS THAT DEPENDING ON THE NUMBER OF SAMPLES CLUSTERS AND GROUND TRUTH CLASSES A COMPLETELY RANDOM LABELING WILL NOT ALWAYS YIELD THE SAME VALUES FOR HOMOGENEITY COMPLETENESS AND HENCE VMEASURE IN PARTICULAR RANDOM LABELING WON’T YIELD ZERO SCORES ESPECIALLY WHEN THE NUMBER OF CLUSTERS IS LARGE  
THIS PROBLEM CAN SAFELY BE IGNORED WHEN THE NUMBER OF SAMPLES IS MORE THAN A THOUSAND AND THE NUMBER OF CLUSTERS IS LESS THAN 10 FOR SMALLER SAMPLE SIZES OR LARGER NUMBER OF CLUSTERS IT IS SAFER TO USE AN ADJUSTED INDEX SUCH AS THE ADJUSTED RAND INDEX ARI  
• THESE METRICS REQUIRE THE KNOWLEDGE OF THE GROUND TRUTH CLASSES WHILE ALMOST NEVER AVAILABLE IN PRACTICE OR REQUIRES MANUAL ASSIGNMENT BY HUMAN ANNOTATORS AS IN THE SUPERVISED LEARNING SETTING

EXAMPLES  
•ADJUSTMENT FOR CHANCE IN CLUSTERING PERFORMANCE EVALUATION ANALYSIS OF THE IMPACT OF THE DATASET SIZE ON THE VALUE OF CLUSTERING MEASURES FOR RANDOM ASSIGNMENTS  
378 CHAPTER 3 USER GUIDE





SCIKITLEARN USER GUIDE RELEASE 0213

MATHEMATICAL FORMULATION

HOMOGENEITY AND COMPLETENESS SCORES ARE FORMALLY GIVEN BY

$$h_1 = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^K \frac{1}{|C_j|} \log \frac{|C_j|}{n}$$

$$h_2 = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^K \frac{1}{|C_j|} \log \frac{|C_j|}{n}$$

$$h_3 = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^K \frac{1}{|C_j|} \log \frac{|C_j|}{n}$$

WHERE  $h_1$  IS THE CONDITIONAL ENTROPY OF THE CLASSES GIVEN THE CLUSTER ASSIGNMENTS AND IS GIVEN BY

$$h_1 = -\sum_{i=1}^n \sum_{j=1}^K \frac{1}{|C_j|} \log \frac{|C_j|}{n}$$

$$h_2 = -\sum_{i=1}^n \sum_{j=1}^K \frac{1}{|C_j|} \log \frac{|C_j|}{n}$$

$$h_3 = -\sum_{i=1}^n \sum_{j=1}^K \frac{1}{|C_j|} \log \frac{|C_j|}{n}$$

$$h_4 = -\sum_{i=1}^n \sum_{j=1}^K \frac{1}{|C_j|} \log \frac{|C_j|}{n}$$

$$h_5 = -\sum_{i=1}^n \sum_{j=1}^K \frac{1}{|C_j|} \log \frac{|C_j|}{n}$$

AND  $h_6$  IS THE ENTROPY OF THE CLASSES AND IS GIVEN BY

$$h_6 = -\sum_{i=1}^n \sum_{j=1}^K \frac{1}{|C_j|} \log \frac{|C_j|}{n}$$

$$h_7 = -\sum_{i=1}^n \sum_{j=1}^K \frac{1}{|C_j|} \log \frac{|C_j|}{n}$$

$$h_8 = -\sum_{i=1}^n \sum_{j=1}^K \frac{1}{|C_j|} \log \frac{|C_j|}{n}$$

$$h_9 = -\sum_{i=1}^n \sum_{j=1}^K \frac{1}{|C_j|} \log \frac{|C_j|}{n}$$

WITH  $n$  THE TOTAL NUMBER OF SAMPLES AND  $n_j$  THE NUMBER OF SAMPLES RESPECTIVELY BELONGING TO CLASS  $j$  AND CLUSTER  $i$  AND FINALLY  $n_{ij}$  THE NUMBER OF SAMPLES FROM CLASS  $j$  ASSIGNED TO CLUSTER  $i$

THE CONDITIONAL ENTROPY OF CLUSTERS GIVEN CLASS  $j$  AND THE ENTROPY OF CLUSTERS  $i$  ARE DEFINED IN A SYM

METRIC MANNER

ROSENBERG AND HIRSCHBERG FURTHER DEFINE VMEASURE AS THE HARMONIC MEAN OF HOMOGENEITY AND COMPLETENESS

$$V = \frac{2 \cdot h_1 \cdot h_2}{h_1 + h_2}$$

$$h_1$$

REFERENCES

• VMEASURE A CONDITIONAL ENTROPYBASED EXTERNAL CLUSTER EVALUATION MEASURE ANDREW ROSENBERG AND JULIA

HIRSCHBERG 2007

FOWLKESMALLOWS SCORES

THE FOWLKESMALLOWS INDEX SKLEARNMETRICSFOWLKESMALLOWSSCORE CAN BE USED WHEN THE GROUND

TRUTH CLASS ASSIGNMENTS OF THE SAMPLES IS KNOWN THE FOWLKESMALLOWS SCORE FMI IS DEFINED AS THE GEOMETRIC MEAN

OF THE PAIRWISE PRECISION AND RECALL

FMI TP /

TP FP TP FN

WHERE TP IS THE NUMBER OF TRUE POSITIVE IE THE NUMBER OF PAIR OF POINTS THAT BELONG TO THE SAME CLUSTERS IN BOTH THE

TRUE LABELS AND THE PREDICTED LABELS FP IS THE NUMBER OF FALSE POSITIVE IE THE NUMBER OF PAIR OF POINTS THAT BELONG

TO THE SAME CLUSTERS IN THE TRUE LABELS AND NOT IN THE PREDICTED LABELS AND FN IS THE NUMBER OF FALSE NEGATIVE IE THE

NUMBER OF PAIR OF POINTS THAT BELONGS IN THE SAME CLUSTERS IN THE PREDICTED LABELS AND NOT IN THE TRUE LABELS

THE SCORE RANGES FROM 0 TO 1 A HIGH VALUE INDICATES A GOOD SIMILARITY BETWEEN TWO CLUSTERS

FROM SKLEARN IMPORT METRICS

LABELSTRUE 0 0 0 1 1 1

LABELSPRED 0 0 1 1 2 2

380 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

METRICSFOWLKESMALLOWSSCORELABELSTRUE LABELSPRED  
047140

ONE CAN PERMUTE 0 AND 1 IN THE PREDICTED LABELS RENAME 2 TO 3 AND GET THE SAME SCORE

LABELSPRED 1 1 0 0 3 3

METRICSFOWLKESMALLOWSSCORELABELSTRUE LABELSPRED  
047140

PERFECT LABELING IS SCORED 10

LABELSPRED LABELSTRUE

METRICSFOWLKESMALLOWSSCORELABELSTRUE LABELSPRED  
10

BAD EG INDEPENDENT LABELINGS HAVE ZERO SCORES

LABELSTRUE 0 1 2 0 3 4 5 1

LABELSPRED 1 1 0 0 2 2 2 2

METRICSFOWLKESMALLOWSSCORELABELSTRUE LABELSPRED  
00

ADVANTAGES

- RANDOM UNIFORM LABEL ASSIGNMENTS HAVE A FMI SCORE CLOSE TO 00 FOR ANY VALUE OF NCLUSTERS AND NSAMPLES WHICH IS NOT THE CASE FOR RAW MUTUAL INFORMATION OR THE VMEASURE FOR INSTANCE
- UPPERBOUNDED AT 1 VALUES CLOSE TO ZERO INDICATE TWO LABEL ASSIGNMENTS THAT ARE LARGELY INDEPENDENT WHILE VALUES CLOSE TO ONE INDICATE SIGNIFICANT AGREEMENT FURTHER VALUES OF EXACTLY 0 INDICATE PURELY INDEPENDENT LABEL ASSIGNMENTS AND A FMI OF EXACTLY 1 INDICATES THAT THE TWO LABEL ASSIGNMENTS ARE EQUAL WITH OR WITHOUT PERMUTATION
- NO ASSUMPTION IS MADE ON THE CLUSTER STRUCTURE CAN BE USED TO COMPARE CLUSTERING ALGORITHMS SUCH AS K MEANS WHICH ASSUMES ISOTROPIC BLOB SHAPES WITH RESULTS OF SPECTRAL CLUSTERING ALGORITHMS WHICH CAN FIND CLUSTER WITH “FOLDED” SHAPES

DRAWBACKS

- CONTRARY TO INERTIA FMIBASED MEASURES REQUIRE THE KNOWLEDGE OF THE GROUND TRUTH CLASSES WHILE ALMOST NEVER AVAILABLE IN PRACTICE OR REQUIRES MANUAL ASSIGNMENT BY HUMAN ANNOTATORS AS IN THE SUPERVISED LEARNING SETTING

REFERENCES

- E B FOWLKES AND C L MALLOWES 1983 “A METHOD FOR COMPARING TWO HIERARCHICAL CLUSTERINGS” JOURNAL OF THE AMERICAN STATISTICAL ASSOCIATION [HTTPWILDFIRESTATUCLAEDUPDFLIBRARYFOWLKESPDF](http://wildfire.statucla.edu/pdf/library/fowlkes.pdf)
- WIKIPEDIA ENTRY FOR THE FOWLKESMALLOWS INDEX

SCIKITLEARN USER GUIDE RELEASE 0213

SILHOUETTE COEFFICIENT

IF THE GROUND TRUTH LABELS ARE NOT KNOWN EVALUATION MUST BE PERFORMED USING THE MODEL ITSELF THE SILHOUETTE COEFFICIENT SKLEARNMETRICSSILHOUETTESCORE IS AN EXAMPLE OF SUCH AN EVALUATION WHERE A HIGHER SILHOUETTE COEFFICIENT SCORE RELATES TO A MODEL WITH BETTER DEFINED CLUSTERS THE SILHOUETTE COEFFICIENT IS DEFINED FOR EACH SAMPLE AND IS COMPOSED OF TWO SCORES

- A THE MEAN DISTANCE BETWEEN A SAMPLE AND ALL OTHER POINTS IN THE SAME CLASS
- B THE MEAN DISTANCE BETWEEN A SAMPLE AND ALL OTHER POINTS IN THE NEXT NEAREST CLUSTER

THE SILHOUETTE COEFFICIENT SFOR A SINGLE SAMPLE IS THEN GIVEN AS

$$s_i = \frac{b_i - a_i}{b_i}$$

THE SILHOUETTE COEFFICIENT FOR A SET OF SAMPLES IS GIVEN AS THE MEAN OF THE SILHOUETTE COEFFICIENT FOR EACH SAMPLE

```
FROM SKLEARN IMPORT METRICS
FROM SKLEARNMETRICS IMPORT PAIRWISEDISTANCES
FROM SKLEARN IMPORT DATASETS
DATASET DATASETSLOADIRIS
X DATASETDATA
Y DATASETTARGET
IN NORMAL USAGE THE SILHOUETTE COEFFICIENT IS APPLIED TO THE RESULTS OF A CLUSTER ANALYSIS
IMPORT NUMPY AS NP
FROM SKLEARNCLUSTER IMPORT KMEANS
KMEANSMODEL KMEANSNCLUSTERS3 RANDOMSTATE1FITX
LABELS KMEANSMODELLABELS
METRICSSILHOUETTESCOREX LABELS METRICEUCLIDEAN
```

055

REFERENCES

- PETER J ROUSSEEUW 1987 “SILHOUETTES A GRAPHICAL AID TO THE INTERPRETATION AND VALIDATION OF CLUSTER ANALYSIS” COMPUTATIONAL AND APPLIED MATHEMATICS 20 53-65 DOI1010160377042787901257

ADVANTAGES

- THE SCORE IS BOUNDED BETWEEN 1 FOR INCORRECT CLUSTERING AND 1 FOR HIGHLY DENSE CLUSTERING SCORES AROUND ZERO INDICATE OVERLAPPING CLUSTERS
- THE SCORE IS HIGHER WHEN CLUSTERS ARE DENSE AND WELL SEPARATED WHICH RELATES TO A STANDARD CONCEPT OF A CLUSTER

DRAWBACKS

- THE SILHOUETTE COEFFICIENT IS GENERALLY HIGHER FOR CONVEX CLUSTERS THAN OTHER CONCEPTS OF CLUSTERS SUCH AS DENSITY BASED CLUSTERS LIKE THOSE OBTAINED THROUGH DBSCAN

382 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES

•SELECTING THE NUMBER OF CLUSTERS WITH SILHOUETTE ANALYSIS ON KMEANS CLUSTERING IN THIS EXAMPLE THE SILHOUETTE ANALYSIS IS USED TO CHOOSE AN OPTIMAL VALUE FOR NCLUSTERS

CALINSKIHARABASZ INDEX

IF THE GROUND TRUTH LABELS ARE NOT KNOWN THE CALINSKIHARABASZ INDEX SKLEARNMETRICS

CALINSKIHARABASZSCORE ALSO KNOWN AS THE VARIANCE RATIO CRITERION CAN BE USED TO EVALUATE THE

MODEL WHERE A HIGHER CALINSKIHARABASZ SCORE RELATES TO A MODEL WITH BETTER DEFINED CLUSTERS

FORCLUSTERS THE CALINSKIHARABASZ SCORE IS GIVEN AS THE RATIO OF THE BETWEENCLUSTERS DISPERSION MEAN AND THE WITHINCLUSTER DISPERSION

TR

TR×-

-1

WHERE IS THE BETWEEN GROUP DISPERSION MATRIX AND IS THE WITHINCLUSTER DISPERSION MATRIX DEFINED BY

Σ

1Σ

∈--

Σ

--

WITHBE THE NUMBER OF POINTS IN OUR DATA BE THE SET OF POINTS IN CLUSTER BE THE CENTER OF CLUSTER BE THE

CENTER OF BE THE NUMBER OF POINTS IN CLUSTER

FROM SKLEARN IMPORT METRICS

FROM SKLEARNMETRICS IMPORT PAIRWISEDISTANCES

FROM SKLEARN IMPORT DATASETS

DATASET DATASETSLOADIRIS

X DATASETDATA

Y DATASETTARGET

IN NORMAL USAGE THE CALINSKIHARABASZ INDEX IS APPLIED TO THE RESULTS OF A CLUSTER ANALYSIS

IMPORT NUMPY AS NP

FROM SKLEARNCLUSTER IMPORT KMEANS

KMEANSMODEL KMEANSNCLUSTERS3 RANDOMSTATE1FITX

LABELS KMEANSMODELLABELS

METRICSCALINSKIHARABASZSCOREX LABELS

56162

ADVANTAGES

- THE SCORE IS HIGHER WHEN CLUSTERS ARE DENSE AND WELL SEPARATED WHICH RELATES TO A STANDARD CONCEPT OF A CLUSTER
- THE SCORE IS FAST TO COMPUTE

SCIKITLEARN USER GUIDE RELEASE 0213

DRAWBACKS

- THE CALINSKI-HARABASZ INDEX IS GENERALLY HIGHER FOR CONVEX CLUSTERS THAN OTHER CONCEPTS OF CLUSTERS SUCH AS DENSITY BASED CLUSTERS LIKE THOSE OBTAINED THROUGH DBSCAN

REFERENCES

- CALINSKI T HARABASZ J 1974 "A DENDRITE METHOD FOR CLUSTER ANALYSIS" COMMUNICATIONS IN STATISTICS THEORY AND METHODS 3 127 DOI10.1080/036109262011560741

DAVIESBOULDIN INDEX

IF THE GROUND TRUTH LABELS ARE NOT KNOWN THE DAVIESBOULDIN INDEX SKLEARNMETRICS

DAVIESBOULDINSCORE CAN BE USED TO EVALUATE THE MODEL WHERE A LOWER DAVIESBOULDIN INDEX RELATES TO A MODEL WITH BETTER SEPARATION BETWEEN THE CLUSTERS

THE INDEX IS DEFINED AS THE AVERAGE SIMILARITY BETWEEN EACH CLUSTER  $\mu_i$  AND ITS MOST SIMILAR ONE  $\mu_j$  IN THE CONTEXT OF THIS INDEX SIMILARITY IS DEFINED AS A MEASURE THAT TRADES OFF

- $\frac{1}{n_i}$  THE AVERAGE DISTANCE BETWEEN EACH POINT OF CLUSTER  $\mu_i$  AND THE CENTROID OF THAT CLUSTER - ALSO KNOWN AS CLUSTER DIAMETER

- $d(\mu_i, \mu_j)$  THE DISTANCE BETWEEN CLUSTER CENTROIDS  $\mu_i$  AND  $\mu_j$

A SIMPLE CHOICE TO CONSTRUCT  $d(\mu_i, \mu_j)$  SO THAT IT IS NONNEGATIVE AND SYMMETRIC IS

$$d(\mu_i, \mu_j) = \frac{1}{2} \left( \frac{1}{n_i} + \frac{1}{n_j} \right) d(\mu_i, \mu_j)$$

THEN THE DAVIESBOULDIN INDEX IS DEFINED AS

$$DBI = \frac{1}{n} \sum_{i=1}^K \left( \frac{1}{n_i} + \frac{1}{n_j} \right) d(\mu_i, \mu_j)$$

$$DBI = \frac{1}{n} \sum_{i=1}^K \left( \frac{1}{n_i} + \frac{1}{n_j} \right) d(\mu_i, \mu_j)$$

$$DBI = \frac{1}{n} \sum_{i=1}^K \left( \frac{1}{n_i} + \frac{1}{n_j} \right) d(\mu_i, \mu_j)$$

$$DBI = \frac{1}{n} \sum_{i=1}^K \left( \frac{1}{n_i} + \frac{1}{n_j} \right) d(\mu_i, \mu_j)$$

ZERO IS THE LOWEST POSSIBLE SCORE VALUES CLOSER TO ZERO INDICATE A BETTER PARTITION

IN NORMAL USAGE THE DAVIESBOULDIN INDEX IS APPLIED TO THE RESULTS OF A CLUSTER ANALYSIS AS FOLLOWS

FROM SKLEARN IMPORT DATASETS

IRIS = DATASETSLOADIRIS

X = IRISDATA

FROM SKLEARNCLUSTER IMPORT KMEANS

FROM SKLEARNMETRICS IMPORT DAVIESBOULDINSCORE

KMEANS = KMEANSNCLUSTERS3 RANDOMSTATE1FITX

LABELS = KMEANSLABELS

DAVIESBOULDINSCOREX LABELS

06619

ADVANTAGES

- THE COMPUTATION OF DAVIESBOULDIN IS SIMPLER THAN THAT OF SILHOUETTE SCORES
- THE INDEX IS COMPUTED ONLY QUANTITIES AND FEATURES INHERENT TO THE DATASET

SCIKITLEARN USER GUIDE RELEASE 0213

DRAWBACKS

- THE DAVIESBOULDIN INDEX IS GENERALLY HIGHER FOR CONVEX CLUSTERS THAN OTHER CONCEPTS OF CLUSTERS SUCH AS DENSITY BASED CLUSTERS LIKE THOSE OBTAINED FROM DBSCAN
- THE USAGE OF CENTROID DISTANCE LIMITS THE DISTANCE METRIC TO EUCLIDEAN SPACE
- A GOOD VALUE REPORTED BY THIS METHOD DOES NOT IMPLY THE BEST INFORMATION RETRIEVAL

REFERENCES

- DAVIES DAVID L BOULDIN DONALD W 1979 "A CLUSTER SEPARATION MEASURE" IEEE TRANSACTIONS ON PATTERN ANALYSIS AND MACHINE INTELLIGENCE PAMI1 2 224227 DOI101109TPAMI19794766909
- HALKIDI MARIA BATISTAKIS YANNIS VAZIRGIANNIS MICHALIS 2001 "ON CLUSTERING VALIDATION TECHNIQUES" JOURNAL OF INTELLIGENT INFORMATION SYSTEMS 1723 107145 DOI101023A1012801612483
- WIKIPEDIA ENTRY FOR DAVIESBOULDIN INDEX

CONTINGENCY MATRIX

CONTINGENCY MATRIX SKLEARNMETRICSCUSTERCONTINGENCYMATRIX REPORTS THE INTERSECTION CARDI NALITY FOR EVERY TRUEPREDICTED CLUSTER PAIR THE CONTINGENCY MATRIX PROVIDES SUFFICIENT STATISTICS FOR ALL CLUSTERING MET RICS WHERE THE SAMPLES ARE INDEPENDENT AND IDENTICALLY DISTRIBUTED AND ONE DOESN'T NEED TO ACCOUNT FOR SOME INSTANCES NOT BEING CLUSTERED

HERE IS AN EXAMPLE

FROM SKLEARNMETRICSCUSTER IMPORT CONTINGENCYMATRIX

X A A A B B B

Y 0 0 1 1 2 2

CONTINGENCYMATRIX X Y

ARRAY2 1 0

0 1 2

THE FIRST ROW OF OUTPUT ARRAY INDICATES THAT THERE ARE THREE SAMPLES WHOSE TRUE CLUSTER IS "A" OF THEM TWO ARE IN PREDICTED CLUSTER 0 ONE IS IN 1 AND NONE IS IN 2 AND THE SECOND ROW INDICATES THAT THERE ARE THREE SAMPLES WHOSE TRUE CLUSTER IS "B" OF THEM NONE IS IN PREDICTED CLUSTER 0 ONE IS IN 1 AND TWO ARE IN 2 ACONFUSION MATRIX FOR CLASSIFICATION IS A SQUARE CONTINGENCY MATRIX WHERE THE ORDER OF ROWS AND COLUMNS CORRESPOND TO A LIST OF CLASSES

ADVANTAGES

- ALLOWS TO EXAMINE THE SPREAD OF EACH TRUE CLUSTER ACROSS PREDICTED CLUSTERS AND VICE VERSA
- THE CONTINGENCY TABLE CALCULATED IS TYPICALLY UTILIZED IN THE CALCULATION OF A SIMILARITY STATISTIC LIKE THE OTHERS LISTED IN THIS DOCUMENT BETWEEN THE TWO CLUSTERINGS

DRAWBACKS

- CONTINGENCY MATRIX IS EASY TO INTERPRET FOR A SMALL NUMBER OF CLUSTERS BUT BECOMES VERY HARD TO INTERPRET FOR A LARGE NUMBER OF CLUSTERS
- 32 UNSUPERVISED LEARNING 385

- IT DOESN'T GIVE A SINGLE METRIC TO USE AS AN OBJECTIVE FOR CLUSTERING OPTIMISATION

REFERENCES

- WIKIPEDIA ENTRY FOR CONTINGENCY MATRIX

324 BICLUSTERING

BICLUSTERING CAN BE PERFORMED WITH THE MODULE SKLEARNCLUSTERBICLUSTER BICLUSTERING ALGORITHMS SIMULTANEOUSLY CLUSTER ROWS AND COLUMNS OF A DATA MATRIX THESE CLUSTERS OF ROWS AND COLUMNS ARE KNOWN AS BICLUSTERS EACH DETERMINES A SUBMATRIX OF THE ORIGINAL DATA MATRIX WITH SOME DESIRED PROPERTIES

FOR INSTANCE GIVEN A MATRIX OF SHAPE 10 10 ONE POSSIBLE BICLUSTER WITH THREE ROWS AND TWO COLUMNS INDUCES A SUBMATRIX OF SHAPE 3 2

```
import numpy as np
data = np.arange(100).reshape(10, 10)
rows = np.array(0, 2, 3)
newaxis = np.newaxis
columns = np.array(1, 2)
data_rows_columns = data[rows, columns]
array([[1, 2],
       [3, 4],
       [5, 6]])
```

FOR VISUALIZATION PURPOSES GIVEN A BICLUSTER THE ROWS AND COLUMNS OF THE DATA MATRIX MAY BE REARRANGED TO MAKE THE BICLUSTER CONTIGUOUS

ALGORITHMS DIFFER IN HOW THEY DEFINE BICLUSTERS SOME OF THE COMMON TYPES INCLUDE

- CONSTANT VALUES CONSTANT ROWS OR CONSTANT COLUMNS
- UNUSUALLY HIGH OR LOW VALUES
- SUBMATRICES WITH LOW VARIANCE
- CORRELATED ROWS OR COLUMNS

ALGORITHMS ALSO DIFFER IN HOW ROWS AND COLUMNS MAY BE ASSIGNED TO BICLUSTERS WHICH LEADS TO DIFFERENT BICLUSTER STRUCTURES BLOCK DIAGONAL OR CHECKERBOARD STRUCTURES OCCUR WHEN ROWS AND COLUMNS ARE DIVIDED INTO PARTITIONS IF EACH ROW AND EACH COLUMN BELONGS TO EXACTLY ONE BICLUSTER THEN REARRANGING THE ROWS AND COLUMNS OF THE DATA MATRIX REVEALS THE BICLUSTERS ON THE DIAGONAL HERE IS AN EXAMPLE OF THIS STRUCTURE WHERE BICLUSTERS HAVE HIGHER AVERAGE VALUES THAN THE OTHER ROWS AND COLUMNS

IN THE CHECKERBOARD CASE EACH ROW BELONGS TO ALL COLUMN CLUSTERS AND EACH COLUMN BELONGS TO ALL ROW CLUSTERS HERE IS AN EXAMPLE OF THIS STRUCTURE WHERE THE VARIANCE OF THE VALUES WITHIN EACH BICLUSTER IS SMALL

AFTER FITTING A MODEL ROW AND COLUMN CLUSTER MEMBERSHIP CAN BE FOUND IN THE ROWS AND COLUMNS ATTRIBUTES ROWSI IS A BINARY VECTOR WITH NONZERO ENTRIES CORRESPONDING TO ROWS THAT BELONG TO BICLUSTER I SIMILARLY COLUMNSI INDICATES WHICH COLUMNS BELONG TO BICLUSTER I

SOME MODELS ALSO HAVE ROWLABELS AND COLUMNLABELS ATTRIBUTES THESE MODELS PARTITION THE ROWS AND COLUMNS SUCH AS IN THE BLOCK DIAGONAL AND CHECKERBOARD BICLUSTER STRUCTURES

NOTE BICLUSTERING HAS MANY OTHER NAMES IN DIFFERENT FIELDS INCLUDING COCLUSTERING TWOMODE CLUSTERING TWOWAY CLUSTERING BLOCK CLUSTERING COUPLED TWOWAY CLUSTERING ETC THE NAMES OF SOME ALGORITHMS SUCH AS THE SPECTRAL COCLUSTERING ALGORITHM REFLECT THESE ALTERNATE NAMES



SCIKITLEARN USER GUIDE RELEASE 0213  
FIG 35 AN EXAMPLE OF BICLUSTERS FORMED BY PARTITIONING ROWS AND COLUMNS  
FIG 36 AN EXAMPLE OF CHECKERBOARD BICLUSTERS  
32 UNSUPERVISED LEARNING 387

SPECTRAL COCLUSTERING

THE SPECTRAL COCLUSTERING ALGORITHM FINDS BICLUSTERS WITH VALUES HIGHER THAN THOSE IN THE CORRESPONDING OTHER ROWS AND COLUMNS EACH ROW AND EACH COLUMN BELONGS TO EXACTLY ONE BICLUSTER SO REARRANGING THE ROWS AND COLUMNS TO MAKE PARTITIONS CONTIGUOUS REVEALS THESE HIGH VALUES ALONG THE DIAGONAL

NOTE THE ALGORITHM TREATS THE INPUT DATA MATRIX AS A BIPARTITE GRAPH THE ROWS AND COLUMNS OF THE MATRIX CORRESPOND TO THE TWO SETS OF VERTICES AND EACH ENTRY CORRESPONDS TO AN EDGE BETWEEN A ROW AND A COLUMN THE ALGORITHM APPROXIMATES THE NORMALIZED CUT OF THIS GRAPH TO FIND HEAVY SUBGRAPHS

MATHEMATICAL FORMULATION

AN APPROXIMATE SOLUTION TO THE OPTIMAL NORMALIZED CUT MAY BE FOUND VIA THE GENERALIZED EIGENVALUE DECOMPOSITION OF THE LAPLACIAN OF THE GRAPH USUALLY THIS WOULD MEAN WORKING DIRECTLY WITH THE LAPLACIAN MATRIX IF THE ORIGINAL DATA MATRIX  $X$  HAS SHAPE  $n \times m$  THE LAPLACIAN MATRIX FOR THE CORRESPONDING BIPARTITE GRAPH HAS SHAPE  $(n+m) \times (n+m)$

HOWEVER IN THIS CASE IT IS POSSIBLE TO WORK DIRECTLY WITH  $X$  WHICH IS SMALLER AND MORE EFFICIENT

THE INPUT MATRIX  $X$  IS PREPROCESSED AS FOLLOWS

$$X' = \begin{bmatrix} X & -12X \\ -12X^T & -12X^T X \end{bmatrix}$$

WHERE  $I$  IS THE DIAGONAL MATRIX WITH ENTRY 1 EQUAL TO  $\sum_j X_{ij}^2$

$D$  IS THE DIAGONAL MATRIX WITH ENTRY 1 EQUAL TO  $\sum_i X_{ij}^2$

$\mathbf{1}$  IS A VECTOR OF ONES

THE SINGULAR VALUE DECOMPOSITION  $U \Sigma V^T$  PROVIDES THE PARTITIONS OF THE ROWS AND COLUMNS OF  $X'$  A SUBSET OF THE LEFT SINGULAR VECTORS GIVES THE ROW PARTITIONS AND A SUBSET OF THE RIGHT SINGULAR VECTORS GIVES THE COLUMN PARTITIONS THE  $\lceil \log_2 n \rceil$  SINGULAR VECTORS STARTING FROM THE SECOND PROVIDE THE DESIRED PARTITIONING INFORMATION THEY ARE USED TO FORM THE MATRIX  $Z$

$$Z = \begin{bmatrix} U & -12U \\ -12U^T & -12U^T U \end{bmatrix}$$

WHERE THE COLUMNS OF  $Z$  ARE  $\ell_2$  NORM 1 AND SIMILARLY FOR  $V$

THEN THE ROWS OF  $Z$  ARE CLUSTERED USING KMEANS THE FIRST  $n$  ROWS LABELS PROVIDE THE ROW PARTITIONING AND THE REMAINING  $n$  COLUMNS LABELS PROVIDE THE COLUMN PARTITIONING

EXAMPLES

- A DEMO OF THE SPECTRAL COCLUSTERING ALGORITHM A SIMPLE EXAMPLE SHOWING HOW TO GENERATE A DATA MATRIX WITH BICLUSTERS AND APPLY THIS METHOD TO IT
- BICLUSTERING DOCUMENTS WITH THE SPECTRAL COCLUSTERING ALGORITHM AN EXAMPLE OF FINDING BICLUSTERS IN THE TWENTY NEWSGROUP DATASET

REFERENCES

- DHILLON INDERJIT S 2001 COCLUSTERING DOCUMENTS AND WORDS USING BIPARTITE SPECTRAL GRAPH PARTITIONING

SPECTRAL BICLUSTERING

THE SPECTRAL BICLUSTERING ALGORITHM ASSUMES THAT THE INPUT DATA MATRIX HAS A HIDDEN CHECKERBOARD STRUCTURE. THE ROWS AND COLUMNS OF A MATRIX WITH THIS STRUCTURE MAY BE PARTITIONED SO THAT THE ENTRIES OF ANY BICLUSTER IN THE CARTESIAN PRODUCT OF ROW CLUSTERS AND COLUMN CLUSTERS ARE APPROXIMATELY CONSTANT. FOR INSTANCE, IF THERE ARE TWO ROW PARTITIONS AND THREE COLUMN PARTITIONS, EACH ROW WILL BELONG TO THREE BICLUSTERS AND EACH COLUMN WILL BELONG TO TWO BICLUSTERS.

THE ALGORITHM PARTITIONS THE ROWS AND COLUMNS OF A MATRIX SO THAT A CORRESPONDING BLOCKWISE CONSTANT CHECKERBOARD MATRIX PROVIDES A GOOD APPROXIMATION TO THE ORIGINAL MATRIX.

MATHEMATICAL FORMULATION

THE INPUT MATRIX  $X$  IS FIRST NORMALIZED TO MAKE THE CHECKERBOARD PATTERN MORE OBVIOUS. THERE ARE THREE POSSIBLE METHODS:

1. INDEPENDENT ROW AND COLUMN NORMALIZATION: AS IN SPECTRAL COCLUSTERING, THIS METHOD MAKES THE ROWS SUM TO A CONSTANT AND THE COLUMNS SUM TO A DIFFERENT CONSTANT.

2. BISTOCHASTIZATION: REPEATED ROW AND COLUMN NORMALIZATION UNTIL CONVERGENCE. THIS METHOD MAKES BOTH ROWS AND COLUMNS SUM TO THE SAME CONSTANT.

3. LOG NORMALIZATION: THE LOG OF THE DATA MATRIX IS COMPUTED  $\log X$ . THEN THE COLUMN MEAN  $\bar{x}_j$ , ROW MEAN  $\bar{y}_i$ , AND OVERALL MEAN  $\bar{\mu}$  OF  $\log X$  ARE COMPUTED. THE FINAL MATRIX IS COMPUTED ACCORDING TO THE FORMULA 
$$\frac{\log X - \bar{x}_j - \bar{y}_i + \bar{\mu}}{\sigma^2}$$

AFTER NORMALIZING, THE FIRST FEW SINGULAR VECTORS ARE COMPUTED, JUST AS IN THE SPECTRAL COCLUSTERING ALGORITHM. IF LOG NORMALIZATION WAS USED, ALL THE SINGULAR VECTORS ARE MEANINGFUL. HOWEVER, IF INDEPENDENT NORMALIZATION OR BISTOCHASTIZATION WERE USED, THE FIRST SINGULAR VECTORS  $v_1$  AND  $v_2$  ARE DISCARDED FROM NOW ON. THE “FIRST” SINGULAR VECTORS REFERS TO  $v_2$  AND  $v_1$  EXCEPT IN THE CASE OF LOG NORMALIZATION.

GIVEN THESE SINGULAR VECTORS, THEY ARE RANKED ACCORDING TO WHICH CAN BE BEST APPROXIMATED BY A PIECEWISE CONSTANT VECTOR. THE APPROXIMATIONS FOR EACH VECTOR ARE FOUND USING ONE-DIMENSIONAL KMEANS AND SCORED USING THE EUCLIDEAN DISTANCE. SOME SUBSET OF THE BEST LEFT AND RIGHT SINGULAR VECTORS ARE SELECTED. NEXT, THE DATA IS PROJECTED TO THIS BEST SUBSET OF SINGULAR VECTORS AND CLUSTERED.

FOR INSTANCE, IF  $U$  SINGULAR VECTORS WERE CALCULATED, THE  $U$  BEST ARE FOUND AS DESCRIBED, WHERE  $X$  LET  $X$  BE THE MATRIX WITH COLUMNS THE  $U$  BEST LEFT SINGULAR VECTORS AND SIMILARLY  $V$  FOR THE RIGHT. TO PARTITION THE ROWS, THE ROWS OF  $X$  ARE PROJECTED TO A  $U$ -DIMENSIONAL SPACE, TREATING THE  $U$  ROWS OF THIS  $U \times U$  MATRIX AS SAMPLES AND CLUSTERING USING KMEANS YIELDS THE ROW LABELS. SIMILARLY, PROJECTING THE COLUMNS TO  $V$  AND CLUSTERING THIS  $U \times U$  MATRIX YIELDS THE COLUMN LABELS.

EXAMPLES

• A DEMO OF THE SPECTRAL BICLUSTERING ALGORITHM: A SIMPLE EXAMPLE SHOWING HOW TO GENERATE A CHECKERBOARD MATRIX AND BICLUSTER IT.

REFERENCES

• KLUGER, YUVAL ET AL. 2003. SPECTRAL BICLUSTERING OF MICROARRAY DATA. COCLUSTERING GENES AND CONDITIONS. 32 UNSUPERVISED LEARNING. 389.

BICLUSTERING EVALUATION

THERE ARE TWO WAYS OF EVALUATING A BICLUSTERING RESULT INTERNAL AND EXTERNAL INTERNAL MEASURES SUCH AS CLUSTER STABILITY RELY ONLY ON THE DATA AND THE RESULT THEMSELVES CURRENTLY THERE ARE NO INTERNAL BICLUSTER MEASURES IN SCIKIT LEARN EXTERNAL MEASURES REFER TO AN EXTERNAL SOURCE OF INFORMATION SUCH AS THE TRUE SOLUTION WHEN WORKING WITH REAL DATA THE TRUE SOLUTION IS USUALLY UNKNOWN BUT BICLUSTERING ARTIFICIAL DATA MAY BE USEFUL FOR EVALUATING ALGORITHMS PRECISELY BECAUSE THE TRUE SOLUTION IS KNOWN

TO COMPARE A SET OF FOUND BICLUSTERS TO THE SET OF TRUE BICLUSTERS TWO SIMILARITY MEASURES ARE NEEDED A SIMILARITY MEASURE FOR INDIVIDUAL BICLUSTERS AND A WAY TO COMBINE THESE INDIVIDUAL SIMILARITIES INTO AN OVERALL SCORE TO COMPARE INDIVIDUAL BICLUSTERS SEVERAL MEASURES HAVE BEEN USED FOR NOW ONLY THE JACCARD INDEX IS IMPLEMENTED

$\frac{|A \cap B|}{|A \cup B|}$

$\frac{|A \cap B|}{|A \cap B|}$

WHERE  $A$  AND  $B$  ARE BICLUSTERS  $|A \cap B|$  IS THE NUMBER OF ELEMENTS IN THEIR INTERSECTION THE JACCARD INDEX ACHIEVES ITS MINIMUM OF 0 WHEN THE BICLUSTERS DO NOT OVERLAP AT ALL AND ITS MAXIMUM OF 1 WHEN THEY ARE IDENTICAL SEVERAL METHODS HAVE BEEN DEVELOPED TO COMPARE TWO SETS OF BICLUSTERS FOR NOW ONLY CONSENSUSSCORE HOCHREITER ET AL 2010 IS AVAILABLE

1 COMPUTE BICLUSTER SIMILARITIES FOR PAIRS OF BICLUSTERS ONE IN EACH SET USING THE JACCARD INDEX OR A SIMILAR MEASURE

2 ASSIGN BICLUSTERS FROM ONE SET TO ANOTHER IN A ONE-TO-ONE FASHION TO MAXIMIZE THE SUM OF THEIR SIMILARITIES THIS STEP IS PERFORMED USING THE HUNGARIAN ALGORITHM

3 THE FINAL SUM OF SIMILARITIES IS DIVIDED BY THE SIZE OF THE LARGER SET

THE MINIMUM CONSENSUS SCORE 0 OCCURS WHEN ALL PAIRS OF BICLUSTERS ARE TOTALLY DISSIMILAR THE MAXIMUM SCORE 1 OCCURS WHEN BOTH SETS ARE IDENTICAL

REFERENCES

- HOCHREITER BODENHOFER ET AL 2010 FABIA FACTOR ANALYSIS FOR BICLUSTER ACQUISITION

325 DECOMPOSING SIGNALS IN COMPONENTS MATRIX FACTORIZATION PROBLEMS

PRINCIPAL COMPONENT ANALYSIS PCA

EXACT PCA AND PROBABILISTIC INTERPRETATION

PCA IS USED TO DECOMPOSE A MULTIVARIATE DATASET IN A SET OF SUCCESSIVE ORTHOGONAL COMPONENTS THAT EXPLAIN A MAXIMUM AMOUNT OF THE VARIANCE IN SCIKITLEARN PCA IS IMPLEMENTED AS A TRANSFORMER OBJECT THAT LEARNS  $k$  COMPONENTS IN ITS FIT METHOD AND CAN BE USED ON NEW DATA TO PROJECT IT ON THESE COMPONENTS

PCA CENTERS BUT DOES NOT SCALE THE INPUT DATA FOR EACH FEATURE BEFORE APPLYING THE SVD THE OPTIONAL PARAMETER `PARAMETERWHITENTRUE` MAKES IT POSSIBLE TO PROJECT THE DATA ONTO THE SINGULAR SPACE WHILE SCALING EACH COMPONENT TO UNIT VARIANCE THIS IS OFTEN USEFUL IF THE MODELS DOWNSTREAM MAKE STRONG ASSUMPTIONS ON THE ISOTROPY OF THE SIGNAL THIS IS FOR EXAMPLE THE CASE FOR SUPPORT VECTOR MACHINES WITH THE RBF KERNEL AND THE KMEANS CLUSTERING ALGORITHM BELOW IS AN EXAMPLE OF THE IRIS DATASET WHICH IS COMPRISED OF 4 FEATURES PROJECTED ON THE 2 DIMENSIONS THAT EXPLAIN MOST VARIANCE

THE PCA OBJECT ALSO PROVIDES A PROBABILISTIC INTERPRETATION OF THE PCA THAT CAN GIVE A LIKELIHOOD OF DATA BASED ON THE AMOUNT OF VARIANCE IT EXPLAINS AS SUCH IT IMPLEMENTS A SCORE METHOD THAT CAN BE USED IN CROSSVALIDATION



EXAMPLES

- COMPARISON OF LDA AND PCA 2D PROJECTION OF IRIS DATASET
- MODEL SELECTION WITH PROBABILISTIC PCA AND FACTOR ANALYSIS FA

INCREMENTAL PCA

THEPCA OBJECT IS VERY USEFUL BUT HAS CERTAIN LIMITATIONS FOR LARGE DATASETS THE BIGGEST LIMITATION IS THAT PCA ONLY SUPPORTS BATCH PROCESSING WHICH MEANS ALL OF THE DATA TO BE PROCESSED MUST FIT IN MAIN MEMORY THE INCREMENTALPCA OBJECT USES A DIFFERENT FORM OF PROCESSING AND ALLOWS FOR PARTIAL COMPUTATIONS WHICH ALMOST EXACTLY MATCH THE RESULTS OFPCA WHILE PROCESSING THE DATA IN A MINIBATCH FASHION INCREMENTALPCA MAKES IT POSSIBLE TO IMPLEMENT OUTOF CORE PRINCIPAL COMPONENT ANALYSIS EITHER BY

- USING ITS PARTIALFIT METHOD ON CHUNKS OF DATA FETCHED SEQUENTIALLY FROM THE LOCAL HARD DRIVE OR A NETWORK DATABASE

- CALLING ITS FIT METHOD ON A MEMORY MAPPED FILE USING NUMPYMEMMAP

INCREMENTALPCA ONLY STORES ESTIMATES OF COMPONENT AND NOISE VARIANCES IN ORDER UPDATE EXPLAINEDVARIANCERATIO INCREMENTALLY THIS IS WHY MEMORY USAGE DEPENDS ON THE NUMBER OF SAMPLES PER BATCH RATHER THAN THE NUMBER OF SAMPLES TO BE PROCESSED IN THE DATASET

AS INPCAINCREMENTALPCA CENTERS BUT DOES NOT SCALE THE INPUT DATA FOR EACH FEATURE BEFORE APPLYING THE SVD

EXAMPLES

- INCREMENTAL PCA

PCA USING RANDOMIZED SVD

IT IS OFTEN INTERESTING TO PROJECT DATA TO A LOWERDIMENSIONAL SPACE THAT PRESERVES MOST OF THE VARIANCE BY DROPPING THE SINGULAR VECTOR OF COMPONENTS ASSOCIATED WITH LOWER SINGULAR VALUES

FOR INSTANCE IF WE WORK WITH 64X64 PIXEL GRAYLEVEL PICTURES FOR FACE RECOGNITION THE DIMENSIONALITY OF THE DATA IS 4096 AND IT IS SLOW TO TRAIN AN RBF SUPPORT VECTOR MACHINE ON SUCH WIDE DATA FURTHERMORE WE KNOW THAT THE INTRINSIC DIMENSIONALITY OF THE DATA IS MUCH LOWER THAN 4096 SINCE ALL PICTURES OF HUMAN FACES LOOK SOMEWHAT ALIKE THE SAMPLES LIE ON A MANIFOLD OF MUCH LOWER DIMENSION SAY AROUND 200 FOR INSTANCE THE PCA ALGORITHM CAN BE USED TO LINEARLY TRANSFORM THE DATA WHILE BOTH REDUCING THE DIMENSIONALITY AND PRESERVE MOST OF THE EXPLAINED VARIANCE AT THE SAME TIME

THE CLASSPCA USED WITH THE OPTIONAL PARAMETER SVDSOLVERRANDOMIZED IS VERY USEFUL IN THAT CASE SINCE WE ARE GOING TO DROP MOST OF THE SINGULAR VECTORS IT IS MUCH MORE EFFICIENT TO LIMIT THE COMPUTATION TO AN APPROXIMATED ESTIMATE OF THE SINGULAR VECTORS WE WILL KEEP TO ACTUALLY PERFORM THE TRANSFORM

FOR INSTANCE THE FOLLOWING SHOWS 16 SAMPLE PORTRAITS CENTERED AROUND 00 FROM THE OLIVETTI DATASET ON THE RIGHT HAND SIDE ARE THE FIRST 16 SINGULAR VECTORS RESHAPED AS PORTRAITS SINCE WE ONLY REQUIRE THE TOP 16 SINGULAR VECTORS OF A DATASET WITH SIZE 400 AND 64x64 4096 THE COMPUTATION TIME IS LESS THAN 1S







SCIKITLEARN USER GUIDE RELEASE 0213

IF WE NOTE  $\mathcal{O}(\text{MAX} \times \text{MAX} \times \text{SAMPLES} \times \text{FEATURES})$  AND  $\mathcal{O}(\text{MIN} \times \text{MIN} \times \text{SAMPLES} \times \text{FEATURES})$  THE TIME COMPLEXITY OF THE RANDOMIZEDPCA IS  $\mathcal{O}(\text{MAX}^2)$

$\text{MAX} \times \text{COMPONENTS}$  INSTEAD OF  $\mathcal{O}(\text{MAX}^2)$

$\text{MAX} \times \text{MIN}$  FOR THE EXACT METHOD IMPLEMENTED IN PCA

THE MEMORY FOOTPRINT OF RANDOMIZED PCA IS ALSO PROPORTIONAL TO  $2 \times \text{MAX} \times \text{COMPONENTS}$  INSTEAD OF  $\text{MAX} \times \text{MIN}$  FOR THE EXACT METHOD

NOTE THE IMPLEMENTATION OF INVERSE TRANSFORM IN PCA WITH SVDSOLVER RANDOMIZED IS NOT THE EXACT

INVERSE TRANSFORM OF TRANSFORM EVEN WHEN `WHITEN=False` DEFAULT

EXAMPLES

- FACES RECOGNITION EXAMPLE USING EIGENFACES AND SVMs
- FACES DATASET DECOMPOSITIONS

REFERENCES

- “FINDING STRUCTURE WITH RANDOMNESS STOCHASTIC ALGORITHMS FOR CONSTRUCTING APPROXIMATE MATRIX DECOMPOSITIONS” HALKO ET AL 2009

KERNEL PCA

KERNELPCA IS AN EXTENSION OF PCA WHICH ACHIEVES NONLINEAR DIMENSIONALITY REDUCTION THROUGH THE USE OF KERNELS SEE PAIRWISE METRICS AFFINITIES AND KERNELS IT HAS MANY APPLICATIONS INCLUDING DENOISING COMPRES SION AND STRUCTURED PREDICTION KERNEL DEPENDENCY ESTIMATION KERNELPCA SUPPORTS BOTH TRANSFORM AND INVERSETRANSFORM

EXAMPLES

- KERNEL PCA

SPARSE PRINCIPAL COMPONENTS ANALYSIS SPARSEPCA AND MINIBATCHSPARSEPCA

SPARSEPCA IS A VARIANT OF PCA WITH THE GOAL OF EXTRACTING THE SET OF SPARSE COMPONENTS THAT BEST RECONSTRUCT THE DATA

MINIBATCH SPARSE PCA MINIBATCHSPARSEPCA IS A VARIANT OF SPARSEPCA THAT IS FASTER BUT LESS ACCURATE THE INCREASED SPEED IS REACHED BY ITERATING OVER SMALL CHUNKS OF THE SET OF FEATURES FOR A GIVEN NUMBER OF ITERATIONS PRINCIPAL COMPONENT ANALYSIS PCA HAS THE DISADVANTAGE THAT THE COMPONENTS EXTRACTED BY THIS METHOD HAVE EXCLU SIVELY DENSE EXPRESSIONS IE THEY HAVE NONZERO COEFFICIENTS WHEN EXPRESSED AS LINEAR COMBINATIONS OF THE ORIGINAL VARIABLES THIS CAN MAKE INTERPRETATION DIFFICULT IN MANY CASES THE REAL UNDERLYING COMPONENTS CAN BE MORE NATURALLY IMAGINED AS SPARSE VECTORS FOR EXAMPLE IN FACE RECOGNITION COMPONENTS MIGHT NATURALLY MAP TO PARTS OF FACES SPARSE PRINCIPAL COMPONENTS YIELDS A MORE PARSIMONIOUS INTERPRETABLE REPRESENTATION CLEARLY EMPHASIZING WHICH OF THE ORIGINAL FEATURES CONTRIBUTE TO THE DIFFERENCES BETWEEN SAMPLES

THE FOLLOWING EXAMPLE ILLUSTRATES 16 COMPONENTS EXTRACTED USING SPARSE PCA FROM THE OLIVETTI FACES DATASET IT CAN BE SEEN HOW THE REGULARIZATION TERM INDUCES MANY ZEROS FURTHERMORE THE NATURAL STRUCTURE OF THE DATA CAUSES THE NONZERO COEFFICIENTS TO BE VERTICALLY ADJACENT THE MODEL DOES NOT ENFORCE THIS MATHEMATICALLY EACH COMPONENT IS

SCIKITLEARN USER GUIDE RELEASE 0213

A VECTOR  $h \in \mathbb{R}^{4096}$  AND THERE IS NO NOTION OF VERTICAL ADJACENCY EXCEPT DURING THE HUMANFRIENDLY VISUALIZATION AS 64X64 PIXEL IMAGES THE FACT THAT THE COMPONENTS SHOWN BELOW APPEAR LOCAL IS THE EFFECT OF THE INHERENT STRUCTURE OF THE DATA WHICH MAKES SUCH LOCAL PATTERNS MINIMIZE RECONSTRUCTION ERROR THERE EXIST SPARSITYINDUCING NORMS THAT TAKE INTO ACCOUNT ADJACENCY AND DIFFERENT KINDS OF STRUCTURE SEE JEN09 FOR A REVIEW OF SUCH METHODS FOR MORE DETAILS ON HOW TO USE SPARSE PCA SEE THE EXAMPLES SECTION BELOW

NOTE THAT THERE ARE MANY DIFFERENT FORMULATIONS FOR THE SPARSE PCA PROBLEM THE ONE IMPLEMENTED HERE IS BASED ONMRL09 THE OPTIMIZATION PROBLEM SOLVED IS A PCA PROBLEM DICTIONARY LEARNING WITH AN  $\ell_1$ PENALTY ON THE COMPONENTS

$$\min_{\mathbf{W}} \arg \min$$

$$\|\mathbf{W}\|_1$$

$$2\|\mathbf{W}\|_2^2$$

$$2\|\mathbf{W}\|_1$$

SUBJECT TO  $\|\mathbf{w}_i\|_2 \leq 1$  FOR ALL  $0 \leq i < n$

THE SPARSITYINDUCING  $\ell_1$ NORM ALSO PREVENTS LEARNING COMPONENTS FROM NOISE WHEN FEW TRAINING SAMPLES ARE AVAILABLE THE DEGREE OF PENALIZATION AND THUS SPARSITY CAN BE ADJUSTED THROUGH THE HYPERPARAMETER ALPHA SMALL VALUES LEAD TO A GENTLY REGULARIZED FACTORIZATION WHILE LARGER VALUES SHRINK MANY COEFFICIENTS TO ZERO

NOTE WHILE IN THE SPIRIT OF AN ONLINE ALGORITHM THE CLASS MINIBATCHSPARSEPCA DOES NOT IMPLEMENT PARTIALFIT BECAUSE THE ALGORITHM IS ONLINE ALONG THE FEATURES DIRECTION NOT THE SAMPLES DIRECTION

EXAMPLES

- FACES DATASET DECOMPOSITIONS

REFERENCES

TRUNCATED SINGULAR VALUE DECOMPOSITION AND LATENT SEMANTIC ANALYSIS

TRUNCATEDSVD IMPLEMENTS A VARIANT OF SINGULAR VALUE DECOMPOSITION SVD THAT ONLY COMPUTES THE  $k$  LARGEST SINGULAR VALUES WHERE  $k$  IS A USER-SPECIFIED PARAMETER

WHEN TRUNCATED SVD IS APPLIED TO TERM-DOCUMENT MATRICES AS RETURNED BY COUNTVECTORIZER OR TFIDFVECTORIZER THIS TRANSFORMATION IS KNOWN AS LATENT SEMANTIC ANALYSIS LSA BECAUSE IT TRANSFORMS SUCH MATRICES TO A “SEMANTIC” SPACE OF LOW DIMENSIONALITY IN PARTICULAR LSA IS KNOWN TO COMBAT THE EFFECTS OF SYNONYMY AND POLYSEMY BOTH OF WHICH ROUGHLY MEAN THERE ARE MULTIPLE MEANINGS PER WORD WHICH CAUSE TERM-DOCUMENT MATRICES TO BE OVERLY SPARSE AND EXHIBIT POOR SIMILARITY UNDER MEASURES SUCH AS COSINE SIMILARITY NOTE LSA IS ALSO KNOWN AS LATENT SEMANTIC INDEXING LSI THOUGH STRICTLY THAT REFERS TO ITS USE IN PERSISTENT INDEXES FOR INFORMATION RETRIEVAL PURPOSES

MATHEMATICALLY TRUNCATED SVD APPLIED TO TRAINING SAMPLES  $X$  PRODUCES A LOW-RANK APPROXIMATION  $\hat{X}$

$$\hat{X} \approx U U^T X$$

AFTER THIS OPERATION  $X_{\text{transformed}}$

$X_{\text{transformed}}$  IS THE TRANSFORMED TRAINING SET WITH  $k$  FEATURES CALLED  $n_{\text{components}}$  IN THE API

TO ALSO TRANSFORM A TEST SET  $X_{\text{test}}$  WE MULTIPLY IT WITH  $U$

$$\hat{X}_{\text{test}} = U X_{\text{test}}$$

NOTE MOST TREATMENTS OF LSA IN THE NATURAL LANGUAGE PROCESSING NLP AND INFORMATION RETRIEVAL IR LITERATURE SWAP THE AXES OF THE MATRIX  $X$  SO THAT IT HAS SHAPE  $n_{\text{features}} \times n_{\text{samples}}$  WE PRESENT LSA IN A DIFFERENT WAY THAT MATCHES THE SCIKITLEARN API BETTER BUT THE SINGULAR VALUES FOUND ARE THE SAME

TRUNCATEDSVD IS VERY SIMILAR TO PCA BUT DIFFERS IN THAT IT WORKS ON SAMPLE MATRICES  $X$  DIRECTLY INSTEAD OF THEIR COVARIANCE MATRICES WHEN THE COLUMNWISE PER-FEATURE MEANS OF  $X$  ARE SUBTRACTED FROM THE FEATURE VALUES TRUNCATED SVD ON THE RESULTING MATRIX IS EQUIVALENT TO PCA IN PRACTICAL TERMS THIS MEANS THAT THE TRUNCATEDSVD TRANSFORMER ACCEPTS SCIPY SPARSE MATRICES WITHOUT THE NEED TO DENSIFY THEM AS DENSIFYING MAY FILL UP MEMORY EVEN FOR MEDIUM-SIZED DOCUMENT COLLECTIONS

WHILE THE TRUNCATEDSVD TRANSFORMER WORKS WITH ANY SPARSE FEATURE MATRIX USING IT ON TF-IDF MATRICES IS RECOMMENDED OVER RAW FREQUENCY COUNTS IN AN LSA-DOCUMENT PROCESSING SETTING IN PARTICULAR SUBLINEAR SCALING AND INVERSE DOCUMENT FREQUENCY SHOULD BE TURNED ON `SUBLINEAR_TF=True` `USE_IDF=True` TO BRING THE FEATURE VALUES CLOSER TO A GAUSSIAN DISTRIBUTION COMPENSATING FOR LSA’S ERRONEOUS ASSUMPTIONS ABOUT TEXTUAL DATA

EXAMPLES

- CLUSTERING TEXT DOCUMENTS USING KMEANS

REFERENCES

• CHRISTOPHER D MANNING PRABHAKAR RAGHAVAN AND HINRICH SCHÜTZE 2008 INTRODUCTION TO INFORMATION RE  
TRIEVAL CAMBRIDGE UNIVERSITY PRESS CHAPTER 18 MATRIX DECOMPOSITIONS LATENT SEMANTIC INDEXING

DICTIONARY LEARNING

SPARSE CODING WITH A PRECOMPUTED DICTIONARY

THE SPARSECODER OBJECT IS AN ESTIMATOR THAT CAN BE USED TO TRANSFORM SIGNALS INTO SPARSE LINEAR COMBINATION OF  
ATOMS FROM A FIXED PRECOMPUTED DICTIONARY SUCH AS A DISCRETE WAVELET BASIS THIS OBJECT THEREFORE DOES NOT IMPLEMENT  
A FIT METHOD THE TRANSFORMATION AMOUNTS TO A SPARSE CODING PROBLEM FINDING A REPRESENTATION OF THE DATA AS A LINEAR  
COMBINATION OF AS FEW DICTIONARY ATOMS AS POSSIBLE ALL VARIATIONS OF DICTIONARY LEARNING IMPLEMENT THE FOLLOWING  
TRANSFORM METHODS CONTROLLABLE VIA THE TRANSFORMMETHOD INITIALIZATION PARAMETER

- ORTHOGONAL MATCHING PURSUIT ORTHOGONAL MATCHING PURSUIT OMP
- LEASTANGLE REGRESSION LEAST ANGLE REGRESSION
- LASSO COMPUTED BY LEASTANGLE REGRESSION
- LASSO USING COORDINATE DESCENT LASSO
- THRESHOLDING

THRESHOLDING IS VERY FAST BUT IT DOES NOT YIELD ACCURATE RECONSTRUCTIONS THEY HAVE BEEN SHOWN USEFUL IN LITERATURE FOR  
CLASSIFICATION TASKS FOR IMAGE RECONSTRUCTION TASKS ORTHOGONAL MATCHING PURSUIT YIELDS THE MOST ACCURATE UNBIASED  
RECONSTRUCTION

THE DICTIONARY LEARNING OBJECTS OFFER VIA THE SPLITCODE PARAMETER THE POSSIBILITY TO SEPARATE THE POSITIVE AND  
NEGATIVE VALUES IN THE RESULTS OF SPARSE CODING THIS IS USEFUL WHEN DICTIONARY LEARNING IS USED FOR EXTRACTING FEATURES  
THAT WILL BE USED FOR SUPERVISED LEARNING BECAUSE IT ALLOWS THE LEARNING ALGORITHM TO ASSIGN DIFFERENT WEIGHTS TO NEGATIVE  
LOADINGS OF A PARTICULAR ATOM FROM TO THE CORRESPONDING POSITIVE LOADING

THE SPLIT CODE FOR A SINGLE SAMPLE HAS LENGTH 2NCOMPONENTS AND IS CONSTRUCTED USING THE FOLLOWING RULE  
FIRST THE REGULAR CODE OF LENGTH NCOMPONENTS IS COMPUTED THEN THE FIRST NCOMPONENTS ENTRIES OF THE  
SPLITCODE ARE FILLED WITH THE POSITIVE PART OF THE REGULAR CODE VECTOR THE SECOND HALF OF THE SPLIT CODE IS FILLED  
WITH THE NEGATIVE PART OF THE CODE VECTOR ONLY WITH A POSITIVE SIGN THEREFORE THE SPLITCODE IS NONNEGATIVE

EXAMPLES

•SPARSE CODING WITH A PRECOMPUTED DICTIONARY

GENERIC DICTIONARY LEARNING

DICTIONARY LEARNING DICTIONARYLEARNING IS A MATRIX FACTORIZATION PROBLEM THAT AMOUNTS TO FINDING A USUALLY  
OVERCOMPLETE DICTIONARY THAT WILL PERFORM WELL AT SPARSELY ENCODING THE FITTED DATA  
REPRESENTING DATA AS SPARSE COMBINATIONS OF ATOMS FROM AN OVERCOMPLETE DICTIONARY IS SUGGESTED TO BE THE WAY THE  
MAMMALIAN PRIMARY VISUAL CORTEX WORKS CONSEQUENTLY DICTIONARY LEARNING APPLIED ON IMAGE PATCHES HAS BEEN SHOWN  
TO GIVE GOOD RESULTS IN IMAGE PROCESSING TASKS SUCH AS IMAGE COMPLETION INPAINTING AND DENOISING AS WELL AS FOR  
SUPERVISED RECOGNITION TASKS

SCIKITLEARN USER GUIDE RELEASE 0213

DICTIONARY LEARNING IS AN OPTIMIZATION PROBLEM SOLVED BY ALTERNATIVELY UPDATING THE SPARSE CODE AS A SOLUTION TO MULTIPLE LASSO PROBLEMS CONSIDERING THE DICTIONARY FIXED AND THEN UPDATING THE DICTIONARY TO BEST FIT THE SPARSE CODE

$\min_{\mathbf{C}} \|\mathbf{C}\|_1$  ARG MIN

$\min_{\mathbf{C}} \|\mathbf{C}\|_1$

$\min_{\mathbf{C}} \|\mathbf{C}\|_1$

$\min_{\mathbf{C}} \|\mathbf{C}\|_1$

SUBJECT TO  $\|\mathbf{C}\|_2 \leq 1$  FOR ALL  $0 \leq i \leq \text{ATOMS}$

AFTER USING SUCH A PROCEDURE TO FIT THE DICTIONARY THE TRANSFORM IS SIMPLY A SPARSE CODING STEP THAT SHARES THE SAME IMPLEMENTATION WITH ALL DICTIONARY LEARNING OBJECTS SEE SPARSE CODING WITH A PRECOMPUTED DICTIONARY

IT IS ALSO POSSIBLE TO CONSTRAIN THE DICTIONARY AND/OR CODE TO BE POSITIVE TO MATCH CONSTRAINTS THAT MAY BE PRESENT IN THE DATA BELOW ARE THE FACES WITH DIFFERENT POSITIVITY CONSTRAINTS APPLIED RED INDICATES NEGATIVE VALUES BLUE INDICATES POSITIVE VALUES AND WHITE REPRESENTS ZEROS

400 CHAPTER 3 USER GUIDE



SCIKITLEARN USER GUIDE RELEASE 0213

THE FOLLOWING IMAGE SHOWS HOW A DICTIONARY LEARNED FROM 4X4 PIXEL IMAGE PATCHES EXTRACTED FROM PART OF THE IMAGE OF A RACCOON FACE LOOKS LIKE

EXAMPLES

- IMAGE DENOISING USING DICTIONARY LEARNING

REFERENCES

- “ONLINE DICTIONARY LEARNING FOR SPARSE CODING” J MAIRAL F BACH J PONCE G SAPIRO 2009

MINIBATCH DICTIONARY LEARNING

MINIBATCHDICTIONARYLEARNING IMPLEMENTS A FASTER BUT LESS ACCURATE VERSION OF THE DICTIONARY LEARNING ALGORITHM THAT IS BETTER SUITED FOR LARGE DATASETS

BY DEFAULT MINIBATCHDICTIONARYLEARNING DIVIDES THE DATA INTO MINIBATCHES AND OPTIMIZES IN AN ONLINE MANNER BY CYCLING OVER THE MINIBATCHES FOR THE SPECIFIED NUMBER OF ITERATIONS HOWEVER AT THE MOMENT IT DOES NOT IMPLEMENT A STOPPING CONDITION

402 CHAPTER 3 USER GUIDE



SCIKITLEARN USER GUIDE RELEASE 0213

THE ESTIMATOR ALSO IMPLEMENTS PARTIALFIT WHICH UPDATES THE DICTIONARY BY ITERATING ONLY ONCE OVER A MINIBATCH THIS CAN BE USED FOR ONLINE LEARNING WHEN THE DATA IS NOT READILY AVAILABLE FROM THE START OR FOR WHEN THE DATA DOES NOT FIT INTO THE MEMORY

CLUSTERING FOR DICTIONARY LEARNING

NOTE THAT WHEN USING DICTIONARY LEARNING TO EXTRACT A REPRESENTATION EG FOR SPARSE CODING CLUSTERING CAN BE A GOOD PROXY TO LEARN THE DICTIONARY FOR INSTANCE THE MINIBATCHKMEANS ESTIMATOR IS COMPUTATIONALLY EFFICIENT AND IMPLEMENTS ONLINE LEARNING WITH A PARTIALFIT METHOD

EXAMPLE ONLINE LEARNING OF A DICTIONARY OF PARTS OF FACES

FACTOR ANALYSIS

IN UNSUPERVISED LEARNING WE ONLY HAVE A DATASET  $X \in \mathbb{R}^{n \times d}$  HOW CAN THIS DATASET BE DESCRIBED MATHEMATICALLY A VERY SIMPLE CONTINUOUS LATENT VARIABLE MODEL FOR  $X$  IS

$$X \approx WH + \epsilon$$

THE VECTOR  $h$  IS CALLED “LATENT” BECAUSE IT IS UNOBSERVED  $\epsilon$  IS CONSIDERED A NOISE TERM DISTRIBUTED ACCORDING TO A GAUSSIAN WITH MEAN 0 AND COVARIANCE  $\sigma^2 I$   $\mu$  IS SOME ARBITRARY OFFSET VECTOR SUCH A MODEL IS CALLED “GENERATIVE” AS IT DESCRIBES HOW  $X$  IS GENERATED FROM  $h$  IF WE USE ALL THE  $w_i$ ’S AS COLUMNS TO FORM A MATRIX  $W$  AND ALL THE  $h_i$ ’S AS COLUMNS OF A MATRIX  $H$  THEN WE CAN WRITE WITH SUITABLY DEFINED MANDE

$$X \approx WH + \epsilon$$

IN OTHER WORDS WE DECOMPOSED MATRIX  $X$

IF  $h$  IS GIVEN THE ABOVE EQUATION AUTOMATICALLY IMPLIES THE FOLLOWING PROBABILISTIC INTERPRETATION

$$p(h) \propto \exp(-\frac{1}{2\sigma^2} h^T h)$$

FOR A COMPLETE PROBABILISTIC MODEL WE ALSO NEED A PRIOR DISTRIBUTION FOR THE LATENT VARIABLE  $h$  THE MOST STRAIGHTFORWARD ASSUMPTION BASED ON THE NICE PROPERTIES OF THE GAUSSIAN DISTRIBUTION IS  $h \sim \mathcal{N}(\mu, \sigma^2 I)$  THIS YIELDS A GAUSSIAN AS THE MARGINAL DISTRIBUTION OF  $X$

$$p(x) \propto \exp(-\frac{1}{2\sigma^2} x^T x)$$

NOW WITHOUT ANY FURTHER ASSUMPTIONS THE IDEA OF HAVING A LATENT VARIABLE  $h$  WOULD BE SUPERFLUOUS -  $X$  CAN BE COMPLETELY MODELLED WITH A MEAN AND A COVARIANCE WE NEED TO IMPOSE SOME MORE SPECIFIC STRUCTURE ON ONE OF THESE TWO PARAMETERS A SIMPLE ADDITIONAL ASSUMPTION REGARDS THE STRUCTURE OF THE ERROR COVARIANCE  $\sigma^2 I$

•  $\sigma^2 I$ ! THIS ASSUMPTION LEADS TO THE PROBABILISTIC MODEL OF PCA

SCIKITLEARN USER GUIDE RELEASE 0213

•  $\Sigma = \text{diag}(\sigma_1^2, \sigma_2^2, \dots, \sigma_p^2)$  THIS MODEL IS CALLED FACTOR ANALYSIS A CLASSICAL STATISTICAL MODEL THE MATRIX  $W$  IS SOMETIMES CALLED THE “FACTOR LOADING MATRIX”

BOTH MODELS ESSENTIALLY ESTIMATE A GAUSSIAN WITH A LOWRANK COVARIANCE MATRIX BECAUSE BOTH MODELS ARE PROBABILISTIC THEY CAN BE INTEGRATED IN MORE COMPLEX MODELS EG MIXTURE OF FACTOR ANALYSERS ONE GETS VERY DIFFERENT MODELS EG FASTICA IF NONGAUSSIAN PRIORS ON THE LATENT VARIABLES ARE ASSUMED

FACTOR ANALYSIS CANPRODUCE SIMILAR COMPONENTS THE COLUMNS OF ITS LOADING MATRIX TO PCA HOWEVER ONE CAN NOT MAKE ANY GENERAL STATEMENTS ABOUT THESE COMPONENTS EG WHETHER THEY ARE ORTHOGONAL

THE MAIN ADVANTAGE FOR FACTOR ANALYSIS OVER PCA IS THAT IT CAN MODEL THE VARIANCE IN EVERY DIRECTION OF THE INPUT SPACE INDEPENDENTLY HETEROSCEDASTIC NOISE

THIS ALLOWS BETTER MODEL SELECTION THAN PROBABILISTIC PCA IN THE PRESENCE OF HETEROSCEDASTIC NOISE

EXAMPLES

•MODEL SELECTION WITH PROBABILISTIC PCA AND FACTOR ANALYSIS FA

404 CHAPTER 3 USER GUIDE



SCIKITLEARN USER GUIDE RELEASE 0213

INDEPENDENT COMPONENT ANALYSIS ICA

INDEPENDENT COMPONENT ANALYSIS SEPARATES A MULTIVARIATE SIGNAL INTO ADDITIVE SUBCOMPONENTS THAT ARE MAXIMALLY INDEPENDENT IT IS IMPLEMENTED IN SCIKITLEARN USING THE FAST ICA ALGORITHM TYPICALLY ICA IS NOT USED FOR REDUCING DIMENSIONALITY BUT FOR SEPARATING SUPERIMPOSED SIGNALS SINCE THE ICA MODEL DOES NOT INCLUDE A NOISE TERM FOR THE MODEL TO BE CORRECT WHITENING MUST BE APPLIED THIS CAN BE DONE INTERNALLY USING THE WHITEN ARGUMENT OR MANUALLY USING ONE OF THE PCA VARIANTS

IT IS CLASSICALLY USED TO SEPARATE MIXED SIGNALS A PROBLEM KNOWN AS BLIND SOURCE SEPARATION AS IN THE EXAMPLE BELOW ICA CAN ALSO BE USED AS YET ANOTHER NON LINEAR DECOMPOSITION THAT FINDS COMPONENTS WITH SOME SPARSITY

406 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES

- BLIND SOURCE SEPARATION USING FASTICA
- FASTICA ON 2D POINT CLOUDS
- FACES DATASET DECOMPOSITIONS

NONNEGATIVE MATRIX FACTORIZATION NMF OR NNMF

NMF WITH THE FROBENIUS NORM

NMF1IS AN ALTERNATIVE APPROACH TO DECOMPOSITION THAT ASSUMES THAT THE DATA AND THE COMPONENTS ARE NONNEGATIVE NMF CAN BE PLUGGED IN INSTEAD OF PCA OR ITS VARIANTS IN THE CASES WHERE THE DATA MATRIX DOES NOT CONTAIN NEGATIVE VALUES IT FINDS A DECOMPOSITION OF SAMPLES INTO TWO MATRICES AND OF NONNEGATIVE ELEMENTS BY OPTIMIZING THE DISTANCE BETWEEN AND THE MATRIX PRODUCT THE MOST WIDELY USED DISTANCE FUNCTION IS THE SQUARED FROBENIUS NORM WHICH IS AN OBVIOUS EXTENSION OF THE EUCLIDEAN NORM TO MATRICES

$\|FRO\|_F^2 = \sum_{i,j} x_{ij}^2$

$\|x\|_2^2 = x^T x$

$\|FRO\|_F^2 = \sum_{i,j} x_{ij}^2$

$\sum_{i,j} x_{ij}^2$

$\sum_{i,j} x_{ij}^2 = \sum_{i,j} x_{ij}^2$

UNLIKEPCA THE REPRESENTATION OF A VECTOR IS OBTAINED IN AN ADDITIVE FASHION BY SUPERIMPOSING THE COMPONENTS WITHOUT SUBTRACTING SUCH ADDITIVE MODELS ARE EFFICIENT FOR REPRESENTING IMAGES AND TEXT

IT HAS BEEN OBSERVED IN HOYER 20042THAT WHEN CAREFULLY CONSTRAINED NMF CAN PRODUCE A PARTSBASED REPRESENTATION OF THE DATASET RESULTING IN INTERPRETABLE MODELS THE FOLLOWING EXAMPLE DISPLAYS 16 SPARSE COMPONENTS FOUND BY NMF FROM THE IMAGES IN THE OLIVETTI FACES DATASET IN COMPARISON WITH THE PCA EIGENFACES

1“LEARNING THE PARTS OF OBJECTS BY NONNEGATIVE MATRIX FACTORIZATION” D LEE S SEUNG 1999

2“NONNEGATIVE MATRIX FACTORIZATION WITH SPARSENESS CONSTRAINTS” P HOYER 2004

THEINIT ATTRIBUTE DETERMINES THE INITIALIZATION METHOD APPLIED WHICH HAS A GREAT IMPACT ON THE PERFORMANCE OF THE METHODNMF IMPLEMENTS THE METHOD NONNEGATIVE DOUBLE SINGULAR VALUE DECOMPOSITION NNDSVD4IS BASED ON TWO SVD PROCESSES ONE APPROXIMATING THE DATA MATRIX THE OTHER APPROXIMATING POSITIVE SECTIONS OF THE RESULTING PARTIAL SVD FACTORS UTILIZING AN ALGEBRAIC PROPERTY OF UNIT RANK MATRICES THE BASIC NNDSVD ALGORITHM IS BETTER FIT FOR SPARSE FACTORIZATION ITS VARIANTS NNDSVDA IN WHICH ALL ZEROS ARE SET EQUAL TO THE MEAN OF ALL ELEMENTS OF THE DATA AND NNDSVDAR IN WHICH THE ZEROS ARE SET TO RANDOM PERTURBATIONS LESS THAN THE MEAN OF THE DATA DIVIDED BY 100 ARE RECOMMENDED IN THE DENSE CASE

NOTE THAT THE MULTIPLICATIVE UPDATE ‘MU’ SOLVER CANNOT UPDATE ZEROS PRESENT IN THE INITIALIZATION SO IT LEADS TO POORER RESULTS WHEN USED JOINTLY WITH THE BASIC NNDSVD ALGORITHM WHICH INTRODUCES A LOT OF ZEROS IN THIS CASE NNDSVDA OR NNDSVDAR SHOULD BE PREFERRED

NMF CAN ALSO BE INITIALIZED WITH CORRECTLY SCALED RANDOM NONNEGATIVE MATRICES BY SETTING INITRANDOM AN INTEGER SEED OR A RANDOMSTATE CAN ALSO BE PASSED TO RANDOMSTATE TO CONTROL REPRODUCIBILITY INNMF L1 AND L2 PRIORS CAN BE ADDED TO THE LOSS FUNCTION IN ORDER TO REGULARIZE THE MODEL THE L2 PRIOR USES THE FROBENIUS NORM WHILE THE L1 PRIOR USES AN ELEMENTWISE L1 NORM AS IN ELASTICNET WE CONTROL THE COMBINATION OF L1 AND L2 WITH THE L1RATIO  $\lambda$  PARAMETER AND THE INTENSITY OF THE REGULARIZATION WITH THE ALPHA  $\lambda$  PARAMETER THEN THE PRIORS TERMS ARE

$$\frac{\lambda}{2} \|W\|_F^2 + \frac{\lambda}{2} \|H\|_F^2$$

$$\frac{\lambda}{2} \|W\|_F^2$$

$$\frac{\lambda}{2} \|H\|_F^2$$

$$\frac{\lambda}{2} \|W\|_F^2$$

$$\frac{\lambda}{2} \|H\|_F^2$$

4“SVD BASED INITIALIZATION A HEAD START FOR NONNEGATIVE MATRIX FACTORIZATION” C BOUTSIDIS E GALLOPOULOS 2008

SCIKITLEARN USER GUIDE RELEASE 0213

AND THE REGULARIZED OBJECTIVE FUNCTION IS

$$\|FRO\|_{\beta}^2 = \|W\|_{\beta}^2 + \|H\|_{\beta}^2 - \beta \sum_{i,j} W_{ij} H_{ij}$$

NMF REGULARIZES BOTH W AND H THE PUBLIC FUNCTION NONNEGATIVEFACTORIZATION ALLOWS A FINER CONTROL THROUGH THE REGULARIZATION ATTRIBUTE AND MAY REGULARIZE ONLY W ONLY H OR BOTH

NMF WITH A BETADIVERGENCE

AS DESCRIBED PREVIOUSLY THE MOST WIDELY USED DISTANCE FUNCTION IS THE SQUARED FROBENIUS NORM WHICH IS AN OBVIOUS EXTENSION OF THE EUCLIDEAN NORM TO MATRICES

$$\|FRO\|_2^2 = \sum_{i,j} W_{ij}^2 + \sum_{i,j} H_{ij}^2$$

OTHER DISTANCE FUNCTIONS CAN BE USED IN NMF AS FOR EXAMPLE THE GENERALIZED KULLBACKLEIBLER KL DIVERGENCE ALSO REFERRED AS IDIVERGENCE

$$\|KL\|_{\beta} = \sum_{i,j} W_{ij} \log \frac{W_{ij}}{H_{ij}} + \sum_{i,j} H_{ij} \log \frac{H_{ij}}{W_{ij}}$$

OR THE ITAKURASAITO IS DIVERGENCE

$$\|IT\|_{\beta} = \sum_{i,j} W_{ij} \log \frac{W_{ij}}{H_{ij}} + \sum_{i,j} H_{ij} \log \frac{H_{ij}}{W_{ij}}$$

THESE THREE DISTANCES ARE SPECIAL CASES OF THE BETADIVERGENCE FAMILY WITH  $\beta = 2, 1, 0$  RESPECTIVELY6 THE BETA DIVERGENCE ARE DEFINED BY

$$\|B\|_{\beta} = \sum_{i,j} W_{ij}^{\beta} + \sum_{i,j} H_{ij}^{\beta}$$

NOTE THAT THIS DEFINITION IS NOT VALID IF  $\beta \in [0, 1]$  YET IT CAN BE CONTINUOUSLY EXTENDED TO THE DEFINITIONS OF  $\|KL\|_{\beta}$  AND  $\|IT\|_{\beta}$  RESPECTIVELY

NMF IMPLEMENTS TWO SOLVERS USING COORDINATE DESCENT ‘CD’5 AND MULTIPLICATIVE UPDATE ‘MU’6 THE ‘MU’ SOLVER CAN OPTIMIZE EVERY BETADIVERGENCE INCLUDING OF COURSE THE FROBENIUS NORM  $\beta = 2$  THE GENERALIZED KULLBACK LEIBLER DIVERGENCE  $\beta = 1$  AND THE ITAKURASAITO DIVERGENCE  $\beta = 0$  NOTE THAT FOR  $\beta \in [1, 2]$  THE ‘MU’ SOLVER IS SIGNIFICANTLY FASTER THAN FOR OTHER VALUES OF  $\beta$  NOTE ALSO THAT WITH A NEGATIVE OR 0 IE ‘ITAKURASAITO’  $\beta$  THE INPUT MATRIX CANNOT CONTAIN ZERO VALUES

THE ‘CD’ SOLVER CAN ONLY OPTIMIZE THE FROBENIUS NORM DUE TO THE UNDERLYING NONCONVEXITY OF NMF THE DIFFERENT SOLVERS MAY CONVERGE TO DIFFERENT MINIMA EVEN WHEN OPTIMIZING THE SAME DISTANCE FUNCTION

NMF IS BEST USED WITH THE FITTRANSFORM METHOD WHICH RETURNS THE MATRIX W THE MATRIX H IS STORED INTO THE FITTED MODEL IN THE COMPONENTS ATTRIBUTE THE METHOD TRANSFORM WILL DECOMPOSE A NEW MATRIX XNEW BASED ON THESE STORED COMPONENTS

```
import numpy as np
X = np.array([1, 2, 1, 3, 12, 4, 1, 5, 08, 6, 1])
from sklearn.decomposition import NMF
model = NMF(n_components=2, init='random', random_state=0)
```

6“ALGORITHMS FOR NONNEGATIVE MATRIX FACTORIZATION WITH THE BETADIVERGENCE” C FEVOTTE J IDIER 2011

5“FAST LOCAL ALGORITHMS FOR LARGE SCALE NONNEGATIVE MATRIX AND TENSOR FACTORIZATIONS” A CICHOCKI A PHAN 2009

32 UNSUPERVISED LEARNING 409

SCIKITLEARN USER GUIDE RELEASE 0213

W MODELFITTRANSFORMX

H MODELCOMPONENTS

XNEW NPARRAY1 0 1 61 1 0 1 4 32 1 0 4

WNEW MODELTRANSFORMXNEW

EXAMPLES

- FACES DATASET DECOMPOSITIONS
- TOPIC EXTRACTION WITH NONNEGATIVE MATRIX FACTORIZATION AND LATENT DIRICHLET ALLOCATION
- BETADIVERGENCE LOSS FUNCTIONS

REFERENCES

LATENT DIRICHLET ALLOCATION LDA

LATENT DIRICHLET ALLOCATION IS A GENERATIVE PROBABILISTIC MODEL FOR COLLECTIONS OF DISCRETE DATASET SUCH AS TEXT CORPORA  
IT IS ALSO A TOPIC MODEL THAT IS USED FOR DISCOVERING ABSTRACT TOPICS FROM A COLLECTION OF DOCUMENTS

THE GRAPHICAL MODEL OF LDA IS A THREELEVEL GENERATIVE MODEL

410 CHAPTER 3 USER GUIDE



NOTE ON NOTATIONS PRESENTED IN THE GRAPHICAL MODEL ABOVE WHICH CAN BE FOUND IN HOFFMAN ET AL 2013

- THE CORPUS IS A COLLECTION OF  $D$  DOCUMENTS
- A DOCUMENT IS A SEQUENCE OF  $N$  WORDS
- THERE ARE  $K$  TOPICS IN THE CORPUS
- THE BOXES REPRESENT REPEATED SAMPLING

IN THE GRAPHICAL MODEL EACH NODE IS A RANDOM VARIABLE AND HAS A ROLE IN THE GENERATIVE PROCESS A SHADED NODE INDICATES AN OBSERVED VARIABLE AND AN UNSHADED NODE INDICATES A HIDDEN LATENT VARIABLE IN THIS CASE WORDS IN THE CORPUS ARE THE ONLY DATA THAT WE OBSERVE THE LATENT VARIABLES DETERMINE THE RANDOM MIXTURE OF TOPICS IN THE CORPUS AND THE DISTRIBUTION OF WORDS IN THE DOCUMENTS THE GOAL OF LDA IS TO USE THE OBSERVED WORDS TO INFER THE HIDDEN TOPIC STRUCTURE

WHEN MODELING TEXT CORPORA THE MODEL ASSUMES THE FOLLOWING GENERATIVE PROCESS FOR A CORPUS WITH  $D$  DOCUMENTS AND  $K$  TOPICS WITH  $N$  CORRESPONDING TO  $N$  COMPONENTS IN THE API

- 1 FOR EACH TOPIC  $k \in \{1, \dots, K\}$  DRAW  $\theta_k \sim \text{DIRICHLET}(\alpha)$  THIS PROVIDES A DISTRIBUTION OVER THE WORDS IE THE PROBABILITY OF A WORD APPEARING IN TOPIC  $k$  CORRESPONDS TO  $\theta_k$
- 2 FOR EACH DOCUMENT  $d \in \{1, \dots, D\}$  DRAW THE TOPIC PROPORTIONS  $\phi_d \sim \text{DIRICHLET}(\beta)$  CORRESPONDS TO  $\phi_d$

3 FOR EACH WORD  $n$  IN DOCUMENT  $d$

- 1 DRAW THE TOPIC ASSIGNMENT  $z_n \sim \text{MULTINOMIAL}(\phi_d)$
- 2 DRAW THE OBSERVED WORD  $w_n \sim \text{MULTINOMIAL}(\theta_{z_n})$

FOR PARAMETER ESTIMATION THE POSTERIOR DISTRIBUTION IS

$$p(\theta, \phi, z | w) \propto \prod_{k=1}^K \theta_k^{\sum_{n=1}^N \mathbb{1}(z_n = k)} \prod_{d=1}^D \phi_d^{\sum_{n=1}^N \mathbb{1}(z_n = k)}$$

SINCE THE POSTERIOR IS INTRACTABLE VARIATIONAL BAYESIAN METHOD USES A SIMPLER DISTRIBUTION  $q(\theta, \phi, z)$  TO APPROXIMATE IT AND THOSE VARIATIONAL PARAMETERS ARE OPTIMIZED TO MAXIMIZE THE EVIDENCE LOWER BOUND ELBO

$$\text{ELBO} = \mathbb{E}_q[\log p(w | \theta, \phi, z)] - \text{KL}(q(\theta, \phi, z) || p(\theta, \phi, z))$$

MAXIMIZING ELBO IS EQUIVALENT TO MINIMIZING THE KULLBACKLEIBLER KL DIVERGENCE BETWEEN  $q(\theta, \phi, z)$  AND THE TRUE POSTERIOR  $p(\theta, \phi, z | w)$

LATENTDIRICHLETALLOCATION IMPLEMENTS THE ONLINE VARIATIONAL BAYES ALGORITHM AND SUPPORTS BOTH ONLINE AND BATCH UPDATE METHODS WHILE THE BATCH METHOD UPDATES VARIATIONAL VARIABLES AFTER EACH FULL PASS THROUGH THE DATA THE ONLINE METHOD UPDATES VARIATIONAL VARIABLES FROM MINIBATCH DATA POINTS

NOTE ALTHOUGH THE ONLINE METHOD IS GUARANTEED TO CONVERGE TO A LOCAL OPTIMUM POINT THE QUALITY OF THE OPTIMUM POINT AND THE SPEED OF CONVERGENCE MAY DEPEND ON MINIBATCH SIZE AND ATTRIBUTES RELATED TO LEARNING RATE SETTING

SCIKITLEARN USER GUIDE RELEASE 0213

WHEN LATENT DIRICHLET ALLOCATION IS APPLIED ON A "DOCUMENTTERM" MATRIX THE MATRIX WILL BE DECOMPOSED INTO A "TOPICTERM" MATRIX AND A "DOCUMENTTOPIC" MATRIX WHILE "TOPICTERM" MATRIX IS STORED AS COMPONENTS IN THE MODEL "DOCUMENTTOPIC" MATRIX CAN BE CALCULATED FROM TRANSFORM METHOD  
LATENT DIRICHLET ALLOCATION ALSO IMPLEMENTS PARTIALFIT METHOD THIS IS USED WHEN DATA CAN BE FETCHED SEQUENTIALLY

EXAMPLES

- TOPIC EXTRACTION WITH NONNEGATIVE MATRIX FACTORIZATION AND LATENT DIRICHLET ALLOCATION

REFERENCES

- "LATENT DIRICHLET ALLOCATION" D BLEI A NG M JORDAN 2003
- "ONLINE LEARNING FOR LATENT DIRICHLET ALLOCATION" M HOFFMAN D BLEI F BACH 2010
- "STOCHASTIC VARIATIONAL INFERENCE" M HOFFMAN D BLEI C WANG J PAISLEY 2013

SEE ALSO DIMENSIONALITY REDUCTION FOR DIMENSIONALITY REDUCTION WITH NEIGHBORHOOD COMPONENTS ANALYSIS

326 COVARIANCE ESTIMATION

MANY STATISTICAL PROBLEMS REQUIRE THE ESTIMATION OF A POPULATION'S COVARIANCE MATRIX WHICH CAN BE SEEN AS AN ESTIMATION OF DATA SET SCATTER PLOT SHAPE MOST OF THE TIME SUCH AN ESTIMATION HAS TO BE DONE ON A SAMPLE WHOSE PROPERTIES SIZE STRUCTURE HOMOGENEITY HAVE A LARGE INFLUENCE ON THE ESTIMATION'S QUALITY THE SKLEARN COVARIANCE PACKAGE PROVIDES TOOLS FOR ACCURATELY ESTIMATING A POPULATION'S COVARIANCE MATRIX UNDER VARIOUS SETTINGS WE ASSUME THAT THE OBSERVATIONS ARE INDEPENDENT AND IDENTICALLY DISTRIBUTED IID

EMPIRICAL COVARIANCE

THE COVARIANCE MATRIX OF A DATA SET IS KNOWN TO BE WELL APPROXIMATED BY THE CLASSICAL MAXIMUM LIKELIHOOD ESTIMATOR OR "EMPIRICAL COVARIANCE" PROVIDED THE NUMBER OF OBSERVATIONS IS LARGE ENOUGH COMPARED TO THE NUMBER OF FEATURES THE VARIABLES DESCRIBING THE OBSERVATIONS MORE PRECISELY THE MAXIMUM LIKELIHOOD ESTIMATOR OF A SAMPLE IS AN UNBIASED ESTIMATOR OF THE CORRESPONDING POPULATION'S COVARIANCE MATRIX

THE EMPIRICAL COVARIANCE MATRIX OF A SAMPLE CAN BE COMPUTED USING THE EMPIRICAL COVARIANCE FUNCTION OF THE PACKAGE OR BY FITTING AN EMPIRICAL COVARIANCE OBJECT TO THE DATA SAMPLE WITH THE EMPIRICAL COVARIANCE FIT METHOD BE CAREFUL THAT RESULTS DEPEND ON WHETHER THE DATA ARE CENTERED SO ONE MAY WANT TO USE THE ASSUMECENTERED PARAMETER ACCURATELY MORE PRECISELY IF ASSUMECENTEREDFALSE THEN THE TEST SET IS SUPPOSED TO HAVE THE SAME MEAN VECTOR AS THE TRAINING SET IF NOT BOTH SHOULD BE CENTERED BY THE USER AND ASSUMECENTEREDTRUE SHOULD BE USED

EXAMPLES

- SEE SHRINKAGE COVARIANCE ESTIMATION LEDOITWOLF VS OAS AND MAX LIKELIHOOD FOR AN EXAMPLE ON HOW TO FIT AN EMPIRICAL COVARIANCE OBJECT TO DATA

SHRUNK COVARIANCE

BASIC SHRINKAGE

DESPITE BEING AN UNBIASED ESTIMATOR OF THE COVARIANCE MATRIX THE MAXIMUM LIKELIHOOD ESTIMATOR IS NOT A GOOD ESTIMATOR OF THE EIGENVALUES OF THE COVARIANCE MATRIX SO THE PRECISION MATRIX OBTAINED FROM ITS INVERSION IS NOT ACCURATE SOMETIMES IT EVEN OCCURS THAT THE EMPIRICAL COVARIANCE MATRIX CANNOT BE INVERTED FOR NUMERICAL REASONS TO AVOID SUCH AN INVERSION PROBLEM A TRANSFORMATION OF THE EMPIRICAL COVARIANCE MATRIX HAS BEEN INTRODUCED THE SHRINKAGE IN SCIKITLEARN THIS TRANSFORMATION WITH A USERDEFINED SHRINKAGE COEFFICIENT CAN BE DIRECTLY APPLIED TO A PRECOMPUTED COVARIANCE WITH THE SHRUNKCOVARIANCE METHOD ALSO A SHRUNK ESTIMATOR OF THE COVARIANCE CAN BE FITTED TO DATA WITH ASHRUNKCOVARIANCE OBJECT AND ITS SHRUNKCOVARIANCEFIT METHOD AGAIN RESULTS DEPEND ON WHETHER THE DATA ARE CENTERED SO ONE MAY WANT TO USE THE ASSUMECENTERED PARAMETER ACCURATELY MATHEMATICALLY THIS SHRINKAGE CONSISTS IN REDUCING THE RATIO BETWEEN THE SMALLEST AND THE LARGEST EIGENVALUES OF THE EMPIRICAL COVARIANCE MATRIX IT CAN BE DONE BY SIMPLY SHIFTING EVERY EIGENVALUE ACCORDING TO A GIVEN OFFSET WHICH IS EQUIVALENT OF FINDING THE L2PENALIZED MAXIMUM LIKELIHOOD ESTIMATOR OF THE COVARIANCE MATRIX IN PRACTICE SHRINKAGE BOILS DOWN TO A SIMPLE A CONVEX TRANSFORMATION  $\hat{\Sigma}_{SHRUNK} = \frac{1}{1+\alpha} \hat{\Sigma} + \frac{\alpha}{1+\alpha} \text{TR}(\hat{\Sigma}) \text{ID}$

CHOOSING THE AMOUNT OF SHRINKAGE  $\alpha$  AMOUNTS TO SETTING A BIASVARIANCE TRADEOFF AND IS DISCUSSED BELOW

- EXAMPLES
- SEE SHRINKAGE COVARIANCE ESTIMATION LEDOITWOLF VS OAS AND MAXLIKELIHOOD FOR AN EXAMPLE ON HOW TO FIT A SHRUNKCOVARIANCE OBJECT TO DATA

LEDOITWOLF SHRINKAGE

IN THEIR 2004 PAPER1 O LEDOIT AND M WOLF PROPOSE A FORMULA TO COMPUTE THE OPTIMAL SHRINKAGE COEFFICIENT  $\alpha$  THAT MINIMIZES THE MEAN SQUARED ERROR BETWEEN THE ESTIMATED AND THE REAL COVARIANCE MATRIX THE LEDOITWOLF ESTIMATOR OF THE COVARIANCE MATRIX CAN BE COMPUTED ON A SAMPLE WITH THE LEDOITWOLF FUNCTION OF THE SKLEARNCOVARIANCE PACKAGE OR IT CAN BE OTHERWISE OBTAINED BY FITTING A LEDOITWOLF OBJECT TO THE SAME SAMPLE

NOTE CASE WHEN POPULATION COVARIANCE MATRIX IS ISOTROPIC

IT IS IMPORTANT TO NOTE THAT WHEN THE NUMBER OF SAMPLES IS MUCH LARGER THAN THE NUMBER OF FEATURES ONE WOULD EXPECT THAT NO SHRINKAGE WOULD BE NECESSARY THE INTUITION BEHIND THIS IS THAT IF THE POPULATION COVARIANCE IS FULL RANK WHEN THE NUMBER OF SAMPLE GROWS THE SAMPLE COVARIANCE WILL ALSO BECOME POSITIVE DEFINITE AS A RESULT NO SHRINKAGE WOULD BE NECESSARY AND THE METHOD SHOULD AUTOMATICALLY DO THIS

THIS HOWEVER IS NOT THE CASE IN THE LEDOITWOLF PROCEDURE WHEN THE POPULATION COVARIANCE HAPPENS TO BE A MULTIPLE OF THE IDENTITY MATRIX IN THIS CASE THE LEDOITWOLF SHRINKAGE ESTIMATE APPROACHES 1 AS THE NUMBER OF SAMPLES INCREASES THIS INDICATES THAT THE OPTIMAL ESTIMATE OF THE COVARIANCE MATRIX IN THE LEDOITWOLF SENSE IS MULTIPLE OF THE IDENTITY SINCE THE POPULATION COVARIANCE IS ALREADY A MULTIPLE OF THE IDENTITY MATRIX THE LEDOITWOLF SOLUTION IS INDEED A REASONABLE ESTIMATE

10 LEDOIT AND M WOLF “A WELLCONDITIONED ESTIMATOR FOR LARGEDIMENSIONAL COVARIANCE MATRICES” JOURNAL OF MULTIVARIATE ANALYSIS 88 ISSUE 2 FEBRUARY 2004 PAGES 365-411

32 UNSUPERVISED LEARNING 413

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES

- SEE SHRINKAGE COVARIANCE ESTIMATION LEDOITWOLF VS OAS AND MAXLIKELIHOOD FOR AN EXAMPLE ON HOW TO FIT A LEDOITWOLF OBJECT TO DATA AND FOR VISUALIZING THE PERFORMANCES OF THE LEDOITWOLF ESTIMATOR IN TERMS OF LIKELIHOOD

REFERENCES

ORACLE APPROXIMATING SHRINKAGE

UNDER THE ASSUMPTION THAT THE DATA ARE GAUSSIAN DISTRIBUTED CHEN ET AL<sup>2</sup> DERIVED A FORMULA AIMED AT CHOOSING A SHRINKAGE COEFFICIENT THAT YIELDS A SMALLER MEAN SQUARED ERROR THAN THE ONE GIVEN BY LEDOIT AND WOLF’S FORMULA THE RESULTING ESTIMATOR IS KNOWN AS THE ORACLE SHRINKAGE APPROXIMATING ESTIMATOR OF THE COVARIANCE

THE OAS ESTIMATOR OF THE COVARIANCE MATRIX CAN BE COMPUTED ON A SAMPLE WITH THE OAS FUNCTION OF THE SKLEARN COVARIANCE PACKAGE OR IT CAN BE OTHERWISE OBTAINED BY FITTING AN OAS OBJECT TO THE SAME SAMPLE

FIG 37 BIASVARIANCE TRADEOFF WHEN SETTING THE SHRINKAGE COMPARING THE CHOICES OF LEDOITWOLF AND OAS ESTIMATORS

REFERENCES

EXAMPLES

- SEE SHRINKAGE COVARIANCE ESTIMATION LEDOITWOLF VS OAS AND MAXLIKELIHOOD FOR AN EXAMPLE ON HOW TO FIT AN OAS OBJECT TO DATA

<sup>2</sup>CHEN ET AL “SHRINKAGE ALGORITHMS FOR MMSE COVARIANCE ESTIMATION” IEEE TRANS ON SIGN PROC V OLUME 58 ISSUE 10 OCTOBER 2010 414 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

- SEE LEDOITWOLF VS OAS ESTIMATION TO VISUALIZE THE MEAN SQUARED ERROR DIFFERENCE BETWEEN A LEDOITWOLF AND ANOAS ESTIMATOR OF THE COVARIANCE

SPARSE INVERSE COVARIANCE

THE MATRIX INVERSE OF THE COVARIANCE MATRIX OFTEN CALLED THE PRECISION MATRIX IS PROPORTIONAL TO THE PARTIAL CORRELATION MATRIX IT GIVES THE PARTIAL INDEPENDENCE RELATIONSHIP IN OTHER WORDS IF TWO FEATURES ARE INDEPENDENT CONDITIONALLY ON THE OTHERS THE CORRESPONDING COEFFICIENT IN THE PRECISION MATRIX WILL BE ZERO THIS IS WHY IT MAKES SENSE TO ESTIMATE A SPARSE PRECISION MATRIX THE ESTIMATION OF THE COVARIANCE MATRIX IS BETTER CONDITIONED BY LEARNING INDEPENDENCE RELATIONS FROM THE DATA THIS IS KNOWN AS COVARIANCE SELECTION

IN THE SMALLSAMPLES SITUATION IN WHICH NSAMPLES IS ON THE ORDER OF NFEATURES OR SMALLER SPARSE INVERSE COVARIANCE ESTIMATORS TEND TO WORK BETTER THAN SHRUNK COVARIANCE ESTIMATORS HOWEVER IN THE OPPOSITE SITUATION OR FOR VERY CORRELATED DATA THEY CAN BE NUMERICALLY UNSTABLE IN ADDITION UNLIKE SHRINKAGE ESTIMATORS SPARSE ESTIMATORS ARE ABLE TO RECOVER OFFDIAGONAL STRUCTURE

THEGRAPHICALLASSO ESTIMATOR USES AN L1 PENALTY TO ENFORCE SPARSITY ON THE PRECISION MATRIX THE HIGHER ITS ALPHA PARAMETER THE MORE SPARSE THE PRECISION MATRIX THE CORRESPONDING GRAPHICALLASSOCV OBJECT USES CROSSVALIDATION TO AUTOMATICALLY SET THE ALPHA PARAMETER

NOTE STRUCTURE RECOVERY

RECOVERING A GRAPHICAL STRUCTURE FROM CORRELATIONS IN THE DATA IS A CHALLENGING THING IF YOU ARE INTERESTED IN SUCH RECOVERY KEEP IN MIND THAT

- RECOVERY IS EASIER FROM A CORRELATION MATRIX THAN A COVARIANCE MATRIX STANDARDIZE YOUR OBSERVATIONS BEFORE RUNNINGGRAPHICALLASSO

- IF THE UNDERLYING GRAPH HAS NODES WITH MUCH MORE CONNECTIONS THAN THE AVERAGE NODE THE ALGORITHM WILL MISS SOME OF THESE CONNECTIONS

32 UNSUPERVISED LEARNING 415

FIG 38 A COMPARISON OF MAXIMUM LIKELIHOOD SHRINKAGE AND SPARSE ESTIMATES OF THE COVARIANCE AND PRECISION MATRIX IN THE VERY SMALL SAMPLES SETTINGS

- IF YOUR NUMBER OF OBSERVATIONS IS NOT LARGE COMPARED TO THE NUMBER OF EDGES IN YOUR UNDERLYING GRAPH YOU WILL NOT RECOVER IT
- EVEN IF YOU ARE IN FAVORABLE RECOVERY CONDITIONS THE ALPHA PARAMETER CHOSEN BY CROSSVALIDATION EG USING THE GRAPHICALASSOCV OBJECT WILL LEAD TO SELECTING TOO MANY EDGES HOWEVER THE RELEVANT EDGES WILL HAVE HEAVIER WEIGHTS THAN THE IRRELEVANT ONES

THE MATHEMATICAL FORMULATION IS THE FOLLOWING

$$\hat{\Sigma} = \underset{\Sigma}{\operatorname{argmin}} \operatorname{tr}(\Sigma) - \log \det \Sigma + \lambda \|\Sigma\|_1$$

WHERE  $\Sigma$  IS THE PRECISION MATRIX TO BE ESTIMATED AND  $\hat{\Sigma}$  IS THE SAMPLE COVARIANCE MATRIX  $\|\Sigma\|_1$  IS THE SUM OF THE ABSOLUTE VALUES OF OFFDIAGONAL COEFFICIENTS OF  $\Sigma$  THE ALGORITHM EMPLOYED TO SOLVE THIS PROBLEM IS THE GLASSO ALGORITHM FROM THE FRIEDMAN 2008 BIOSTATISTICS PAPER IT IS THE SAME ALGORITHM AS IN THE R GLASSO PACKAGE

EXAMPLES

- SPARSE INVERSE COVARIANCE ESTIMATION EXAMPLE ON SYNTHETIC DATA SHOWING SOME RECOVERY OF A STRUCTURE AND COMPARING TO OTHER COVARIANCE ESTIMATORS
- VISUALIZING THE STOCK MARKET STRUCTURE EXAMPLE ON REAL STOCK MARKET DATA FINDING WHICH SYMBOLS ARE MOST LINKED

REFERENCES

- FRIEDMAN ET AL “SPARSE INVERSE COVARIANCE ESTIMATION WITH THE GRAPHICAL LASSO” BIOSTATISTICS 9 PP 432 2008

SCIKITLEARN USER GUIDE RELEASE 0213

ROBUST COVARIANCE ESTIMATION

REAL DATA SETS ARE OFTEN SUBJECT TO MEASUREMENT OR RECORDING ERRORS REGULAR BUT UNCOMMON OBSERVATIONS MAY ALSO APPEAR FOR A VARIETY OF REASONS OBSERVATIONS WHICH ARE VERY UNCOMMON ARE CALLED OUTLIERS THE EMPIRICAL COVARIANCE ESTIMATOR AND THE SHRUNK COVARIANCE ESTIMATORS PRESENTED ABOVE ARE VERY SENSITIVE TO THE PRESENCE OF OUTLIERS IN THE DATA THEREFORE ONE SHOULD USE ROBUST COVARIANCE ESTIMATORS TO ESTIMATE THE COVARIANCE OF ITS REAL DATA SETS ALTERNATIVELY ROBUST COVARIANCE ESTIMATORS CAN BE USED TO PERFORM OUTLIER DETECTION AND DISCARD DOWNWEIGHT SOME OBSERVATIONS ACCORDING TO FURTHER PROCESSING OF THE DATA

THE SKLEARN COVARIANCE PACKAGE IMPLEMENTS A ROBUST ESTIMATOR OF COVARIANCE THE MINIMUM COVARIANCE DETERMINANT

MINIMUM COVARIANCE DETERMINANT

THE MINIMUM COVARIANCE DETERMINANT ESTIMATOR IS A ROBUST ESTIMATOR OF A DATA SET'S COVARIANCE INTRODUCED BY PJ ROUSSEEUW IN 3 THE IDEA IS TO FIND A GIVEN PROPORTION  $h$  OF "GOOD" OBSERVATIONS WHICH ARE NOT OUTLIERS AND COMPUTE THEIR EMPIRICAL COVARIANCE MATRIX THIS EMPIRICAL COVARIANCE MATRIX IS THEN RESCALED TO COMPENSATE THE PERFORMED SELECTION OF OBSERVATIONS "CONSISTENCY STEP" HAVING COMPUTED THE MINIMUM COVARIANCE DETERMINANT ESTIMATOR ONE CAN GIVE WEIGHTS TO OBSERVATIONS ACCORDING TO THEIR MAHALANOBIS DISTANCE LEADING TO A REWEIGHTED ESTIMATE OF THE COVARIANCE MATRIX OF THE DATA SET "REWEIGHTING STEP"

ROUSSEEUW AND VAN DRIESSEN 4 DEVELOPED THE FASTMCD ALGORITHM IN ORDER TO COMPUTE THE MINIMUM COVARIANCE DETERMINANT THIS ALGORITHM IS USED IN SCIKITLEARN WHEN FITTING AN MCD OBJECT TO DATA THE FASTMCD ALGORITHM ALSO COMPUTES A ROBUST ESTIMATE OF THE DATA SET LOCATION AT THE SAME TIME

RAW ESTIMATES CAN BE ACCESSED AS RAWLOCATION AND RAWCOVARIANCE ATTRIBUTES OF A MINCOVD ET ROBUST COVARIANCE ESTIMATOR OBJECT

REFERENCES

EXAMPLES

- SEE ROBUST VS EMPIRICAL COVARIANCE ESTIMATE FOR AN EXAMPLE ON HOW TO FIT A MINCOVD ET OBJECT TO DATA AND SEE HOW THE ESTIMATE REMAINS ACCURATE DESPITE THE PRESENCE OF OUTLIERS
- SEE ROBUST COVARIANCE ESTIMATION AND MAHALANOBIS DISTANCES RELEVANCE TO VISUALIZE THE DIFFERENCE BETWEEN EMPIRICAL COVARIANCE AND MINCOVD ET COVARIANCE ESTIMATORS IN TERMS OF MAHALANOBIS DISTANCE SO WE GET A BETTER ESTIMATE OF THE PRECISION MATRIX TOO

3P J ROUSSEEUW LEAST MEDIAN OF SQUARES REGRESSION J AM STAT ASS 79 871 1984

4A FAST ALGORITHM FOR THE MINIMUM COVARIANCE DETERMINANT ESTIMATOR 1999 AMERICAN STATISTICAL ASSOCIATION AND THE AMERICAN QUALITY TECHNOMETRICS

32 UNSUPERVISED LEARNING 417

SCIKITLEARN USER GUIDE RELEASE 0213

INFLUENCE OF OUTLIERS ON LOCATION AND COVARIANCE

ESTIMATESSEPARATING INLIERS FROM OUTLIERS USING A MAHALANOBIS DISTANCE

327 NOVELTY AND OUTLIER DETECTION

MANY APPLICATIONS REQUIRE BEING ABLE TO DECIDE WHETHER A NEW OBSERVATION BELONGS TO THE SAME DISTRIBUTION AS EXISTING OBSERVATIONS IT IS AN INLIER OR SHOULD BE CONSIDERED AS DIFFERENT IT IS AN OUTLIER OFTEN THIS ABILITY IS USED TO CLEAN REAL DATA SETS TWO IMPORTANT DISTINCTIONS MUST BE MADE

OUTLIER DETECTION THE TRAINING DATA CONTAINS OUTLIERS WHICH ARE DEFINED AS OBSERVATIONS THAT ARE FAR FROM THE OTHERS OUTLIER DETECTION ESTIMATORS THUS TRY TO FIT THE REGIONS WHERE THE TRAINING DATA IS THE MOST CONCENTRATED IGNORING THE DEVIANT OBSERVATIONS

NOVELTY DETECTION THE TRAINING DATA IS NOT POLLUTED BY OUTLIERS AND WE ARE INTERESTED IN DETECTING WHETHER ANEW OBSERVATION IS AN OUTLIER IN THIS CONTEXT AN OUTLIER IS ALSO CALLED A NOVELTY

OUTLIER DETECTION AND NOVELTY DETECTION ARE BOTH USED FOR ANOMALY DETECTION WHERE ONE IS INTERESTED IN DETECTING ABNORMAL OR UNUSUAL OBSERVATIONS OUTLIER DETECTION IS THEN ALSO KNOWN AS UNSUPERVISED ANOMALY DETECTION AND NOVELTY DETECTION AS SEMISUPERVISED ANOMALY DETECTION IN THE CONTEXT OF OUTLIER DETECTION THE OUTLIERSANOMALIES CANNOT FORM A DENSE CLUSTER AS AVAILABLE ESTIMATORS ASSUME THAT THE OUTLIERSANOMALIES ARE LOCATED IN LOW DENSITY REGIONS ON THE CONTRARY IN THE CONTEXT OF NOVELTY DETECTION NOVELTIESANOMALIES CAN FORM A DENSE CLUSTER AS LONG AS THEY ARE IN A LOW DENSITY REGION OF THE TRAINING DATA CONSIDERED AS NORMAL IN THIS CONTEXT

THE SCIKITLEARN PROJECT PROVIDES A SET OF MACHINE LEARNING TOOLS THAT CAN BE USED BOTH FOR NOVELTY OR OUTLIER DETECTION THIS STRATEGY IS IMPLEMENTED WITH OBJECTS LEARNING IN AN UNSUPERVISED WAY FROM THE DATA

ESTIMATORFITXTRAIN

NEW OBSERVATIONS CAN THEN BE SORTED AS INLIERS OR OUTLIERS WITH A PREDICT METHOD

ESTIMATORPREDICTXTEST

INLIERS ARE LABELED 1 WHILE OUTLIERS ARE LABELED 1 THE PREDICT METHOD MAKES USE OF A THRESHOLD ON THE RAW SCORING FUNCTION COMPUTED BY THE ESTIMATOR THIS SCORING FUNCTION IS ACCESSIBLE THROUGH THE SCORESAMPLES METHOD WHILE THE THRESHOLD CAN BE CONTROLLED BY THE CONTAMINATION PARAMETER

THEDECISIONFUNCTION METHOD IS ALSO DEFINED FROM THE SCORING FUNCTION IN SUCH A WAY THAT NEGATIVE VALUES ARE OUTLIERS AND NONNEGATIVE ONES ARE INLIERS

ESTIMATORDECISIONFUNCTIONXTEST

418 CHAPTER 3 USER GUIDE



SCIKITLEARN USER GUIDE RELEASE 0213

NOTE THAT NEIGHBORSLOCALOUTLIERFACTOR DOES NOT SUPPORT PREDICT DECISIONFUNCTION AND SCORESAMPLERESAMPLES METHODS BY DEFAULT BUT ONLY A FITPREDICT METHOD AS THIS ESTIMATOR WAS ORIGINALLY MEANT TO BE APPLIED FOR OUTLIER DETECTION THE SCORES OF ABNORMALITY OF THE TRAINING SAMPLES ARE ACCESSIBLE THROUGH THE NEGATIVEOUTLIERFACTOR ATTRIBUTE  
IF YOU REALLY WANT TO USE NEIGHBORSLOCALOUTLIERFACTOR FOR NOVELTY DETECTION IE PREDICT LABELS OR COMPUTE THE SCORE OF ABNORMALITY OF NEW UNSEEN DATA YOU CAN INSTANTIATE THE ESTIMATOR WITH THE NOVELTY PARAMETER SET TO TRUE BEFORE FITTING THE ESTIMATOR IN THIS CASE FITPREDICT IS NOT AVAILABLE

WARNING NOVELTY DETECTION WITH LOCAL OUTLIER FACTOR  
WHENNOVELTY IS SET TOTRUE BE AWARE THAT YOU MUST ONLY USE PREDICT DECISIONFUNCTION ANDSCORESAMPLERESAMPLES ON NEW UNSEEN DATA AND NOT ON THE TRAINING SAMPLES AS THIS WOULD LEAD TO WRONG RESULTS THE SCORES OF ABNORMALITY OF THE TRAINING SAMPLES ARE ALWAYS ACCESSIBLE THROUGH THE NEGATIVEOUTLIERFACTOR ATTRIBUTE

THE BEHAVIOR OF NEIGHBORSLOCALOUTLIERFACTOR IS SUMMARIZED IN THE FOLLOWING TABLE

METHOD	OUTLIER DETECTION	NOVELTY DETECTION
FITPREDICT	OK	NOT AVAILABLE
PREDICT	NOT AVAILABLE	USE ONLY ON NEW DATA
DECISIONFUNCTION	NOT AVAILABLE	USE ONLY ON NEW DATA
SCORESAMPLERESAMPLES	USENEGATIVEOUTLIERFACTOR	USE ONLY ON NEW DATA

OVERVIEW OF OUTLIER DETECTION METHODS  
A COMPARISON OF THE OUTLIER DETECTION ALGORITHMS IN SCIKITLEARN LOCAL OUTLIER FACTOR LOF DOES NOT SHOW A DECISION BOUNDARY IN BLACK AS IT HAS NO PREDICT METHOD TO BE APPLIED ON NEW DATA WHEN IT IS USED FOR OUTLIER DETECTION ENSEMBLEISOLATIONFOREST ANDNEIGHBORSLOCALOUTLIERFACTOR PERFORM REASONABLY WELL ON THE DATA SETS CONSIDERED HERE THE SVMONECLASSSSVM IS KNOWN TO BE SENSITIVE TO OUTLIERS AND THUS DOES NOT PERFORM VERY WELL FOR OUTLIER DETECTION FINALLY COVARIANCEELLIPTICENVELOPE ASSUMES THE DATA IS GAUSSIAN AND LEARNS AN ELLIPSE FOR MORE DETAILS ON THE DIFFERENT ESTIMATORS REFER TO THE EXAMPLE COMPARING ANOMALY DETECTION ALGORITHMS FOR OUTLIER DETECTION ON TOY DATASETS AND THE SECTIONS HEREUNDER

EXAMPLES  
• SEE COMPARING ANOMALY DETECTION ALGORITHMS FOR OUTLIER DETECTION ON TOY DATASETS FOR A COMPARISON OF THE SVMONECLASSSSVM THEENSEMBLEISOLATIONFOREST THENEIGHBORSLOCALOUTLIERFACTOR ANDCOVARIANCEELLIPTICENVELOPE

NOVELTY DETECTION  
CONSIDER A DATA SET OF  $n$  OBSERVATIONS FROM THE SAME DISTRIBUTION DESCRIBED BY  $d$  FEATURES CONSIDER NOW THAT WE ADD ONE MORE OBSERVATION TO THAT DATA SET IS THE NEW OBSERVATION SO DIFFERENT FROM THE OTHERS THAT WE CAN DOUBT IT IS REGULAR IE DOES IT COME FROM THE SAME DISTRIBUTION OR ON THE CONTRARY IS IT SO SIMILAR TO THE OTHER THAT WE CANNOT DISTINGUISH IT FROM THE ORIGINAL OBSERVATIONS THIS IS THE QUESTION ADDRESSED BY THE NOVELTY DETECTION TOOLS AND METHODS IN GENERAL IT IS ABOUT TO LEARN A ROUGH CLOSE FRONTIER DELIMITING THE CONTOUR OF THE INITIAL OBSERVATIONS DISTRIBUTION PLOTTED IN EMBEDDING  $d$  DIMENSIONAL SPACE THEN IF FURTHER OBSERVATIONS LAY WITHIN THE FRONTIERDELIMITED SUBSPACE



SCIKITLEARN USER GUIDE RELEASE 0213

THEY ARE CONSIDERED AS COMING FROM THE SAME POPULATION THAN THE INITIAL OBSERVATIONS OTHERWISE IF THEY LAY OUTSIDE THE FRONTIER WE CAN SAY THAT THEY ARE ABNORMAL WITH A GIVEN CONFIDENCE IN OUR ASSESSMENT

THE ONECLASS SVM HAS BEEN INTRODUCED BY SCHÖLKOPF ET AL FOR THAT PURPOSE AND IMPLEMENTED IN THE SUPPORT VECTOR MACHINES MODULE IN THE SVMONECLASSSSVM OBJECT IT REQUIRES THE CHOICE OF A KERNEL AND A SCALAR PARAMETER TO DEFINE A FRONTIER THE RBF KERNEL IS USUALLY CHOSEN ALTHOUGH THERE EXISTS NO EXACT FORMULA OR ALGORITHM TO SET ITS BANDWIDTH PARAMETER THIS IS THE DEFAULT IN THE SCIKITLEARN IMPLEMENTATION THE  $\gamma$ PARAMETER ALSO KNOWN AS THE MARGIN OF THE ONECLASS SVM CORRESPONDS TO THE PROBABILITY OF FINDING A NEW BUT REGULAR OBSERVATION OUTSIDE THE FRONTIER

REFERENCES

- ESTIMATING THE SUPPORT OF A HIGHDIMENSIONAL DISTRIBUTION SCHÖLKOPF BERNHARD ET AL NEURAL COMPUTATION 137 2001 14431471

EXAMPLES

- SEE ONECLASS SVM WITH NONLINEAR KERNEL RBF FOR VISUALIZING THE FRONTIER LEARNED AROUND SOME DATA BY A SVMONECLASSSSVM OBJECT
- SPECIES DISTRIBUTION MODELING

OUTLIER DETECTION

OUTLIER DETECTION IS SIMILAR TO NOVELTY DETECTION IN THE SENSE THAT THE GOAL IS TO SEPARATE A CORE OF REGULAR OBSERVATIONS FROM SOME POLLUTING ONES CALLED OUTLIERS YET IN THE CASE OF OUTLIER DETECTION WE DON'T HAVE A CLEAN DATA SET REPRESENTING THE POPULATION OF REGULAR OBSERVATIONS THAT CAN BE USED TO TRAIN ANY TOOL

SCIKITLEARN USER GUIDE RELEASE 0213

FITTING AN ELLIPTIC ENVELOPE

ONE COMMON WAY OF PERFORMING OUTLIER DETECTION IS TO ASSUME THAT THE REGULAR DATA COME FROM A KNOWN DISTRIBUTION  
EG DATA ARE GAUSSIAN DISTRIBUTED FROM THIS ASSUMPTION WE GENERALLY TRY TO DEFINE THE “SHAPE” OF THE DATA AND CAN  
DEFINE OUTLYING OBSERVATIONS AS OBSERVATIONS WHICH STAND FAR ENOUGH FROM THE FIT SHAPE  
THE SCIKITLEARN PROVIDES AN OBJECT COVARIANCEELLIPTICENVELOPE THAT FITS A ROBUST COVARIANCE ESTIMATE TO  
THE DATA AND THUS FITS AN ELLIPSE TO THE CENTRAL DATA POINTS IGNORING POINTS OUTSIDE THE CENTRAL MODE  
FOR INSTANCE ASSUMING THAT THE INLIER DATA ARE GAUSSIAN DISTRIBUTED IT WILL ESTIMATE THE INLIER LOCATION AND COVARIANCE  
IN A ROBUST WAY IE WITHOUT BEING INFLUENCED BY OUTLIERS THE MAHALANOBIS DISTANCES OBTAINED FROM THIS ESTIMATE IS  
USED TO DERIVE A MEASURE OF OUTLYINGNESS THIS STRATEGY IS ILLUSTRATED BELOW

EXAMPLES

- SEE ROBUST COVARIANCE ESTIMATION AND MAHALANOBIS DISTANCES RELEVANCE FOR AN ILLUSTRATION OF THE DIF  
ERENCE BETWEEN USING A STANDARD COVARIANCEEMPIRICALCOVARIANCE OR A ROBUST ESTIMATE  
COVARIANCEMINCOVDET OF LOCATION AND COVARIANCE TO ASSESS THE DEGREE OF OUTLYINGNESS OF AN OB  
SERVATION

REFERENCES

- ROUSSEEUW PJ VAN DRIESSEN K “A FAST ALGORITHM FOR THE MINIMUM COVARIANCE DETERMINANT ESTIMATOR”  
TECHNOMETRICS 41:3 212 1999

ISOLATION FOREST

ONE EFFICIENT WAY OF PERFORMING OUTLIER DETECTION IN HIGHDIMENSIONAL DATASETS IS TO USE RANDOM FORESTS THE ENSEMBLEISOLATIONFOREST ‘ISOLATES’ OBSERVATIONS BY RANDOMLY SELECTING A FEATURE AND THEN RANDOMLY SELECTING A SPLIT VALUE BETWEEN THE MAXIMUM AND MINIMUM VALUES OF THE SELECTED FEATURE SINCE RECURSIVE PARTITIONING CAN BE REPRESENTED BY A TREE STRUCTURE THE NUMBER OF SPLITTINGS REQUIRED TO ISOLATE A SAMPLE IS EQUIVALENT TO THE PATH LENGTH FROM THE ROOT NODE TO THE TERMINATING NODE THIS PATH LENGTH AVERAGED OVER A FOREST OF SUCH RANDOM TREES IS A MEASURE OF NORMALITY AND OUR DECISION FUNCTION RANDOM PARTITIONING PRODUCES NOTICEABLY SHORTER PATHS FOR ANOMALIES HENCE WHEN A FOREST OF RANDOM TREES COLLECTIVELY PRODUCE SHORTER PATH LENGTHS FOR PARTICULAR SAMPLES THEY ARE HIGHLY LIKELY TO BE ANOMALIES THE IMPLEMENTATION OF ENSEMBLEISOLATIONFOREST IS BASED ON AN ENSEMBLE OF TREE EXTRATREEREgressor FOLLOWING ISOLATION FOREST ORIGINAL PAPER THE MAXIMUM DEPTH OF EACH TREE IS SET TO $\lceil \log_2 n \rceil$  WHERE  $n$  IS THE NUMBER OF SAMPLES USED TO BUILD THE TREE SEE LIU ET AL 2008 FOR MORE DETAILS THIS ALGORITHM IS ILLUSTRATED BELOW THE ENSEMBLEISOLATIONFOREST SUPPORTSWARMSTARTTRUE WHICH ALLOWS YOU TO ADD MORE TREES TO AN ALREADY FITTED MODEL

```
FROM SKLEARNENSEMBLE IMPORT ISOLATIONFOREST
IMPORT NUMPY AS NP
X = NPARRAY(1 1 2 1 3 2 0 0 20 50 3 5)
CLF = ISOLATIONFOREST(NESTIMATORS=10, WARMSTART=True)
CLF.fit(X) # FIT 10 TREES
CLF.set_params(n_estimators=20) # ADD 10 MORE TREES
CLF.fit(X) # FIT THE ADDED TREES
```

- EXAMPLES
- SEE ISOLATIONFOREST EXAMPLE FOR AN ILLUSTRATION OF THE USE OF ISOLATIONFOREST
- 32 UNSUPERVISED LEARNING 423

SCIKITLEARN USER GUIDE RELEASE 0213

• SEE COMPARING ANOMALY DETECTION ALGORITHMS FOR OUTLIER DETECTION ON TOY DATASETS FOR A COMPARISON OF ENSEMBLEISOLATIONFOREST WITHNEIGHBORSLOCALOUTLIERFACTOR SVMONECLASSSVM  
TUNED TO PERFORM LIKE AN OUTLIER DETECTION METHOD AND A COVARIANCEBASED OUTLIER DETECTION WITH COVARIANCEELLIPTICENVELOPE

REFERENCES

• LIU FEI TONY TING KAI MING AND ZHOU ZHIHUA “ISOLATION FOREST” DATA MINING 2008 ICDM’08 EIGHTH IEEE INTERNATIONAL CONFERENCE ON

LOCAL OUTLIER FACTOR

ANOTHER EFFICIENT WAY TO PERFORM OUTLIER DETECTION ON MODERATELY HIGH DIMENSIONAL DATASETS IS TO USE THE LOCAL OUTLIER FACTOR LOF ALGORITHM

THE NEIGHBORS LOCAL OUTLIER FACTOR LOF ALGORITHM COMPUTES A SCORE CALLED LOCAL OUTLIER FACTOR REFLECTING THE DEGREE OF ABNORMALITY OF THE OBSERVATIONS IT MEASURES THE LOCAL DENSITY DEVIATION OF A GIVEN DATA POINT WITH RESPECT TO ITS NEIGHBORS THE IDEA IS TO DETECT THE SAMPLES THAT HAVE A SUBSTANTIALLY LOWER DENSITY THAN THEIR NEIGHBORS IN PRACTICE THE LOCAL DENSITY IS OBTAINED FROM THE K NEAREST NEIGHBORS THE LOF SCORE OF AN OBSERVATION IS EQUAL TO THE RATIO OF THE AVERAGE LOCAL DENSITY OF HIS K NEAREST NEIGHBORS AND ITS OWN LOCAL DENSITY A NORMAL INSTANCE IS EXPECTED TO HAVE A LOCAL DENSITY SIMILAR TO THAT OF ITS NEIGHBORS WHILE ABNORMAL DATA ARE EXPECTED TO HAVE MUCH SMALLER LOCAL DENSITY

THE NUMBER K OF NEIGHBORS CONSIDERED ALIAS PARAMETER NNEIGHBORS IS TYPICALLY CHOSEN 1 GREATER THAN THE MINIMUM NUMBER OF OBJECTS A CLUSTER HAS TO CONTAIN SO THAT OTHER OBJECTS CAN BE LOCAL OUTLIERS RELATIVE TO THIS CLUSTER AND 2 SMALLER THAN THE MAXIMUM NUMBER OF CLOSE BY OBJECTS THAT CAN POTENTIALLY BE LOCAL OUTLIERS IN PRACTICE SUCH INFORMATION IS GENERALLY NOT AVAILABLE AND TAKING NNEIGHBORS=20 APPEARS TO WORK WELL IN GENERAL WHEN THE PROPORTION OF OUTLIERS IS HIGH IE GREATER THAN 10 AS IN THE EXAMPLE BELOW NNEIGHBORS SHOULD BE GREATER NNEIGHBORS=35 IN THE EXAMPLE BELOW

THE STRENGTH OF THE LOF ALGORITHM IS THAT IT TAKES BOTH LOCAL AND GLOBAL PROPERTIES OF DATASETS INTO CONSIDERATION IT CAN PERFORM WELL EVEN IN DATASETS WHERE ABNORMAL SAMPLES HAVE DIFFERENT UNDERLYING DENSITIES THE QUESTION IS NOT HOW ISOLATED THE SAMPLE IS BUT HOW ISOLATED IT IS WITH RESPECT TO THE SURROUNDING NEIGHBORHOOD

WHEN APPLYING LOF FOR OUTLIER DETECTION THERE ARE NO PREDICT DECISIONFUNCTION AND SCORES SAMPLES METHODS BUT ONLY A FITPREDICT METHOD THE SCORES OF ABNORMALITY OF THE TRAINING SAMPLES ARE ACCESSIBLE THROUGH THE NEGATIVEOUTLIERFACTOR ATTRIBUTE NOTE THAT PREDICT DECISIONFUNCTION AND SCORES SAMPLES CAN BE USED ON NEW UNSEEN DATA WHEN LOF IS APPLIED FOR NOVELTY DETECTION IE WHEN THE NOVELTY PARAMETER IS SET TO TRUE SEE NOVELTY DETECTION WITH LOCAL OUTLIER FACTOR

THIS STRATEGY IS ILLUSTRATED BELOW

EXAMPLES

• SEE OUTLIER DETECTION WITH LOCAL OUTLIER FACTOR LOF FOR AN ILLUSTRATION OF THE USE OF NEIGHBORS LOCALOUTLIERFACTOR

• SEE COMPARING ANOMALY DETECTION ALGORITHMS FOR OUTLIER DETECTION ON TOY DATASETS FOR A COMPARISON WITH OTHER ANOMALY DETECTION METHODS

424 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

REFERENCES

- BREUNIG KRIEGEL NG AND SANDER 2000 LOF IDENTIFYING DENSITYBASED LOCAL OUTLIERS PROC ACM SIGMOD

NOVELTY DETECTION WITH LOCAL OUTLIER FACTOR

TO USE NEIGHBORSLOCALOUTLIERFACTOR FOR NOVELTY DETECTION IE PREDICT LABELS OR COMPUTE THE SCORE OF ABNORMALITY OF NEW UNSEEN DATA YOU NEED TO INSTANTIATE THE ESTIMATOR WITH THE NOVELTY PARAMETER SET TO TRUE BEFORE FITTING THE ESTIMATOR

LOF = LOCALOUTLIERFACTOR(NOVELTY = TRUE

LOFFITXTRAIN

NOTE THAT FITPREDICT IS NOT AVAILABLE IN THIS CASE

WARNING NOVELTY DETECTION WITH LOCAL OUTLIER FACTOR'

WHEN NOVELTY IS SET TO TRUE BE AWARE THAT YOU MUST ONLY USE PREDICT DECISIONFUNCTION

AND SCORES SAMPLES ON NEW UNSEEN DATA AND NOT ON THE TRAINING SAMPLES AS THIS WOULD LEAD TO

WRONG RESULTS THE SCORES OF ABNORMALITY OF THE TRAINING SAMPLES ARE ALWAYS ACCESSIBLE THROUGH THE

NEGATIVEOUTLIERFACTOR ATTRIBUTE

NOVELTY DETECTION WITH LOCAL OUTLIER FACTOR IS ILLUSTRATED BELOW

32 UNSUPERVISED LEARNING 425

328 DENSITY ESTIMATION

DENSITY ESTIMATION WALKS THE LINE BETWEEN UNSUPERVISED LEARNING FEATURE ENGINEERING AND DATA MODELING SOME OF THE MOST POPULAR AND USEFUL DENSITY ESTIMATION TECHNIQUES ARE MIXTURE MODELS SUCH AS GAUSSIAN MIXTURES SKLEARN MIXTUREGAUSSIANMIXTURE AND NEIGHBORBASED APPROACHES SUCH AS THE KERNEL DENSITY ESTIMATE SKLEARN NEIGHBORSKERNELDENSITY GAUSSIAN MIXTURES ARE DISCUSSED MORE FULLY IN THE CONTEXT OF CLUSTERING BECAUSE THE TECHNIQUE IS ALSO USEFUL AS AN UNSUPERVISED CLUSTERING SCHEME DENSITY ESTIMATION IS A VERY SIMPLE CONCEPT AND MOST PEOPLE ARE ALREADY FAMILIAR WITH ONE COMMON DENSITY ESTIMATION TECHNIQUE THE HISTOGRAM DENSITY ESTIMATION HISTOGRAMS A HISTOGRAM IS A SIMPLE VISUALIZATION OF DATA WHERE BINS ARE DEFINED AND THE NUMBER OF DATA POINTS WITHIN EACH BIN IS TALLIED AN EXAMPLE OF A HISTOGRAM CAN BE SEEN IN THE UPPERLEFT PANEL OF THE FOLLOWING FIGURE



SCIKITLEARN USER GUIDE RELEASE 0213

A MAJOR PROBLEM WITH HISTOGRAMS HOWEVER IS THAT THE CHOICE OF BINNING CAN HAVE A DISPROPORTIONATE EFFECT ON THE RESULTING VISUALIZATION CONSIDER THE UPPERRIGHT PANEL OF THE ABOVE FIGURE IT SHOWS A HISTOGRAM OVER THE SAME DATA WITH THE BINS SHIFTED RIGHT THE RESULTS OF THE TWO VISUALIZATIONS LOOK ENTIRELY DIFFERENT AND MIGHT LEAD TO DIFFERENT INTERPRETATIONS OF THE DATA

INTUITIVELY ONE CAN ALSO THINK OF A HISTOGRAM AS A STACK OF BLOCKS ONE BLOCK PER POINT BY STACKING THE BLOCKS IN THE APPROPRIATE GRID SPACE WE RECOVER THE HISTOGRAM BUT WHAT IF INSTEAD OF STACKING THE BLOCKS ON A REGULAR GRID WE CENTER EACH BLOCK ON THE POINT IT REPRESENTS AND SUM THE TOTAL HEIGHT AT EACH LOCATION THIS IDEA LEADS TO THE LOWERLEFT VISUALIZATION IT IS PERHAPS NOT AS CLEAN AS A HISTOGRAM BUT THE FACT THAT THE DATA DRIVE THE BLOCK LOCATIONS MEAN THAT IT IS A MUCH BETTER REPRESENTATION OF THE UNDERLYING DATA

THIS VISUALIZATION IS AN EXAMPLE OF A KERNEL DENSITY ESTIMATION IN THIS CASE WITH A TOPHAT KERNEL IE A SQUARE BLOCK AT EACH POINT WE CAN RECOVER A SMOOTHER DISTRIBUTION BY USING A SMOOTHER KERNEL THE BOTTOMRIGHT PLOT SHOWS A GAUSSIAN KERNEL DENSITY ESTIMATE IN WHICH EACH POINT CONTRIBUTES A GAUSSIAN CURVE TO THE TOTAL THE RESULT IS A SMOOTH DENSITY ESTIMATE WHICH IS DERIVED FROM THE DATA AND FUNCTIONS AS A POWERFUL NONPARAMETRIC MODEL OF THE DISTRIBUTION OF POINTS

KERNEL DENSITY ESTIMATION

KERNEL DENSITY ESTIMATION IN SCIKITLEARN IS IMPLEMENTED IN THE SKLEARNNEIGHBORSKERNELDENSITY ESTIMATOR WHICH USES THE BALL TREE OR KD TREE FOR EFFICIENT QUERIES SEE NEAREST NEIGHBORS FOR A DISCUSSION OF THESE THOUGH THE ABOVE EXAMPLE USES A 1D DATA SET FOR SIMPLICITY KERNEL DENSITY ESTIMATION CAN BE PERFORMED IN ANY NUMBER OF DIMENSIONS THOUGH IN PRACTICE THE CURSE OF DIMENSIONALITY CAUSES ITS PERFORMANCE TO DEGRADE IN HIGH DIMENSIONS IN THE FOLLOWING FIGURE 100 POINTS ARE DRAWN FROM A BIMODAL DISTRIBUTION AND THE KERNEL DENSITY ESTIMATES ARE SHOWN FOR THREE CHOICES OF KERNELS

IT’S CLEAR HOW THE KERNEL SHAPE AFFECTS THE SMOOTHNESS OF THE RESULTING DISTRIBUTION THE SCIKITLEARN KERNEL DENSITY ESTIMATOR CAN BE USED AS FOLLOWS

```
FROM SKLEARNNEIGHBORSKDE IMPORT KERNELDENSITY
IMPORT NUMPY AS NP
X = NPARRAY(1 1 2 1 3 2 1 1 2 1 3 2)
KDE = KERNELDENSITY(KERNELGAUSSIAN, BANDWIDTH=0.2, FIT_X=
KDE.SAMPLES)
ARRAY(0.41075698 0.41075698 0.41076071 0.41075698 0.41075698
0.41076071)
```

HERE WE HAVE USED KERNELGAUSSIAN AS SEEN ABOVE MATHEMATICALLY A KERNEL IS A POSITIVE FUNCTION  $K(x, x')$  WHICH IS CONTROLLED BY THE BANDWIDTH PARAMETER  $h$  GIVEN THIS KERNEL FORM THE DENSITY ESTIMATE AT A POINT  $x$  WITHIN A GROUP OF POINTS  $\{x_1, \dots, x_n\}$  IS GIVEN BY

$$\hat{p}(x) = \frac{1}{n} \sum_{i=1}^n K(x, x_i) / h$$

THE BANDWIDTH HERE ACTS AS A SMOOTHING PARAMETER CONTROLLING THE TRADEOFF BETWEEN BIAS AND VARIANCE IN THE RESULT A LARGE BANDWIDTH LEADS TO A VERY SMOOTH IE HIGHBIAS DENSITY DISTRIBUTION A SMALL BANDWIDTH LEADS TO AN UNSMOOTH IE HIGHVARIANCE DENSITY DISTRIBUTION

SKLEARNNEIGHBORSKERNELDENSITY IMPLEMENTS SEVERAL COMMON KERNEL FORMS WHICH ARE SHOWN IN THE FOLLOWING FIGURE

SCIKITLEARN USER GUIDE RELEASE 0213  
THE FORM OF THESE KERNELS IS AS FOLLOWS  
• GAUSSIAN KERNEL KERNEL GAUSSIAN

$$K(x,y) = \frac{1}{\sqrt{2\pi}h} \exp\left(-\frac{\|x-y\|^2}{2h^2}\right)$$
  
• TOPHAT KERNEL KERNEL TOPHAT

$$K(x,y) = \frac{1}{h} \max\left(0, 1 - \frac{\|x-y\|}{h}\right)$$
  
• EPANECHNIKOV KERNEL KERNEL EPANECHNIKOV

$$K(x,y) = \frac{1}{2h} \exp\left(-\frac{\|x-y\|}{h}\right)$$
  
• EXPONENTIAL KERNEL KERNEL EXPONENTIAL

$$K(x,y) = \frac{1}{h} \exp\left(-\frac{\|x-y\|}{h}\right)$$
  
• LINEAR KERNEL KERNEL LINEAR

$$K(x,y) = \frac{1}{2h} \exp\left(-\frac{\|x-y\|}{h}\right)$$
  
• COSINE KERNEL KERNEL COSINE

$$K(x,y) = \frac{1}{2h} \exp\left(-\frac{\|x-y\|}{h}\right)$$
  
THE KERNEL DENSITY ESTIMATOR CAN BE USED WITH ANY OF THE VALID DISTANCE METRICS SEE SKLEARNNEIGHBORS  
DISTANCEMETRIC FOR A LIST OF AVAILABLE METRICS THOUGH THE RESULTS ARE PROPERLY NORMALIZED ONLY FOR THE EUCLIDEAN  
METRIC ONE PARTICULARLY USEFUL METRIC IS THE HAVERSINE DISTANCE WHICH MEASURES THE ANGULAR DISTANCE BETWEEN POINTS  
ON A SPHERE HERE IS AN EXAMPLE OF USING A KERNEL DENSITY ESTIMATE FOR A VISUALIZATION OF GEOSPATIAL DATA IN THIS CASE  
THE DISTRIBUTION OF OBSERVATIONS OF TWO DIFFERENT SPECIES ON THE SOUTH AMERICAN CONTINENT  
32 UNSUPERVISED LEARNING 429

SCIKITLEARN USER GUIDE RELEASE 0213

ONE OTHER USEFUL APPLICATION OF KERNEL DENSITY ESTIMATION IS TO LEARN A NONPARAMETRIC GENERATIVE MODEL OF A DATASET IN ORDER TO EFFICIENTLY DRAW NEW SAMPLES FROM THIS GENERATIVE MODEL HERE IS AN EXAMPLE OF USING THIS PROCESS TO CREATE A NEW SET OF HANDWRITTEN DIGITS USING A GAUSSIAN KERNEL LEARNED ON A PCA PROJECTION OF THE DATA  
THE “NEW” DATA CONSISTS OF LINEAR COMBINATIONS OF THE INPUT DATA WITH WEIGHTS PROBABILISTICALLY DRAWN GIVEN THE KDE MODEL

430 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES

- SIMPLE 1D KERNEL DENSITY ESTIMATION COMPUTATION OF SIMPLE KERNEL DENSITY ESTIMATES IN ONE DIMENSION
- KERNEL DENSITY ESTIMATION AN EXAMPLE OF USING KERNEL DENSITY ESTIMATION TO LEARN A GENERATIVE MODEL OF THE HANDWRITTEN DIGITS DATA AND DRAWING NEW SAMPLES FROM THIS MODEL
- KERNEL DENSITY ESTIMATE OF SPECIES DISTRIBUTIONS AN EXAMPLE OF KERNEL DENSITY ESTIMATION USING THE HAVER SINE DISTANCE METRIC TO VISUALIZE GEOSPATIAL DATA

329 NEURAL NETWORK MODELS UNSUPERVISED

RESTRICTED BOLTZMANN MACHINES

RESTRICTED BOLTZMANN MACHINES RBM ARE UNSUPERVISED NONLINEAR FEATURE LEARNERS BASED ON A PROBABILISTIC MODEL THE FEATURES EXTRACTED BY AN RBM OR A HIERARCHY OF RBMS OFTEN GIVE GOOD RESULTS WHEN FED INTO A LINEAR CLASSIFIER SUCH AS A LINEAR SVM OR A PERCEPTRON

THE MODEL MAKES ASSUMPTIONS REGARDING THE DISTRIBUTION OF INPUTS AT THE MOMENT SCIKITLEARN ONLY PROVIDES BERNOULLIRBM WHICH ASSUMES THE INPUTS ARE EITHER BINARY VALUES OR VALUES BETWEEN 0 AND 1 EACH ENCODING THE PROBABILITY THAT THE SPECIFIC FEATURE WOULD BE TURNED ON

THE RBM TRIES TO MAXIMIZE THE LIKELIHOOD OF THE DATA USING A PARTICULAR GRAPHICAL MODEL THE PARAMETER LEARNING ALGORITHM USED STOCHASTIC MAXIMUM LIKELIHOOD PREVENTS THE REPRESENTATIONS FROM STRAYING FAR FROM THE INPUT DATA WHICH MAKES THEM CAPTURE INTERESTING REGULARITIES BUT MAKES THE MODEL LESS USEFUL FOR SMALL DATASETS AND USUALLY NOT USEFUL FOR DENSITY ESTIMATION

THE METHOD GAINED POPULARITY FOR INITIALIZING DEEP NEURAL NETWORKS WITH THE WEIGHTS OF INDEPENDENT RBMS THIS METHOD IS KNOWN AS UNSUPERVISED PRETRAINING

EXAMPLES

- RESTRICTED BOLTZMANN MACHINE FEATURES FOR DIGIT CLASSIFICATION
- GRAPHICAL MODEL AND PARAMETRIZATION
- THE GRAPHICAL MODEL OF AN RBM IS A FULLYCONNECTED BIPARTITE GRAPH
- 32 UNSUPERVISED LEARNING 431

SCIKITLEARN USER GUIDE RELEASE 0213

THE NODES ARE RANDOM VARIABLES WHOSE STATES DEPEND ON THE STATE OF THE OTHER NODES THEY ARE CONNECTED TO THE MODEL IS THEREFORE PARAMETERIZED BY THE WEIGHTS OF THE CONNECTIONS AS WELL AS ONE INTERCEPT BIAS TERM FOR EACH VISIBLE AND HIDDEN UNIT OMITTED FROM THE IMAGE FOR SIMPLICITY

THE ENERGY FUNCTION MEASURES THE QUALITY OF A JOINT ASSIGNMENT

$E(\mathbf{v}, \mathbf{h}) = -\sum_i v_i \sum_j w_{ij} h_j - \sum_i b_i v_i - \sum_j c_j h_j$

IN THE FORMULA ABOVE  $\mathbf{b}$  AND  $\mathbf{c}$  ARE THE INTERCEPT VECTORS FOR THE VISIBLE AND HIDDEN LAYERS RESPECTIVELY THE JOINT

432 CHAPTER 3 USER GUIDE

PROBABILITY OF THE MODEL IS DEFINED IN TERMS OF THE ENERGY

$$E_{\text{model}} = -\sum_i \sum_j w_{ij} x_i y_j - \sum_i b_i x_i - \sum_j c_j y_j$$

where

THE WORD RESTRICTED REFERS TO THE BIPARTITE STRUCTURE OF THE MODEL WHICH PROHIBITS DIRECT INTERACTION BETWEEN HIDDEN UNITS OR BETWEEN VISIBLE UNITS THIS MEANS THAT THE FOLLOWING CONDITIONAL INDEPENDENCIES ARE ASSUMED

$$p(y_j | x_i, y_{-j}) = p(y_j | x_i)$$

$$p(x_i | y_j, x_{-i}) = p(x_i | y_j)$$

THE BIPARTITE STRUCTURE ALLOWS FOR THE USE OF EFFICIENT BLOCK GIBBS SAMPLING FOR INFERENCE

BERNOULLI RESTRICTED BOLTZMANN MACHINES

IN THEBERNOULLIRBM ALL UNITS ARE BINARY STOCHASTIC UNITS THIS MEANS THAT THE INPUT DATA SHOULD EITHER BE BINARY OR REALVALUED BETWEEN 0 AND 1 SIGNIFYING THE PROBABILITY THAT THE VISIBLE UNIT WOULD TURN ON OR OFF THIS IS A GOOD MODEL FOR CHARACTER RECOGNITION WHERE THE INTEREST IS ON WHICH PIXELS ARE ACTIVE AND WHICH AREN'T FOR IMAGES OF NATURAL SCENES IT NO LONGER FITS BECAUSE OF BACKGROUND DEPTH AND THE TENDENCY OF NEIGHBOURING PIXELS TO TAKE THE SAME VALUES THE CONDITIONAL PROBABILITY DISTRIBUTION OF EACH UNIT IS GIVEN BY THE LOGISTIC SIGMOID ACTIVATION FUNCTION OF THE INPUT IT RECEIVES

$$p(x_i = 1 | \mathbf{h}) = \sigma(\sum_j w_{ij} h_j + b_i)$$

$$p(h_j = 1 | \mathbf{x}) = \sigma(\sum_i w_{ji} x_i + c_j)$$

$$p(x_i = 1 | \mathbf{h}) = \sigma(\sum_j w_{ij} h_j + b_i)$$

$$p(h_j = 1 | \mathbf{x}) = \sigma(\sum_i w_{ji} x_i + c_j)$$

WHERE  $\sigma$  IS THE LOGISTIC SIGMOID FUNCTION

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

$$1 - \sigma(z) = \sigma(-z)$$

STOCHASTIC MAXIMUM LIKELIHOOD LEARNING

THE TRAINING ALGORITHM IMPLEMENTED IN BERNOLLIIRBM IS KNOWN AS STOCHASTIC MAXIMUM LIKELIHOOD SML OR PERSISTENT CONTRASTIVE DIVERGENCE PCD OPTIMIZING MAXIMUM LIKELIHOOD DIRECTLY IS INFEASIBLE BECAUSE OF THE FORM OF THE DATA LIKELIHOOD

$$\log p(\mathbf{x}, \mathbf{h}) = \log \sum_{\mathbf{h}} p(\mathbf{x}, \mathbf{h})$$

$$h_j - \sum_i w_{ji} x_i - c_j - \log \sum_i \exp(h_j - \sum_i w_{ji} x_i - c_j)$$

$$\frac{\partial}{\partial w_{ji}} \log p(\mathbf{x}, \mathbf{h}) = h_j - \sum_i w_{ji} x_i - c_j$$

FOR SIMPLICITY THE EQUATION ABOVE IS WRITTEN FOR A SINGLE TRAINING EXAMPLE THE GRADIENT WITH RESPECT TO THE WEIGHTS IS FORMED OF TWO TERMS CORRESPONDING TO THE ONES ABOVE THEY ARE USUALLY KNOWN AS THE POSITIVE GRADIENT AND THE NEGATIVE GRADIENT BECAUSE OF THEIR RESPECTIVE SIGNS IN THIS IMPLEMENTATION THE GRADIENTS ARE ESTIMATED OVER MINIBATCHES OF SAMPLES

IN MAXIMIZING THE LOGLIKELIHOOD THE POSITIVE GRADIENT MAKES THE MODEL PREFER HIDDEN STATES THAT ARE COMPATIBLE WITH THE OBSERVED TRAINING DATA BECAUSE OF THE BIPARTITE STRUCTURE OF RBMS IT CAN BE COMPUTED EFFICIENTLY THE NEGATIVE GRADIENT HOWEVER IS INTRACTABLE ITS GOAL IS TO LOWER THE ENERGY OF JOINT STATES THAT THE MODEL PREFERS THEREFORE MAKING IT STAY TRUE TO THE DATA IT CAN BE APPROXIMATED BY MARKOV CHAIN MONTE CARLO USING BLOCK GIBBS SAMPLING BY ITERATIVELY SAMPLING EACH OF  $\mathbf{x}$  AND  $\mathbf{h}$  GIVEN THE OTHER UNTIL THE CHAIN MIXES SAMPLES GENERATED IN THIS WAY ARE SOMETIMES REFERRED AS FANTASY PARTICLES THIS IS INEFFICIENT AND IT IS DIFFICULT TO DETERMINE WHETHER THE MARKOV CHAIN MIXES THE CONTRASTIVE DIVERGENCE METHOD SUGGESTS TO STOP THE CHAIN AFTER A SMALL NUMBER OF ITERATIONS  $k$  USUALLY EVEN 1 THIS METHOD IS FAST AND HAS LOW VARIANCE BUT THE SAMPLES ARE FAR FROM THE MODEL DISTRIBUTION

SCIKITLEARN USER GUIDE RELEASE 0213

PERSISTENT CONTRASTIVE DIVERGENCE ADDRESSES THIS INSTEAD OF STARTING A NEW CHAIN EACH TIME THE GRADIENT IS NEEDED AND PERFORMING ONLY ONE GIBBS SAMPLING STEP IN PCD WE KEEP A NUMBER OF CHAINS FANTASY PARTICLES THAT ARE UPDATED □ GIBBS STEPS AFTER EACH WEIGHT UPDATE THIS ALLOWS THE PARTICLES TO EXPLORE THE SPACE MORE THOROUGHLY REFERENCES

- “A FAST LEARNING ALGORITHM FOR DEEP BELIEF NETS” G HINTON S OSINDERO Y W TEH 2006
- “TRAINING RESTRICTED BOLTZMANN MACHINES USING APPROXIMATIONS TO THE LIKELIHOOD GRADIENT” T TIELEMAN 2008

33 MODEL SELECTION AND EVALUATION

331 CROSSVALIDATION EVALUATING ESTIMATOR PERFORMANCE

LEARNING THE PARAMETERS OF A PREDICTION FUNCTION AND TESTING IT ON THE SAME DATA IS A METHODOLOGICAL MISTAKE A MODEL THAT WOULD JUST REPEAT THE LABELS OF THE SAMPLES THAT IT HAS JUST SEEN WOULD HAVE A PERFECT SCORE BUT WOULD FAIL TO PREDICT ANYTHING USEFUL ON YETUNSEEN DATA THIS SITUATION IS CALLED OVERFITTING TO AVOID IT IT IS COMMON PRACTICE WHEN PERFORMING A SUPERVISED MACHINE LEARNING EXPERIMENT TO HOLD OUT PART OF THE AVAILABLE DATA AS A TEST SETXTEST YTEST NOTE THAT THE WORD “EXPERIMENT” IS NOT INTENDED TO DENOTE ACADEMIC USE ONLY BECAUSE EVEN IN COMMERCIAL SETTINGS MACHINE LEARNING USUALLY STARTS OUT EXPERIMENTALLY HERE IS A FLOWCHART OF TYPICAL CROSS VALIDATION WORKFLOW IN MODEL TRAINING THE BEST PARAMETERS CAN BE DETERMINED BY GRID SEARCH TECHNIQUES IN SCIKITLEARN A RANDOM SPLIT INTO TRAINING AND TEST SETS CAN BE QUICKLY COMPUTED WITH THE TRAINTESTSPLIT HELPER FUNCTION LET’S LOAD THE IRIS DATA SET TO FIT A LINEAR SUPPORT VECTOR MACHINE ON IT

```
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn import datasets
from sklearn import svm
iris = datasets.load_iris()
iris_data_shape, iris_target_shape
```

150, 4, 150
WE CAN NOW QUICKLY SAMPLE A TRAINING SET WHILE HOLDING OUT 40 OF THE DATA FOR TESTING EVALUATING OUR CLASSIFIER
434 CHAPTER 3 USER GUIDE



```
SCIKITLEARN USER GUIDE RELEASE 0213
XTRAIN XTEST YTRAIN YTEST  TRAINTESTSPLIT
IRISDATA IRISTARGET TESTSIZE04 RANDOMSTATE0
XTRAINSHAPE YTRAINSHAPE
90 4 90
XTESTSHAPE YTESTSHAPE
60 4 60
CLF  SVMSVCKERNELLINEAR C1FITXTRAIN YTRAIN
CLFSCOREXTEST YTEST
096
```

WHEN EVALUATING DIFFERENT SETTINGS “HYPERPARAMETERS” FOR ESTIMATORS SUCH AS THE CSETTING THAT MUST BE MANUALLY SET FOR AN SVM THERE IS STILL A RISK OF OVERFITTING ON THE TEST SET BECAUSE THE PARAMETERS CAN BE TWEAKED UNTIL THE ESTIMATOR PERFORMS OPTIMALLY THIS WAY KNOWLEDGE ABOUT THE TEST SET CAN “LEAK” INTO THE MODEL AND EVALUATION METRICS NO LONGER REPORT ON GENERALIZATION PERFORMANCE TO SOLVE THIS PROBLEM YET ANOTHER PART OF THE DATASET CAN BE HELD OUT AS A SO CALLED “VALIDATION SET” TRAINING PROCEEDS ON THE TRAINING SET AFTER WHICH EVALUATION IS DONE ON THE VALIDATION SET AND WHEN THE EXPERIMENT SEEMS TO BE SUCCESSFUL FINAL EVALUATION CAN BE DONE ON THE TEST SET HOWEVER BY PARTITIONING THE AVAILABLE DATA INTO THREE SETS WE DRASTICALLY REDUCE THE NUMBER OF SAMPLES WHICH CAN BE USED FOR LEARNING THE MODEL AND THE RESULTS CAN DEPEND ON A PARTICULAR RANDOM CHOICE FOR THE PAIR OF TRAIN VALIDATION SETS

A SOLUTION TO THIS PROBLEM IS A PROCEDURE CALLED CROSSVALIDATION CV FOR SHORT A TEST SET SHOULD STILL BE HELD OUT FOR FINAL EVALUATION BUT THE VALIDATION SET IS NO LONGER NEEDED WHEN DOING CV IN THE BASIC APPROACH CALLED KFOLD CV THE TRAINING SET IS SPLIT INTO KSMALLER SETS OTHER APPROACHES ARE DESCRIBED BELOW BUT GENERALLY FOLLOW THE SAME PRINCIPLES THE FOLLOWING PROCEDURE IS FOLLOWED FOR EACH OF THE K“FOLDS”

- A MODEL IS TRAINED USING  $\frac{n-1}{k}$  OF THE FOLDS AS TRAINING DATA
  - THE RESULTING MODEL IS VALIDATED ON THE REMAINING PART OF THE DATA IE IT IS USED AS A TEST SET TO COMPUTE A PERFORMANCE MEASURE SUCH AS ACCURACY
- THE PERFORMANCE MEASURE REPORTED BY KFOLD CROSSVALIDATION IS THEN THE AVERAGE OF THE VALUES COMPUTED IN THE LOOP THIS APPROACH CAN BE COMPUTATIONALLY EXPENSIVE BUT DOES NOT WASTE TOO MUCH DATA AS IS THE CASE WHEN FIXING AN ARBITRARY VALIDATION SET WHICH IS A MAJOR ADVANTAGE IN PROBLEMS SUCH AS INVERSE INFERENCE WHERE THE NUMBER OF SAMPLES IS VERY SMALL

SCIKITLEARN USER GUIDE RELEASE 0213

COMPUTING CROSSVALIDATED METRICS

THE SIMPLEST WAY TO USE CROSSVALIDATION IS TO CALL THE CROSSVALSCORE HELPER FUNCTION ON THE ESTIMATOR AND THE DATASET

THE FOLLOWING EXAMPLE DEMONSTRATES HOW TO ESTIMATE THE ACCURACY OF A LINEAR KERNEL SUPPORT VECTOR MACHINE ON THE IRIS DATASET BY SPLITTING THE DATA FITTING A MODEL AND COMPUTING THE SCORE 5 CONSECUTIVE TIMES WITH DIFFERENT SPLITS EACH TIME

```
FROM SKLEARNMODELSELECTION IMPORT CROSSVALSCORE
CLF SVMSVCKERNELLINEAR C1
SCORES CROSSVALSCORECLF IRISDATA IRISTARGET CV5
SCORES
ARRAY096 1 096 096 1
```

THE MEAN SCORE AND THE 95 CONFIDENCE INTERVAL OF THE SCORE ESTIMATE ARE HENCE GIVEN BY

```
PRINTACCURACY 02F02F SCORESMEAN SCORESSTD 2
ACCURACY 098 003
```

BY DEFAULT THE SCORE COMPUTED AT EACH CV ITERATION IS THE SCORE METHOD OF THE ESTIMATOR IT IS POSSIBLE TO CHANGE THIS BY USING THE SCORING PARAMETER

```
FROM SKLEARN IMPORT METRICS
SCORES CROSSVALSCORE
CLF IRISDATA IRISTARGET CV5 SCORINGF1MACRO
SCORES
ARRAY096 1 096 096 1
```

SEETHE SCORING PARAMETER DEFINING MODEL EVALUATION RULES FOR DETAILS IN THE CASE OF THE IRIS DATASET THE SAMPLES ARE BALANCED ACROSS TARGET CLASSES HENCE THE ACCURACY AND THE F1SCORE ARE ALMOST EQUAL

WHEN THECVARGUMENT IS AN INTEGER CROSSVALSCORE USES THEKFOLD ORSTRATIFIEDKFOLD STRATEGIES BY DEFAULT THE LATTER BEING USED IF THE ESTIMATOR DERIVES FROM CLASSIFIERMIXIN

IT IS ALSO POSSIBLE TO USE OTHER CROSS VALIDATION STRATEGIES BY PASSING A CROSS VALIDATION ITERATOR INSTEAD FOR INSTANCE

```
FROM SKLEARNMODELSELECTION IMPORT SHUFFLESPLIT
NSAMPLES IRISDATASHAPE0
CV SHUFFLESPLITNSPLITS5 TESTSIZE03 RANDOMSTATE0
CROSSVALSCORECLF IRISDATA IRISTARGET CVCV
ARRAY0977 0977 1 0955 1
```

ANOTHER OPTION IS TO USE AN ITERABLE YIELDING TRAIN TEST SPLITS AS ARRAYS OF INDICES FOR EXAMPLE

```
DEF CUSTOMCV2FOLDSEX
N XSHAPE0
I 1
WHILE I 2
IDX NPARANGEN I 1 2 N I 2 DTYPEINT
YIELD IDX IDX
I 1
```

```
CUSTOMCV CUSTOMCV2FOLDSSIRISDATA
CROSSVALSCORECLF IRISDATA IRISTARGET CVCUSTOMCV
ARRAY1 0973
```

436 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

DATA TRANSFORMATION WITH HELD OUT DATA

JUST AS IT IS IMPORTANT TO TEST A PREDICTOR ON DATA HELDOUT FROM TRAINING PREPROCESSING SUCH AS STANDARDIZATION  
FEATURE SELECTION ETC AND SIMILAR DATA TRANSFORMATIONS SIMILARLY SHOULD BE LEARNT FROM A TRAINING SET AND APPLIED  
TO HELDOUT DATA FOR PREDICTION

FROM SKLEARN IMPORT PREPROCESSING

XTRAIN XTEST YTRAIN YTEST TRAINTESTSPLIT

IRISDATA IRISTARGET TESTSIZE04 RANDOMSTATE0

SCALER PREPROCESSINGSTANDARDSCALERFITXTRAIN

XTRAINTRANSFORMED SCALERTRANSFORMXTRAIN

CLF SVMSCVC1FITXTRAINTRANSFORMED YTRAIN

XTESTTRANSFORMED SCALERTRANSFORMXTEST

CLFSCOREXTESTTRANSFORMED YTEST

09333

APIPELNE MAKES IT EASIER TO COMPOSE ESTIMATORS PROVIDING THIS BEHAVIOR UNDER CROSSVALIDATION

FROM SKLEARNPIPELINE IMPORT MAKEPIPELINE

CLF MAKEPIPELINEPREPROCESSINGSTANDARDSCALER SVMSCVC1

CROSSVALSCORECLF IRISDATA IRISTARGET CVCV

ARRAY0977 0933 0955 0933 0977

SEPIPELINES AND COMPOSITE ESTIMATORS

THE CROSSVALIDATE FUNCTION AND MULTIPLE METRIC EVALUATION

THECROSSVALIDATE FUNCTION DIFFERS FROM CROSSVALSCORE IN TWO WAYS

- IT ALLOWS SPECIFYING MULTIPLE METRICS FOR EVALUATION
- IT RETURNS A DICT CONTAINING FITTIMES SCORETIMES AND OPTIONALLY TRAINING SCORES AS WELL AS FITTED ESTIMATORS IN  
ADDITION TO THE TEST SCORE

FOR SINGLE METRIC EVALUATION WHERE THE SCORING PARAMETER IS A STRING CALLABLE OR NONE THE KEYS WILL BE  
TESTSCORE FITTIME SCORETIME

AND FOR MULTIPLE METRIC EVALUATION THE RETURN VALUE IS A DICT WITH THE FOLLOWING KEYS

TESTSCORER1NAME TESTSCORER2NAME TESTSCORER FITTIME  
SCORETIME

RETURNTRAINSCORE IS SET TOFALSE BY DEFAULT TO SAVE COMPUTATION TIME TO EVALUATE THE SCORES ON THE TRAINING  
SET AS WELL YOU NEED TO BE SET TO TRUE

YOU MAY ALSO RETAIN THE ESTIMATOR FITTED ON EACH TRAINING SET BY SETTING RETURNESTIMATORTRUE

THE MULTIPLE METRICS CAN BE SPECIFIED EITHER AS A LIST TUPLE OR SET OF PREDEFINED SCORER NAMES

FROM SKLEARNMODELSELECTION IMPORT CROSSVALIDATE

FROM SKLEARNMETRICS IMPORT RECALLSCORE

SCORING PRECISIONMACRO RECALLMACRO

CLF SVMSCVKERNELLINEAR C1 RANDOMSTATE0

SCORES CROSSVALIDATECLF IRISDATA IRISTARGET SCORINGSCORING

CV5

SORTEDSCORESKEYS

FITTIME SCORETIME TESTPRECISIONMACRO TESTRECALLMACRO

SCORESTESTRECALLMACRO

ARRAY096 1 096 096 1

33 MODEL SELECTION AND EVALUATION 437

SCIKITLEARN USER GUIDE RELEASE 0213

OR AS A DICT MAPPING SCORER NAME TO A PREDEFINED OR CUSTOM SCORING FUNCTION

FROM SKLEARNMETRICSSCORER IMPORT MAKESCORER

SCORING PRECMACRO PRECISIONMACRO

RECMACRO MAKESCORERRECALLSCORE AVERAGEMACRO

SCORES CROSSVALIDATECLF IRISDATA IRISTARGET SCORINGSCORING

CV5 RETURNTRAINSCORE TRUE

SORTEDSCORESKEYS

FITTIME SCORETIME TESTPRECMACRO TESTRECMACRO

TRAINPRECMACRO TRAINRECMACRO

SCORESTRAINRECMACRO

ARRAY097 097 099 098 098

HERE IS AN EXAMPLE OF CROSSVALIDATE USING A SINGLE METRIC

SCORES CROSSVALIDATECLF IRISDATA IRISTARGET

SCORINGPRECISIONMACRO CV5

RETURNESTIMATOR TRUE

SORTEDSCORESKEYS

ESTIMATOR FITTIME SCORETIME TESTSCORE

OBTAINING PREDICTIONS BY CROSSVALIDATION

THE FUNCTION CROSSVALPREDICT HAS A SIMILAR INTERFACE TO CROSSVALSCORE BUT RETURNS FOR EACH ELEMENT IN THE INPUT THE PREDICTION THAT WAS OBTAINED FOR THAT ELEMENT WHEN IT WAS IN THE TEST SET ONLY CROSSVALIDATION STRATEGIES THAT ASSIGN ALL ELEMENTS TO A TEST SET EXACTLY ONCE CAN BE USED OTHERWISE AN EXCEPTION IS RAISED

WARNING NOTE ON INAPPROPRIATE USAGE OF CROSSVALPREDICT

THE RESULT OF CROSSVALPREDICT MAY BE DIFFERENT FROM THOSE OBTAINED USING CROSSVALSCORE AS THE ELEMENTS ARE GROUPED IN DIFFERENT WAYS THE FUNCTION CROSSVALSCORE TAKES AN AVERAGE OVER CROSSVALIDATION FOLDS WHEREAS CROSSVALPREDICT SIMPLY RETURNS THE LABELS OR PROBABILITIES FROM SEVERAL DISTINCT MODELS UNDISTINGUISHED THUS CROSSVALPREDICT IS NOT AN APPROPRIATE MEASURE OF GENERALISATION ERROR

THE FUNCTION CROSSVALPREDICT IS APPROPRIATE FOR

- VISUALIZATION OF PREDICTIONS OBTAINED FROM DIFFERENT MODELS
- MODEL BLENDING WHEN PREDICTIONS OF ONE SUPERVISED ESTIMATOR ARE USED TO TRAIN ANOTHER ESTIMATOR IN ENSEMBLE METHODS

THE AVAILABLE CROSS VALIDATION ITERATORS ARE INTRODUCED IN THE FOLLOWING SECTION

EXAMPLES

- RECEIVER OPERATING CHARACTERISTIC ROC WITH CROSS VALIDATION
- RECURSIVE FEATURE ELIMINATION WITH CROSSVALIDATION
- PARAMETER ESTIMATION USING GRID SEARCH WITH CROSSVALIDATION
- SAMPLE PIPELINE FOR TEXT FEATURE EXTRACTION AND EVALUATION
- PLOTING CROSSVALIDATED PREDICTIONS
- NESTED VERSUS NONNESTED CROSSVALIDATION

438 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

CROSS VALIDATION ITERATORS

THE FOLLOWING SECTIONS LIST UTILITIES TO GENERATE INDICES THAT CAN BE USED TO GENERATE DATASET SPLITS ACCORDING TO DIFFERENT CROSS VALIDATION STRATEGIES

CROSSVALIDATION ITERATORS FOR IID DATA

ASSUMING THAT SOME DATA IS INDEPENDENT AND IDENTICALLY DISTRIBUTED IID IS MAKING THE ASSUMPTION THAT ALL SAMPLES STEM FROM THE SAME GENERATIVE PROCESS AND THAT THE GENERATIVE PROCESS IS ASSUMED TO HAVE NO MEMORY OF PAST GENERATED SAMPLES

THE FOLLOWING CROSSVALIDATORS CAN BE USED IN SUCH CASES

NOTE

WHILE IID DATA IS A COMMON ASSUMPTION IN MACHINE LEARNING THEORY IT RARELY HOLDS IN PRACTICE IF ONE KNOWS THAT THE SAMPLES HAVE BEEN GENERATED USING A TIMEDEPENDENT PROCESS IT’S SAFER TO USE A TIMESERIES AWARE CROSSVALIDATION SCHEME SIMILARLY IF WE KNOW THAT THE GENERATIVE PROCESS HAS A GROUP STRUCTURE SAMPLES FROM COLLECTED FROM DIFFERENT SUBJECTS EXPERIMENTS MEASUREMENT DEVICES IT SAFER TO USE GROUPWISE CROSSVALIDATION

KFOLD

KFOLD DIVIDES ALL THE SAMPLES IN  $n$  GROUPS OF SAMPLES CALLED FOLDS IF  $n \leq 1$  THIS IS EQUIVALENT TO THE LEAVE ONE OUT STRATEGY OF EQUAL SIZES IF POSSIBLE THE PREDICTION FUNCTION IS LEARNED USING  $n-1$  FOLDS AND THE FOLD LEFT OUT IS USED FOR TEST

EXAMPLE OF 2FOLD CROSSVALIDATION ON A DATASET WITH 4 SAMPLES

```
import numpy as np
from sklearn.model_selection import KFold
X = A B C D
kf = KFold(n_splits=2)
for train, test in kf.split(X):
    print(s, 'train test')
2 3 0 1
0 1 2 3
```

HERE IS A VISUALIZATION OF THE CROSSVALIDATION BEHAVIOR NOTE THAT KFOLD IS NOT AFFECTED BY CLASSES OR GROUPS

33 MODEL SELECTION AND EVALUATION 439

SCIKITLEARN USER GUIDE RELEASE 0213

EACH FOLD IS CONSTITUTED BY TWO ARRAYS THE FIRST ONE IS RELATED TO THE TRAINING SET AND THE SECOND ONE TO THE TEST SET  
THUS ONE CAN CREATE THE TRAININGTEST SETS USING NUMPY INDEXING

```
X NPARRAY0 0 1 1 1 1 2 2
Y NPARRAY0 1 0 1
XTRAIN XTEST YTRAIN YTEST XTRAIN XTEST YTRAIN YTEST
```

REPEATED KFOLD

REPEATEDKFOLD REPEATS KFOLD N TIMES IT CAN BE USED WHEN ONE REQUIRES TO RUN KFOLD N TIMES PRODUCING  
DIFFERENT SPLITS IN EACH REPETITION

EXAMPLE OF 2FOLD KFOLD REPEATED 2 TIMES

```
IMPORT NUMPY AS NP
FROM SKLEARNMODELSELECTION IMPORT REPEATEDKFOLD
X NPARRAY1 2 3 4 1 2 3 4
RANDOMSTATE 12883823
RKF REPEATEDKFOLDNSPLITS2 NREPEATS2 RANDOMSTATERRANDOMSTATE
FOR TRAIN TEST INRKFSPLITX
PRINTS S TRAIN TEST
```

```
2 3 0 1
0 1 2 3
0 2 1 3
1 3 0 2
```

SIMILARLYREPEATEDSTRATIFIEDKFOLD REPEATS STRATIFIED KFOLD N TIMES WITH DIFFERENT RANDOMIZATION IN EACH  
REPETITION

LEAVE ONE OUT LOO

LEAVEONEOUT OR LOO IS A SIMPLE CROSSVALIDATION EACH LEARNING SET IS CREATED BY TAKING ALL THE SAMPLES EXCEPT  
ONE THE TEST SET BEING THE SAMPLE LEFT OUT THUS FOR  $\frac{1}{n}$ SAMPLES WE HAVE  $\frac{1}{n}$ DIFFERENT TRAINING SETS AND  $\frac{1}{n}$ DIFFERENT TESTS  
SET THIS CROSSVALIDATION PROCEDURE DOES NOT WASTE MUCH DATA AS ONLY ONE SAMPLE IS REMOVED FROM THE TRAINING SET

```
FROM SKLEARNMODELSELECTION IMPORT LEAVEONEOUT
X 1 2 3 4
LOO LEAVEONEOUT
FOR TRAIN TEST INLOOSPLITX
PRINTS S TRAIN TEST
```

```
1 2 3 0
0 2 3 1
0 1 3 2
0 1 2 3
```

POTENTIAL USERS OF LOO FOR MODEL SELECTION SHOULD WEIGH A FEW KNOWN CAVEATS WHEN COMPARED WITH  $\frac{1}{n}$ FOLD CROSS  
VALIDATION ONE BUILDS  $\frac{1}{n}$ MODELS FROM  $\frac{1}{n}$ SAMPLES INSTEAD OF  $\frac{1}{n}$ MODELS WHERE  $\frac{1}{n}$  MOREOVER EACH IS TRAINED ON  $\frac{1}{n}-1$   
SAMPLES RATHER THAN  $\frac{1}{n}-1$  IN BOTH WAYS ASSUMING  $\frac{1}{n}$ IS NOT TOO LARGE AND  $\frac{1}{n}$  LOO IS MORE COMPUTATIONALLY  
EXPENSIVE THAN  $\frac{1}{n}$ FOLD CROSS VALIDATION

IN TERMS OF ACCURACY LOO OFTEN RESULTS IN HIGH VARIANCE AS AN ESTIMATOR FOR THE TEST ERROR INTUITIVELY SINCE  $\frac{1}{n}-1$ OF  
THE $\frac{1}{n}$ SAMPLES ARE USED TO BUILD EACH MODEL MODELS CONSTRUCTED FROM FOLDS ARE VIRTUALLY IDENTICAL TO EACH OTHER AND TO  
THE MODEL BUILT FROM THE ENTIRE TRAINING SET

440 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

HOWEVER IF THE LEARNING CURVE IS STEEP FOR THE TRAINING SIZE IN QUESTION THEN 5 OR 10 FOLD CROSS VALIDATION CAN OVERESTIMATE THE GENERALIZATION ERROR

AS A GENERAL RULE MOST AUTHORS AND EMPIRICAL EVIDENCE SUGGEST THAT 5 OR 10 FOLD CROSS VALIDATION SHOULD BE PREFERRED TO LOO

REFERENCES

- [HTTPWWWFAQSORGFAQSAIFAQNEURALNETSPART3SECTION12HTML](http://www.faq.sorgfa.org/faq/neuralnets/part3/section12.html)
- T HASTIE R TIBSHIRANI J FRIEDMAN THE ELEMENTS OF STATISTICAL LEARNING SPRINGER 2009
- L BREIMAN P SPECTOR SUBMODEL SELECTION AND EVALUATION IN REGRESSION THE XRANDOM CASE INTERNATIONAL STATISTICAL REVIEW 1992
- R KOHAVI A STUDY OF CROSSVALIDATION AND BOOTSTRAP FOR ACCURACY ESTIMATION AND MODEL SELECTION INTL JNT CONF AI
- R BHARAT RAO G FUNG R ROSALES ON THE DANGERS OF CROSSVALIDATION AN EXPERIMENTAL EVALUATION SIAM 2008
- G JAMES D WITTEN T HASTIE R TIBSHIRANI AN INTRODUCTION TO STATISTICAL LEARNING SPRINGER 2013

LEAVE P OUT LPO

LEAVEOUT IS VERY SIMILAR TO LEAVEONEOUT AS IT CREATES ALL THE POSSIBLE TRAININGTEST SETS BY REMOVING  $\frac{1}{n}$  SAMPLES FROM THE COMPLETE SET FOR  $\frac{1}{n}$  SAMPLES THIS PRODUCES  $\frac{1}{n}$  TRAINTEST PAIRS UNLIKE LEAVEONEOUT ANDKFOLD THE TEST SETS WILL OVERLAP FOR  $\frac{1}{n}$

EXAMPLE OF LEAVE2OUT ON A DATASET WITH 4 SAMPLES

```
FROM SKLEARNMODELSELECTION IMPORT LEAVEPOUT
X NPONES4
LPO LEAVEPOUTP2
FOR TRAIN TEST INLPOSPLITX
PRINTS S TRAIN TEST
2 3 0 1
1 3 0 2
1 2 0 3
0 3 1 2
0 2 1 3
0 1 2 3
```

RANDOM PERMUTATIONS CROSSVALIDATION AKA SHUFFLE SPLIT

SHUFFLESPLIT

THESHUFFLESPLIT ITERATOR WILL GENERATE A USER DEFINED NUMBER OF INDEPENDENT TRAIN TEST DATASET SPLITS SAMPLES ARE FIRST SHUFFLED AND THEN SPLIT INTO A PAIR OF TRAIN AND TEST SETS

IT IS POSSIBLE TO CONTROL THE RANDOMNESS FOR REPRODUCIBILITY OF THE RESULTS BY EXPLICITLY SEEDING THE RANDOMSTATE PSEUDO RANDOM NUMBER GENERATOR

HERE IS A USAGE EXAMPLE

33 MODEL SELECTION AND EVALUATION 441

SCIKITLEARN USER GUIDE RELEASE 0213

```
FROM SKLEARNMODELSELECTION IMPORT SHUFFLESPLIT
X NPARANGE10
SS SHUFFLESPLITNSPLITS5 TESTSIZE025
RANDOMSTATE0
FOR TRAININDEX TESTINDEX INSSSPLITX
PRINTS S TRAININDEX TESTINDEX
9 1 6 7 3 0 5 2 8 4
2 9 8 0 6 7 4 3 5 1
4 5 1 0 6 9 7 2 3 8
2 7 5 8 0 3 4 6 1 9
4 1 0 6 8 9 3 5 2 7
```

HERE IS A VISUALIZATION OF THE CROSSVALIDATION BEHAVIOR NOTE THAT SHUFFLESPLIT IS NOT AFFECTED BY CLASSES OR GROUPS

SHUFFLESPLIT IS THUS A GOOD ALTERNATIVE TO KFOLD CROSS VALIDATION THAT ALLOWS A FINER CONTROL ON THE NUMBER OF ITERATIONS AND THE PROPORTION OF SAMPLES ON EACH SIDE OF THE TRAIN TEST SPLIT

CROSSVALIDATION ITERATORS WITH STRATIFICATION BASED ON CLASS LABELS

SOME CLASSIFICATION PROBLEMS CAN EXHIBIT A LARGE IMBALANCE IN THE DISTRIBUTION OF THE TARGET CLASSES FOR INSTANCE THERE COULD BE SEVERAL TIMES MORE NEGATIVE SAMPLES THAN POSITIVE SAMPLES IN SUCH CASES IT IS RECOMMENDED TO USE STRATIFIED SAMPLING AS IMPLEMENTED IN STRATIFIEDKFOLD ANDSTRATIFIEDSHUFFLESPLIT TO ENSURE THAT RELATIVE CLASS FREQUENCIES IS APPROXIMATELY PRESERVED IN EACH TRAIN AND VALIDATION FOLD

STRATIFIED KFOLD

STRATIFIEDKFOLD IS A VARIATION OF KFOLD WHICH RETURNS STRATIFIED FOLDS EACH SET CONTAINS APPROXIMATELY THE SAME PERCENTAGE OF SAMPLES OF EACH TARGET CLASS AS THE COMPLETE SET

EXAMPLE OF STRATIFIED 3FOLD CROSSVALIDATION ON A DATASET WITH 10 SAMPLES FROM TWO SLIGHTLY UNBALANCED CLASSES

```
FROM SKLEARNMODELSELECTION IMPORT STRATIFIEDKFOLD
X NPONES10
Y 0 0 0 0 1 1 1 1 1 1
SKF STRATIFIEDKFOLDNSPLITS3
FOR TRAIN TEST INSKFSPLITX Y
```

442 CHAPTER 3 USER GUIDE



SCIKITLEARN USER GUIDE RELEASE 0213

PRINTS S TRAIN TEST

2 3 6 7 8 9 0 1 4 5

0 1 3 4 5 8 9 2 6 7

0 1 2 4 5 6 7 3 8 9

HERE IS A VISUALIZATION OF THE CROSSVALIDATION BEHAVIOR

REPEATEDSTRATIFIEDKFOLD CAN BE USED TO REPEAT STRATIFIED KFOLD N TIMES WITH DIFFERENT RANDOMIZATION IN EACH REPETITION

STRATIFIED SHUFFLE SPLIT

STRATIFIEDSHUFFLESPLIT IS A VARIATION OF SHUFFLESPLIT WHICH RETURNS STRATIFIED SPLITS IEWHICH CREATES SPLITS BY PRESERVING THE SAME PERCENTAGE FOR EACH TARGET CLASS AS IN THE COMPLETE SET

HERE IS A VISUALIZATION OF THE CROSSVALIDATION BEHAVIOR

CROSSVALIDATION ITERATORS FOR GROUPED DATA

THE IID ASSUMPTION IS BROKEN IF THE UNDERLYING GENERATIVE PROCESS YIELD GROUPS OF DEPENDENT SAMPLES

33 MODEL SELECTION AND EVALUATION 443

SCIKITLEARN USER GUIDE RELEASE 0213

SUCH A GROUPING OF DATA IS DOMAIN SPECIFIC AN EXAMPLE WOULD BE WHEN THERE IS MEDICAL DATA COLLECTED FROM MULTIPLE PATIENTS WITH MULTIPLE SAMPLES TAKEN FROM EACH PATIENT AND SUCH DATA IS LIKELY TO BE DEPENDENT ON THE INDIVIDUAL GROUP IN OUR EXAMPLE THE PATIENT ID FOR EACH SAMPLE WILL BE ITS GROUP IDENTIFIER

IN THIS CASE WE WOULD LIKE TO KNOW IF A MODEL TRAINED ON A PARTICULAR SET OF GROUPS GENERALIZES WELL TO THE UNSEEN GROUPS TO MEASURE THIS WE NEED TO ENSURE THAT ALL THE SAMPLES IN THE VALIDATION FOLD COME FROM GROUPS THAT ARE NOT REPRESENTED AT ALL IN THE PAIRED TRAINING FOLD

THE FOLLOWING CROSSVALIDATION SPLITTERS CAN BE USED TO DO THAT THE GROUPING IDENTIFIER FOR THE SAMPLES IS SPECIFIED VIA THEGROUPS PARAMETER

GROUP KFOLD

GROUPKFOLD IS A VARIATION OF KFOLD WHICH ENSURES THAT THE SAME GROUP IS NOT REPRESENTED IN BOTH TESTING AND TRAINING SETS FOR EXAMPLE IF THE DATA IS OBTAINED FROM DIFFERENT SUBJECTS WITH SEVERAL SAMPLES PERSUBJECT AND IF THE MODEL IS FLEXIBLE ENOUGH TO LEARN FROM HIGHLY PERSON SPECIFIC FEATURES IT COULD FAIL TO GENERALIZE TO NEW SUBJECTS GROUPKFOLD MAKES IT POSSIBLE TO DETECT THIS KIND OF OVERFITTING SITUATIONS

IMAGINE YOU HAVE THREE SUBJECTS EACH WITH AN ASSOCIATED NUMBER FROM 1 TO 3

```
FROM SKLEARNMODELSELECTION IMPORT GROUPKFOLD
X 01 02 22 24 23 455 58 88 9 10
Y  A B B B C C C D D D
GROUPS 1 1 1 2 2 2 3 3 3 3
GKF  GROUPKFOLDNSPLITS3
FOR TRAIN TEST INGKFSPPLITX Y GROUPSGROUPS
PRINTS S  TRAIN TEST
0 1 2 3 4 5 6 7 8 9
0 1 2 6 7 8 9 3 4 5
3 4 5 6 7 8 9 0 1 2
```

EACH SUBJECT IS IN A DIFFERENT TESTING FOLD AND THE SAME SUBJECT IS NEVER IN BOTH TESTING AND TRAINING NOTICE THAT THE FOLDS DO NOT HAVE EXACTLY THE SAME SIZE DUE TO THE IMBALANCE IN THE DATA

HERE IS A VISUALIZATION OF THE CROSSVALIDATION BEHAVIOR

444 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

LEAVE ONE GROUP OUT

LEAVEONEGROUPOUT IS A CROSSVALIDATION SCHEME WHICH HOLDS OUT THE SAMPLES ACCORDING TO A THIRDPARTY PROVIDED ARRAY OF INTEGER GROUPS THIS GROUP INFORMATION CAN BE USED TO ENCODE ARBITRARY DOMAIN SPECIFIC PREDEFINED CROSS VALIDATION FOLDS

EACH TRAINING SET IS THUS CONSTITUTED BY ALL THE SAMPLES EXCEPT THE ONES RELATED TO A SPECIFIC GROUP

FOR EXAMPLE IN THE CASES OF MULTIPLE EXPERIMENTS LEAVEONEGROUPOUT CAN BE USED TO CREATE A CROSSVALIDATION BASED ON THE DIFFERENT EXPERIMENTS WE CREATE A TRAINING SET USING THE SAMPLES OF ALL THE EXPERIMENTS EXCEPT ONE

```
FROM SKLEARNMODELSELECTION IMPORT LEAVEONEGROUPOUT
X 1 5 10 50 60 70 80
Y 0 1 1 2 2 2 2
GROUPS 1 1 2 2 3 3 3
LOGO LEAVEONEGROUPOUT
FOR TRAIN TEST INLOGOSPLITX Y GROUPSGROUPS
PRINTS S TRAIN TEST
2 3 4 5 6 0 1
0 1 4 5 6 2 3
0 1 2 3 4 5 6
```

ANOTHER COMMON APPLICATION IS TO USE TIME INFORMATION FOR INSTANCE THE GROUPS COULD BE THE YEAR OF COLLECTION OF THE SAMPLES AND THUS ALLOW FOR CROSSVALIDATION AGAINST TIMEBASED SPLITS

LEAVE P GROUPS OUT

LEAVEPGROUPSOUT IS SIMILAR AS LEAVEONEGROUPOUT BUT REMOVES SAMPLES RELATED TO n GROUPS FOR EACH TRAINING TEST SET

EXAMPLE OF LEAVE2GROUP OUT

```
FROM SKLEARNMODELSELECTION IMPORT LEAVEPGROUPSOUT
X NPARANGE6
Y 1 1 1 2 2 2
GROUPS 1 1 2 2 3 3
LPGO LEAVEPGROUPSOUTNGROUPS2
FOR TRAIN TEST INLPGOSPLITX Y GROUPSGROUPS
PRINTS S TRAIN TEST
4 5 0 1 2 3
2 3 0 1 4 5
0 1 2 3 4 5
```

GROUP SHUFFLE SPLIT

THEGROUPSHUFFLESPLIT ITERATOR BEHAVES AS A COMBINATION OF SHUFFLESPLIT ANDLEAVEPGROUPSOUT AND GENERATES A SEQUENCE OF RANDOMIZED PARTITIONS IN WHICH A SUBSET OF GROUPS ARE HELD OUT FOR EACH SPLIT

HERE IS A USAGE EXAMPLE

```
FROM SKLEARNMODELSELECTION IMPORT GROUPSHUFFLESPLIT
X 01 02 22 24 23 455 58 0001
Y A B B B C C C A
```

33 MODEL SELECTION AND EVALUATION 445

SCIKITLEARN USER GUIDE RELEASE 0213

GROUPS 1 1 2 2 3 3 4 4  
GSS GROUPSHUFFLESPLITNSPLITS4 TESTSIZE05 RANDOMSTATE0  
FOR TRAIN TEST INGSSSPLITX Y GROUPSGROUPS  
PRINTS S TRAIN TEST

0 1 2 3 4 5 6 7  
2 3 6 7 0 1 4 5  
2 3 4 5 0 1 6 7  
4 5 6 7 0 1 2 3

HERE IS A VISUALIZATION OF THE CROSSVALIDATION BEHAVIOR  
THIS CLASS IS USEFUL WHEN THE BEHAVIOR OF LEAVEPGROUPSOUT IS DESIRED BUT THE NUMBER OF GROUPS IS LARGE ENOUGH  
THAT GENERATING ALL POSSIBLE PARTITIONS WITH  $\frac{1}{n}$  GROUPS WITHHELD WOULD BE PROHIBITIVELY EXPENSIVE IN SUCH A SCE  
NARIOGROUPSHUFFLESPLIT PROVIDES A RANDOM SAMPLE WITH REPLACEMENT OF THE TRAIN TEST SPLITS GENERATED BY  
LEAVEPGROUPSOUT  
PREDEFINED FOLDSPLITS VALIDATIONSETS  
FOR SOME DATASETS A PREDEFINED SPLIT OF THE DATA INTO TRAINING AND VALIDATION FOLD OR INTO SEVERAL CROSSVALIDATION FOLDS  
ALREADY EXISTS USING PREDEFINEDSPLIT IT IS POSSIBLE TO USE THESE FOLDS EG WHEN SEARCHING FOR HYPERPARAMETERS  
FOR EXAMPLE WHEN USING A VALIDATION SET SET THE TESTFOLD TO 0 FOR ALL SAMPLES THAT ARE PART OF THE VALIDATION SET  
AND TO 1 FOR ALL OTHER SAMPLES  
CROSS VALIDATION OF TIME SERIES DATA  
TIME SERIES DATA IS CHARACTERISED BY THE CORRELATION BETWEEN OBSERVATIONS THAT ARE NEAR IN TIME AUTOCORRELATION HOW  
EVER CLASSICAL CROSSVALIDATION TECHNIQUES SUCH AS KFOLD ANDSHUFFLESPLIT ASSUME THE SAMPLES ARE INDEPENDENT  
AND IDENTICALLY DISTRIBUTED AND WOULD RESULT IN UNREASONABLE CORRELATION BETWEEN TRAINING AND TESTING INSTANCES YIELD  
ING POOR ESTIMATES OF GENERALISATION ERROR ON TIME SERIES DATA THEREFORE IT IS VERY IMPORTANT TO EVALUATE OUR MODEL  
FOR TIME SERIES DATA ON THE “FUTURE” OBSERVATIONS LEAST LIKE THOSE THAT ARE USED TO TRAIN THE MODEL TO ACHIEVE THIS ONE  
SOLUTION IS PROVIDED BY TIMESERIESSPLIT  
446 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

TIME SERIES SPLIT

TIMESERIESSPLIT IS A VARIATION OF KFOLD WHICH RETURNS FIRST  $\lfloor \frac{n}{k} \rfloor$  FOLDS AS TRAIN SET AND THE  $\lfloor \frac{n}{k} \rfloor + 1$  TH FOLD AS TEST SET

NOTE THAT UNLIKE STANDARD CROSSVALIDATION METHODS SUCCESSIVE TRAINING SETS ARE SUPERSETS OF THOSE THAT COME BEFORE THEM ALSO IT ADDS ALL SURPLUS DATA TO THE FIRST TRAINING PARTITION WHICH IS ALWAYS USED TO TRAIN THE MODEL

THIS CLASS CAN BE USED TO CROSSVALIDATE TIME SERIES DATA SAMPLES THAT ARE OBSERVED AT FIXED TIME INTERVALS

EXAMPLE OF 3SPLIT TIME SERIES CROSSVALIDATION ON A DATASET WITH 6 SAMPLES

```
FROM SKLEARNMODELSELECTION IMPORT TIMESERIESSPLIT
X NPARRAY1 2 3 4 1 2 3 4 1 2 3 4
Y NPARRAY1 2 3 4 5 6
TSCV TIMESERIESSPLITNSPLITS3
PRINTTSCV
TIMESERIESSPLITMAXTRAINSIZENONE NSPLITS3
FOR TRAIN TEST INTSCVSPLITX
PRINTS S TRAIN TEST
0 1 2 3
0 1 2 3 4
0 1 2 3 4 5
```

HERE IS A VISUALIZATION OF THE CROSSVALIDATION BEHAVIOR

A NOTE ON SHUFFLING

IF THE DATA ORDERING IS NOT ARBITRARY EG SAMPLES WITH THE SAME CLASS LABEL ARE CONTIGUOUS SHUFFLING IT FIRST MAY BE ESSENTIAL TO GET A MEANINGFUL CROSS VALIDATION RESULT HOWEVER THE OPPOSITE MAY BE TRUE IF THE SAMPLES ARE NOT INDEPENDENTLY AND IDENTICALLY DISTRIBUTED FOR EXAMPLE IF SAMPLES CORRESPOND TO NEWS ARTICLES AND ARE ORDERED BY THEIR TIME OF PUBLICATION THEN SHUFFLING THE DATA WILL LIKELY LEAD TO A MODEL THAT IS OVERFIT AND AN INFLATED VALIDATION SCORE IT WILL BE TESTED ON SAMPLES THAT ARE ARTIFICIALLY SIMILAR CLOSE IN TIME TO TRAINING SAMPLES

SOME CROSS VALIDATION ITERATORS SUCH AS KFOLD HAVE AN INBUILT OPTION TO SHUFFLE THE DATA INDICES BEFORE SPLITTING THEM

NOTE THAT

- THIS CONSUMES LESS MEMORY THAN SHUFFLING THE DATA DIRECTLY
- BY DEFAULT NO SHUFFLING OCCURS INCLUDING FOR THE STRATIFIED K FOLD CROSS VALIDATION PERFORMED BY SPECIFYING CVSOMEINTEGER TOCROSSVALSCORE GRID SEARCH ETC KEEP IN MIND THAT TRAINTESTSPLIT STILL RETURNS A RANDOM SPLIT

33 MODEL SELECTION AND EVALUATION 447

SCIKITLEARN USER GUIDE RELEASE 0213

- THERANDOMSTATE PARAMETER DEFAULTS TO NONE MEANING THAT THE SHUFFLING WILL BE DIFFERENT EVERY TIME
- KFOLD SHUFFLETRUE IS ITERATED HOWEVER GRIDSEARCHCV WILL USE THE SAME SHUFFLING FOR EACH SET OF PARAMETERS VALIDATED BY A SINGLE CALL TO ITS FIT METHOD
- TO GET IDENTICAL RESULTS FOR EACH SPLIT SET RANDOMSTATE TO AN INTEGER

CROSS VALIDATION AND MODEL SELECTION

CROSS VALIDATION ITERATORS CAN ALSO BE USED TO DIRECTLY PERFORM MODEL SELECTION USING GRID SEARCH FOR THE OPTIMAL HYPERPARAMETERS OF THE MODEL THIS IS THE TOPIC OF THE NEXT SECTION TUNING THE HYPERPARAMETERS OF AN ESTIMATOR

332 TUNING THE HYPERPARAMETERS OF AN ESTIMATOR

HYPERPARAMETERS ARE PARAMETERS THAT ARE NOT DIRECTLY LEARNT WITHIN ESTIMATORS IN SCIKITLEARN THEY ARE PASSED AS ARGUMENTS TO THE CONSTRUCTOR OF THE ESTIMATOR CLASSES TYPICAL EXAMPLES INCLUDE CKERNEL ANDGAMMA FOR SUPPORT VECTOR CLASSIFIER ALPHA FOR LASSO ETC

IT IS POSSIBLE AND RECOMMENDED TO SEARCH THE HYPERPARAMETER SPACE FOR THE BEST CROSS VALIDATION SCORE

ANY PARAMETER PROVIDED WHEN CONSTRUCTING AN ESTIMATOR MAY BE OPTIMIZED IN THIS MANNER SPECIFICALLY TO FIND THE NAMES AND CURRENT VALUES FOR ALL PARAMETERS FOR A GIVEN ESTIMATOR USE

ESTIMATORGETPARAMS

A SEARCH CONSISTS OF

- AN ESTIMATOR REGRESSOR OR CLASSIFIER SUCH AS SKLEARN SVM SVC
- A PARAMETER SPACE
- A METHOD FOR SEARCHING OR SAMPLING CANDIDATES
- A CROSSVALIDATION SCHEME AND
- ASCORE FUNCTION

SOME MODELS ALLOW FOR SPECIALIZED EFFICIENT PARAMETER SEARCH STRATEGIES OUTLINED BELOW TWO GENERIC APPROACHES TO SAMPLING SEARCH CANDIDATES ARE PROVIDED IN SCIKITLEARN FOR GIVEN VALUES GRIDSEARCHCV EXHAUSTIVELY CONSIDERS ALL PARAMETER COMBINATIONS WHILE RANDOMIZEDSEARCHCV CAN SAMPLE A GIVEN NUMBER OF CANDIDATES FROM A PARAMETER SPACE WITH A SPECIFIED DISTRIBUTION AFTER DESCRIBING THESE TOOLS WE DETAIL BEST PRACTICE APPLICABLE TO BOTH APPROACHES NOTE THAT IT IS COMMON THAT A SMALL SUBSET OF THOSE PARAMETERS CAN HAVE A LARGE IMPACT ON THE PREDICTIVE OR COMPUTATION PERFORMANCE OF THE MODEL WHILE OTHERS CAN BE LEFT TO THEIR DEFAULT VALUES IT IS RECOMMENDED TO READ THE DOCSTRING OF THE ESTIMATOR CLASS TO GET A FINER UNDERSTANDING OF THEIR EXPECTED BEHAVIOR POSSIBLY BY READING THE ENCLOSED REFERENCE TO THE LITERATURE

EXHAUSTIVE GRID SEARCH

THE GRID SEARCH PROVIDED BY GRIDSEARCHCV EXHAUSTIVELY GENERATES CANDIDATES FROM A GRID OF PARAMETER VALUES SPECIFIED WITH THE PARAMGRID PARAMETER FOR INSTANCE THE FOLLOWING PARAMGRID

PARAMGRID

C 1 10 100 1000 KERNEL LINEAR

C 1 10 100 1000 GAMMA 0001 00001 KERNEL RBF

SCIKITLEARN USER GUIDE RELEASE 0213  
SPECIFIES THAT TWO GRIDS SHOULD BE EXPLORED ONE WITH A LINEAR KERNEL AND C VALUES IN 1 10 100 1000 AND THE SECOND ONE WITH AN RBF KERNEL AND THE CROSSPRODUCT OF C VALUES RANGING IN 1 10 100 1000 AND GAMMA VALUES IN 0001 00001

THEGRIDSEARCHCV INSTANCE IMPLEMENTS THE USUAL ESTIMATOR API WHEN “FITTING” IT ON A DATASET ALL THE POSSIBLE COMBINATIONS OF PARAMETER VALUES ARE EVALUATED AND THE BEST COMBINATION IS RETAINED

EXAMPLES

- SEE PARAMETER ESTIMATION USING GRID SEARCH WITH CROSSVALIDATION FOR AN EXAMPLE OF GRID SEARCH COMPUTATION ON THE DIGITS DATASET
  - SEE SAMPLE PIPELINE FOR TEXT FEATURE EXTRACTION AND EVALUATION FOR AN EXAMPLE OF GRID SEARCH COUPLING PARAMETERS FROM A TEXT DOCUMENTS FEATURE EXTRACTOR NGRAM COUNT VECTORIZER AND TFIDF TRANSFORMER WITH A CLASSIFIER HERE A LINEAR SVM TRAINED WITH SGD WITH EITHER ELASTIC NET OR L2 PENALTY USING A PIPELINE PIPELINE INSTANCE
  - SEE NESTED VERSUS NONNESTED CROSSVALIDATION FOR AN EXAMPLE OF GRID SEARCH WITHIN A CROSS VALIDATION LOOP ON THE IRIS DATASET THIS IS THE BEST PRACTICE FOR EVALUATING THE PERFORMANCE OF A MODEL WITH GRID SEARCH
  - SEE DEMONSTRATION OF MULTIMETRIC EVALUATION ON CROSSVALSCORE AND GRIDSEARCHCV FOR AN EXAMPLE OF GRIDSEARCHCV BEING USED TO EVALUATE MULTIPLE METRICS SIMULTANEOUSLY
  - SEE BALANCE MODEL COMPLEXITY AND CROSSVALIDATED SCORE FOR AN EXAMPLE OF USING REFITCALLABLE IN TERFACE INGRIDSEARCHCV THE EXAMPLE SHOWS HOW THIS INTERFACE ADDS CERTAIN AMOUNT OF FLEXIBILITY IN IDENTIFYING THE “BEST” ESTIMATOR THIS INTERFACE CAN ALSO BE USED IN MULTIPLE METRICS EVALUATION
- RANDOMIZED PARAMETER OPTIMIZATION
- WHILE USING A GRID OF PARAMETER SETTINGS IS CURRENTLY THE MOST WIDELY USED METHOD FOR PARAMETER OPTIMIZATION OTHER SEARCH METHODS HAVE MORE FAVOURABLE PROPERTIES RANDOMIZEDSEARCHCV IMPLEMENTS A RANDOMIZED SEARCH OVER PARAMETERS WHERE EACH SETTING IS SAMPLED FROM A DISTRIBUTION OVER POSSIBLE PARAMETER VALUES THIS HAS TWO MAIN BENEFITS OVER AN EXHAUSTIVE SEARCH
- A BUDGET CAN BE CHOSEN INDEPENDENT OF THE NUMBER OF PARAMETERS AND POSSIBLE VALUES
  - ADDING PARAMETERS THAT DO NOT INFLUENCE THE PERFORMANCE DOES NOT DECREASE EFFICIENCY
- SPECIFYING HOW PARAMETERS SHOULD BE SAMPLED IS DONE USING A DICTIONARY VERY SIMILAR TO SPECIFYING PARAMETERS FOR GRIDSEARCHCV ADDITIONALLY A COMPUTATION BUDGET BEING THE NUMBER OF SAMPLED CANDIDATES OR SAMPLING ITERATIONS IS SPECIFIED USING THE NITER PARAMETER FOR EACH PARAMETER EITHER A DISTRIBUTION OVER POSSIBLE VALUES OR A LIST OF DISCRETE CHOICES WHICH WILL BE SAMPLED UNIFORMLY CAN BE SPECIFIED
- C SCIPYSTATSEXPONSCALE100 GAMMA SCIPYSTATSEXPONSCALE1
- KERNEL RBF CLASSWEIGHTBALANCED NONE

THIS EXAMPLE USES THE SCIPYSTATS MODULE WHICH CONTAINS MANY USEFUL DISTRIBUTIONS FOR SAMPLING PARAMETERS SUCH ASEXPON GAMMA UNIFORM ORRANDINT IN PRINCIPLE ANY FUNCTION CAN BE PASSED THAT PROVIDES A RVS RANDOM VARIATE SAMPLE METHOD TO SAMPLE A VALUE A CALL TO THE RVS FUNCTION SHOULD PROVIDE INDEPENDENT RANDOM SAMPLES FROM POSSIBLE PARAMETER VALUES ON CONSECUTIVE CALLS

WARNING THE DISTRIBUTIONS IN SCIPYSTATS PRIOR TO VERSION SCIPY 016 DO NOT ALLOW SPECIFYING A RANDOM STATE INSTEAD THEY USE THE GLOBAL NUMPY RANDOM STATE THAT CAN BE SEEDED VIA NPRANDOM SEED OR SET USING NPRANDOMSETSTATE HOWEVER BEGINNING SCIKITLEARN 018 THE SKLEARN

33 MODEL SELECTION AND EVALUATION 449

SCIKITLEARN USER GUIDE RELEASE 0213

MODELSELECTION MODULE SETS THE RANDOM STATE PROVIDED BY THE USER IF SCIPY 016 IS ALSO AVAILABLE

FOR CONTINUOUS PARAMETERS SUCH AS CABOVE IT IS IMPORTANT TO SPECIFY A CONTINUOUS DISTRIBUTION TO TAKE FULL ADVANTAGE OF THE RANDOMIZATION THIS WAY INCREASING NITER WILL ALWAYS LEAD TO A FINER SEARCH

EXAMPLES

- COMPARING RANDOMIZED SEARCH AND GRID SEARCH FOR HYPERPARAMETER ESTIMATION COMPARES THE USAGE AND EFFICIENCY OF RANDOMIZED SEARCH AND GRID SEARCH

REFERENCES

- BERGSTRA J AND BENGIO Y RANDOM SEARCH FOR HYPERPARAMETER OPTIMIZATION THE JOURNAL OF MACHINE LEARNING RESEARCH 2012

TIPS FOR PARAMETER SEARCH

SPECIFYING AN OBJECTIVE METRIC

BY DEFAULT PARAMETER SEARCH USES THE SCORE FUNCTION OF THE ESTIMATOR TO EVALUATE A PARAMETER SETTING THESE ARE THE SKLEARNMETRICSACCURACYScore FOR CLASSIFICATION AND SKLEARNMETRICSR2Score FOR REGRESSION FOR SOME APPLICATIONS OTHER SCORING FUNCTIONS ARE BETTER SUITED FOR EXAMPLE IN UNBALANCED CLASSIFICATION THE ACCURACY SCORE IS OFTEN UNINFORMATIVE AN ALTERNATIVE SCORING FUNCTION CAN BE SPECIFIED VIA THE SCORING PARAMETER TO GRIDSEARCHCV RANDOMIZEDSEARCHCV AND MANY OF THE SPECIALIZED CROSSVALIDATION TOOLS DESCRIBED BELOW SEE THE SCORING PARAMETER DEFINING MODEL EVALUATION RULES FOR MORE DETAILS

SPECIFYING MULTIPLE METRICS FOR EVALUATION

GRIDSEARCHCV ANDRANDOMIZEDSEARCHCV ALLOW SPECIFYING MULTIPLE METRICS FOR THE SCORING PARAMETER MULTIMETRIC SCORING CAN EITHER BE SPECIFIED AS A LIST OF STRINGS OF PREDEFINED SCORE NAMES OR A DICT MAPPING THE SCORER NAME TO THE SCORER FUNCTION ANDOR THE PREDEFINED SCORER NAMES SEE USING MULTIPLE METRIC EVALUATION FOR MORE DETAILS

WHEN SPECIFYING MULTIPLE METRICS THE REFIT PARAMETER MUST BE SET TO THE METRIC STRING FOR WHICH THE BESTPARAMS WILL BE FOUND AND USED TO BUILD THE BESTESTIMATOR ON THE WHOLE DATASET IF THE SEARCH SHOULD NOT BE REFIT SET REFITFALSE LEAVING REFIT TO THE DEFAULT VALUE NONE WILL RESULT IN AN ERROR WHEN USING MULTIPLE METRICS

SEEDEMONSTRATION OF MULTIMETRIC EVALUATION ON CROSSVALSCORE AND GRIDSEARCHCV FOR AN EXAMPLE USAGE

COMPOSITE ESTIMATORS AND PARAMETER SPACES

PIPELINE CHAINING ESTIMATORS DESCRIBES BUILDING COMPOSITE ESTIMATORS WHOSE PARAMETER SPACE CAN BE SEARCHED WITH THESE TOOLS

450 CHAPTER 3 USER GUIDE



SCIKITLEARN USER GUIDE RELEASE 0213

MODEL SELECTION DEVELOPMENT AND EVALUATION

MODEL SELECTION BY EVALUATING VARIOUS PARAMETER SETTINGS CAN BE SEEN AS A WAY TO USE THE LABELED DATA TO “TRAIN” THE PARAMETERS OF THE GRID

WHEN EVALUATING THE RESULTING MODEL IT IS IMPORTANT TO DO IT ON HELDOUT SAMPLES THAT WERE NOT SEEN DURING THE GRID SEARCH PROCESS IT IS RECOMMENDED TO SPLIT THE DATA INTO A DEVELOPMENT SET TO BE FED TO THE GRIDSEARCHCV INSTANCE AND AN EVALUATION SET TO COMPUTE PERFORMANCE METRICS

THIS CAN BE DONE BY USING THE TRAINTESTSPLIT UTILITY FUNCTION

PARALLELISM

GRIDSEARCHCV ANDRANDOMIZEDSEARCHCV EVALUATE EACH PARAMETER SETTING INDEPENDENTLY COMPUTATIONS CAN BE RUN IN PARALLEL IF YOUR OS SUPPORTS IT BY USING THE KEYWORD NJOBS1 SEE FUNCTION SIGNATURE FOR MORE DETAILS

ROBUSTNESS TO FAILURE

SOME PARAMETER SETTINGS MAY RESULT IN A FAILURE TO FIT ONE OR MORE FOLDS OF THE DATA BY DEFAULT THIS WILL CAUSE THE ENTIRE SEARCH TO FAIL EVEN IF SOME PARAMETER SETTINGS COULD BE FULLY EVALUATED SETTING ERRORSORE0 ORNP NAN WILL MAKE THE PROCEDURE ROBUST TO SUCH FAILURE ISSUING A WARNING AND SETTING THE SCORE FOR THAT FOLD TO 0 OR NAN BUT COMPLETING THE SEARCH

ALTERNATIVES TO BRUTE FORCE PARAMETER SEARCH

MODEL SPECIFIC CROSSVALIDATION

SOME MODELS CAN FIT DATA FOR A RANGE OF VALUES OF SOME PARAMETER ALMOST AS EFFICIENTLY AS FITTING THE ESTIMATOR FOR A SINGLE VALUE OF THE PARAMETER THIS FEATURE CAN BE LEVERAGED TO PERFORM A MORE EFFICIENT CROSSVALIDATION USED FOR MODEL SELECTION OF THIS PARAMETER

THE MOST COMMON PARAMETER AMENABLE TO THIS STRATEGY IS THE PARAMETER ENCODING THE STRENGTH OF THE REGULARIZER IN THIS CASE WE SAY THAT WE COMPUTE THE REGULARIZATION PATH OF THE ESTIMATOR

HERE IS THE LIST OF SUCH MODELS

LINEARMODELELASTICNETCV L1RATIO EPS ELASTIC NET MODEL WITH ITERATIVE FITTING ALONG A REGULARIZATION PATH

LINEARMODELLARSCV FITINTERCEPT CROSSVALIDATED LEAST ANGLE REGRESSION MODEL

LINEARMODELLASSOCV EPS NALPHAS LASSO LINEAR MODEL WITH ITERATIVE FITTING ALONG A REGULARIZATION PATH

LINEARMODELLASSOLARSCV FITINTERCEPT CROSSVALIDATED LASSO USING THE LARS ALGORITHM

LINEARMODELLOGISTICREGRESSIONCV CS LOGISTIC REGRESSION CV AKA LOGIT MAXENT CLASSIFIER

LINEARMODELMULTITASKELASTICNETCV MULTITASK L1L2 ELASTICNET WITH BUILTIN CROSSVALIDATION

LINEARMODELMULTITASKLASSOCV EPS MULTITASK LASSO MODEL TRAINED WITH L1L2 MIXEDNORM AS REGULARIZER

LINEARMODELORTHOGONALMATCHINGPURSUITCV CROSSVALIDATED ORTHOGONAL MATCHING PURSUIT MODEL

OMP

LINEARMODELRIDGECV ALPHAS RIDGE REGRESSION WITH BUILTIN CROSSVALIDATION

CONTINUED ON NEXT PAGE

33 MODEL SELECTION AND EVALUATION 451

LINEARMODEL RIDGECLASSIFIER CV ALPHAS

RIDGE CLASSIFIER WITH BUILTIN CROSSVALIDATION

SKLEARNLINEARMODEL ELASTICNET CV

CLASS SKLEARNLINEARMODEL ELASTICNET CV L1RATIO 0.5 EPS 0.001 NALPHAS 100 AL

PHAS NONE FIT INTERCEPT TRUE NORMALIZE FALSE

PRECOMPUTE 'AUTO' MAXITER 1000 TOL 0.0001

CV 'WARN' COPY X TRUE VERBOSE 0 NJOBS NONE

POSITIVE FALSE RANDOM STATE NONE SELEC

TION 'CYCLIC'

ELASTIC NET MODEL WITH ITERATIVE FITTING ALONG A REGULARIZATION PATH

SEE GLOSSARY ENTRY FOR CROSSVALIDATION ESTIMATOR

READ MORE IN THE USER GUIDE

PARAMETERS

L1RATIO FLOAT OR ARRAY OF FLOATS OPTIONAL FLOAT BETWEEN 0 AND 1 PASSED TO ELASTICNET SCAL

ING BETWEEN L1 AND L2 PENALTIES FOR L1RATIO 0 THE PENALTY IS AN L2 PENALTY FOR

L1RATIO 1 IT IS AN L1 PENALTY FOR 0 L1RATIO 1 THE PENALTY IS A COM

BINATION OF L1 AND L2 THIS PARAMETER CAN BE A LIST IN WHICH CASE THE DIFFERENT VALUES ARE

TESTED BY CROSSVALIDATION AND THE ONE GIVING THE BEST PREDICTION SCORE IS USED NOTE THAT A

GOOD CHOICE OF LIST OF VALUES FOR L1RATIO IS OFTEN TO PUT MORE VALUES CLOSE TO 1 IE LASSO

AND LESS CLOSE TO 0 IE RIDGE AS IN 1 5 7 9 95 99 1

EPS FLOAT OPTIONAL LENGTH OF THE PATH EPS 1E3 MEANS THAT ALPHA MIN ALPHA MAX

1E3

NALPHAS INT OPTIONAL NUMBER OF ALPHAS ALONG THE REGULARIZATION PATH USED FOR EACH L1RATIO

ALPHAS NUMPY ARRAY OPTIONAL LIST OF ALPHAS WHERE TO COMPUTE THE MODELS IF NONE ALPHAS ARE

SET AUTOMATICALLY

FIT INTERCEPT BOOLEAN WHETHER TO CALCULATE THE INTERCEPT FOR THIS MODEL IF SET TO FALSE NO

INTERCEPT WILL BE USED IN CALCULATIONS EG DATA IS EXPECTED TO BE ALREADY CENTERED

NORMALIZE BOOLEAN OPTIONAL DEFAULT FALSE THIS PARAMETER IS IGNORED WHEN

FIT INTERCEPT IS SET TO FALSE IF TRUE THE REGRESSORS X WILL BE NORMALIZED BE

FORE REGRESSION BY SUBTRACTING THE MEAN AND DIVIDING BY THE L2 NORM IF YOU WISH TO

STANDARDIZE PLEASE USE SKLEARN PREPROCESSING STANDARD SCALER BEFORE

CALLING FIT ON AN ESTIMATOR WITH NORMALIZE FALSE

PRECOMPUTE TRUE FALSE 'AUTO' ARRAY LIKE WHETHER TO USE A PRECOMPUTED GRAM MATRIX TO

SPEED UP CALCULATIONS IF SET TO AUTO LET US DECIDE THE GRAM MATRIX CAN ALSO BE PASSED

AS ARGUMENT

MAXITER INT OPTIONAL THE MAXIMUM NUMBER OF ITERATIONS

TOL FLOAT OPTIONAL THE TOLERANCE FOR THE OPTIMIZATION IF THE UPDATES ARE SMALLER THAN TOL THE

OPTIMIZATION CODE CHECKS THE DUAL GAP FOR OPTIMALITY AND CONTINUES UNTIL IT IS SMALLER THAN

TOL

CV INT CROSSVALIDATION GENERATOR OR AN ITERABLE OPTIONAL DETERMINES THE CROSSVALIDATION SPLIT

TING STRATEGY POSSIBLE INPUTS FOR CV ARE

• NONE TO USE THE DEFAULT 3FOLD CROSSVALIDATION

SCIKITLEARN USER GUIDE RELEASE 0213

- INTEGER TO SPECIFY THE NUMBER OF FOLDS
- CV SPLITTER
- AN ITERABLE YIELDING TRAIN TEST SPLITS AS ARRAYS OF INDICES

FOR INTEGERS NONE INPUTS K FOLD IS USED

REFER USER GUIDE FOR THE VARIOUS CROSS VALIDATION STRATEGIES THAT CAN BE USED HERE  
CHANGED IN VERSION 020 CV DEFAULT VALUE IF NONE WILL CHANGE FROM 3 FOLD TO 5 FOLD IN  
V022

COPYX BOOLEAN OPTIONAL DEFAULT TRUE IF TRUE X WILL BE COPIED ELSE IT MAY BE OVERWRITTEN  
VERBOSE BOOL OR INTEGER AMOUNT OF VERBOSITY

NJOBS INT OR NONE OPTIONAL DEFAULT NONE NUMBER OF CPUS TO USE DURING THE CROSS VALIDA  
TION NONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT 1 MEANS USING  
ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

POSITIVE BOOL OPTIONAL WHEN SET TO TRUE FORCES THE COEFFICIENTS TO BE POSITIVE

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULT NONE THE SEED OF THE  
PSEUDO RANDOM NUMBER GENERATOR THAT SELECTS A RANDOM FEATURE TO UPDATE IF INT RAN  
DOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE  
RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS  
THE RANDOMSTATE INSTANCE USED BY NRANDOM USED WHEN SELECTION 'RANDOM'  
SELECTION STR DEFAULT 'CYCLIC' IF SET TO 'RANDOM' A RANDOM COEFFICIENT IS UPDATED EVERY ITERATION  
RATHER THAN LOOPING OVER FEATURES SEQUENTIALLY BY DEFAULT THIS SETTING TO 'RANDOM' OFTEN  
LEADS TO SIGNIFICANTLY FASTER CONVERGENCE ESPECIALLY WHEN TOL IS HIGHER THAN 1E4

ATTRIBUTES

ALPHA FLOAT THE AMOUNT OF PENALIZATION CHOSEN BY CROSS VALIDATION

L1RATIO FLOAT THE COMPROMISE BETWEEN L1 AND L2 PENALIZATION CHOSEN BY CROSS VALIDATION

COEF ARRAY SHAPE NFEATURES NTARGETS NFEATURES PARAMETER VECTOR W IN THE COST FUNC  
TION FORMULA

INTERCEPT FLOAT ARRAY SHAPE NTARGETS NFEATURES INDEPENDENT TERM IN THE DECISION FUNC  
TION

MSEPATH ARRAY SHAPE NL1RATIO NALPHA NFOLDS MEAN SQUARE ERROR FOR THE TEST SET ON  
EACH FOLD VARYING L1RATIO AND ALPHA

ALPHAS NUMPY ARRAY SHAPE NALPHAS OR NL1RATIO NALPHAS THE GRID OF ALPHAS USED FOR  
FITTING FOR EACH L1RATIO

NITER INT NUMBER OF ITERATIONS RUN BY THE COORDINATE DESCENT SOLVER TO REACH THE SPECIFIED  
TOLERANCE FOR THE OPTIMAL ALPHA

SEE ALSO

ENETPATH

ELASTICNET

NOTES

FOR AN EXAMPLE SEE EXAMPLES LINEAR MODEL PLOT LASSO MODEL SELECTION PY  
33 MODEL SELECTION AND EVALUATION 453

SCIKITLEARN USER GUIDE RELEASE 0213

TO AVOID UNNECESSARY MEMORY DUPLICATION THE X ARGUMENT OF THE FIT METHOD SHOULD BE DIRECTLY PASSED AS A FORTRANCONTIGUOUS NUMPY ARRAY

THE PARAMETER L1RATIO CORRESPONDS TO ALPHA IN THE GLMNET R PACKAGE WHILE ALPHA CORRESPONDS TO THE LAMBDA PARAMETER IN GLMNET MORE SPECIFICALLY THE OPTIMIZATION OBJECTIVE IS

1 2NSAMPLES Y XW22

ALPHA L1RATIO W1

05ALPHA1 L1RATIO W22

IF YOU ARE INTERESTED IN CONTROLLING THE L1 AND L2 PENALTY SEPARATELY KEEP IN MIND THAT THIS IS EQUIVALENT TO

AL1 BL2

FOR

ALPHA A B ANDL1RATIO A A B

EXAMPLES

FROM SKLEARNLINEARMODEL IMPORT ELASTICNETCV

FROM SKLEARNDATASETS IMPORT MAKEREGRESSION

X Y MAKEREGRESSIONNFEATURES2 RANDOMSTATE0

REGR ELASTICNETCVCV5 RANDOMSTATE0

REGRFITX Y

ELASTICNETCVALPHASNONE COPYXTRUE CV5 EPS0001 FITINTERCEPTTRUE

L1RATIO05 MAXITER1000 NALPHAS100 NJOBSNONE

NORMALIZEFALSE POSITIVEFALSE PRECOMPUTEAUTO RANDOMSTATE0

SELECTIONCYCLIC TOL00001 VERBOSE0

PRINTREGRALPHA

0199

PRINTREGRINTERCEPT

0398

PRINTREGRPREDICT0 0

0398

METHODS

FITSELF X Y FIT LINEAR MODEL WITH COORDINATE DESCENT

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PATH X Y L1RATIO EPS NALPHAS COMPUTE ELASTIC NET PATH WITH COORDINATE DESCENT

PREDICT SELF X PREDICT USING THE LINEAR MODEL

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELF L1RATIO05 EPS0001 NALPHAS100 ALPHASNONE FITINTERCEPTTRUE NORMALIZEFALSE PRECOMPUTE'AUTO' MAXITER1000 TOL00001 CV'WARN' COPYXTRUE VERBOSE0 NJOBSNONE POSITIVEFALSE RANDOMSTATENONE SELECTION'CYCLIC'

FITSELFXY

FIT LINEAR MODEL WITH COORDINATE DESCENT

454 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

FIT IS ON GRID OF ALPHAS AND BEST ALPHA ESTIMATED BY CROSSVALIDATION

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA PASS DIRECTLY AS FORTRAN  
CONTIGUOUS DATA TO AVOID UNNECESSARY MEMORY DUPLICATION IF Y IS MONOOUTPUT X CAN  
BE SPARSE

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NTARGETS TARGET VALUES

GETPARAMS SELFDEEPTTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

STATICPATHXYL1RATIO05 EPS0001 NALPHAS100 ALPHASNONE PRECOMPUTE'AUTO'  
XYNONE COPYXTRUE COEFINITNONE VERBOSEFALSE RETURNNITERFALSE POSI  
TIVEFALSE CHECKINPUTTRUE PARAMS

COMPUTE ELASTIC NET PATH WITH COORDINATE DESCENT

THE ELASTIC NET OPTIMIZATION FUNCTION VARIES FOR MONO AND MULTIOUTPUTS

FOR MONOOUTPUT TASKS IT IS

1 2NSAMPLES Y XW22  
ALPHA L1RATIO W1  
05ALPHA1 L1RATIO W22

FOR MULTIOUTPUT TASKS IT IS

1 2 NSAMPLES Y XWFRO2  
ALPHA L1RATIO W21  
05ALPHA1 L1RATIO WFRO2

WHERE

W21 SUMI SQRTSUMJ WIJ2  
IE THE SUM OF NORM OF EACH ROW

READ MORE IN THE USER GUIDE

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA PASS DIRECTLY AS FORTRAN  
CONTIGUOUS DATA TO AVOID UNNECESSARY MEMORY DUPLICATION IF YIS MONOOUTPUT THEN X  
CAN BE SPARSE

YNDARRAY SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TARGET VALUES

L1RATIO FLOAT OPTIONAL FLOAT BETWEEN 0 AND 1 PASSED TO ELASTIC NET SCALING BETWEEN L1 AND  
L2 PENALTIES L1RATIO1 CORRESPONDS TO THE LASSO

EPS FLOAT LENGTH OF THE PATH EPS1E3 MEANS THAT ALPHAMIN ALPHAMAX  
1E3

NALPHAS INT OPTIONAL NUMBER OF ALPHAS ALONG THE REGULARIZATION PATH

33 MODEL SELECTION AND EVALUATION 455

SCIKITLEARN USER GUIDE RELEASE 0213

ALPHAS NDARRAY OPTIONAL LIST OF ALPHAS WHERE TO COMPUTE THE MODELS IF NONE ALPHAS ARE SET AUTOMATICALLY

PRECOMPUTE TRUE FALSE 'AUTO' ARRAYLIKE WHETHER TO USE A PRECOMPUTED GRAM MATRIX TO SPEED UP CALCULATIONS IF SET TO AUTO LET US DECIDE THE GRAM MATRIX CAN ALSO BE PASSED AS ARGUMENT

XYARRAYLIKE OPTIONAL XY NPDOTXT Y THAT CAN BE PRECOMPUTED IT IS USEFUL ONLY WHEN THE GRAM MATRIX IS PRECOMPUTED

COPYX BOOLEAN OPTIONAL DEFAULT TRUE IF TRUE X WILL BE COPIED ELSE IT MAY BE OVERWRITTEN

COEFINIT ARRAY SHAPE NFEATURES NONE THE INITIAL VALUES OF THE COEFFICIENTS

VERBOSE BOOL OR INTEGER AMOUNT OF VERBOSITY

RETURNNNITER BOOL WHETHER TO RETURN THE NUMBER OF ITERATIONS OR NOT

POSITIVE BOOL DEFAULT FALSE IF SET TO TRUE FORCES COEFFICIENTS TO BE POSITIVE ONLY ALLOWED WHENYNDIM 1

CHECKINPUT BOOL DEFAULT TRUE SKIP INPUT VALIDATION CHECKS INCLUDING THE GRAM MATRIX WHEN PROVIDED ASSUMING THERE ARE HANDLED BY THE CALLER WHEN CHECKINPUTFALSE

PARAMS KWARGS KEYWORD ARGUMENTS PASSED TO THE COORDINATE DESCENT SOLVER

RETURNS

ALPHAS ARRAY SHAPE NALPHAS THE ALPHAS ALONG THE PATH WHERE MODELS ARE COMPUTED

COEFS ARRAY SHAPE NFEATURES NALPHAS OR NOUTPUTS NFEATURES NALPHAS COEFFICIENTS ALONG THE PATH

DUALGAPS ARRAY SHAPE NALPHAS THE DUAL GAPS AT THE END OF THE OPTIMIZATION FOR EACH ALPHA

NITERS ARRAYLIKE SHAPE NALPHAS THE NUMBER OF ITERATIONS TAKEN BY THE COORDINATE DESCENT OPTIMIZER TO REACH THE SPECIFIED TOLERANCE FOR EACH ALPHA IS RETURNED WHEN RETURNNNITER IS SET TO TRUE

SEE ALSO

MULTITASKELASTICNET

MULTITASKELASTICNETCV

ELASTICNET

ELASTICNETCV

NOTES

FOR AN EXAMPLE SEE EXAMPLESLINEARMODELPLOTLASSOCOORDINATEDDESCENTPATHPY

PREDICTSELF

PREDICT USING THE LINEAR MODEL

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES SAMPLES

RETURNS

456 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

CARRAY SHAPE NSAMPLES RETURNS PREDICTED VALUES

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{U}{V}$  WHERE U IS THE RESIDUAL SUM OF SQUARES  $Y_{TRUE} - Y_{PRED}$  2SUM AND V IS THE TOTAL SUM OF SQUARES  $Y_{TRUE} - Y_{TRUEMEAN}$  2SUM THE BEST POSSIBLE SCORE IS 10 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 00

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES FOR SOME ESTIMATORS THIS MAY BE A PRECOMPUTED KERNEL MATRIX INSTEAD SHAPE NSAMPLES NSAMPLESFITTED WHERE NSAMPLESFITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE VALUES FOR X SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT R2 OF SELF PREDICTX WRT Y

NOTES

THE R2 SCORE USED WHEN CALLING SCORE ON A REGRESSOR WILL USE MULTIOUTPUTUNIFORMAVERAGE FROM VERSION 023 TO KEEP CONSISTENT WITH METRICSR2SCORE THIS WILL INFLUENCE THE SCORE METHOD OF ALL THE MULTIOUTPUT REGRESSORS EXCEPT FOR MULTIOUTPUTMULTIOUTPUTREGRESSOR TO SPECIFY THE DEFAULT VALUE MANUALLY AND AVOID THE WARNING PLEASE EITHER CALL METRICSR2SCORE DIRECTLY OR MAKE A CUSTOM SCORER WITH METRICSMAKESCORER THE BUILTIN SCORER R2 USES MULTIOUTPUTUNIFORMAVERAGE

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

SKLEARNLINEARMODEL LARSCV

CLASSSKLEARNLINEARMODEL LARSCVFITINTERCEPTTRUE VERBOSEFALSE MAXITER500

NORMALIZETRUE PRECOMPUTE'AUTO'

CV'WARN' MAXNALPHAS1000 NJOBSNONE

EPS2220446049250313E16 COPYXTRUE POSITIVEFALSE

CROSSVALIDATED LEAST ANGLE REGRESSION MODEL

SEE GLOSSARY ENTRY FOR CROSSVALIDATION ESTIMATOR

READ MORE IN THE USER GUIDE

PARAMETERS

FITINTERCEPT BOOLEAN WHETHER TO CALCULATE THE INTERCEPT FOR THIS MODEL IF SET TO FALSE NO INTERCEPT WILL BE USED IN CALCULATIONS EG DATA IS EXPECTED TO BE ALREADY CENTERED

33 MODEL SELECTION AND EVALUATION 457

SCIKITLEARN USER GUIDE RELEASE 0213

VERBOSE BOOLEAN OR INTEGER OPTIONAL SETS THE VERBOSITY AMOUNT

MAXITER INTEGER OPTIONAL MAXIMUM NUMBER OF ITERATIONS TO PERFORM

NORMALIZE BOOLEAN OPTIONAL DEFAULT TRUE THIS PARAMETER IS IGNORED WHEN FITINTERCEPT IS SET TO FALSE IF TRUE THE REGRESSORS X WILL BE NORMALIZED BEFORE REGRESSION BY SUBTRACTING THE MEAN AND DIVIDING BY THE L2NORM IF YOU WISH TO STANDARDIZE PLEASE USE SKLEARNPREPROCESSINGSTANDARDSCALER BEFORE CALLING FIT ON AN ESTIMATOR WITHNORMALIZEFALSE

PRECOMPUTE TRUE FALSE 'AUTO' ARRAYLIKE WHETHER TO USE A PRECOMPUTED GRAM MATRIX TO SPEED UP CALCULATIONS IF SET TO AUTO LET US DECIDE THE GRAM MATRIX CANNOT BE PASSED AS ARGUMENT SINCE WE WILL USE ONLY SUBSETS OF X

CVINT CROSSVALIDATION GENERATOR OR AN ITERABLE OPTIONAL DETERMINES THE CROSSVALIDATION SPLITTING STRATEGY POSSIBLE INPUTS FOR CV ARE

- NONE TO USE THE DEFAULT 3FOLD CROSSVALIDATION
- INTEGER TO SPECIFY THE NUMBER OF FOLDS
- CV SPLITTER
- AN ITERABLE YIELDING TRAIN TEST SPLITS AS ARRAYS OF INDICES

FOR INTEGERNONE INPUTS KFOLD IS USED

REFER USER GUIDE FOR THE VARIOUS CROSSVALIDATION STRATEGIES THAT CAN BE USED HERE

CHANGED IN VERSION 020 CVDEFAULT VALUE IF NONE WILL CHANGE FROM 3FOLD TO 5FOLD IN V022

MAXNALPHAS INTEGER OPTIONAL THE MAXIMUM NUMBER OF POINTS ON THE PATH USED TO COMPUTE THE RESIDUALS IN THE CROSSVALIDATION

NJOBS INT OR NONE OPTIONAL DEFAULTNONE NUMBER OF CPUS TO USE DURING THE CROSS VALIDATIONNONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

EPS FLOAT OPTIONAL THE MACHINEPRECISION REGULARIZATION IN THE COMPUTATION OF THE CHOLESKY DIAGONAL FACTORS INCREASE THIS FOR VERY ILLCONDITIONED SYSTEMS

COPYX BOOLEAN OPTIONAL DEFAULT TRUE IF TRUE X WILL BE COPIED ELSE IT MAY BE OVERWRITTEN

POSITIVE BOOLEAN DEFAULTFALSE RESTRICT COEFFICIENTS TO BE 0 BE AWARE THAT YOU MIGHT WANT TO REMOVE FITINTERCEPT WHICH IS SET TRUE BY DEFAULT

DEPRECATED SINCE VERSION 020 THE OPTION IS BROKEN AND DEPRECATED IT WILL BE REMOVED IN V022

ATTRIBUTES

COEF ARRAY SHAPE NFEATURES PARAMETER VECTOR W IN THE FORMULATION FORMULA

INTERCEPT FLOAT INDEPENDENT TERM IN DECISION FUNCTION

COEFPATH ARRAY SHAPE NFEATURES NALPHAS THE VARYING VALUES OF THE COEFFICIENTS ALONG THE PATH

ALPHA FLOAT THE ESTIMATED REGULARIZATION PARAMETER ALPHA

ALPHAS ARRAY SHAPE NALPHAS THE DIFFERENT VALUES OF ALPHA ALONG THE PATH

CVALPHAS ARRAY SHAPE NCVALPHAS ALL THE VALUES OF ALPHA ALONG THE PATH FOR THE DIFFERENT FOLDS

458 CHAPTER 3 USER GUIDE



SCIKITLEARN USER GUIDE RELEASE 0213

MSEPATH ARRAY SHAPE NFOLDS NCVALPHAS THE MEAN SQUARE ERROR ON LEFTOUT FOR EACH FOLD  
ALONG THE PATH ALPHA VALUES GIVEN BY CVALPHAS  
NITER ARRAYLIKE OR INT THE NUMBER OF ITERATIONS RUN BY LARS WITH THE OPTIMAL ALPHA  
SEE ALSO  
LARSPATH LASSOLARS LASSOLARSCV

EXAMPLES

```
FROM SKLEARNLINEARMODEL IMPORT LARSCV
FROM SKLEARNDATASETS IMPORT MAKEREGRESSION
X Y MAKEREGRESSIONNSAMPLES200 NOISE40 RANDOMSTATE0
REG LARSCVCV5FITX Y
REGSCOREX Y
09996
REGALPHA
00254
REGPREDICTX1
ARRAY1540842
```

METHODS

FITSELF X Y FIT THE MODEL USING X Y AS TRAINING DATA  
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR  
PREDICT SELF X PREDICT USING THE LINEAR MODEL  
SCORE SELF X Y SAMPLEWEIGHT RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PRE  
DICTION  
SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR  
INIT SELF FITINTERCEPTTRUE VERBOSEFALSE MAXITER500 NORMALIZETRUE PRECOM  
PUTE'AUTO' CV'WARN' MAXNALPHAS1000 NJOBSNONE EPS2220446049250313E  
16COPYXTRUE POSITIVEFALSE  
FITSELFXY  
FIT THE MODEL USING X Y AS TRAINING DATA  
PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA  
YARRAYLIKE SHAPE NSAMPLES TARGET VALUES  
RETURNS  
SELF OBJECT RETURNS AN INSTANCE OF SELF  
GETPARAMS SELFDEEPTTRUE  
GET PARAMETERS FOR THIS ESTIMATOR  
PARAMETERS  
DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBOBJECTS THAT ARE ESTIMATORS  
RETURNS  
33 MODEL SELECTION AND EVALUATION 459



SCIKITLEARN USER GUIDE RELEASE 0213  
SKLEARNLINEARMODEL LASSOCV  
CLASSSSKLEARNLINEARMODEL LASSOCVEPS0001 NALPHAS100 ALPHASNONE FITINTERCEPTTRUE  
NORMALIZEFALSE PRECOMPUTE'AUTO' MAXITER1000  
TOL00001 COPYXTRUE CV'WARN' VERBOSEFALSE  
NJOBSNONE POSITIVEFALSE RANDOMSTATENONE SELECTION'CYCLIC'  
LASSO LINEAR MODEL WITH ITERATIVE FITTING ALONG A REGULARIZATION PATH  
SEE GLOSSARY ENTRY FOR CROSSVALIDATION ESTIMATOR  
THE BEST MODEL IS SELECTED BY CROSSVALIDATION  
THE OPTIMIZATION OBJECTIVE FOR LASSO IS  
1 2 NSAMPLES Y XW22 ALPHA W1  
READ MORE IN THE USER GUIDE  
PARAMETERS  
EPS FLOAT OPTIONAL LENGTH OF THE PATH EPS1E3 MEANS THATALPHAMIN ALPHAMAX  
1E3  
NALPHAS INT OPTIONAL NUMBER OF ALPHAS ALONG THE REGULARIZATION PATH  
ALPHAS NUMPY ARRAY OPTIONAL LIST OF ALPHAS WHERE TO COMPUTE THE MODELS IF NONE ALPHAS  
ARE SET AUTOMATICALLY  
FITINTERCEPT BOOLEAN DEFAULT TRUE WHETHER TO CALCULATE THE INTERCEPT FOR THIS MODEL IF SET TO  
FALSE NO INTERCEPT WILL BE USED IN CALCULATIONS EG DATA IS EXPECTED TO BE ALREADY CENTERED  
NORMALIZE BOOLEAN OPTIONAL DEFAULT FALSE THIS PARAMETER IS IGNORED WHEN  
FITINTERCEPT IS SET TO FALSE IF TRUE THE REGRESSORS X WILL BE NORMALIZED BEFORE  
REGRESSION BY SUBTRACTING THE MEAN AND DIVIDING BY THE L2NORM IF YOU WISH TO  
STANDARDIZE PLEASE USE SKLEARNPREPROCESSINGSTANDARDSCALER BEFORE  
CALLINGFIT ON AN ESTIMATOR WITH NORMALIZEFALSE  
PRECOMPUTE TRUE FALSE 'AUTO' ARRAYLIKE WHETHER TO USE A PRECOMPUTED GRAM MATRIX TO  
SPEED UP CALCULATIONS IF SET TO AUTO LET US DECIDE THE GRAM MATRIX CAN ALSO BE PASSED  
AS ARGUMENT  
MAXITER INT OPTIONAL THE MAXIMUM NUMBER OF ITERATIONS  
TOLFLOAT OPTIONAL THE TOLERANCE FOR THE OPTIMIZATION IF THE UPDATES ARE SMALLER THAN TOL THE  
OPTIMIZATION CODE CHECKS THE DUAL GAP FOR OPTIMALITY AND CONTINUES UNTIL IT IS SMALLER THAN  
TOL  
COPYX BOOLEAN OPTIONAL DEFAULT TRUE IF TRUE X WILL BE COPIED ELSE IT MAY BE OVERWRITTEN  
CVINT CROSSVALIDATION GENERATOR OR AN ITERABLE OPTIONAL DETERMINES THE CROSSVALIDATION SPLIT  
TING STRATEGY POSSIBLE INPUTS FOR CV ARE  
• NONE TO USE THE DEFAULT 3FOLD CROSSVALIDATION  
• INTEGER TO SPECIFY THE NUMBER OF FOLDS  
• CV SPLITTER  
• AN ITERABLE YIELDING TRAIN TEST SPLITS AS ARRAYS OF INDICES  
FOR INTEGERNONE INPUTS KFOLD IS USED  
REFER USER GUIDE FOR THE VARIOUS CROSSVALIDATION STRATEGIES THAT CAN BE USED HERE  
33 MODEL SELECTION AND EVALUATION 461

SCIKITLEARN USER GUIDE RELEASE 0213

CHANGED IN VERSION 020 CVDEFAULT VALUE IF NONE WILL CHANGE FROM 3FOLD TO 5FOLD IN V022

VERBOSE BOOL OR INTEGER AMOUNT OF VERBOSITY

NJOBS INT OR NONE OPTIONAL DEFAULTNONE NUMBER OF CPUS TO USE DURING THE CROSS VALIDATIONNONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

POSITIVE BOOL OPTIONAL IF POSITIVE RESTRICT REGRESSION COEFFICIENTS TO BE POSITIVE

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULT NONE THE SEED OF THE PSEUDO RANDOM NUMBER GENERATOR THAT SELECTS A RANDOM FEATURE TO UPDATE IF INT RANDOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NPRANDOM USED WHEN SELECTION 'RANDOM'

SELECTION STR DEFAULT 'CYCLIC' IF SET TO 'RANDOM' A RANDOM COEFFICIENT IS UPDATED EVERY ITERATION RATHER THAN LOOPING OVER FEATURES SEQUENTIALLY BY DEFAULT THIS SETTING TO 'RANDOM' OFTEN LEADS TO SIGNIFICANTLY FASTER CONVERGENCE ESPECIALLY WHEN TOL IS HIGHER THAN 1E4

ATTRIBUTES

ALPHA FLOAT THE AMOUNT OF PENALIZATION CHOSEN BY CROSS VALIDATION

COEF ARRAY SHAPE NFEATURES NTARGETS NFEATURES PARAMETER VECTOR W IN THE COST FUNCTION FORMULA

INTERCEPT FLOAT ARRAY SHAPE NTARGETS INDEPENDENT TERM IN DECISION FUNCTION

MSEPATH ARRAY SHAPE NALPHAS NFOLDS MEAN SQUARE ERROR FOR THE TEST SET ON EACH FOLD VARYING ALPHA

ALPHAS NUMPY ARRAY SHAPE NALPHAS THE GRID OF ALPHAS USED FOR FITTING

DUALGAP NDARRAY SHAPE THE DUAL GAP AT THE END OF THE OPTIMIZATION FOR THE OPTIMAL ALPHA ALPHA

NITER INT NUMBER OF ITERATIONS RUN BY THE COORDINATE DESCENT SOLVER TO REACH THE SPECIFIED TOLERANCE FOR THE OPTIMAL ALPHA

SEE ALSO

LARSPATH

LASSOPATH

LASSOLARS

LASSO

LASSOLARSCV

NOTES

FOR AN EXAMPLE SEE EXAMPLESLINEARMODELPLOTLASSOMODELSELECTIONPY

TO AVOID UNNECESSARY MEMORY DUPLICATION THE X ARGUMENT OF THE FIT METHOD SHOULD BE DIRECTLY PASSED AS A FORTRANCONTIGUOUS NUMPY ARRAY

462 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES

FROM SKLEARNLINEARMODEL IMPORT LASSOCV  
FROM SKLEARNDATASETS IMPORT MAKEREGRESSION

X Y MAKEREGRESSIONNOISE4 RANDOMSTATE0

REG LASSOCVCV5 RANDOMSTATE0FITX Y

REGSCOREX Y

09993

REGPREDICTX1

ARRAY784951

METHODS

FITSELF X Y FIT LINEAR MODEL WITH COORDINATE DESCENT

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PATH X Y EPS NALPHAS ALPHAS COMPUTE LASSO PATH WITH COORDINATE DESCENT

PREDICT SELF X PREDICT USING THE LINEAR MODEL

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PRE  
DICTION

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELF EPS0001 NALPHAS100 ALPHASNONE FITINTERCEPTTRUE NORMALIZEFALSE PRE

COMPUTE'AUTO' MAXITER1000 TOL00001 COPYXTRUE CV'WARN' VERBOSEFALSE

NJOBSNONE POSITIVEFALSE RANDOMSTATENONE SELECTION'CYCLIC'

FITSELFXY

FIT LINEAR MODEL WITH COORDINATE DESCENT

FIT IS ON GRID OF ALPHAS AND BEST ALPHA ESTIMATED BY CROSSVALIDATION

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA PASS DIRECTLY AS FORTRAN

CONTIGUOUS DATA TO AVOID UNNECESSARY MEMORY DUPLICATION IF Y IS MONOOUTPUT X CAN  
BE SPARSE

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NTARGETS TARGET VALUES

GETPARAMS SELFDEEPTTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

STATICPATHXYEPS0001 NALPHAS100 ALPHASNONE PRECOMPUTE'AUTO' XYNONE

COPYXTRUE COEFINITNONE VERBOSEFALSE RETURNNITERFALSE POSITIVEFALSE

PARAMS

COMPUTE LASSO PATH WITH COORDINATE DESCENT

THE LASSO OPTIMIZATION FUNCTION VARIES FOR MONO AND MULTIOUTPUTS

FOR MONOOUTPUT TASKS IT IS

33 MODEL SELECTION AND EVALUATION 463

SCIKITLEARN USER GUIDE RELEASE 0213

1 2 NSAMPLES Y XW22 ALPHA W1  
FOR MULTIOUTPUT TASKS IT IS

1 2 NSAMPLES Y XW2FRO ALPHA W21  
WHERE  
W21 SUMI SQRTSUMJ WIJ2  
IE THE SUM OF NORM OF EACH ROW  
READ MORE IN THE USER GUIDE

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES TRAINING DATA PASS DIRECTLY  
AS FORTRANCONTIGUOUS DATA TO AVOID UNNECESSARY MEMORY DUPLICATION IF YIS MONOOUTPUT  
THENXCAN BE SPARSE

YNDARRAY SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TARGET VALUES  
EPS FLOAT OPTIONAL LENGTH OF THE PATH EPS1E3 MEANS THAT ALPHAMIN  
ALPHAMAX 1E3  
NALPHAS INT OPTIONAL NUMBER OF ALPHAS ALONG THE REGULARIZATION PATH  
ALPHAS NDARRAY OPTIONAL LIST OF ALPHAS WHERE TO COMPUTE THE MODELS IF NONE ALPHAS ARE  
SET AUTOMATICALLY

PRECOMPUTE TRUE FALSE 'AUTO' ARRAYLIKE WHETHER TO USE A PRECOMPUTED GRAM MATRIX  
TO SPEED UP CALCULATIONS IF SET TO AUTO LET US DECIDE THE GRAM MATRIX CAN ALSO BE  
PASSED AS ARGUMENT

XYARRAYLIKE OPTIONAL XY NPDOTXT Y THAT CAN BE PRECOMPUTED IT IS USEFUL ONLY  
WHEN THE GRAM MATRIX IS PRECOMPUTED

COPYX BOOLEAN OPTIONAL DEFAULT TRUE IF TRUE X WILL BE COPIED ELSE IT MAY BE OVER  
WRITTEN

COEFINIT ARRAY SHAPE NFEATURES NONE THE INITIAL VALUES OF THE COEFFICIENTS  
VERBOSE BOOL OR INTEGER AMOUNT OF VERBOSITY

RETURNNNITER BOOL WHETHER TO RETURN THE NUMBER OF ITERATIONS OR NOT  
POSITIVE BOOL DEFAULT FALSE IF SET TO TRUE FORCES COEFFICIENTS TO BE POSITIVE ONLY ALLOWED  
WHENYNDIM 1

PARAMS KWARGS KEYWORD ARGUMENTS PASSED TO THE COORDINATE DESCENT SOLVER

RETURNS

ALPHAS ARRAY SHAPE NALPHAS THE ALPHAS ALONG THE PATH WHERE MODELS ARE COMPUTED

COEFS ARRAY SHAPE NFEATURES NALPHAS OR NOUTPUTS NFEATURES NALPHAS COEFFICIENTS  
ALONG THE PATH

DUALGAPS ARRAY SHAPE NALPHAS THE DUAL GAPS AT THE END OF THE OPTIMIZATION FOR EACH  
ALPHA

NITERS ARRAYLIKE SHAPE NALPHAS THE NUMBER OF ITERATIONS TAKEN BY THE COORDINATE DE  
SCENT OPTIMIZER TO REACH THE SPECIFIED TOLERANCE FOR EACH ALPHA

464 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

SEE ALSO

LARSPATH

LASSO

LASSOLARS

LASSOCV

LASSOLARSCV

SKLEARNDECOMPOSITIONSPARSEENCODE

NOTES

FOR AN EXAMPLE SEE EXAMPLESLINEARMODELPLOTLASSOCOORDINATEDESCENTPATHPY

TO AVOID UNNECESSARY MEMORY DUPLICATION THE X ARGUMENT OF THE FIT METHOD SHOULD BE DIRECTLY PASSED AS A

FORTRANCONTIGUOUS NUMPY ARRAY

NOTE THAT IN CERTAIN CASES THE LARS SOLVER MAY BE SIGNIFICANTLY FASTER TO IMPLEMENT THIS FUNCTIONALITY IN

PARTICULAR LINEAR INTERPOLATION CAN BE USED TO RETRIEVE MODEL COEFFICIENTS BETWEEN THE VALUES OUTPUT BY

LARSPATH

EXAMPLES

COMPARING LASSOPATH AND LARSPATH WITH INTERPOLATION

X NPARRAY1 2 31 23 54 43T

Y NPARRAY1 2 31

USE LASSOPATH TO COMPUTE A COEFFICIENT PATH

COEFPATH LASSOPATHX Y ALPHAS5 1 5

PRINTCOEFPATH

0 0 046874778

02159048 04425765 023689075

NOW USE LARSPATH AND 1D LINEAR INTERPOLATION TO COMPUTE THE

SAME PATH

FROM SKLEARNLINEARMODEL IMPORT LARSPATH

ALPHAS ACTIVE COEFPATHLARS LARSPATHX Y METHODLASSO

FROM SCIPY IMPORT INTERPOLATE

COEFPATHCONTINUOUS INTERPOLATEINTERP1DALPHAS1

COEFPATHLARS 1

PRINTCOEFPATHCONTINUOUS5 1 5

0 0 046915237

02159048 04425765 023668876

PREDICTSELF

PREDICT USING THE LINEAR MODEL

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES SAMPLES

RETURNS

CARRAY SHAPE NSAMPLES RETURNS PREDICTED VALUES

33 MODEL SELECTION AND EVALUATION 465

SCIKITLEARN USER GUIDE RELEASE 0213

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{U}{V}$  WHERE U IS THE RESIDUAL SUM OF SQUARES  $\sum (y_{true} - y_{pred})^2$  AND V IS THE TOTAL SUM OF SQUARES  $\sum (y_{true} - y_{true\_mean})^2$  THE BEST POSSIBLE SCORE IS 1.0 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 0.0

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES FOR SOME ESTIMATORS THIS MAY BE A PRECOMPUTED KERNEL MATRIX INSTEAD SHAPE NSAMPLES NSAMPLESFITTED WHERE NSAMPLESFITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE VALUES FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT R2 OF SELF-PREDICT X WRT Y

NOTES

THE R2 SCORE USED WHEN CALLING SCORE ON A REGRESSOR WILL USE MULTIOUTPUTUNIFORMAVERAGE FROM VERSION 0.23 TO KEEP CONSISTENT WITH METRICS.R2SCORE THIS WILL INFLUENCE THE SCORE METHOD OF ALL THE MULTIOUTPUT REGRESSORS EXCEPT FOR MULTIOUTPUTMULTIOUTPUTREGRESSOR TO SPECIFY THE DEFAULT VALUE MANUALLY AND AVOID THE WARNING PLEASE EITHER CALL METRICS.R2SCORE DIRECTLY OR MAKE A CUSTOM SCORER WITH METRICS.MAKESCORER THE BUILT-IN SCORER R2 USES MULTIOUTPUTUNIFORMAVERAGE

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENT.PARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

EXAMPLES USING SKLEARN.LINEAR.MODEL.LASSOCV

- CROSSVALIDATION ON DIABETES DATASET EXERCISE
- FEATURE SELECTION USING SELECTFROMMODEL AND LASSOCV
- LASSO MODEL SELECTION CROSSVALIDATION AIC BIC

SKLEARN.LINEAR.MODEL.LASSO.LARSCV

CLASS SKLEARN.LINEAR.MODEL.LASSO.LARSCV FIT\_INTERCEPT=True VERBOSE=False MAX\_ITER=500

NORMALIZE=True PRECOMPUTE='AUTO'

CV='WARN' MAX\_N\_ALPHAS=1000 N\_JOBS=None

EPS=2.220446049250313E-16 COPY\_X=True POSITIVE=False

CROSSVALIDATED LASSO USING THE LARS ALGORITHM

466 CHAPTER 3 USER GUIDE



SCIKITLEARN USER GUIDE RELEASE 0213  
SEE GLOSSARY ENTRY FOR CROSSVALIDATION ESTIMATOR  
THE OPTIMIZATION OBJECTIVE FOR LASSO IS  
 $\frac{1}{2} \sum_{i=1}^n \text{NSAMPLES } Y_i - \sum_{j=1}^p \text{XW}_{22} \text{ ALPHA } W_{1j}$   
READ MORE IN THE USER GUIDE  
PARAMETERS  
FITINTERCEPT BOOLEAN WHETHER TO CALCULATE THE INTERCEPT FOR THIS MODEL IF SET TO FALSE NO  
INTERCEPT WILL BE USED IN CALCULATIONS EG DATA IS EXPECTED TO BE ALREADY CENTERED  
VERBOSE BOOLEAN OR INTEGER OPTIONAL SETS THE VERBOSITY AMOUNT  
MAXITER INTEGER OPTIONAL MAXIMUM NUMBER OF ITERATIONS TO PERFORM  
NORMALIZE BOOLEAN OPTIONAL DEFAULT TRUE THIS PARAMETER IS IGNORED WHEN FITINTERCEPT  
IS SET TO FALSE IF TRUE THE REGRESSORS X WILL BE NORMALIZED BEFORE REGRESSION BY SUB  
TRACTING THE MEAN AND DIVIDING BY THE L2NORM IF YOU WISH TO STANDARDIZE PLEASE USE  
SKLEARNPREPROCESSINGSTANDARDSCALER BEFORE CALLING FIT ON AN ESTIMATOR  
WITHNORMALIZEFALSE  
PRECOMPUTE TRUE FALSE 'AUTO' WHETHER TO USE A PRECOMPUTED GRAM MATRIX TO SPEED UP  
CALCULATIONS IF SET TO AUTO LET US DECIDE THE GRAM MATRIX CANNOT BE PASSED AS ARGUMENT  
SINCE WE WILL USE ONLY SUBSETS OF X  
CVINT CROSSVALIDATION GENERATOR OR AN ITERABLE OPTIONAL DETERMINES THE CROSSVALIDATION SPLIT  
TING STRATEGY POSSIBLE INPUTS FOR CV ARE  
• NONE TO USE THE DEFAULT 3FOLD CROSSVALIDATION  
• INTEGER TO SPECIFY THE NUMBER OF FOLDS  
• CV SPLITTER  
• AN ITERABLE YIELDING TRAIN TEST SPLITS AS ARRAYS OF INDICES  
FOR INTEGERNONE INPUTS KFOLD IS USED  
REFER USER GUIDE FOR THE VARIOUS CROSSVALIDATION STRATEGIES THAT CAN BE USED HERE  
CHANGED IN VERSION 020 CVDEFAULT VALUE IF NONE WILL CHANGE FROM 3FOLD TO 5FOLD IN  
V022  
MAXNALPHAS INTEGER OPTIONAL THE MAXIMUM NUMBER OF POINTS ON THE PATH USED TO COMPUTE  
THE RESIDUALS IN THE CROSSVALIDATION  
NJOBS INT OR NONE OPTIONAL DEFAULTNONE NUMBER OF CPUS TO USE DURING THE CROSS VALIDA  
TIONNONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING  
ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS  
EPS FLOAT OPTIONAL THE MACHINEPRECISION REGULARIZATION IN THE COMPUTATION OF THE CHOLESKY  
DIAGONAL FACTORS INCREASE THIS FOR VERY ILLCONDITIONED SYSTEMS  
COPYX BOOLEAN OPTIONAL DEFAULT TRUE IF TRUE X WILL BE COPIED ELSE IT MAY BE OVERWRITTEN  
POSITIVE BOOLEAN DEFAULTFALSE RESTRICT COEFFICIENTS TO BE  $\geq 0$  BE AWARE THAT YOU MIGHT  
WANT TO REMOVE FITINTERCEPT WHICH IS SET TRUE BY DEFAULT UNDER THE POSITIVE RESTRICTION THE  
MODEL COEFFICIENTS DO NOT CONVERGE TO THE ORDINARYLEASTSQUARES SOLUTION FOR SMALL VALUES  
OF ALPHA ONLY COEFFICIENTS UP TO THE SMALLEST ALPHA VALUE  $\text{ALPHASALPHAS} = 0$   
MIN WHEN FITPATHTRUE REACHED BY THE STEPWISE LARSLASSO ALGORITHM ARE TYPICALLY IN  
CONGRUENCE WITH THE SOLUTION OF THE COORDINATE DESCENT LASSO ESTIMATOR AS A CONSEQUENCE  
33 MODEL SELECTION AND EVALUATION 467

SCIKITLEARN USER GUIDE RELEASE 0213

USING LASSOLARSCV ONLY MAKES SENSE FOR PROBLEMS WHERE A SPARSE SOLUTION IS EXPECTED  
AND/OR REACHED

ATTRIBUTES

COEF ARRAY SHAPE NFEATURES PARAMETER VECTOR W IN THE FORMULATION FORMULA

INTERCEPT FLOAT INDEPENDENT TERM IN DECISION FUNCTION

COEFPATH ARRAY SHAPE NFEATURES NALPHAS THE VARYING VALUES OF THE COEFFICIENTS ALONG THE  
PATH

ALPHA FLOAT THE ESTIMATED REGULARIZATION PARAMETER ALPHA

ALPHAS ARRAY SHAPE NALPHAS THE DIFFERENT VALUES OF ALPHA ALONG THE PATH

CVALPHAS ARRAY SHAPE NCVALPHAS ALL THE VALUES OF ALPHA ALONG THE PATH FOR THE DIFFERENT  
FOLDS

MSEPATH ARRAY SHAPE NFOLDS NCVALPHAS THE MEAN SQUARE ERROR ON LEFTOUT FOR EACH FOLD  
ALONG THE PATH ALPHA VALUES GIVEN BY CVALPHAS

NITER ARRAYLIKE OR INT THE NUMBER OF ITERATIONS RUN BY LARS WITH THE OPTIMAL ALPHA

SEE ALSO

LARSPATH LASSOLARS LARSCV LASSOCV

NOTES

THE OBJECT SOLVES THE SAME PROBLEM AS THE LASSOCV OBJECT HOWEVER UNLIKE THE LASSOCV IT FIND THE RELEVANT  
ALPHAS VALUES BY ITSELF IN GENERAL BECAUSE OF THIS PROPERTY IT WILL BE MORE STABLE HOWEVER IT IS MORE FRAGILE TO  
HEAVILY MULTICOLLINEAR DATASETS

IT IS MORE EFFICIENT THAN THE LASSOCV IF ONLY A SMALL NUMBER OF FEATURES ARE SELECTED COMPARED TO THE TOTAL  
NUMBER FOR INSTANCE IF THERE ARE VERY FEW SAMPLES COMPARED TO THE NUMBER OF FEATURES

EXAMPLES

```
FROM SKLEARNLINEARMODEL IMPORT LASSOLARSCV
FROM SKLEARNDATASETS IMPORT MAKEREGRESSION
X Y MAKEREGRESSIONNOISE40 RANDOMSTATE0
REG LASSOLARSCVCV5FITX Y
REGSCOREX Y
09992
REGALPHA
00484
REGPREDICTX1
ARRAY778723
```

METHODS

FITSELF X Y FIT THE MODEL USING X Y AS TRAINING DATA

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PREDICT SELF X PREDICT USING THE LINEAR MODEL

CONTINUED ON NEXT PAGE

468 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

TABLE 35 – CONTINUED FROM PREVIOUS PAGE

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELF FIT INTERCEPT TRUE VERBOSE FALSE MAX ITER 500 NORMALIZE TRUE PRECOMPUTE 'AUTO' CV 'WARN' MAX NALPHAS 1000 NJOBS NONE EPS 2.220446049250313E-16 COPY X TRUE POSITIVE FALSE

FIT SELF X Y

FIT THE MODEL USING X Y AS TRAINING DATA

PARAMETERS

X ARRAY LIKE SHAPE NSAMPLES NFEATURES TRAINING DATA

Y ARRAY LIKE SHAPE NSAMPLES TARGET VALUES

RETURNS

SELF OBJECT RETURNS AN INSTANCE OF SELF

GETPARAMS SELF DEEP TRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICT SELF X

PREDICT USING THE LINEAR MODEL

PARAMETERS

X ARRAY LIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES SAMPLES

RETURNS

C ARRAY SHAPE NSAMPLES RETURNS PREDICTED VALUES

SCORE SELF X Y SAMPLEWEIGHT NONE

RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{U}{V}$  WHERE U IS THE RESIDUAL SUM OF SQUARES Y TRUE Y PRED 2 SUM AND V IS THE TOTAL SUM OF SQUARES Y TRUE Y TRUE MEAN 2 SUM THE BEST POSSIBLE SCORE IS 10 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 00

PARAMETERS

X ARRAY LIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES FOR SOME ESTIMATORS THIS MAY BE A PRECOMPUTED KERNEL MATRIX INSTEAD SHAPE NSAMPLES NSAMPLES FITTED WHERE NSAMPLES FITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR

Y ARRAY LIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE VALUES FOR X

SAMPLEWEIGHT ARRAY LIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT R2 OF SELF PREDICT X WRT Y

33 MODEL SELECTION AND EVALUATION 469

SCIKITLEARN USER GUIDE RELEASE 0213

NOTES

THE R2 SCORE USED WHEN CALLING SCORE ON A REGRESSOR WILL USE MULTIOUTPUTUNIFORMAVERAGE FROM VERSION 023 TO KEEP CONSISTENT WITH METRICSR2SCORE THIS WILL INFLUENCE THE SCORE METHOD OF ALL THE MULTIOUTPUT REGRESSORS EXCEPT FOR MULTIOUTPUTMULTIOUTPUTPUTREGRESSOR TO SPECIFY THE DEFAULT VALUE MANUALLY AND AVOID THE WARNING PLEASE EITHER CALL METRICSR2SCORE DIRECTLY OR MAKE A CUSTOM SCORER WITH METRICSMAKESCORER THE BUILTIN SCORER R2 USES MULTIOUTPUTUNIFORMAVERAGE

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

EXAMPLES USING SKLEARNLINEARMODELLASSOLARSCV

•LASSO MODEL SELECTION CROSSVALIDATION AIC BIC

SKLEARNLINEARMODEL LOGISTICREGRESSIONCV

CLASSSSKLEARNLINEARMODEL LOGISTICREGRESSIONCV CS10 FITINTERCEPTTRUE CV'WARN'

DUALFALSE PENALTY'L2' SCOR

INGNONE SOLVER'LBFGS' TOL00001

MAXITER100 CLASSWEIGHTNONE

NJOBSNONE VERBOSE0 REFITTRUE IN

TERCEPTSCALING10 MULTICLASS'WARN'

RANDOMSTATENONE L1RATIOSNONE

LOGISTIC REGRESSION CV AKA LOGIT MAXENT CLASSIFIER

SEE GLOSSARY ENTRY FOR CROSSVALIDATION ESTIMATOR

THIS CLASS IMPLEMENTS LOGISTIC REGRESSION USING LIBLINEAR NEWTONCG SAG OF LBFGS OPTIMIZER THE NEWTONCG SAG AND LBFGS SOLVERS SUPPORT ONLY L2 REGULARIZATION WITH PRIMAL FORMULATION THE LIBLINEAR SOLVER SUPPORTS BOTH L1 AND L2 REGULARIZATION WITH A DUAL FORMULATION ONLY FOR THE L2 PENALTY ELASTICNET PENALTY IS ONLY SUPPORTED BY THE SAGA SOLVER

FOR THE GRID OF CSVALUES AND L1RATIOS VALUES THE BEST HYPERPARAMETER IS SELECTED BY THE CROSSVALIDATOR

STRATIFIEDKFOLD BUT IT CAN BE CHANGED USING THE CVPARAMETER THE 'NEWTONCG' 'SAG' 'SAGA' AND 'LBFGS'

SOLVERS CAN WARMSTART THE COEFFICIENTS SEE GLOSSARY

READ MORE IN THE USER GUIDE

PARAMETERS

CSLIST OF FLOATS OR INT OPTIONAL DEFAULT10 EACH OF THE VALUES IN CS DESCRIBES THE INVERSE OF REGULARIZATION STRENGTH IF CS IS AS AN INT THEN A GRID OF CS VALUES ARE CHOSEN IN A LOGARITH MIC SCALE BETWEEN 1E4 AND 1E4 LIKE IN SUPPORT VECTOR MACHINES SMALLER VALUES SPECIFY STRONGER REGULARIZATION

FITINTERCEPT BOOL OPTIONAL DEFAULTTRUE SPECIFIES IF A CONSTANT AKA BIAS OR INTERCEPT

SHOULD BE ADDED TO THE DECISION FUNCTION

470 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

CVINT OR CROSSVALIDATION GENERATOR OPTIONAL DEFAULTNONE THE DEFAULT CROSSVALIDATION GENERATOR USED IS STRATIFIED KFOLDS IF AN INTEGER IS PROVIDED THEN IT IS THE NUMBER OF FOLDS USED SEE THE MODULE SKLEARNMODELSELECTION MODULE FOR THE LIST OF POSSIBLE CROSSVALIDATION OBJECTS

CHANGED IN VERSION 020 CVDEFAULT VALUE IF NONE WILL CHANGE FROM 3FOLD TO 5FOLD IN V022

DUAL BOOL OPTIONAL DEFAULTFALSE DUAL OR PRIMAL FORMULATION DUAL FORMULATION IS ONLY IMPLEMENTED FOR L2 PENALTY WITH LIBLINEAR SOLVER PREFER DUALFALSE WHEN NSAMPLES NFEATURES

PENALTY STR 'L1' 'L2' OR 'ELASTICNET' OPTIONAL DEFAULT'L2' USED TO SPECIFY THE NORM USED IN THE PENALIZATION THE 'NEWTONCG' 'SAG' AND 'LBFGS' SOLVERS SUPPORT ONLY L2 PENALTIES 'ELASTICNET' IS ONLY SUPPORTED BY THE 'SAGA' SOLVER

SCORING STRING CALLABLE OR NONE OPTIONAL DEFAULTNONE A STRING SEE MODEL EVALUATION DOCUMENTATION OR A SCORER CALLABLE OBJECT FUNCTION WITH SIGNATURE SCORERESTIMATOR X Y FOR A LIST OF SCORING FUNCTIONS THAT CAN BE USED LOOK AT SKLEARNMETRICS THE DEFAULT SCORING OPTION USED IS 'ACCURACY'

SOLVER STR 'NEWTONCG' 'LBFGS' 'LIBLINEAR' 'SAG' 'SAGA' OPTIONAL DEFAULT'LBFGS' ALGORITHM TO USE IN THE OPTIMIZATION PROBLEM

- FOR SMALL DATASETS 'LIBLINEAR' IS A GOOD CHOICE WHEREAS 'SAG' AND 'SAGA' ARE FASTER FOR LARGE ONES
- FOR MULTICLASS PROBLEMS ONLY 'NEWTONCG' 'SAG' 'SAGA' AND 'LBFGS' HANDLE MULTINOMIAL LOSS 'LIBLINEAR' IS LIMITED TO ONEVERUSREST SCHEMES
- 'NEWTONCG' 'LBFGS' AND 'SAG' ONLY HANDLE L2 PENALTY WHEREAS 'LIBLINEAR' AND 'SAGA' HANDLE L1 PENALTY
- 'LIBLINEAR' MIGHT BE SLOWER IN LOGISTICREGRESSIONCV BECAUSE IT DOES NOT HANDLE WARM STARTING

NOTE THAT 'SAG' AND 'SAGA' FAST CONVERGENCE IS ONLY GUARANTEED ON FEATURES WITH APPROXIMATELY THE SAME SCALE YOU CAN PREPROCESS THE DATA WITH A SCALER FROM SKLEARNPREPROCESSING

NEW IN VERSION 017 STOCHASTIC AVERAGE GRADIENT DESCENT SOLVER

NEW IN VERSION 019 SAGA SOLVER

TOLFLOAT OPTIONAL DEFAULT1E4 TOLERANCE FOR STOPPING CRITERIA

MAXITER INT OPTIONAL DEFAULT100 MAXIMUM NUMBER OF ITERATIONS OF THE OPTIMIZATION ALGORITHM

CLASSWEIGHT DICT OR 'BALANCED' OPTIONAL DEFAULTNONE WEIGHTS ASSOCIATED WITH CLASSES IN THE FORMCLASSLABEL WEIGHT IF NOT GIVEN ALL CLASSES ARE SUPPOSED TO HAVE WEIGHT ONE

THE "BALANCED" MODE USES THE VALUES OF Y TO AUTOMATICALLY ADJUST WEIGHTS INVERSELY PROPORTIONAL TO CLASS FREQUENCIES IN THE INPUT DATA AS NSAMPLES NCLASSES NP

BINCOUNTY

NOTE THAT THESE WEIGHTS WILL BE MULTIPLIED WITH SAMPLEWEIGHT PASSED THROUGH THE FIT METHOD IF SAMPLEWEIGHT IS SPECIFIED

NEW IN VERSION 017 CLASSWEIGHT 'BALANCED'

33 MODEL SELECTION AND EVALUATION 471

SCIKITLEARN USER GUIDE RELEASE 0213

NJOBS INT OR NONE OPTIONAL DEFAULTNONE NUMBER OF CPU CORES USED DURING THE CROSS  
VALIDATION LOOP NONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1  
MEANS USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

VERBOSE INT OPTIONAL DEFAULT0 FOR THE 'LIBLINEAR' 'SAG' AND 'LBFGS' SOLVERS SET VERBOSE TO  
ANY POSITIVE NUMBER FOR VERBOSITY

REFIT BOOL OPTIONAL DEFAULTTRUE IF SET TO TRUE THE SCORES ARE AVERAGED ACROSS ALL FOLDS AND  
THE COEFS AND THE C THAT CORRESPONDS TO THE BEST SCORE IS TAKEN AND A FINAL REFIT IS DONE USING  
THESE PARAMETERS OTHERWISE THE COEFS INTERCEPTS AND C THAT CORRESPOND TO THE BEST SCORES  
ACROSS FOLDS ARE AVERAGED

INTERCEPTSCALING FLOAT OPTIONAL DEFAULT1 USEFUL ONLY WHEN THE SOLVER 'LIBLINEAR' IS USED  
AND SELFITINTERCEPT IS SET TO TRUE IN THIS CASE X BECOMES X SELFINTERCEPTSCALING  
IE A "SYNTHETIC" FEATURE WITH CONSTANT VALUE EQUAL TO INTERCEPTSCALING IS AP  
PENDE TO THE INSTANCE VECTOR THE INTERCEPT BECOMES INTERCEPTSCALING  
SYNTHETICFEATUREWEIGHT

NOTE THE SYNTHETIC FEATURE WEIGHT IS SUBJECT TO L1L2 REGULARIZATION AS ALL OTHER FEATURES TO  
LESSEN THE EFFECT OF REGULARIZATION ON SYNTHETIC FEATURE WEIGHT AND THEREFORE ON THE INTERCEPT  
INTERCEPTSCALING HAS TO BE INCREASED

MULTICLASS STR 'OVR' 'MULTINOMIAL' 'AUTO' OPTIONAL DEFAULT'OVR' IF THE OPTION CHOSEN  
IS 'OVR' THEN A BINARY PROBLEM IS FIT FOR EACH LABEL FOR 'MULTINOMIAL' THE LOSS MINIMISED  
IS THE MULTINOMIAL LOSS FIT ACROSS THE ENTIRE PROBABILITY DISTRIBUTION EVEN WHEN THE DATA IS  
BINARY 'MULTINOMIAL' IS UNAVAILABLE WHEN SOLVER'LIBLINEAR' 'AUTO' SELECTS 'OVR' IF THE DATA  
IS BINARY OR IF SOLVER'LIBLINEAR' AND OTHERWISE SELECTS 'MULTINOMIAL'

NEW IN VERSION 018 STOCHASTIC AVERAGE GRADIENT DESCENT SOLVER FOR 'MULTINOMIAL' CASE  
CHANGED IN VERSION 020 DEFAULT WILL CHANGE FROM 'OVR' TO 'AUTO' IN 022

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RAN  
DOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE  
RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS  
THE RANDOMSTATE INSTANCE USED BY NPRANDOM

L1RATIOS LIST OF FLOAT OR NONE OPTIONAL DEFAULTNONE THE LIST OF ELASTICNET MIXING PARAME  
TER WITH0 L1RATIO 1 ONLY USED IF PENALTYELASTICNET A VALUE OF  
0 IS EQUIVALENT TO USING PENALTYL2 WHILE 1 IS EQUIVALENT TO USING PENALTYL1  
FOR0 L1RATIO 1 THE PENALTY IS A COMBINATION OF L1 AND L2

ATTRIBUTES

CLASSES ARRAY SHAPE NCLASSES A LIST OF CLASS LABELS KNOWN TO THE CLASSIFIER

COEF ARRAY SHAPE 1 NFEATURES OR NCLASSES NFEATURES COEFFICIENT OF THE FEATURES IN THE  
DECISION FUNCTION

COEF IS OF SHAPE 1 NFEATURES WHEN THE GIVEN PROBLEM IS BINARY

INTERCEPT ARRAY SHAPE 1 OR NCLASSES INTERCEPT AKA BIAS ADDED TO THE DECISION FUNC  
TION

IFFITINTERCEPT IS SET TO FALSE THE INTERCEPT IS SET TO ZERO INTERCEPT IS OF  
SHAPE1 WHEN THE PROBLEM IS BINARY

CS ARRAY SHAPE NCS ARRAY OF C IE INVERSE OF REGULARIZATION PARAMETER VALUES USED FOR  
CROSSVALIDATION

L1RATIOS ARRAY SHAPE NL1RATIOS ARRAY OF L1RATIOS USED FOR CROSSVALIDATION IF NO  
L1RATIO IS USED IE PENALTY IS NOT 'ELASTICNET' THIS IS SET TO NONE

472 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

COEFSPATHS ARRAY SHAPE NFOLDS NCS NFEATURES OR NFOLDS NCS NFEATURES 1 DICT WITH CLASSES AS THE KEYS AND THE PATH OF COEFFICIENTS OBTAINED DURING CROSSVALIDATING ACROSS EACH FOLD AND THEN ACROSS EACH CS AFTER DOING AN OVR FOR THE CORRESPONDING CLASS AS VALUES IF THE 'MULTICLASS' OPTION IS SET TO 'MULTINOMIAL' THEN THE COEFSPATHS ARE THE COEFFICIENTS CORRESPONDING TO EACH CLASS EACH DICT VALUE HAS SHAPE NFOLDS NCS NFEATURES ORNFOLDS NCS NFEATURES 1 DEPENDING ON WHETHER THE INTERCEPT IS FIT OR NOT IF PENALTYELASTICNET THE SHAPE ISNFOLDS NCS NL1RATIOS NFEATURES ORNFOLDS NCS NL1RATIOS NFEATURES 1 SCORES DICT DICT WITH CLASSES AS THE KEYS AND THE VALUES AS THE GRID OF SCORES OBTAINED DURING CROSSVALIDATING EACH FOLD AFTER DOING AN OVR FOR THE CORRESPONDING CLASS IF THE 'MULTICLASS' OPTION GIVEN IS 'MULTINOMIAL' THEN THE SAME SCORES ARE REPEATED ACROSS ALL CLASSES SINCE THIS IS THE MULTINOMIAL CLASS EACH DICT VALUE HAS SHAPE NFOLDS NCS ORNFOLDS NCS NL1RATIOS IFPENALTYELASTICNET CARRAY SHAPE NCLASSES OR NCLASSES 1 ARRAY OF C THAT MAPS TO THE BEST SCORES ACROSS EVERY CLASS IF REFIT IS SET TO FALSE THEN FOR EACH CLASS THE BEST C IS THE AVERAGE OF THE C'S THAT CORRESPOND TO THE BEST SCORES FOR EACH FOLD CIS OF SHAPENCLASSES WHEN THE PROBLEM IS BINARY L1RATIO ARRAY SHAPE NCLASSES OR NCLASSES 1 ARRAY OF L1RATIO THAT MAPS TO THE BEST SCORES ACROSS EVERY CLASS IF REFIT IS SET TO FALSE THEN FOR EACH CLASS THE BEST L1RATIO IS THE AVERAGE OF THE L1RATIO'S THAT CORRESPOND TO THE BEST SCORES FOR EACH FOLD L1RATIO IS OF SHAPENCLASSES WHEN THE PROBLEM IS BINARY NITER ARRAY SHAPE NCLASSES NFOLDS NCS OR 1 NFOLDS NCS ACTUAL NUMBER OF ITERATIONS FOR ALL CLASSES FOLDS AND CS IN THE BINARY OR MULTINOMIAL CASES THE FIRST DIMENSION IS EQUAL TO 1 IF PENALTYELASTICNET THE SHAPE IS NCLASSES NFOLDS NCS NL1RATIOS OR1 NFOLDS NCS NL1RATIOS SEE ALSO LOGISTICREGRESSION EXAMPLES FROM SKLEARNDATASETS IMPORT LOADIRIS FROM SKLEARNLINEARMODEL IMPORT LOGISTICREGRESSIONCV X Y LOADIRISRETURNXY TRUE CLF LOGISTICREGRESSIONCVCV5 RANDOMSTATE0 MULTICLASSMULTINOMIALFITX Y CLFPREDICTX2 ARRAY0 0 CLFPREDICTPROBAX2 SHAPE 2 3 CLFSCOREX Y 098 METHODS DECISIONFUNCTION SELF X PREDICT CONFIDENCE SCORES FOR SAMPLES DENSIFY SELF CONVERT COEFFICIENT MATRIX TO DENSE ARRAY FORMAT CONTINUED ON NEXT PAGE 33 MODEL SELECTION AND EVALUATION 473

SCIKITLEARN USER GUIDE RELEASE 0213

TABLE 36 – CONTINUED FROM PREVIOUS PAGE

FITSELF X Y SAMPLEWEIGHT FIT THE MODEL ACCORDING TO THE GIVEN TRAINING DATA

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PREDICT SELF X PREDICT CLASS LABELS FOR SAMPLES IN X

PREDICTLOGPROBA SELF X LOG OF PROBABILITY ESTIMATES

PREDICTPROBA SELF X PROBABILITY ESTIMATES

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE SCORE USING THE SCORING OPTION ON THE GIVEN TEST DATA AND LABELS

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

SPARSIFY SELF CONVERT COEFFICIENT MATRIX TO SPARSE FORMAT

INIT SELFCS10 FITINTERCEPTTRUE CV'WARN' DUALFALSE PENALTY'L2' SCORINGNONE SOLVER'LBFGS' TOL00001 MAXITER100 CLASSWEIGHTNONE NJOBSNONE VER BOSE0 REFITTRUE INTERCEPTSCALING10 MULTICLASS'WARN' RANDOMSTATENONE L1RATIOSNONE

DECISIONFUNCTION SELF X

PREDICT CONFIDENCE SCORES FOR SAMPLES

THE CONFIDENCE SCORE FOR A SAMPLE IS THE SIGNED DISTANCE OF THAT SAMPLE TO THE HYPERPLANE

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES SAMPLES

RETURNS

ARRAY SHAPENSAMPLES IF NCLASSES > 2 ELSE NSAMPLES NCLASSES CONFIDENCE

SCORES PER SAMPLE CLASS COMBINATION IN THE BINARY CASE CONFIDENCE SCORE FOR SELFCLASSES1 WHERE 0 MEANS THIS CLASS WOULD BE PREDICTED

DENSIFYSELF

CONVERT COEFFICIENT MATRIX TO DENSE ARRAY FORMAT

CONVERTS THE COEF MEMBER BACK TO A NUMPYNDARRAY THIS IS THE DEFAULT FORMAT OF COEF AND IS REQUIRED FOR FITTING SO CALLING THIS METHOD IS ONLY REQUIRED ON MODELS THAT HAVE PREVIOUSLY BEEN SPARSIFIED OTHERWISE IT IS A NOOP

RETURNS

SELF ESTIMATOR

FITSELFXYSAMPLEWEIGHTNONE

FIT THE MODEL ACCORDING TO THE GIVEN TRAINING DATA

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES TRAINING VECTOR WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

YARRAYLIKE SHAPE NSAMPLES TARGET VECTOR RELATIVE TO X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL ARRAY OF WEIGHTS THAT ARE ASSIGNED TO INDIVIDUAL SAMPLES IF NOT PROVIDED THEN EACH SAMPLE IS GIVEN UNIT WEIGHT

RETURNS

SELF OBJECT

GETPARAMS SELFDEEPTREE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

474 CHAPTER 3 USER GUIDE



SCIKITLEARN USER GUIDE RELEASE 0213

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELF

PREDICT CLASS LABELS FOR SAMPLES IN X

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES SAMPLES

RETURNS

CARRAY SHAPE NSAMPLES PREDICTED CLASS LABEL PER SAMPLE

PREDICTLOGPROBA SELF

LOG OF PROBABILITY ESTIMATES

THE RETURNED ESTIMATES FOR ALL CLASSES ARE ORDERED BY THE LABEL OF CLASSES

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES

RETURNS

TARRAYLIKE SHAPE NSAMPLES NCLASSES RETURNS THE LOGPROBABILITY OF THE SAMPLE FOR EACH CLASS IN THE MODEL WHERE CLASSES ARE ORDERED AS THEY ARE IN SELFCLASSES

PREDICTPROBA SELF

PROBABILITY ESTIMATES

THE RETURNED ESTIMATES FOR ALL CLASSES ARE ORDERED BY THE LABEL OF CLASSES

FOR A MULTICLASS PROBLEM IF MULTICLASS IS SET TO BE "MULTINOMIAL" THE SOFTMAX FUNCTION IS USED TO FIND THE PREDICTED PROBABILITY OF EACH CLASS ELSE USE A ONEVSREST APPROACH IE CALCULATE THE PROBABILITY OF EACH CLASS ASSUMING IT TO BE POSITIVE USING THE LOGISTIC FUNCTION AND NORMALIZE THESE VALUES ACROSS ALL THE CLASSES

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES

RETURNS

TARRAYLIKE SHAPE NSAMPLES NCLASSES RETURNS THE PROBABILITY OF THE SAMPLE FOR EACH CLASS IN THE MODEL WHERE CLASSES ARE ORDERED AS THEY ARE IN SELFCLASSES

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE SCORE USING THE SCORING OPTION ON THE GIVEN TEST DATA AND LABELS

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES

YARRAYLIKE SHAPE NSAMPLES TRUE LABELS FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT SCORE OF SELF PREDICTX WRT Y

33 MODEL SELECTION AND EVALUATION 475

SCIKITLEARN USER GUIDE RELEASE 0213

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

SPARSIFY SELF

CONVERT COEFFICIENT MATRIX TO SPARSE FORMAT

CONVERTS THE COEF MEMBER TO A SCIPYSPARSE MATRIX WHICH FOR L1REGULARIZED MODELS CAN BE MUCH MORE MEMORY AND STORAGEEFFICIENT THAN THE USUAL NUMPYNDARRAY REPRESENTATION

THEINTERCEPT MEMBER IS NOT CONVERTED

RETURNS

SELF ESTIMATOR

NOTES

FOR NONSPARSE MODELS IE WHEN THERE ARE NOT MANY ZEROS IN COEF THIS MAY ACTUALLY INCREASE MEMORY USAGE SO USE THIS METHOD WITH CARE A RULE OF THUMB IS THAT THE NUMBER OF ZERO ELEMENTS WHICH CAN BE COMPUTED WITH COEF\_0SUM\_ MUST BE MORE THAN 50 FOR THIS TO PROVIDE SIGNIFICANT BENEFITS

AFTER CALLING THIS METHOD FURTHER FITTING WITH THE PARTIALFIT METHOD IF ANY WILL NOT WORK UNTIL YOU CALL DENSIFY

SKLEARNLINEARMODEL MULTITASKELASTICNETCV

CLASSSSKLEARNLINEARMODEL MULTITASKELASTICNETCV L1RATIO0.05 EPS0.0001 NALPHAS100

ALPHASNONE FITINTERCEPTTRUE

NORMALIZEFALSE MAXITER1000

TOL0.00001 CV'WARN' COPYXTRUE

VERBOSE0 NJOBSNONE RAN

DOMSTATENONE SELECTION'CYCLIC'

MULTITASK L1L2 ELASTICNET WITH BUILTIN CROSSVALIDATION

SEE GLOSSARY ENTRY FOR CROSSVALIDATION ESTIMATOR

THE OPTIMIZATION OBJECTIVE FOR MULTITASKELASTICNET IS

$\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^m \|y_{ij} - x_{ij}^T W\|^2$

ALPHA L1RATIO W21

0.05ALPHA1 L1RATIO W21

WHERE

$W_{21} = \sum_i \sqrt{\sum_j W_{ij}^2}$

IE THE SUM OF NORM OF EACH ROW

READ MORE IN THE USER GUIDE

PARAMETERS

476 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

L1RATIO FLOAT OR ARRAY OF FLOATS THE ELASTICNET MIXING PARAMETER WITH 0 L1RATIO 1 FOR L1RATIO 1 THE PENALTY IS AN L1L2 PENALTY FOR L1RATIO 0 IT IS AN L2 PENALTY FOR 0 L1RATIO 1 THE PENALTY IS A COMBINATION OF L1L2 AND L2 THIS PARAMETER CAN BE A LIST IN WHICH CASE THE DIFFERENT VALUES ARE TESTED BY CROSSVALIDATION AND THE ONE GIVING THE BEST PREDICTION SCORE IS USED NOTE THAT A GOOD CHOICE OF LIST OF VALUES FOR L1RATIO IS OFTEN TO PUT MORE VALUES CLOSE TO 1 IE LASSO AND LESS CLOSE TO 0 IE RIDGE AS IN 1 5 7 9 95 99 1

EPS FLOAT OPTIONAL LENGTH OF THE PATH EPS1E3 MEANS THATALPHAMIN ALPHAMAX 1E3

NALPHAS INT OPTIONAL NUMBER OF ALPHAS ALONG THE REGULARIZATION PATH  
ALPHAS ARRAYLIKE OPTIONAL LIST OF ALPHAS WHERE TO COMPUTE THE MODELS IF NOT PROVIDED SET AUTOMATICALLY

FITINTERCEPT BOOLEAN WHETHER TO CALCULATE THE INTERCEPT FOR THIS MODEL IF SET TO FALSE NO INTERCEPT WILL BE USED IN CALCULATIONS EG DATA IS EXPECTED TO BE ALREADY CENTERED  
NORMALIZE BOOLEAN OPTIONAL DEFAULT FALSE THIS PARAMETER IS IGNORED WHEN  
FITINTERCEPT IS SET TO FALSE IF TRUE THE REGRESSORS X WILL BE NORMALIZED BE  
FORE REGRESSION BY SUBTRACTING THE MEAN AND DIVIDING BY THE L2NORM IF YOU WISH TO  
STANDARDIZE PLEASE USE SKLEARNPREPROCESSINGSTANDARDSCALER BEFORE  
CALLINGFIT ON AN ESTIMATOR WITH NORMALIZEFALSE

MAXITER INT OPTIONAL THE MAXIMUM NUMBER OF ITERATIONS  
TOLFLOAT OPTIONAL THE TOLERANCE FOR THE OPTIMIZATION IF THE UPDATES ARE SMALLER THAN TOL THE OPTIMIZATION CODE CHECKS THE DUAL GAP FOR OPTIMALITY AND CONTINUES UNTIL IT IS SMALLER THAN TOL

CVINT CROSSVALIDATION GENERATOR OR AN ITERABLE OPTIONAL DETERMINES THE CROSSVALIDATION SPLITTING STRATEGY POSSIBLE INPUTS FOR CV ARE

- NONE TO USE THE DEFAULT 3FOLD CROSSVALIDATION
- INTEGER TO SPECIFY THE NUMBER OF FOLDS
- CV SPLITTER
- AN ITERABLE YIELDING TRAIN TEST SPLITS AS ARRAYS OF INDICES

FOR INTEGERNONE INPUTS KFOLD IS USED  
REFER USER GUIDE FOR THE VARIOUS CROSSVALIDATION STRATEGIES THAT CAN BE USED HERE  
CHANGED IN VERSION 020 CVDEFAULT VALUE IF NONE WILL CHANGE FROM 3FOLD TO 5FOLD IN V022

COPYX BOOLEAN OPTIONAL DEFAULT TRUE IF TRUE X WILL BE COPIED ELSE IT MAY BE OVERWRITTEN  
VERBOSE BOOL OR INTEGER AMOUNT OF VERBOSITY  
NJOBS INT OR NONE OPTIONAL DEFAULTNONE NUMBER OF CPUS TO USE DURING THE CROSS VALI  
DATION NOTE THAT THIS IS USED ONLY IF MULTIPLE VALUES FOR L1RATIO ARE GIVEN NONE MEANS 1  
UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL PROCESSORS SEE  
GLOSSARY FOR MORE DETAILS

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULT NONE THE SEED OF THE  
PSEUDO RANDOM NUMBER GENERATOR THAT SELECTS A RANDOM FEATURE TO UPDATE IF INT RAN  
DOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE  
RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS  
THE RANDOMSTATE INSTANCE USED BY NPRANDOM USED WHEN SELECTION 'RANDOM'

SCIKITLEARN USER GUIDE RELEASE 0213

SELECTION STR DEFAULT 'CYCLIC' IF SET TO 'RANDOM' A RANDOM COEFFICIENT IS UPDATED EVERY ITERATION  
RATHER THAN LOOPING OVER FEATURES SEQUENTIALLY BY DEFAULT THIS SETTING TO 'RANDOM' OFTEN  
LEADS TO SIGNIFICANTLY FASTER CONVERGENCE ESPECIALLY WHEN TOL IS HIGHER THAN 1E4

ATTRIBUTES

INTERCEPT ARRAY SHAPE NTASKS INDEPENDENT TERM IN DECISION FUNCTION

COEF ARRAY SHAPE NTASKS NFEATURES PARAMETER VECTOR W IN THE COST FUNCTION FORMULA  
NOTE THATCOEF STORES THE TRANSPOSE OF WWT

ALPHA FLOAT THE AMOUNT OF PENALIZATION CHOSEN BY CROSS VALIDATION

MSEPATH ARRAY SHAPE NALPHAS NFOLDS OR NL1RATIO NALPHAS NFOLDS MEAN SQUARE  
ERROR FOR THE TEST SET ON EACH FOLD VARYING ALPHA

ALPHAS NUMPY ARRAY SHAPE NALPHAS OR NL1RATIO NALPHAS THE GRID OF ALPHAS USED FOR  
FITTING FOR EACH L1RATIO

L1RATIO FLOAT BEST L1RATIO OBTAINED BY CROSSVALIDATION

NITER INT NUMBER OF ITERATIONS RUN BY THE COORDINATE DESCENT SOLVER TO REACH THE SPECIFIED  
TOLERANCE FOR THE OPTIMAL ALPHA

SEE ALSO

MULTITASKELASTICNET

ELASTICNETCV

MULTITASKLASSOCV

NOTES

THE ALGORITHM USED TO FIT THE MODEL IS COORDINATE DESCENT  
TO AVOID UNNECESSARY MEMORY DUPLICATION THE X ARGUMENT OF THE FIT METHOD SHOULD BE DIRECTLY PASSED AS A  
FORTRANCONTIGUOUS NUMPY ARRAY

EXAMPLES

```
FROM SKLEARN IMPORT LINEARMODEL
CLF LINEARMODELMULTITASKELASTICNETCVCV3
CLFFIT00 1 1 2 2
0 0 1 1 2 2
```

MULTITASKELASTICNETCVALPHASNONE COPYXTRUE CV3 EPS0001  
FITINTERCEPTTRUE L1RATIO05 MAXITER1000 NALPHAS100  
NJOBSNONE NORMALIZEFALSE RANDOMSTATENONE SELECTIONCYCLIC  
TOL00001 VERBOSE0  
PRINTCLFCOEF  
052875032 046958558  
052875032 046958558  
PRINTCLFINTERCEPT  
000166409 000166409

478 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

METHODS

FITSELF X Y FIT LINEAR MODEL WITH COORDINATE DESCENT

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PATH X Y L1RATIO EPS NALPHAS COMPUTE ELASTIC NET PATH WITH COORDINATE DESCENT

PREDICT SELF X PREDICT USING THE LINEAR MODEL

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELF L1RATIO05 EPS0001 NALPHAS100 ALPHASNONE FITINTERCEPTTRUE NORMALIZEFALSE MAXITER1000 TOL00001 CV'WARN' COPYXTRUE VERBOSE0 NJOBSNONE RANDOMSTATENONE SELECTION'CYCLIC'

FITSELFXY

FIT LINEAR MODEL WITH COORDINATE DESCENT

FIT IS ON GRID OF ALPHAS AND BEST ALPHA ESTIMATED BY CROSSVALIDATION

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA PASS DIRECTLY AS FORTRAN CONTIGUOUS DATA TO AVOID UNNECESSARY MEMORY DUPLICATION IF Y IS MONOOUTPUT X CAN BE SPARSE

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NTARGETS TARGET VALUES

GETPARAMS SELFDEEPTTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

STATICPATHXY L1RATIO05 EPS0001 NALPHAS100 ALPHASNONE PRECOMPUTE'AUTO' XYNONE COPYXTRUE COEFINITNONE VERBOSEFALSE RETURNNNITERFALSE POSITIVEFALSE CHECKINPUTTRUE PARAMS

COMPUTE ELASTIC NET PATH WITH COORDINATE DESCENT

THE ELASTIC NET OPTIMIZATION FUNCTION VARIES FOR MONO AND MULTIOUTPUTS

FOR MONOOUTPUT TASKS IT IS

1 2NSAMPLES Y XW22

ALPHA L1RATIO W1

05ALPHA1 L1RATIO W22

FOR MULTIOUTPUT TASKS IT IS

1 2 NSAMPLES Y XWFRO2

ALPHA L1RATIO W21

05ALPHA1 L1RATIO WFRO2

WHERE

33 MODEL SELECTION AND EVALUATION 479

SCIKITLEARN USER GUIDE RELEASE 0213

$W_{21} = \sum_j \sqrt{w_{ij}^2}$   
IE THE SUM OF NORM OF EACH ROW  
READ MORE IN THE USER GUIDE

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA PASS DIRECTLY AS FORTRAN  
CONTIGUOUS DATA TO AVOID UNNECESSARY MEMORY DUPLICATION IF Y IS MONOOUTPUT THEN X  
CAN BE SPARSE

YNDARRAY SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TARGET VALUES

L1RATIO FLOAT OPTIONAL FLOAT BETWEEN 0 AND 1 PASSED TO ELASTIC NET SCALING BETWEEN L1 AND  
L2 PENALTIES L1RATIO1 CORRESPONDS TO THE LASSO

EPS FLOAT LENGTH OF THE PATH EPS1E3 MEANS THAT ALPHAMIN ALPHAMAX  
1E3

NALPHAS INT OPTIONAL NUMBER OF ALPHAS ALONG THE REGULARIZATION PATH  
ALPHAS NDARRAY OPTIONAL LIST OF ALPHAS WHERE TO COMPUTE THE MODELS IF NONE ALPHAS ARE  
SET AUTOMATICALLY

PRECOMPUTE TRUE FALSE 'AUTO' ARRAYLIKE WHETHER TO USE A PRECOMPUTED GRAM MATRIX  
TO SPEED UP CALCULATIONS IF SET TO AUTO LET US DECIDE THE GRAM MATRIX CAN ALSO BE  
PASSED AS ARGUMENT

XYARRAYLIKE OPTIONAL XY NPDOTXT Y THAT CAN BE PRECOMPUTED IT IS USEFUL ONLY  
WHEN THE GRAM MATRIX IS PRECOMPUTED

COPYX BOOLEAN OPTIONAL DEFAULT TRUE IF TRUE X WILL BE COPIED ELSE IT MAY BE OVER  
WRITTEN

COEFINIT ARRAY SHAPE NFEATURES NONE THE INITIAL VALUES OF THE COEFFICIENTS

VERBOSE BOOL OR INTEGER AMOUNT OF VERBOSITY

RETURNNITER BOOL WHETHER TO RETURN THE NUMBER OF ITERATIONS OR NOT

POSITIVE BOOL DEFAULT FALSE IF SET TO TRUE FORCES COEFFICIENTS TO BE POSITIVE ONLY ALLOWED  
WHEN YNDIM 1

CHECKINPUT BOOL DEFAULT TRUE SKIP INPUT VALIDATION CHECKS INCLUDING THE GRAM MATRIX  
WHEN PROVIDED ASSUMING THERE ARE HANDLED BY THE CALLER WHEN CHECKINPUTFALSE

PARAMS KWARGS KEYWORD ARGUMENTS PASSED TO THE COORDINATE DESCENT SOLVER

RETURNS

ALPHAS ARRAY SHAPE NALPHAS THE ALPHAS ALONG THE PATH WHERE MODELS ARE COMPUTED

COEFS ARRAY SHAPE NFEATURES NALPHAS OR NOUTPUTS NFEATURES NALPHAS COEFFICIENTS  
ALONG THE PATH

DUALGAPS ARRAY SHAPE NALPHAS THE DUAL GAPS AT THE END OF THE OPTIMIZATION FOR EACH  
ALPHA

NITERS ARRAYLIKE SHAPE NALPHAS THE NUMBER OF ITERATIONS TAKEN BY THE COORDINATE  
DESCENT OPTIMIZER TO REACH THE SPECIFIED TOLERANCE FOR EACH ALPHA IS RETURNED WHEN  
RETURNNITER IS SET TO TRUE

SEE ALSO

480 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

MULTITASKELASTICNET

MULTITASKELASTICNETCV

ELASTICNET

ELASTICNETCV

NOTES

FOR AN EXAMPLE SEE EXAMPLESLINEARMODELPLOTLASSOCOORDINATEDDESCENTPATHPY

PREDICTSELF

PREDICT USING THE LINEAR MODEL

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES SAMPLES

RETURNS

CARRAY SHAPE NSAMPLES RETURNS PREDICTED VALUES

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{U}{V}$  WHERE U IS THE RESIDUAL SUM OF SQUARES  $Y_{TRUE} - Y_{PRED}$

2SUM AND V IS THE TOTAL SUM OF SQUARES  $Y_{TRUE} - Y_{TRUEMEAN}$  2SUM THE BEST POSSIBLE SCORE

IS 10 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS

PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 00

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES FOR SOME ESTIMATORS THIS MAY

BE A PRECOMPUTED KERNEL MATRIX INSTEAD SHAPE NSAMPLES NSAMPLESFITTED WHERE

NSAMPLESFITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE VALUES FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT R2 OF SELF-PREDICTX WRT Y

NOTES

THE R2 SCORE USED WHEN CALLING SCORE ON A REGRESSOR WILL USE MULTIOUTPUTUNIFORMAVERAGE

FROM VERSION 023 TO KEEP CONSISTENT WITH METRICSR2SCORE THIS WILL INFLUENCE THE SCORE

METHOD OF ALL THE MULTIOUTPUT REGRESSORS EXCEPT FOR MULTIOUTPUTMULTIOUTPUTREGRESSOR

TO SPECIFY THE DEFAULT VALUE MANUALLY AND AVOID THE WARNING PLEASE EITHER CALL METRICSR2SCORE

DIRECTLY OR MAKE A CUSTOM SCORER WITH METRICSMAKESCORER THE BUILTIN SCORER R2 USES

MULTIOUTPUTUNIFORMAVERAGE

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE

PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT

OF A NESTED OBJECT

33 MODEL SELECTION AND EVALUATION 481

SCIKITLEARN USER GUIDE RELEASE 0213

RETURNS  
SELF

SKLEARNLINEARMODEL MULTITASKLASSOCV  
CLASSSSKLEARNLINEARMODEL MULTITASKLASSOCV EPS0001 NALPHAS100 ALPHASNONE  
FITINTERCEPTTRUE NORMALIZEFALSE  
MAXITER1000 TOL00001 COPYXTRUE  
CV'WARN' VERBOSEFALSE NJOBSNONE  
RANDOMSTATENONE SELECTION'CYCLIC'

MULTITASK LASSO MODEL TRAINED WITH L1L2 MIXEDNORM AS REGULARIZER  
SEE GLOSSARY ENTRY FOR CROSSVALIDATION ESTIMATOR  
THE OPTIMIZATION OBJECTIVE FOR MULTITASKLASSO IS  
1 2 NSAMPLES Y XWFRO2 ALPHA W21  
WHERE  
W21 SUMI SQRTSUMJ WIJ2  
IE THE SUM OF NORM OF EACH ROW  
READ MORE IN THE USER GUIDE

PARAMETERS  
EPS FLOAT OPTIONAL LENGTH OF THE PATH EPS1E3 MEANS THATALPHAMIN ALPHAMAX  
1E3  
NALPHAS INT OPTIONAL NUMBER OF ALPHAS ALONG THE REGULARIZATION PATH  
ALPHAS ARRAYLIKE OPTIONAL LIST OF ALPHAS WHERE TO COMPUTE THE MODELS IF NOT PROVIDED SET  
AUTOMATICALLY  
FITINTERCEPT BOOLEAN WHETHER TO CALCULATE THE INTERCEPT FOR THIS MODEL IF SET TO FALSE NO  
INTERCEPT WILL BE USED IN CALCULATIONS EG DATA IS EXPECTED TO BE ALREADY CENTERED  
NORMALIZE BOOLEAN OPTIONAL DEFAULT FALSE THIS PARAMETER IS IGNORED WHEN  
FITINTERCEPT IS SET TO FALSE IF TRUE THE REGRESSORS X WILL BE NORMALIZED BE  
FORE REGRESSION BY SUBTRACTING THE MEAN AND DIVIDING BY THE L2NORM IF YOU WISH TO  
STANDARDIZE PLEASE USE SKLEARNPREPROCESSINGSTANDARDSCALER BEFORE  
CALLINGFIT ON AN ESTIMATOR WITH NORMALIZEFALSE  
MAXITER INT OPTIONAL THE MAXIMUM NUMBER OF ITERATIONS  
TOLFLOAT OPTIONAL THE TOLERANCE FOR THE OPTIMIZATION IF THE UPDATES ARE SMALLER THAN TOL THE  
OPTIMIZATION CODE CHECKS THE DUAL GAP FOR OPTIMALITY AND CONTINUES UNTIL IT IS SMALLER THAN  
TOL  
COPYX BOOLEAN OPTIONAL DEFAULT TRUE IF TRUE X WILL BE COPIED ELSE IT MAY BE OVERWRITTEN  
CVINT CROSSVALIDATION GENERATOR OR AN ITERABLE OPTIONAL DETERMINES THE CROSSVALIDATION SPLIT  
TING STRATEGY POSSIBLE INPUTS FOR CV ARE

- NONE TO USE THE DEFAULT 3FOLD CROSSVALIDATION
- INTEGER TO SPECIFY THE NUMBER OF FOLDS
- CV SPLITTER

482 CHAPTER 3 USER GUIDE



SCIKITLEARN USER GUIDE RELEASE 0213

• AN ITERABLE YIELDING TRAIN TEST SPLITS AS ARRAYS OF INDICES

FOR INTEGERNONE INPUTS KFOLD IS USED

REFER USER GUIDE FOR THE VARIOUS CROSSVALIDATION STRATEGIES THAT CAN BE USED HERE

CHANGED IN VERSION 020 CVDEFAULT VALUE IF NONE WILL CHANGE FROM 3FOLD TO 5FOLD IN

V022

VERBOSE BOOL OR INTEGER AMOUNT OF VERBOSITY

NJOBS INT OR NONE OPTIONAL DEFAULTNONE NUMBER OF CPUS TO USE DURING THE CROSS VALI

DATION NOTE THAT THIS IS USED ONLY IF MULTIPLE VALUES FOR L1RATIO ARE GIVEN NONE MEANS 1

UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL PROCESSORS SEE

GLOSSARY FOR MORE DETAILS

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULT NONE THE SEED OF THE

PSEUDO RANDOM NUMBER GENERATOR THAT SELECTS A RANDOM FEATURE TO UPDATE IF INT RAN

DOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE

RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS

THE RANDOMSTATE INSTANCE USED BY NPRANDOM USED WHEN SELECTION 'RANDOM'

SELECTION STR DEFAULT 'CYCLIC' IF SET TO 'RANDOM' A RANDOM COEFFICIENT IS UPDATED EVERY ITERATION

RATHER THAN LOOPING OVER FEATURES SEQUENTIALLY BY DEFAULT THIS SETTING TO 'RANDOM' OFTEN

LEADS TO SIGNIFICANTLY FASTER CONVERGENCE ESPECIALLY WHEN TOL IS HIGHER THAN 1E4

ATTRIBUTES

INTERCEPT ARRAY SHAPE NTASKS INDEPENDENT TERM IN DECISION FUNCTION

COEF ARRAY SHAPE NTASKS NFEATURES PARAMETER VECTOR W IN THE COST FUNCTION FORMULA

NOTE THATCOEF STORES THE TRANSPOSE OF WWT

ALPHA FLOAT THE AMOUNT OF PENALIZATION CHOSEN BY CROSS VALIDATION

MSEPATH ARRAY SHAPE NALPHAS NFOLDS MEAN SQUARE ERROR FOR THE TEST SET ON EACH FOLD

VARYING ALPHA

ALPHAS NUMPY ARRAY SHAPE NALPHAS THE GRID OF ALPHAS USED FOR FITTING

NITER INT NUMBER OF ITERATIONS RUN BY THE COORDINATE DESCENT SOLVER TO REACH THE SPECIFIED

TOLERANCE FOR THE OPTIMAL ALPHA

SEE ALSO

MULTITASKELASTICNET

ELASTICNETCV

MULTITASKELASTICNETCV

NOTES

THE ALGORITHM USED TO FIT THE MODEL IS COORDINATE DESCENT

TO AVOID UNNECESSARY MEMORY DUPLICATION THE X ARGUMENT OF THE FIT METHOD SHOULD BE DIRECTLY PASSED AS A

FORTRANCONTIGUOUS NUMPY ARRAY

33 MODEL SELECTION AND EVALUATION 483

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES

```
FROM SKLEARNLINEARMODEL IMPORT MULTITASKLASSOCV
FROM SKLEARNDATASETS IMPORT MAKEREGRESSION
FROM SKLEARNMETRICS IMPORT R2SCORE
X Y MAKEREGRESSIONNTARGETS2 NOISE4 RANDOMSTATE0
REG MULTITASKLASSOCVCV5 RANDOMSTATE0FITX Y
R2SCOREY REGPREDICTX
09994
REGALPHA
05713
REGPREDICTX1
ARRAY1537971 949015
```

METHODS

FITSELF X Y FIT LINEAR MODEL WITH COORDINATE DESCENT

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PATH X Y EPS NALPHAS ALPHAS COMPUTE LASSO PATH WITH COORDINATE DESCENT

PREDICT SELF X PREDICT USING THE LINEAR MODEL

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PRE  
DICTION

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELFEPS0001 NALPHAS100 ALPHASNONE FITINTERCEPTTRUE NORMALIZEFALSE

MAXITER1000 TOL00001 COPYXTRUE CV'WARN' VERBOSEFALSE NJOBSNONE RAN  
DOMSTATENONE SELECTION'CYCLIC'

FITSELFXY

FIT LINEAR MODEL WITH COORDINATE DESCENT

FIT IS ON GRID OF ALPHAS AND BEST ALPHA ESTIMATED BY CROSSVALIDATION

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA PASS DIRECTLY AS FORTRAN  
CONTIGUOUS DATA TO AVOID UNNECESSARY MEMORY DUPLICATION IF Y IS MONOOUTPUT X CAN  
BE SPARSE

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NTARGETS TARGET VALUES

GETPARAMS SELFDEEPTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

STATICPATHXYEPS0001 NALPHAS100 ALPHASNONE PRECOMPUTE'AUTO' XYNONE

COPYXTRUE COEFINITNONE VERBOSEFALSE RETURNNNITERFALSE POSITIVEFALSE

PARAMS

COMPUTE LASSO PATH WITH COORDINATE DESCENT

484 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213  
 THE LASSO OPTIMIZATION FUNCTION VARIES FOR MONO AND MULTIOUTPUTS  
 FOR MONOOUTPUT TASKS IT IS  
 1 2 NSAMPLES Y XW22 ALPHA W1  
 FOR MULTIOUTPUT TASKS IT IS  
 1 2 NSAMPLES Y XW2FRO ALPHA W21  
 WHERE  
 W21 SUMI SQRTSUMJ WIJ2  
 IE THE SUM OF NORM OF EACH ROW  
 READ MORE IN THE USER GUIDE  
 PARAMETERS  
 XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES TRAINING DATA PASS DIRECTLY  
 AS FORTRANCONTIGUOUS DATA TO AVOID UNNECESSARY MEMORY DUPLICATION IF YIS MONOOUTPUT  
 THENXCAN BE SPARSE  
 YNDARRAY SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TARGET VALUES  
 EPS FLOAT OPTIONAL LENGTH OF THE PATH EPS1E3 MEANS THAT ALPHAMIN  
 ALPHAMAX 1E3  
 NALPHAS INT OPTIONAL NUMBER OF ALPHAS ALONG THE REGULARIZATION PATH  
 ALPHAS NDARRAY OPTIONAL LIST OF ALPHAS WHERE TO COMPUTE THE MODELS IF NONE ALPHAS ARE  
 SET AUTOMATICALLY  
 PRECOMPUTE TRUE FALSE 'AUTO' ARRAYLIKE WHETHER TO USE A PRECOMPUTED GRAM MATRIX  
 TO SPEED UP CALCULATIONS IF SET TO AUTO LET US DECIDE THE GRAM MATRIX CAN ALSO BE  
 PASSED AS ARGUMENT  
 XYARRAYLIKE OPTIONAL XY NPDOTXT Y THAT CAN BE PRECOMPUTED IT IS USEFUL ONLY  
 WHEN THE GRAM MATRIX IS PRECOMPUTED  
 COPYX BOOLEAN OPTIONAL DEFAULT TRUE IF TRUE X WILL BE COPIED ELSE IT MAY BE OVER  
 WRITTEN  
 COEFINIT ARRAY SHAPE NFEATURES NONE THE INITIAL VALUES OF THE COEFFICIENTS  
 VERBOSE BOOL OR INTEGER AMOUNT OF VERBOSITY  
 RETURNNITER BOOL WHETHER TO RETURN THE NUMBER OF ITERATIONS OR NOT  
 POSITIVE BOOL DEFAULT FALSE IF SET TO TRUE FORCES COEFFICIENTS TO BE POSITIVE ONLY ALLOWED  
 WHENYNDIM 1  
 PARAMS KWARGS KEYWORD ARGUMENTS PASSED TO THE COORDINATE DESCENT SOLVER  
 RETURNS  
 ALPHAS ARRAY SHAPE NALPHAS THE ALPHAS ALONG THE PATH WHERE MODELS ARE COMPUTED  
 COEFS ARRAY SHAPE NFEATURES NALPHAS OR NOUTPUTS NFEATURES NALPHAS COEFFICIENTS  
 ALONG THE PATH  
 DUALGAPS ARRAY SHAPE NALPHAS THE DUAL GAPS AT THE END OF THE OPTIMIZATION FOR EACH  
 ALPHA  
 33 MODEL SELECTION AND EVALUATION 485

SCIKITLEARN USER GUIDE RELEASE 0213

NITERS ARRAYLIKE SHAPE NALPHAS THE NUMBER OF ITERATIONS TAKEN BY THE COORDINATE DESCENT OPTIMIZER TO REACH THE SPECIFIED TOLERANCE FOR EACH ALPHA

SEE ALSO

LARSPATH

LASSO

LASSOLARS

LASSOCV

LASSOLARSCV

SKLEARNDECOMPOSITIONSPARSEENCODER

NOTES

FOR AN EXAMPLE SEE EXAMPLESLINEARMODELPLOTLASSOCOORDINATEDDESCENTPATHPY

TO AVOID UNNECESSARY MEMORY DUPLICATION THE X ARGUMENT OF THE FIT METHOD SHOULD BE DIRECTLY PASSED AS A FORTRANCONTIGUOUS NUMPY ARRAY

NOTE THAT IN CERTAIN CASES THE LARS SOLVER MAY BE SIGNIFICANTLY FASTER TO IMPLEMENT THIS FUNCTIONALITY IN PARTICULAR LINEAR INTERPOLATION CAN BE USED TO RETRIEVE MODEL COEFFICIENTS BETWEEN THE VALUES OUTPUT BY LARSPATH

EXAMPLES

COMPARING LASSOPATH AND LARSPATH WITH INTERPOLATION

X NPARRAY1 2 31 23 54 43T

Y NPARRAY1 2 31

USE LASSOPATH TO COMPUTE A COEFFICIENT PATH

COEFPATH LASSOPATHX Y ALPHAS5 1 5

PRINTCOEFPATH

0 0 046874778

02159048 04425765 023689075

NOW USE LARSPATH AND 1D LINEAR INTERPOLATION TO COMPUTE THE SAME PATH

FROM SKLEARNLINEARMODEL IMPORT LARSPATH

ALPHAS ACTIVE COEFPATHLARS LARSPATHX Y METHODLASSO

FROM SCIPY IMPORT INTERPOLATE

COEFPATHCONTINUOUS INTERPOLATEINTERP1DALPHAS1

COEFPATHLARS 1

PRINTCOEFPATHCONTINUOUS5 1 5

0 0 046915237

02159048 04425765 023668876

PREDICTSELF

PREDICT USING THE LINEAR MODEL

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES SAMPLES

RETURNS

486 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

CARRAY SHAPE NSAMPLES RETURNS PREDICTED VALUES

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{U}{V}$  WHERE U IS THE RESIDUAL SUM OF SQUARES  $Y_{TRUE} - Y_{PRED}$   $2SUM$  AND V IS THE TOTAL SUM OF SQUARES  $Y_{TRUE} - Y_{TRUEMEAN}$   $2SUM$  THE BEST POSSIBLE SCORE IS 10 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 00

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES FOR SOME ESTIMATORS THIS MAY BE A PRECOMPUTED KERNEL MATRIX INSTEAD SHAPE NSAMPLES NSAMPLESFITTED WHERE NSAMPLESFITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE VALUES FOR X SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT R2 OF SELF PREDICTX WRT Y

NOTES

THE R2 SCORE USED WHEN CALLING SCORE ON A REGRESSOR WILL USE MULTIOUTPUTUNIFORMAVERAGE FROM VERSION 023 TO KEEP CONSISTENT WITH METRICSR2SCORE THIS WILL INFLUENCE THE SCORE METHOD OF ALL THE MULTIOUTPUT REGRESSORS EXCEPT FOR MULTIOUTPUTMULTIOUTPUTREGRESSOR TO SPECIFY THE DEFAULT VALUE MANUALLY AND AVOID THE WARNING PLEASE EITHER CALL METRICSR2SCORE DIRECTLY OR MAKE A CUSTOM SCORER WITH METRICSMAKESCORER THE BUILTIN SCORER R2 USES MULTIOUTPUTUNIFORMAVERAGE

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

SKLEARNLINEARMODEL ORTHOGONALMATCHINGPURSUITCV

CLASSSSKLEARNLINEARMODEL ORTHOGONALMATCHINGPURSUITCV COPYTRUE FITINTERCEPTTRUE NORMALIZETRUE

MAXITERNONE CV'WARN'

NJOBSNONE VERBOSEFALSE

CROSSVALIDATED ORTHOGONAL MATCHING PURSUIT MODEL OMP

SEE GLOSSARY ENTRY FOR CROSSVALIDATION ESTIMATOR

READ MORE IN THE USER GUIDE

PARAMETERS

COPY BOOL OPTIONAL WHETHER THE DESIGN MATRIX X MUST BE COPIED BY THE ALGORITHM A FALSE VALUE IS ONLY HELPFUL IF X IS ALREADY FORTRANORDERED OTHERWISE A COPY IS MADE ANYWAY

33 MODEL SELECTION AND EVALUATION 487

SCIKITLEARN USER GUIDE RELEASE 0213

FITINTERCEPT BOOLEAN OPTIONAL WHETHER TO CALCULATE THE INTERCEPT FOR THIS MODEL IF SET TO FALSE NO INTERCEPT WILL BE USED IN CALCULATIONS EG DATA IS EXPECTED TO BE ALREADY CENTERED

NORMALIZE BOOLEAN OPTIONAL DEFAULT TRUE THIS PARAMETER IS IGNORED WHEN FITINTERCEPT IS SET TO FALSE IF TRUE THE REGRESSORS X WILL BE NORMALIZED BEFORE REGRESSION BY SUBTRACTING THE MEAN AND DIVIDING BY THE L2NORM IF YOU WISH TO STANDARDIZE PLEASE USE SKLEARNPREPROCESSINGSTANDARDSCALER BEFORE CALLING FIT ON AN ESTIMATOR

WITHNORMALIZEFALSE

MAXITER INTEGER OPTIONAL MAXIMUM NUMBERS OF ITERATIONS TO PERFORM THEREFORE MAXIMUM FEATURES TO INCLUDE 10 OF NFEATURES BUT AT LEAST 5 IF AVAILABLE

CVINT CROSSVALIDATION GENERATOR OR AN ITERABLE OPTIONAL DETERMINES THE CROSSVALIDATION SPLITTING STRATEGY POSSIBLE INPUTS FOR CV ARE

- NONE TO USE THE DEFAULT 3FOLD CROSSVALIDATION
- INTEGER TO SPECIFY THE NUMBER OF FOLDS
- CV SPLITTER
- AN ITERABLE YIELDING TRAIN TEST SPLITS AS ARRAYS OF INDICES

FOR INTEGERNONE INPUTS KFOLD IS USED

REFER USER GUIDE FOR THE VARIOUS CROSSVALIDATION STRATEGIES THAT CAN BE USED HERE

CHANGED IN VERSION 020 CVDEFAULT VALUE IF NONE WILL CHANGE FROM 3FOLD TO 5FOLD IN V022

NJOBS INT OR NONE OPTIONAL DEFAULTNONE NUMBER OF CPUS TO USE DURING THE CROSS VALIDATIONNONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

VERBOSE BOOLEAN OR INTEGER OPTIONAL SETS THE VERBOSITY AMOUNT

ATTRIBUTES

INTERCEPT FLOAT OR ARRAY SHAPE NTARGETS INDEPENDENT TERM IN DECISION FUNCTION

COEF ARRAY SHAPE NFEATURES OR NTARGETS NFEATURES PARAMETER VECTOR W IN THE PROBLEM FORMULATION

NNONZEROCOEF INT ESTIMATED NUMBER OF NONZERO COEFFICIENTS GIVING THE BEST MEAN SQUARED ERROR OVER THE CROSSVALIDATION FOLDS

NITER INT OR ARRAYLIKE NUMBER OF ACTIVE FEATURES ACROSS EVERY TARGET FOR THE MODEL REFIT WITH THE BEST HYPERPARAMETERS GOT BY CROSSVALIDATING ACROSS ALL FOLDS

SEE ALSO

ORTHOGONALMP

ORTHOGONALMPGRAM

LARSPATH

LARS

LASSOLARS

ORTHOGONALMATCHINGPURSUIT

LARSCV

LASSOLARSCV

488 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

DECOMPOSITIONSPARSEENCODE

EXAMPLES

FROM SKLEARNLINEARMODEL IMPORT ORTHOGONALMATCHINGPURSUITCV

FROM SKLEARNDATASETS IMPORT MAKEREGRESSION

X Y MAKEREGRESSIONNFEATURES100 NINFORMATIVE10

NOISE4 RANDOMSTATE0

REG ORTHOGONALMATCHINGPURSUITCVCV5FITX Y

REGSCOREX Y

09991

REGNNONZEROCOEFS

10

REGPREDICTX1

ARRAY783854

METHODS

FITSELF X Y FIT THE MODEL USING X Y AS TRAINING DATA

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PREDICT SELF X PREDICT USING THE LINEAR MODEL

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PRE

DICTION

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELF COPYTRUE FITINTERCEPTTRUE NORMALIZETRUE MAXITERNONE CV'WARN'

NJOBSNONE VERBOSEFALSE

FITSELFXY

FIT THE MODEL USING X Y AS TRAINING DATA

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA

YARRAYLIKE SHAPE NSAMPLES TARGET VALUES WILL BE CAST TO X'S DTYPE IF NECESSARY

RETURNS

SELF OBJECT RETURNS AN INSTANCE OF SELF

GETPARAMS SELFDEEPTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELF

PREDICT USING THE LINEAR MODEL

PARAMETERS

33 MODEL SELECTION AND EVALUATION 489

SCIKITLEARN USER GUIDE RELEASE 0213

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES SAMPLES

RETURNS

CARRAY SHAPE NSAMPLES RETURNS PREDICTED VALUES

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{U}{V}$  WHERE U IS THE RESIDUAL SUM OF SQUARES  $Y_{TRUE} - Y_{PRED}$   $2 \sum$  AND V IS THE TOTAL SUM OF SQUARES  $Y_{TRUE} - Y_{TRUEMEAN}$   $2 \sum$  THE BEST POSSIBLE SCORE IS 10 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 00

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES FOR SOME ESTIMATORS THIS MAY BE A PRECOMPUTED KERNEL MATRIX INSTEAD SHAPE NSAMPLES NSAMPLESFITTED WHERE NSAMPLESFITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE VALUES FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT R2 OF SELF PREDICT X WRT Y

NOTES

THE R2 SCORE USED WHEN CALLING SCORE ON A REGRESSOR WILL USE MULTIOUTPUTUNIFORMAVERAGE FROM VERSION 023 TO KEEP CONSISTENT WITH METRICSR2SCORE THIS WILL INFLUENCE THE SCORE METHOD OF ALL THE MULTIOUTPUT REGRESSORS EXCEPT FOR MULTIOUTPUTMULTIOUTPUTREGRESSOR TO SPECIFY THE DEFAULT VALUE MANUALLY AND AVOID THE WARNING PLEASE EITHER CALL METRICSR2SCORE DIRECTLY OR MAKE A CUSTOM SCORER WITH METRICSMAKESCORER THE BUILTIN SCORER R2 USES MULTIOUTPUTUNIFORMAVERAGE

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

EXAMPLES USING SKLEARNLINEARMODEL ORTHOGONAL MATCHING PURSUIT CV

- ORTHOGONAL MATCHING PURSUIT

SKLEARNLINEARMODEL RIDGECV

CLASS SKLEARNLINEARMODEL RIDGECV ALPHAS 01 10 100 FIT INTERCEPT TRUE NORMAL

IZE FALSE SCORING NONE CV NONE GCV MODE NONE

STORE CV VALUES FALSE

RIDGE REGRESSION WITH BUILTIN CROSS VALIDATION

490 CHAPTER 3 USER GUIDE



SCIKITLEARN USER GUIDE RELEASE 0213  
SEE GLOSSARY ENTRY FOR CROSSVALIDATION ESTIMATOR  
BY DEFAULT IT PERFORMS GENERALIZED CROSSVALIDATION WHICH IS A FORM OF EFFICIENT LEAVEONEOUT CROSS  
VALIDATION  
READ MORE IN THE USER GUIDE  
PARAMETERS  
ALPHAS NUMPY ARRAY OF SHAPE NALPHAS ARRAY OF ALPHA VALUES TO TRY REGULARIZATION STRENGTH  
MUST BE A POSITIVE FLOAT REGULARIZATION IMPROVES THE CONDITIONING OF THE PROBLEM AND RE  
DUCES THE VARIANCE OF THE ESTIMATES LARGER VALUES SPECIFY STRONGER REGULARIZATION ALPHA  
CORRESPONDS TO C1 IN OTHER LINEAR MODELS SUCH AS LOGISTICREGRESSION OR LINEARSVC IF  
USING GENERALIZED CROSSVALIDATION ALPHAS MUST BE POSITIVE  
FITINTERCEPT BOOLEAN WHETHER TO CALCULATE THE INTERCEPT FOR THIS MODEL IF SET TO FALSE NO  
INTERCEPT WILL BE USED IN CALCULATIONS EG DATA IS EXPECTED TO BE ALREADY CENTERED  
NORMALIZE BOOLEAN OPTIONAL DEFAULT FALSE THIS PARAMETER IS IGNORED WHEN  
FITINTERCEPT IS SET TO FALSE IF TRUE THE REGRESSORS X WILL BE NORMALIZED BE  
FORE REGRESSION BY SUBTRACTING THE MEAN AND DIVIDING BY THE L2NORM IF YOU WISH TO  
STANDARDIZE PLEASE USE SKLEARNPREPROCESSINGSTANDARDSCALER BEFORE  
CALLINGFIT ON AN ESTIMATOR WITH NORMALIZEFALSE  
SCORING STRING CALLABLE OR NONE OPTIONAL DEFAULT NONE A STRING SEE MODEL EVALUATION DOC  
UMENTATION OR A SCORER CALLABLE OBJECT FUNCTION WITH SIGNATURE SCORERESTIMATOR  
X Y IF NONE THE NEGATIVE MEAN SQUARED ERROR IF CV IS 'AUTO' OR NONE IE WHEN USING  
GENERALIZED CROSSVALIDATION AND R2 SCORE OTHERWISE  
CVINT CROSSVALIDATION GENERATOR OR AN ITERABLE OPTIONAL DETERMINES THE CROSSVALIDATION SPLIT  
TING STRATEGY POSSIBLE INPUTS FOR CV ARE  
• NONE TO USE THE EFFICIENT LEAVEONEOUT CROSSVALIDATION ALSO KNOWN AS GENERALIZED  
CROSSVALIDATION  
• INTEGER TO SPECIFY THE NUMBER OF FOLDS  
• CV SPLITTER  
• AN ITERABLE YIELDING TRAIN TEST SPLITS AS ARRAYS OF INDICES  
FOR INTEGERNONE INPUTS IF YIS BINARY OR MULTICLASS SKLEARNMODELSELECTION  
STRATIFIEDKFOLD IS USED ELSE SKLEARNMODELSELECTIONKFOLD IS USED  
REFER USER GUIDE FOR THE VARIOUS CROSSVALIDATION STRATEGIES THAT CAN BE USED HERE  
GCVMODE NONE 'AUTO' 'SVD' 'EIGEN' OPTIONAL FLAG INDICATING WHICH STRATEGY TO USE WHEN  
PERFORMING GENERALIZED CROSSVALIDATION OPTIONS ARE  
AUTO USE SVD IFNSAMPLES NFEATURES OTHERWISE USE EIGEN  
SVD FORCE USE OF SINGULAR VALUE DECOMPOSITION OF X WHEN X IS  
DENSE EIGENVALUE DECOMPOSITION OF XTX WHEN X ISSPARSE  
EIGEN FORCE COMPUTATION VIA EIGENDECOMPOSITION OF XXT  
THE 'AUTO' MODE IS THE DEFAULT AND IS INTENDED TO PICK THE CHEAPER OPTION OF THE TWO DEPEND  
ING ON THE SHAPE OF THE TRAINING DATA  
STORECVVALUES BOOLEAN DEFAULTFALSE FLAG INDICATING IF THE CROSSVALIDATION VALUES CORRE  
SPONDING TO EACH ALPHA SHOULD BE STORED IN THE CVVALUES ATTRIBUTE SEE BELOW THIS  
FLAG IS ONLY COMPATIBLE WITH CVNONE IE USING GENERALIZED CROSSVALIDATION  
ATTRIBUTES  
33 MODEL SELECTION AND EVALUATION 491

SCIKITLEARN USER GUIDE RELEASE 0213

CVVALUES ARRAY SHAPE NSAMPLES NALPHAS OR SHAPE NSAMPLES NTARGETS NALPHAS  
OPTIONAL CROSSVALIDATION VALUES FOR EACH ALPHA IF STORECVVALUESTRUE AND  
CVNONE AFTERFIT HAS BEEN CALLED THIS ATTRIBUTE WILL CONTAIN THE MEAN SQUARED  
ERRORS BY DEFAULT OR THE VALUES OF THE LOSSSCOREFUNC FUNCTION IF PROVIDED IN THE  
CONSTRUCTOR

COEF ARRAY SHAPE NFEATURES OR NTARGETS NFEATURES WEIGHT VECTORS  
INTERCEPT FLOAT ARRAY SHAPE NTARGETS INDEPENDENT TERM IN DECISION FUNCTION SET TO 00  
IFFITINTERCEPT FALSE

ALPHA FLOAT ESTIMATED REGULARIZATION PARAMETER

SEE ALSO

RIDGE RIDGE REGRESSION  
RIDGECLASSIFIER RIDGE CLASSIFIER  
RIDGECLASSIFIERCV RIDGE CLASSIFIER WITH BUILTIN CROSS VALIDATION

EXAMPLES

FROM SKLEARNDATASETS IMPORT LOADDIABETES  
FROM SKLEARNLINEARMODEL IMPORT RIDGECV  
X Y LOADDIABETESRETURNXY TRUE  
CLF RIDGECVALPHAS1E3 1E2 1E1 1FITX Y  
CLFSCOREX Y  
05166

METHODS

FITSELF X Y SAMPLEWEIGHT FIT RIDGE REGRESSION MODEL  
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR  
PREDICT SELF X PREDICT USING THE LINEAR MODEL  
SCORE SELF X Y SAMPLEWEIGHT RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PRE  
DICTION

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR  
INIT SELFALPHAS01 10100 FITINTERCEPTTRUE NORMALIZEFALSE SCORINGNONE CVNONE  
GCVMODENONE STORECVVALUESFALSE  
FITSELFXYSAMPLEWEIGHTNONE  
FIT RIDGE REGRESSION MODEL

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA IF USING GCV WILL BE CAST TO  
FLOAT64 IF NECESSARY  
YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NTARGETS TARGET VALUES WILL BE CAST TO  
X'S DTYPE IF NECESSARY  
SAMPLEWEIGHT FLOAT OR ARRAYLIKE OF SHAPE NSAMPLES SAMPLE WEIGHT

RETURNS

492 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

SELF OBJECT  
NOTES

WHEN SAMPLEWEIGHT IS PROVIDED THE SELECTED HYPERPARAMETER MAY DEPEND ON WHETHER WE USE GENERALIZED CROSSVALIDATION CVNONE OR CV'AUTO' OR ANOTHER FORM OF CROSSVALIDATION BECAUSE ONLY GENERALIZED CROSSVALIDATION TAKES THE SAMPLE WEIGHTS INTO ACCOUNT WHEN COMPUTING THE VALIDATION SCORE

GETPARAMS SELFDEEPTURE  
GET PARAMETERS FOR THIS ESTIMATOR  
PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS  
PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELF  
PREDICT USING THE LINEAR MODEL  
PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES SAMPLES  
RETURNS

CARRAY SHAPE NSAMPLES RETURNS PREDICTED VALUES

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{U}{V}$  WHERE U IS THE RESIDUAL SUM OF SQUARES YTRUE YPRED

2SUM AND V IS THE TOTAL SUM OF SQUARES YTRUE YTRUEMEAN 2SUM THE BEST POSSIBLE SCORE

IS 10 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 00

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES FOR SOME ESTIMATORS THIS MAY

BE A PRECOMPUTED KERNEL MATRIX INSTEAD SHAPE NSAMPLES NSAMPLESFITTED WHERE

NSAMPLESFITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE VALUES FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT R2 OF SELF PREDICTX WRT Y

NOTES

THE R2 SCORE USED WHEN CALLING SCORE ON A REGRESSOR WILL USE MULTIOUTPUTUNIFORMAVERAGE

FROM VERSION 023 TO KEEP CONSISTENT WITH METRICSR2SCORE THIS WILL INFLUENCE THE SCORE

METHOD OF ALL THE MULTIOUTPUT REGRESSORS EXCEPT FOR MULTIOUTPUTMULTIOUTPUTREGRESSOR

TO SPECIFY THE DEFAULT VALUE MANUALLY AND AVOID THE WARNING PLEASE EITHER CALL METRICSR2SCORE

DIRECTLY OR MAKE A CUSTOM SCORER WITH METRICSMAKESCORER THE BUILTIN SCORER R2 USES

MULTIOUTPUTUNIFORMAVERAGE

33 MODEL SELECTION AND EVALUATION 493

SCIKITLEARN USER GUIDE RELEASE 0213

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

EXAMPLES USING SKLEARNLINEARMODELRIDGE CV

- FACE COMPLETION WITH A MULTIOUTPUT ESTIMATORS
- EFFECT OF TRANSFORMING THE TARGETS IN REGRESSION MODEL

SKLEARNLINEARMODEL RIDGECLASSIFIER CV

CLASSSSKLEARNLINEARMODEL RIDGECLASSIFIER CV ALPHAS01 10100 FITINTERCEPTTRUE

NORMALIZEFALSE SCORINGNONE CVNONE

CLASSWEIGHTNONE STORECVVALUESFALSE

RIDGE CLASSIFIER WITH BUILTIN CROSSVALIDATION

SEE GLOSSARY ENTRY FOR CROSSVALIDATION ESTIMATOR

BY DEFAULT IT PERFORMS GENERALIZED CROSSVALIDATION WHICH IS A FORM OF EFFICIENT LEAVEONEOUT CROSS VALIDATION CURRENTLY ONLY THE NFEATURES NSAMPLES CASE IS HANDLED EFFICIENTLY

READ MORE IN THE USER GUIDE

PARAMETERS

ALPHAS NUMPY ARRAY OF SHAPE NALPHAS ARRAY OF ALPHA VALUES TO TRY REGULARIZATION STRENGTH

MUST BE A POSITIVE FLOAT REGULARIZATION IMPROVES THE CONDITIONING OF THE PROBLEM AND REDUCES THE VARIANCE OF THE ESTIMATES LARGER VALUES SPECIFY STRONGER REGULARIZATION ALPHA CORRESPONDS TO C1 IN OTHER LINEAR MODELS SUCH AS LOGISTICREGRESSION OR LINEARSVC

FITINTERCEPT BOOLEAN WHETHER TO CALCULATE THE INTERCEPT FOR THIS MODEL IF SET TO FALSE NO INTERCEPT WILL BE USED IN CALCULATIONS EG DATA IS EXPECTED TO BE ALREADY CENTERED

NORMALIZE BOOLEAN OPTIONAL DEFAULT FALSE THIS PARAMETER IS IGNORED WHEN

FITINTERCEPT IS SET TO FALSE IF TRUE THE REGRESSORS X WILL BE NORMALIZED BEFORE REGRESSION BY SUBTRACTING THE MEAN AND DIVIDING BY THE L2NORM IF YOU WISH TO STANDARDIZE PLEASE USE SKLEARNPREPROCESSINGSTANDARDSCALER BEFORE

CALLINGFIT ON AN ESTIMATOR WITH NORMALIZEFALSE

SCORING STRING CALLABLE OR NONE OPTIONAL DEFAULT NONE A STRING SEE MODEL EVALUATION DOCUMENTATION OR A SCORER CALLABLE OBJECT FUNCTION WITH SIGNATURE SCORERESTIMATOR

X Y

CVINT CROSSVALIDATION GENERATOR OR AN ITERABLE OPTIONAL DETERMINES THE CROSSVALIDATION SPLITTING STRATEGY POSSIBLE INPUTS FOR CV ARE

- NONE TO USE THE EFFICIENT LEAVEONEOUT CROSSVALIDATION
- INTEGER TO SPECIFY THE NUMBER OF FOLDS
- CV SPLITTER
- AN ITERABLE YIELDING TRAIN TEST SPLITS AS ARRAYS OF INDICES

494 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

REFER USER GUIDE FOR THE VARIOUS CROSSVALIDATION STRATEGIES THAT CAN BE USED HERE

CLASSWEIGHT DICT OR 'BALANCED' OPTIONAL WEIGHTS ASSOCIATED WITH CLASSES IN THE FORM

CLASSLABEL WEIGHT IF NOT GIVEN ALL CLASSES ARE SUPPOSED TO HAVE WEIGHT ONE

THE "BALANCED" MODE USES THE VALUES OF Y TO AUTOMATICALLY ADJUST WEIGHTS INVERSELY PRO

PORTIONAL TO CLASS FREQUENCIES IN THE INPUT DATA AS NSAMPLES NCLASSES NP

BINCOUNTY

STORECVVALUES BOOLEAN DEFAULTFALSE FLAG INDICATING IF THE CROSSVALIDATION VALUES CORRE

SPONDING TO EACH ALPHA SHOULD BE STORED IN THE CVVALUES ATTRIBUTE SEE BELOW THIS

FLAG IS ONLY COMPATIBLE WITH CVNONE IE USING GENERALIZED CROSSVALIDATION

ATTRIBUTES

CVVALUES ARRAY SHAPE NSAMPLES NTARGETS NALPHAS OPTIONAL CROSSVALIDATION VALUES

FOR EACH ALPHA IF STORECVVALUETRUE ANDCVNONE AFTERFIT HAS BEEN

CALLED THIS ATTRIBUTE WILL CONTAIN THE MEAN SQUARED ERRORS BY DEFAULT OR THE VALUES OF THE

LOSSSCOREFUNC FUNCTION IF PROVIDED IN THE CONSTRUCTOR

COEF ARRAY SHAPE 1 NFEATURES OR NTARGETS NFEATURES COEFFICIENT OF THE FEATURES IN THE

DECISION FUNCTION

COEF IS OF SHAPE 1 NFEATURES WHEN THE GIVEN PROBLEM IS BINARY

INTERCEPT FLOAT ARRAY SHAPE NTARGETS INDEPENDENT TERM IN DECISION FUNCTION SET TO 00

IFFITINTERCEPT FALSE

ALPHA FLOAT ESTIMATED REGULARIZATION PARAMETER

SEE ALSO

RIDGE RIDGE REGRESSION

RIDGECLASSIFIER RIDGE CLASSIFIER

RIDGECV RIDGE REGRESSION WITH BUILTIN CROSS VALIDATION

NOTES

FOR MULTICLASS CLASSIFICATION NCLASS CLASSIFIERS ARE TRAINED IN A ONEVERSUSALL APPROACH CONCRETELY THIS IS

IMPLEMENTED BY TAKING ADVANTAGE OF THE MULTIVARIATE RESPONSE SUPPORT IN RIDGE

EXAMPLES

```
FROM SKLEARNDATASETS IMPORT LOADBREASTCANCER
FROM SKLEARNLINEARMODEL IMPORT RIDGECLASSIFIERCV
X Y LOADBREASTCANCERRETURNXY TRUE
CLF RIDGECLASSIFIERCVALPHAS1E3 1E2 1E1 1FITX Y
CLFSCOREX Y
```

09630

METHODS

33 MODEL SELECTION AND EVALUATION 495

SCIKITLEARN USER GUIDE RELEASE 0213

DECISIONFUNCTION SELF X PREDICT CONFIDENCE SCORES FOR SAMPLES

FITSELF X Y SAMPLEWEIGHT FIT THE RIDGE CLASSIFIER

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PREDICT SELF X PREDICT CLASS LABELS FOR SAMPLES IN X

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELFALPHAS01 10100 FITINTERCEPTTRUE NORMALIZEFALSE SCORINGNONE CVNONE

CLASSWEIGHTNONE STORECVVALUESFALSE

DECISIONFUNCTION SELF X

PREDICT CONFIDENCE SCORES FOR SAMPLES

THE CONFIDENCE SCORE FOR A SAMPLE IS THE SIGNED DISTANCE OF THAT SAMPLE TO THE HYPERPLANE

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES SAMPLES

RETURNS

ARRAY SHAPENSAMPLES IF NCLASSES > 2 ELSE NSAMPLES NCLASSES CONFIDENCE

SCORES PER SAMPLE CLASS COMBINATION IN THE BINARY CASE CONFIDENCE SCORE FOR

SELFCLASSES1 WHERE 0 MEANS THIS CLASS WOULD BE PREDICTED

FITSELFXYSAMPLEWEIGHTNONE

FIT THE RIDGE CLASSIFIER

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING VECTORS WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES WHEN USING GCV WILL BE CAST TO FLOAT64 IF NECESSARY

YARRAYLIKE SHAPE NSAMPLES TARGET VALUES WILL BE CAST TO X'S DTYPE IF NECESSARY

SAMPLEWEIGHT FLOAT OR NUMPY ARRAY OF SHAPE NSAMPLES SAMPLE WEIGHT

RETURNS

SELF OBJECT

GETPARAMS SELFDEEPTREE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELF X

PREDICT CLASS LABELS FOR SAMPLES IN X

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES SAMPLES

RETURNS

496 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

CARRAY SHAPE NSAMPLES PREDICTED CLASS LABEL PER SAMPLE

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS

IN MULTILABEL CLASSIFICATION THIS IS THE SUBSET ACCURACY WHICH IS A HARSH METRIC SINCE YOU REQUIRE FOR EACH SAMPLE THAT EACH LABEL SET BE CORRECTLY PREDICTED

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE LABELS FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT MEAN ACCURACY OF SELF PREDICTX WRT Y

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

INFORMATION CRITERION

SOME MODELS CAN OFFER AN INFORMATION THEORETIC CLOSED FORM FORMULA OF THE OPTIMAL ESTIMATE OF THE REGULARIZATION PARAMETER BY COMPUTING A SINGLE REGULARIZATION PATH INSTEAD OF SEVERAL WHEN USING CROSS VALIDATION

HERE IS THE LIST OF MODELS BENEFITING FROM THE AKAIKE INFORMATION CRITERION AIC OR THE BAYESIAN INFORMATION CRITERION BIC FOR AUTOMATED MODEL SELECTION

LINEAR MODEL LASSO LARSIC CRITERION LASSO MODEL FIT WITH LARS USING BIC OR AIC FOR MODEL SELECTION

SKLEARN LINEAR MODEL LASSO LARSIC

CLASS SKLEARN LINEAR MODEL LASSO LARSIC CRITERION 'AIC' FIT INTERCEPT TRUE VERBOSE FALSE

NORMALIZE TRUE PRECOMPUTE 'AUTO' MAX ITER 500

EPS 2.220446049250313E-16 COPY X TRUE POSITIVE FALSE

LASSO MODEL FIT WITH LARS USING BIC OR AIC FOR MODEL SELECTION

THE OPTIMIZATION OBJECTIVE FOR LASSO IS

$\frac{1}{2} \|y - Xw\|_2^2 + \alpha \|w\|_1$

AIC IS THE AKAIKE INFORMATION CRITERION AND BIC IS THE BAYES INFORMATION CRITERION SUCH CRITERIA ARE USEFUL TO SELECT THE VALUE OF THE REGULARIZATION PARAMETER BY MAKING A TRADEOFF BETWEEN THE GOODNESS OF FIT AND THE COMPLEXITY OF THE MODEL A GOOD MODEL SHOULD EXPLAIN WELL THE DATA WHILE BEING SIMPLE

READ MORE IN THE USER GUIDE

33 MODEL SELECTION AND EVALUATION 497

SCIKITLEARN USER GUIDE RELEASE 0213

PARAMETERS

CRITERION 'BIC' 'AIC' THE TYPE OF CRITERION TO USE  
FITINTERCEPT BOOLEAN WHETHER TO CALCULATE THE INTERCEPT FOR THIS MODEL IF SET TO FALSE NO INTERCEPT WILL BE USED IN CALCULATIONS EG DATA IS EXPECTED TO BE ALREADY CENTERED  
VERBOSE BOOLEAN OR INTEGER OPTIONAL SETS THE VERBOSITY AMOUNT  
NORMALIZE BOOLEAN OPTIONAL DEFAULT TRUE THIS PARAMETER IS IGNORED WHEN FITINTERCEPT IS SET TO FALSE IF TRUE THE REGRESSORS X WILL BE NORMALIZED BEFORE REGRESSION BY SUBTRACTING THE MEAN AND DIVIDING BY THE L2NORM IF YOU WISH TO STANDARDIZE PLEASE USE SKLEARNPREPROCESSINGSTANDARDSCALER BEFORE CALLING FIT ON AN ESTIMATOR  
WITHNORMALIZEFALSE  
PRECOMPUTE TRUE FALSE 'AUTO' ARRAYLIKE WHETHER TO USE A PRECOMPUTED GRAM MATRIX TO SPEED UP CALCULATIONS IF SET TO AUTO LET US DECIDE THE GRAM MATRIX CAN ALSO BE PASSED AS ARGUMENT  
MAXITER INTEGER OPTIONAL MAXIMUM NUMBER OF ITERATIONS TO PERFORM CAN BE USED FOR EARLY STOPPING  
EPS FLOAT OPTIONAL THE MACHINEPRECISION REGULARIZATION IN THE COMPUTATION OF THE CHOLESKY DIAGONAL FACTORS INCREASE THIS FOR VERY ILLCONDITIONED SYSTEMS UNLIKE THE TOL PARAMETER IN SOME ITERATIVE OPTIMIZATIONBASED ALGORITHMS THIS PARAMETER DOES NOT CONTROL THE TOLERANCE OF THE OPTIMIZATION  
COPYX BOOLEAN OPTIONAL DEFAULT TRUE IF TRUE X WILL BE COPIED ELSE IT MAY BE OVERWRITTEN  
POSITIVE BOOLEAN DEFAULTFALSE RESTRICT COEFFICIENTS TO BE 0 BE AWARE THAT YOU MIGHT WANT TO REMOVE FITINTERCEPT WHICH IS SET TRUE BY DEFAULT UNDER THE POSITIVE RESTRICTION THE MODEL COEFFICIENTS DO NOT CONVERGE TO THE ORDINARYLEASTSQUARES SOLUTION FOR SMALL VALUES OF ALPHA ONLY COEFFICIENTS UP TO THE SMALLEST ALPHA VALUE ALPHASALPHAS 0  
MIN WHEN FITPATHTRUE REACHED BY THE STEPWISE LARSLASSO ALGORITHM ARE TYPICALLY IN CONGRUENCE WITH THE SOLUTION OF THE COORDINATE DESCENT LASSO ESTIMATOR AS A CONSEQUENCE USING LASSOLARSIC ONLY MAKES SENSE FOR PROBLEMS WHERE A SPARSE SOLUTION IS EXPECTED ANDOR REACHED

ATTRIBUTES

COEF ARRAY SHAPE NFEATURES PARAMETER VECTOR W IN THE FORMULATION FORMULA  
INTERCEPT FLOAT INDEPENDENT TERM IN DECISION FUNCTION  
ALPHA FLOAT THE ALPHA PARAMETER CHOSEN BY THE INFORMATION CRITERION  
NITER INT NUMBER OF ITERATIONS RUN BY LARSPATH TO FIND THE GRID OF ALPHAS  
CRITERION ARRAY SHAPE NALPHAS THE VALUE OF THE INFORMATION CRITERIA 'AIC' 'BIC' ACROSS ALL ALPHAS THE ALPHA WHICH HAS THE SMALLEST INFORMATION CRITERION IS CHOSEN THIS VALUE IS LARGER BY A FACTOR OF NSAMPLES COMPARED TO EQNS 215 AND 216 IN ZOU ET AL 2007  
SEE ALSO

LARSPATH LASSOLARS LASSOLARSCV

NOTES

THE ESTIMATION OF THE NUMBER OF DEGREES OF FREEDOM IS GIVEN BY  
498 CHAPTER 3 USER GUIDE



SCIKITLEARN USER GUIDE RELEASE 0213  
“ON THE DEGREES OF FREEDOM OF THE LASSO” HUI ZOU TREVOR HASTIE AND ROBERT TIBSHIRANI ANN STATIST V OLUME  
35 NUMBER 5 2007 21732192  
HTTPSENWIKIPEDIAORGWIKIAKAIKINFORMATIONCRITERION HTTPSENWIKIPEDIAORGWIKIBAYESIAN  
INFORMATIONCRITERION  
EXAMPLES  
FROM SKLEARN IMPORT LINEARMODEL  
REG LINEARMODELLASSOLARSICCRITERIONBIC  
REGFIT1 1 0 0 1 1 11111 0 11111  
  
LASSOLARSICCOPYXTRUE CRITERIONBIC EPS FITINTERCEPTTRUE  
MAXITER500 NORMALIZETRUE POSITIVEFALSE PRECOMPUTEAUTO  
VERBOSEFALSE  
PRINTREGCOEF  
0 111  
METHODS  
FITSELF X Y COPYX FIT THE MODEL USING X Y AS TRAINING DATA  
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR  
PREDICT SELF X PREDICT USING THE LINEAR MODEL  
SCORE SELF X Y SAMPLEWEIGHT RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PRE  
DICTION  
SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR  
INIT SELF CRITERION‘AIC’ FITINTERCEPTTRUE VERBOSEFALSE NORMALIZETRUE PRECOM  
PUTE‘AUTO’ MAXITER500 EPS2220446049250313E16 COPYXTRUE POSITIVEFALSE  
FITSELFXYCOPYXNONE  
FIT THE MODEL USING X Y AS TRAINING DATA  
PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA  
YARRAYLIKE SHAPE NSAMPLES TARGET VALUES WILL BE CAST TO X’S DTYPE IF NECESSARY  
COPYX BOOLEAN OPTIONAL DEFAULT NONE IF PROVIDED THIS PARAMETER WILL OVERRIDE THE  
CHOICE OF COPYX MADE AT INSTANCE CREATION IF TRUE X WILL BE COPIED ELSE IT MAY  
BE OVERWRITTEN  
RETURNS  
SELF OBJECT RETURNS AN INSTANCE OF SELF  
GETPARAMS SELFDEEPTTRUE  
GET PARAMETERS FOR THIS ESTIMATOR  
PARAMETERS  
DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBOBJECTS THAT ARE ESTIMATORS  
RETURNS  
33 MODEL SELECTION AND EVALUATION 499

```

SCIKITLEARN USER GUIDE RELEASE 0213
PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES
PREDICTSELF
PREDICT USING THE LINEAR MODEL
PARAMETERS
XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES SAMPLES
RETURNS
CARRAY SHAPE NSAMPLES RETURNS PREDICTED VALUES
SCORESELFXYSAMPLEWEIGHTNONE
RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION
THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{U}{V}$  WHERE U IS THE RESIDUAL SUM OF SQUARES  $Y_{TRUE} - Y_{PRED}$ 
2SUM AND V IS THE TOTAL SUM OF SQUARES  $Y_{TRUE} - Y_{TRUEMEAN}$  2SUM THE BEST POSSIBLE SCORE
IS 10 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS
PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 00
PARAMETERS
XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES FOR SOME ESTIMATORS THIS MAY
BE A PRECOMPUTED KERNEL MATRIX INSTEAD SHAPE NSAMPLES NSAMPLESFITTED WHERE
NSAMPLESFITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR
YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE VALUES FOR X
SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS
RETURNS
SCORE FLOAT R2 OF SELF PREDICTX WRT Y
NOTES
THE R2 SCORE USED WHEN CALLING SCORE ON A REGRESSOR WILL USE MULTIOUTPUTUNIFORMAVERAGE
FROM VERSION 023 TO KEEP CONSISTENT WITH METRICSR2SCORE THIS WILL INFLUENCE THE SCORE
METHOD OF ALL THE MULTIOUTPUT REGRESSORS EXCEPT FOR MULTIOUTPUTMULTIOUTPUTREGRESSOR
TO SPECIFY THE DEFAULT VALUE MANUALLY AND AVOID THE WARNING PLEASE EITHER CALL METRICSR2SCORE
DIRECTLY OR MAKE A CUSTOM SCORER WITH METRICSMAKESCORER THE BUILTIN SCORER R2 USES
MULTIOUTPUTUNIFORMAVERAGE
SETPARAMS SELFPARAMS
SET THE PARAMETERS OF THIS ESTIMATOR
THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE
PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT
OF A NESTED OBJECT
RETURNS
SELF
EXAMPLES USING SKLEARNLINEARMODELLASSOLARSIC
•LASSO MODEL SELECTION CROSSVALIDATION AIC BIC
500 CHAPTER 3 USER GUIDE

```

SCIKITLEARN USER GUIDE RELEASE 0213

OUT OF BAG ESTIMATES

WHEN USING ENSEMBLE METHODS BASE UPON BAGGING IE GENERATING NEW TRAINING SETS USING SAMPLING WITH REPLACEMENT PART OF THE TRAINING SET REMAINS UNUSED FOR EACH CLASSIFIER IN THE ENSEMBLE A DIFFERENT PART OF THE TRAINING SET IS LEFT OUT

THIS LEFT OUT PORTION CAN BE USED TO ESTIMATE THE GENERALIZATION ERROR WITHOUT HAVING TO RELY ON A SEPARATE VALIDATION SET THIS ESTIMATE COMES “FOR FREE” AS NO ADDITIONAL DATA IS NEEDED AND CAN BE USED FOR MODEL SELECTION

THIS IS CURRENTLY IMPLEMENTED IN THE FOLLOWING CLASSES

ENSEMBLERANDOMFORESTCLASSIFIER A RANDOM FOREST CLASSIFIER

ENSEMBLERANDOMFORESTREGRESSOR A RANDOM FOREST REGRESSOR

ENSEMBLEEXTRATREESCLASSIFIER AN EXTRATREES CLASSIFIER

ENSEMBLEEXTRATREESREGRESSOR NESTIMATORS

AN EXTRATREES REGRESSOR

ENSEMBLEGRADIENTBOOSTINGCLASSIFIER LOSS

GRADIENT BOOSTING FOR CLASSIFICATION

ENSEMBLEGRADIENTBOOSTINGREGRESSOR LOSS

GRADIENT BOOSTING FOR REGRESSION

SKLEARNENSEMBLE RANDOMFORESTCLASSIFIER

CLASSSSKLEARNENSEMBLE RANDOMFORESTCLASSIFIER NESTIMATORS’WARN’ CRITE

RION’GINI’ MAXDEPTHNONE

MINSAMPLESSPLIT2 MINSAMPLESLEAF1

MINWEIGHTFRACTIONLEAF00

MAXFEATURES’AUTO’

MAXLEAFNODESNONE

MINIMPURITYDECREASE00

MINIMPURITYSPLITNONE BOOTSTRAPTRUE

OOBSCOREFALSE NJOBSNONE

RANDOMSTATENONE VERBOSE0

WARMSTARTFALSE CLASSWEIGHTNONE

A RANDOM FOREST CLASSIFIER

A RANDOM FOREST IS A META ESTIMATOR THAT FITS A NUMBER OF DECISION TREE CLASSIFIERS ON VARIOUS SUBSAMPLES OF THE DATASET AND USES AVERAGING TO IMPROVE THE PREDICTIVE ACCURACY AND CONTROL OVERFITTING THE SUBSAMPLE SIZE IS ALWAYS THE SAME AS THE ORIGINAL INPUT SAMPLE SIZE BUT THE SAMPLES ARE DRAWN WITH REPLACEMENT IF

BOOTSTRAPTRUE DEFAULT

READ MORE IN THE USER GUIDE

PARAMETERS

NESTIMATORS INTEGER OPTIONAL DEFAULT10 THE NUMBER OF TREES IN THE FOREST

CHANGED IN VERSION 020 THE DEFAULT VALUE OF NESTIMATORS WILL CHANGE FROM 10 IN

VERSION 020 TO 100 IN VERSION 022

CRITERION STRING OPTIONAL DEFAULT”GINI” THE FUNCTION TO MEASURE THE QUALITY OF A SPLIT SUP

PORTED CRITERIA ARE “GINI” FOR THE GINI IMPURITY AND “ENTROPY” FOR THE INFORMATION GAIN NOTE

THIS PARAMETER IS TREESPECIFIC

MAXDEPTH INTEGER OR NONE OPTIONAL DEFAULTNONE THE MAXIMUM DEPTH OF THE TREE IF

NONE THEN NODES ARE EXPANDED UNTIL ALL LEAVES ARE PURE OR UNTIL ALL LEAVES CONTAIN LESS THAN

33 MODEL SELECTION AND EVALUATION 501

SCIKITLEARN USER GUIDE RELEASE 0213

MINSAMPLESSPLIT SAMPLES

MINSAMPLESSPLIT INT FLOAT OPTIONAL DEFAULT2 THE MINIMUM NUMBER OF SAMPLES REQUIRED TO SPLIT AN INTERNAL NODE

- IF INT THEN CONSIDER MINSAMPLESSPLIT AS THE MINIMUM NUMBER
- IF FLOAT THEN MINSAMPLESSPLIT IS A FRACTION AND CEILMINSAMPLESSPLIT

NSAMPLES ARE THE MINIMUM NUMBER OF SAMPLES FOR EACH SPLIT

CHANGED IN VERSION 018 ADDED FLOAT VALUES FOR FRACTIONS

MINSAMPLESLEAF INT FLOAT OPTIONAL DEFAULT1 THE MINIMUM NUMBER OF SAMPLES REQUIRED

TO BE AT A LEAF NODE A SPLIT POINT AT ANY DEPTH WILL ONLY BE CONSIDERED IF IT LEAVES AT LEAST MINSAMPLESLEAF TRAINING SAMPLES IN EACH OF THE LEFT AND RIGHT BRANCHES THIS MAY

HAVE THE EFFECT OF SMOOTHING THE MODEL ESPECIALLY IN REGRESSION

- IF INT THEN CONSIDER MINSAMPLESLEAF AS THE MINIMUM NUMBER
- IF FLOAT THEN MINSAMPLESLEAF IS A FRACTION AND CEILMINSAMPLESLEAF

NSAMPLES ARE THE MINIMUM NUMBER OF SAMPLES FOR EACH NODE

CHANGED IN VERSION 018 ADDED FLOAT VALUES FOR FRACTIONS

MINWEIGHTFRACTIONLEAF FLOAT OPTIONAL DEFAULT0 THE MINIMUM WEIGHTED FRACTION OF THE

SUM TOTAL OF WEIGHTS OF ALL THE INPUT SAMPLES REQUIRED TO BE AT A LEAF NODE SAMPLES HAVE EQUAL WEIGHT WHEN SAMPLEWEIGHT IS NOT PROVIDED

MAXFEATURES INT FLOAT STRING OR NONE OPTIONAL DEFAULT" AUTO" THE NUMBER OF FEATURES TO CONSIDER WHEN LOOKING FOR THE BEST SPLIT

- IF INT THEN CONSIDER MAXFEATURES FEATURES AT EACH SPLIT
- IF FLOAT THEN MAXFEATURES IS A FRACTION AND INTMAXFEATURES

NFEATURES FEATURES ARE CONSIDERED AT EACH SPLIT

- IF "AUTO" THEN MAXFEATURESSQRTNFEATURES
- IF "SQRT" THEN MAXFEATURESSQRTNFEATURES SAME AS "AUTO"
- IF "LOG2" THEN MAXFEATURESLOG2NFEATURES
- IF NONE THEN MAXFEATURESNFEATURES

NOTE THE SEARCH FOR A SPLIT DOES NOT STOP UNTIL AT LEAST ONE VALID PARTITION OF THE NODE SAMPLES IS FOUND EVEN IF IT REQUIRES TO EFFECTIVELY INSPECT MORE THAN MAXFEATURES FEATURES

MAXLEAFNODES INT OR NONE OPTIONAL DEFAULTNONE GROW TREES WITH MAXLEAFNODES

IN BESTFIRST FASHION BEST NODES ARE DEFINED AS RELATIVE REDUCTION IN IMPURITY IF NONE THEN UNLIMITED NUMBER OF LEAF NODES

MINIMPURITYDECREASE FLOAT OPTIONAL DEFAULT0 A NODE WILL BE SPLIT IF THIS SPLIT INDUCES

A DECREASE OF THE IMPURITY GREATER THAN OR EQUAL TO THIS VALUE

THE WEIGHTED IMPURITY DECREASE EQUATION IS THE FOLLOWING

NT N IMPURITY NTR NT RIGHTIMPURITY

NTL NT LEFTIMPURITY

WHERE N IS THE TOTAL NUMBER OF SAMPLES NT IS THE NUMBER OF SAMPLES AT THE CURRENT NODE

NTL IS THE NUMBER OF SAMPLES IN THE LEFT CHILD AND NTR IS THE NUMBER OF SAMPLES IN

THE RIGHT CHILD

NNTNTR ANDNTL ALL REFER TO THE WEIGHTED SUM IF SAMPLEWEIGHT IS PASSED

502 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

NEW IN VERSION 019

MINIMPURITYSPLIT FLOAT DEFAULT1E7 THRESHOLD FOR EARLY STOPPING IN TREE GROWTH A NODE WILL SPLIT IF ITS IMPURITY IS ABOVE THE THRESHOLD OTHERWISE IT IS A LEAF  
DEPRECATED SINCE VERSION 019 MINIMPURITYSPLIT HAS BEEN DEPRECATED IN FAVOR OF MINIMPURITYDECREASE IN 019 THE DEFAULT VALUE OF MINIMPURITYSPLIT WILL CHANGE FROM 1E7 TO 0 IN 023 AND IT WILL BE REMOVED IN 025 USE MINIMPURITYDECREASE INSTEAD

BOOTSTRAP BOOLEAN OPTIONAL DEFAULTTRUE WHETHER BOOTSTRAP SAMPLES ARE USED WHEN BUILDING TREES IF FALSE THE WHOLE DATSET IS USED TO BUILD EACH TREE

OOBSCORE BOOL DEFAULTFALSE WHETHER TO USE OUTFBAG SAMPLES TO ESTIMATE THE GENERALIZATION ACCURACY

NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF JOBS TO RUN IN PARALLEL FOR BOTH FIT ANDPREDICT NONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT  
1MEANS USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RAN DOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NPRANDOM

VERBOSE INT OPTIONAL DEFAULT0 CONTROLS THE VERBOSITY WHEN FITTING AND PREDICTING

WARMSTART BOOL OPTIONAL DEFAULTFALSE WHEN SET TO TRUE REUSE THE SOLUTION OF THE PREVIOUS CALL TO FIT AND ADD MORE ESTIMATORS TO THE ENSEMBLE OTHERWISE JUST FIT A WHOLE NEW FOREST SEE THE GLOSSARY

CLASSWEIGHT DICT LIST OF DICTS “BALANCED” “BALANCEDSUBSAMPLE” OR NONE OPTIONAL DE FAULTNONE WEIGHTS ASSOCIATED WITH CLASSES IN THE FORM CLASSLABEL WEIGHT  
IF NOT GIVEN ALL CLASSES ARE SUPPOSED TO HAVE WEIGHT ONE FOR MULTIOUTPUT PROBLEMS A LIST OF DICTS CAN BE PROVIDED IN THE SAME ORDER AS THE COLUMNS OF Y  
NOTE THAT FOR MULTIOUTPUT INCLUDING MULTILABEL WEIGHTS SHOULD BE DEFINED FOR EACH CLASS OF EVERY COLUMN IN ITS OWN DICT FOR EXAMPLE FOR FOURCLASS MULTILABEL CLASSIFICATION WEIGHTS SHOULD BE 0 1 1 1 0 1 1 5 0 1 1 1 0 1 1 1 INSTEAD OF 11 25  
31 41  
THE “BALANCED” MODE USES THE VALUES OF Y TO AUTOMATICALLY ADJUST WEIGHTS INVERSELY PROPORTIONAL TO CLASS FREQUENCIES IN THE INPUT DATA AS NSAMPLES NCLASSES NP  
BINCOUNTY

THE “BALANCEDSUBSAMPLE” MODE IS THE SAME AS “BALANCED” EXCEPT THAT WEIGHTS ARE COM PUTED BASED ON THE BOOTSTRAP SAMPLE FOR EVERY TREE GROWN  
FOR MULTIOUTPUT THE WEIGHTS OF EACH COLUMN OF Y WILL BE MULTIPLIED  
NOTE THAT THESE WEIGHTS WILL BE MULTIPLIED WITH SAMPLEWEIGHT PASSED THROUGH THE FIT METHOD IF SAMPLEWEIGHT IS SPECIFIED

ATTRIBUTES

ESTIMATORS LIST OF DECISIONTREECLASSIFIER THE COLLECTION OF FITTED SUBESTIMATORS

CLASSES ARRAY OF SHAPE NCLASSES OR A LIST OF SUCH ARRAYS THE CLASSES LABELS SINGLE OUTPUT PROBLEM OR A LIST OF ARRAYS OF CLASS LABELS MULTIOUTPUT PROBLEM

NCLASSES INT OR LIST THE NUMBER OF CLASSES SINGLE OUTPUT PROBLEM OR A LIST CONTAINING THE NUMBER OF CLASSES FOR EACH OUTPUT MULTIOUTPUT PROBLEM

33 MODEL SELECTION AND EVALUATION 503

SCIKITLEARN USER GUIDE RELEASE 0213

NFEATURES INT THE NUMBER OF FEATURES WHEN FIT IS PERFORMED

NOUTPUTS INT THE NUMBER OF OUTPUTS WHEN FIT IS PERFORMED

FEATUREIMPORTANCES ARRAY OF SHAPE NFEATURES RETURN THE FEATURE IMPORTANCES

THE HIGHER THE MORE IMPORTANT THE FEATURE

OOBSCORE FLOAT SCORE OF THE TRAINING DATASET OBTAINED USING AN OUTFBAG ESTIMATE

OOBDECISIONFUNCTION ARRAY OF SHAPE NSAMPLES NCLASSES DECISION FUNCTION COM

PUTED WITH OUTFBAG ESTIMATE ON THE TRAINING SET IF NESTIMATORS IS SMALL IT MIGHT

BE POSSIBLE THAT A DATA POINT WAS NEVER LEFT OUT DURING THE BOOTSTRAP IN THIS CASE

OOBDECISIONFUNCTION MIGHT CONTAIN NAN

SEE ALSO

DECISIONTREECLASSIFIER EXTRATREESCLASSIFIER

NOTES

THE DEFAULT VALUES FOR THE PARAMETERS CONTROLLING THE SIZE OF THE TREES EG MAXDEPTH

MINSAMPLESLEAF ETC LEAD TO FULLY GROWN AND UNPRUNED TREES WHICH CAN POTENTIALLY BE VERY LARGE ON

SOME DATA SETS TO REDUCE MEMORY CONSUMPTION THE COMPLEXITY AND SIZE OF THE TREES SHOULD BE CONTROLLED BY

SETTING THOSE PARAMETER VALUES

THE FEATURES ARE ALWAYS RANDOMLY PERMUTED AT EACH SPLIT THEREFORE THE BEST FOUND SPLIT MAY VARY EVEN WITH

THE SAME TRAINING DATA MAXFEATURESNFEATURES ANDBOOTSTRAPFALSE IF THE IMPROVEMENT OF THE

CRITERION IS IDENTICAL FOR SEVERAL SPLITS ENUMERATED DURING THE SEARCH OF THE BEST SPLIT TO OBTAIN A DETERMINISTIC

BEHAVIOUR DURING FITTING RANDOMSTATE HAS TO BE FIXED

REFERENCES

R45F14345C0001

EXAMPLES

FROM SKLEARNENSEMBLE IMPORT RANDOMFORESTCLASSIFIER

FROM SKLEARNDATASETS IMPORT MAKECLASSIFICATION

X Y MAKECLASSIFICATIONNSAMPLES1000 NFEATURES4

NINFORMATIVE2 NREDUNDANT0

RANDOMSTATE0 SHUFFLE FALSE

CLF RANDOMFORESTCLASSIFIERNESTIMATORS100 MAXDEPTH2

RANDOMSTATE0

CLFFITX Y

RANDOMFORESTCLASSIFIERBOOTSTRAPTRUE CLASSWEIGHTNONE CRITERIONGINI

MAXDEPTH2 MAXFEATURESAUTO MAXLEAFNODESNONE

MINIMPURITYDECREASE00 MINIMPURITYSPLITNONE

MINSAMPLESLEAF1 MINSAMPLESSPLIT2

MINWEIGHTFRACTIONLEAF00 NESTIMATORS100 NJOBSNONE

OOBSCOREFALSE RANDOMSTATE0 VERBOSE0 WARMSTARTFALSE

PRINTCLFFEATUREIMPORTANCES

014205973 076664038 00282433 006305659

PRINTCLFPREDICT0 0 0 0

1

504 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

METHODS

APPLY SELF X APPLY TREES IN THE FOREST TO X RETURN LEAF INDICES  
DECISIONPATH SELF X RETURN THE DECISION PATH IN THE FOREST  
FITSELF X Y SAMPLEWEIGHT BUILD A FOREST OF TREES FROM THE TRAINING SET X Y  
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR  
PREDICT SELF X PREDICT CLASS FOR X  
PREDICTLOGPROBA SELF X PREDICT CLASS LOGPROBABILITIES FOR X  
PREDICTPROBA SELF X PREDICT CLASS PROBABILITIES FOR X  
SCORE SELF X Y SAMPLEWEIGHT RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS  
SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR  
INIT SELFNESTIMATORS'WARN' CRITERION'GINI' MAXDEPTHNONE MINSAMPLESSPLIT2  
MINSAMPLESLEAF1 MINWEIGHTFRACTIONLEAF00 MAXFEATURES'AUTO'  
MAXLEAFNODESNONE MINIMPURITYDECREASE00 MINIMPURITYSPLITNONE  
BOOTSTRAPTRUE OOBSCOREFALSE NJOBSNONE RANDOMSTATENONE VERBOSE0  
WARMSTARTFALSE CLASSWEIGHTNONE  
APPLYSELF X  
APPLY TREES IN THE FOREST TO X RETURN LEAF INDICES  
PARAMETERS  
XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTERNALLY ITS DTYPE WILL BE CONVERTED TO DTYPE NPFLOAT32 IF A SPARSE MATRIX IS PROVIDED IT WILL BE CONVERTED INTO A SPARSE CSR MATRIX  
RETURNS  
XLEAVES ARRAYLIKE SHAPE NSAMPLES NESTIMATORS FOR EACH DATAPOINT X IN X AND FOR EACH TREE IN THE FOREST RETURN THE INDEX OF THE LEAF X ENDS UP IN  
DECISIONPATH SELF X  
RETURN THE DECISION PATH IN THE FOREST  
NEW IN VERSION 018  
PARAMETERS  
XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTERNALLY ITS DTYPE WILL BE CONVERTED TO DTYPE NPFLOAT32 IF A SPARSE MATRIX IS PROVIDED IT WILL BE CONVERTED INTO A SPARSE CSR MATRIX  
RETURNS  
INDICATOR SPARSE CSR ARRAY SHAPE NSAMPLES NNODES RETURN A NODE INDICATOR MATRIX WHERE NON ZERO ELEMENTS INDICATES THAT THE SAMPLES GOES THROUGH THE NODES  
NNODESPTR ARRAY OF SIZE NESTIMATORS 1 THE COLUMNS FROM INDICATOR  
TORNNODESPTR INNNODESPTR I1 GIVES THE INDICATOR VALUE FOR THE ITH ESTIMATOR  
FEATUREIMPORTANCES  
RETURN THE FEATURE IMPORTANCES THE HIGHER THE MORE IMPORTANT THE FEATURE  
RETURNS  
33 MODEL SELECTION AND EVALUATION 505

SCIKITLEARN USER GUIDE RELEASE 0213

FEATUREIMPORTANCES ARRAY SHAPE NFEATURES THE VALUES OF THIS ARRAY SUM TO 1 UNLESS ALL TREES ARE SINGLE NODE TREES CONSISTING OF ONLY THE ROOT NODE IN WHICH CASE IT WILL BE AN ARRAY OF ZEROS

FITSELFXYSAMPLEWEIGHTNONE

BUILD A FOREST OF TREES FROM THE TRAINING SET X Y

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX OF SHAPE NSAMPLES NFEATURES THE TRAINING INPUT SAMPLES INTERNALLY ITS DTYPE WILL BE CONVERTED TO DTYPENPFLOAT32 IF A SPARSE MATRIX IS PROVIDED IT WILL BE CONVERTED INTO A SPARSE CSCMATRIX

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS THE TARGET VALUES CLASS LABELS IN CLASSIFICATION REAL NUMBERS IN REGRESSION

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OR NONE SAMPLE WEIGHTS IF NONE THEN SAMPLES ARE EQUALLY WEIGHTED SPLITS THAT WOULD CREATE CHILD NODES WITH NET ZERO OR NEGATIVE WEIGHT ARE IGNORED WHILE SEARCHING FOR A SPLIT IN EACH NODE IN THE CASE OF CLASSIFICATION SPLITS ARE ALSO IGNORED IF THEY WOULD RESULT IN ANY SINGLE CLASS CARRYING A NEGATIVE WEIGHT IN EITHER CHILD NODE

RETURNS

SELF OBJECT

GETPARAMS SELFDEEPTREE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELFXY

PREDICT CLASS FOR X

THE PREDICTED CLASS OF AN INPUT SAMPLE IS A VOTE BY THE TREES IN THE FOREST WEIGHTED BY THEIR PROBABILITY ESTIMATES THAT IS THE PREDICTED CLASS IS THE ONE WITH HIGHEST MEAN PROBABILITY ESTIMATE ACROSS THE TREES

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX OF SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTERNALLY ITS DTYPE WILL BE CONVERTED TO DTYPENPFLOAT32 IF A SPARSE MATRIX IS PROVIDED IT WILL BE CONVERTED INTO A SPARSE CSRMATRIX

RETURNS

YARRAY OF SHAPE NSAMPLES OR NSAMPLES NOUTPUTS THE PREDICTED CLASSES

PREDICTLOGPROBA SELFXY

PREDICT CLASS LOGPROBABILITIES FOR X

THE PREDICTED CLASS LOGPROBABILITIES OF AN INPUT SAMPLE IS COMPUTED AS THE LOG OF THE MEAN PREDICTED CLASS PROBABILITIES OF THE TREES IN THE FOREST

PARAMETERS

506 CHAPTER 3 USER GUIDE



SCIKITLEARN USER GUIDE RELEASE 0213

XARRAYLIKE OR SPARSE MATRIX OF SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTERNALLY ITS DTYPE WILL BE CONVERTED TO DYPENPFLOAT32 IF A SPARSE MATRIX IS PROVIDED IT WILL BE CONVERTED INTO A SPARSE CSRMATRIX

RETURNS

PARRAY OF SHAPE NSAMPLES NCLASSES OR A LIST OF NOUTPUTS SUCH ARRAYS IF NOUTPUTS

1 THE CLASS PROBABILITIES OF THE INPUT SAMPLES THE ORDER OF THE CLASSES CORRESPONDS TO THAT IN THE ATTRIBUTE CLASSES

PREDICTPROBA SELF

PREDICT CLASS PROBABILITIES FOR X

THE PREDICTED CLASS PROBABILITIES OF AN INPUT SAMPLE ARE COMPUTED AS THE MEAN PREDICTED CLASS PROBABILITIES OF THE TREES IN THE FOREST THE CLASS PROBABILITY OF A SINGLE TREE IS THE FRACTION OF SAMPLES OF THE SAME CLASS IN A LEAF

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX OF SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTERNALLY ITS DTYPE WILL BE CONVERTED TO DYPENPFLOAT32 IF A SPARSE MATRIX IS PROVIDED IT WILL BE CONVERTED INTO A SPARSE CSRMATRIX

RETURNS

PARRAY OF SHAPE NSAMPLES NCLASSES OR A LIST OF NOUTPUTS SUCH ARRAYS IF NOUTPUTS

1 THE CLASS PROBABILITIES OF THE INPUT SAMPLES THE ORDER OF THE CLASSES CORRESPONDS TO THAT IN THE ATTRIBUTE CLASSES

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS

IN MULTILABEL CLASSIFICATION THIS IS THE SUBSET ACCURACY WHICH IS A HARSH METRIC SINCE YOU REQUIRE FOR EACH SAMPLE THAT EACH LABEL SET BE CORRECTLY PREDICTED

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE LABELS FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT MEAN ACCURACY OF SELF PREDICTX WRT Y

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

EXAMPLES USING SKLEARNENSEMBLERANDOMFORESTCLASSIFIER

- COMPARISON OF CALIBRATION OF CLASSIFIERS
- PROBABILITY CALIBRATION FOR 3CLASS CLASSIFICATION

33 MODEL SELECTION AND EVALUATION 507

SCIKITLEARN USER GUIDE RELEASE 0213

- CLASSIFIER COMPARISON
- INDUCTIVE CLUSTERING
- PLOT CLASS PROBABILITIES CALCULATED BY THE VOTINGCLASSIFIER
- OOB ERRORS FOR RANDOM FORESTS
- FEATURE TRANSFORMATIONS WITH ENSEMBLES OF TREES
- PLOT THE DECISION SURFACES OF ENSEMBLES OF TREES ON THE IRIS DATASET
- COMPARING RANDOMIZED SEARCH AND GRID SEARCH FOR HYPERPARAMETER ESTIMATION
- CLASSIFICATION OF TEXT DOCUMENTS USING SPARSE FEATURES

SKLEARNENSEMBLE RANDOMFORESTREGRESSOR  
CLASSSSKLEARNENSEMBLE RANDOMFORESTREGRESSOR NESTIMATORS'WARN' CRITERION'MSE' MAXDEPTHNONE  
MINSAMPLESSPLIT2 MINSAMPLESLEAF1  
MINWEIGHTFRACTIONLEAF00  
MAXFEATURES'AUTO' MAXLEAFNODESNONE  
MINIMPURITYDECREASE00  
MINIMPURITYSPLITNONE BOOTSTRAPTRUE  
OOBSCOREFALSE NJOBSNONE  
RANDOMSTATENONE VERBOSE0  
WARMSTARTFALSE  
A RANDOM FOREST REGRESSOR

A RANDOM FOREST IS A META ESTIMATOR THAT FITS A NUMBER OF CLASSIFYING DECISION TREES ON VARIOUS SUBSAMPLES OF THE DATASET AND USES AVERAGING TO IMPROVE THE PREDICTIVE ACCURACY AND CONTROL OVERFITTING THE SUBSAMPLE SIZE IS ALWAYS THE SAME AS THE ORIGINAL INPUT SAMPLE SIZE BUT THE SAMPLES ARE DRAWN WITH REPLACEMENT IF BOOTSTRAPTRUE DEFAULT  
READ MORE IN THE USER GUIDE  
PARAMETERS

NESTIMATORS INTEGER OPTIONAL DEFAULT10 THE NUMBER OF TREES IN THE FOREST  
CHANGED IN VERSION 020 THE DEFAULT VALUE OF NESTIMATORS WILL CHANGE FROM 10 IN  
VERSION 020 TO 100 IN VERSION 022  
CRITERION STRING OPTIONAL DEFAULT"mse" THE FUNCTION TO MEASURE THE QUALITY OF A SPLIT SUPPORTED CRITERIA ARE "mse" FOR THE MEAN SQUARED ERROR WHICH IS EQUAL TO VARIANCE REDUCTION AS FEATURE SELECTION CRITERION AND "mae" FOR THE MEAN ABSOLUTE ERROR  
NEW IN VERSION 018 MEAN ABSOLUTE ERROR MAE CRITERION  
MAXDEPTH INTEGER OR NONE OPTIONAL DEFAULTNONE THE MAXIMUM DEPTH OF THE TREE IF NONE THEN NODES ARE EXPANDED UNTIL ALL LEAVES ARE PURE OR UNTIL ALL LEAVES CONTAIN LESS THAN MINSAMPLESSPLIT SAMPLES  
MINSAMPLESSPLIT INT FLOAT OPTIONAL DEFAULT2 THE MINIMUM NUMBER OF SAMPLES REQUIRED TO SPLIT AN INTERNAL NODE  
• IF INT THEN CONSIDER MINSAMPLESSPLIT AS THE MINIMUM NUMBER  
• IF FLOAT THEN MINSAMPLESSPLIT IS A FRACTION AND CEILMINSAMPLESSPLIT  
NSAMPLES ARE THE MINIMUM NUMBER OF SAMPLES FOR EACH SPLIT  
CHANGED IN VERSION 018 ADDED FLOAT VALUES FOR FRACTIONS  
508 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

MINSAMPLESLEAF INT FLOAT OPTIONAL DEFAULT1 THE MINIMUM NUMBER OF SAMPLES REQUIRED TO BE AT A LEAF NODE A SPLIT POINT AT ANY DEPTH WILL ONLY BE CONSIDERED IF IT LEAVES AT LEAST MINSAMPLESLEAF TRAINING SAMPLES IN EACH OF THE LEFT AND RIGHT BRANCHES THIS MAY HAVE THE EFFECT OF SMOOTHING THE MODEL ESPECIALLY IN REGRESSION

- IF INT THEN CONSIDER MINSAMPLESLEAF AS THE MINIMUM NUMBER
- IF FLOAT THEN MINSAMPLESLEAF IS A FRACTION AND CEILMINSAMPLESLEAF NSAMPLES ARE THE MINIMUM NUMBER OF SAMPLES FOR EACH NODE

CHANGED IN VERSION 018 ADDED FLOAT VALUES FOR FRACTIONS

MINWEIGHTFRACTIONLEAF FLOAT OPTIONAL DEFAULT0 THE MINIMUM WEIGHTED FRACTION OF THE SUM TOTAL OF WEIGHTS OF ALL THE INPUT SAMPLES REQUIRED TO BE AT A LEAF NODE SAMPLES HAVE EQUAL WEIGHT WHEN SAMPLEWEIGHT IS NOT PROVIDED

MAXFEATURES INT FLOAT STRING OR NONE OPTIONAL DEFAULT" AUTO" THE NUMBER OF FEATURES TO CONSIDER WHEN LOOKING FOR THE BEST SPLIT

- IF INT THEN CONSIDER MAXFEATURES FEATURES AT EACH SPLIT
- IF FLOAT THEN MAXFEATURES IS A FRACTION AND INTMAXFEATURES NFEATURES FEATURES ARE CONSIDERED AT EACH SPLIT
- IF "AUTO" THEN MAXFEATURESNFEATURES
- IF "SQRT" THEN MAXFEATURESSQRTNFEATURES
- IF "LOG2" THEN MAXFEATURESLOG2NFEATURES
- IF NONE THEN MAXFEATURESNFEATURES

NOTE THE SEARCH FOR A SPLIT DOES NOT STOP UNTIL AT LEAST ONE VALID PARTITION OF THE NODE SAMPLES IS FOUND EVEN IF IT REQUIRES TO EFFECTIVELY INSPECT MORE THAN MAXFEATURES FEATURES

MAXLEAFNODES INT OR NONE OPTIONAL DEFAULTNONE GROW TREES WITH MAXLEAFNODES

IN BESTFIRST FASHION BEST NODES ARE DEFINED AS RELATIVE REDUCTION IN IMPURITY IF NONE THEN UNLIMITED NUMBER OF LEAF NODES

MINIMPURITYDECREASE FLOAT OPTIONAL DEFAULT0 A NODE WILL BE SPLIT IF THIS SPLIT INDUCES

A DECREASE OF THE IMPURITY GREATER THAN OR EQUAL TO THIS VALUE

THE WEIGHTED IMPURITY DECREASE EQUATION IS THE FOLLOWING

$$NT \cdot N \cdot \text{IMPURITY} - NTR \cdot NT \cdot \text{RIGHTIMPURITY}$$

$$- NTL \cdot NT \cdot \text{LEFTIMPURITY}$$

WHERE N IS THE TOTAL NUMBER OF SAMPLES NT IS THE NUMBER OF SAMPLES AT THE CURRENT NODE

NTL IS THE NUMBER OF SAMPLES IN THE LEFT CHILD AND NTR IS THE NUMBER OF SAMPLES IN

THE RIGHT CHILD

$$NNT - NTR \text{ AND } NNTL \text{ ALL REFER TO THE WEIGHTED SUM IF SAMPLEWEIGHT IS PASSED}$$

NEW IN VERSION 019

MINIMPURITYSPLIT FLOAT DEFAULT1E7 THRESHOLD FOR EARLY STOPPING IN TREE GROWTH A NODE

WILL SPLIT IF ITS IMPURITY IS ABOVE THE THRESHOLD OTHERWISE IT IS A LEAF

DEPRECATED SINCE VERSION 019 MINIMPURITYSPLIT HAS BEEN DEPRECATED IN FAVOR OF

MINIMPURITYDECREASE IN 019 THE DEFAULT VALUE OF MINIMPURITYSPLIT

WILL CHANGE FROM 1E7 TO 0 IN 023 AND IT WILL BE REMOVED IN 025 USE

MINIMPURITYDECREASE INSTEAD

33 MODEL SELECTION AND EVALUATION 509

SCIKITLEARN USER GUIDE RELEASE 0213

BOOTSTRAP BOOLEAN OPTIONAL DEFAULTTRUE WHETHER BOOTSTRAP SAMPLES ARE USED WHEN BUILDING TREES IF FALSE THE WHOLE DATSET IS USED TO BUILD EACH TREE

OOBSCORE BOOL OPTIONAL DEFAULTFALSE WHETHER TO USE OUTOFBAG SAMPLES TO ESTIMATE THE R2 ON UNSEEN DATA

NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF JOBS TO RUN IN PARALLEL FOR BOTH FIT ANDPREDICT NONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT

1MEANS USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RAN

DOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE

RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NPRANDOM

VERBOSE INT OPTIONAL DEFAULT0 CONTROLS THE VERBOSITY WHEN FITTING AND PREDICTING

WARMSTART BOOL OPTIONAL DEFAULTFALSE WHEN SET TO TRUE REUSE THE SOLUTION OF THE PREVIOUS CALL TO FIT AND ADD MORE ESTIMATORS TO THE ENSEMBLE OTHERWISE JUST FIT A WHOLE NEW FOREST SEE THE GLOSSARY

ATTRIBUTES

ESTIMATORS LIST OF DECISIONTREEREgressor THE COLLECTION OF FITTED SUBESTIMATORS

FEATUREIMPORTANCES ARRAY OF SHAPE NFEATURES RETURN THE FEATURE IMPORTANCES

THE HIGHER THE MORE IMPORTANT THE FEATURE

NFEATURES INT THE NUMBER OF FEATURES WHEN FIT IS PERFORMED

NOUTPUTS INT THE NUMBER OF OUTPUTS WHEN FIT IS PERFORMED

OOBSCORE FLOAT SCORE OF THE TRAINING DATASET OBTAINED USING AN OUTOFBAG ESTIMATE

OOBPREDICTION ARRAY OF SHAPE NSAMPLES PREDICTION COMPUTED WITH OUTOFBAG ESTIMATE ON THE TRAINING SET

SEE ALSO

DECISIONTREEREgressor EXTRATREESRegressor

NOTES

THE DEFAULT VALUES FOR THE PARAMETERS CONTROLLING THE SIZE OF THE TREES EG MAXDEPTH

MINSAMPLESLEAF ETC LEAD TO FULLY GROWN AND UNPRUNED TREES WHICH CAN POTENTIALLY BE VERY LARGE ON

SOME DATA SETS TO REDUCE MEMORY CONSUMPTION THE COMPLEXITY AND SIZE OF THE TREES SHOULD BE CONTROLLED BY SETTING THOSE PARAMETER VALUES

THE FEATURES ARE ALWAYS RANDOMLY PERMUTED AT EACH SPLIT THEREFORE THE BEST FOUND SPLIT MAY VARY EVEN WITH

THE SAME TRAINING DATA MAXFEATURESNFEATURES ANDBOOTSTRAPFALSE IF THE IMPROVEMENT OF THE

CRITERION IS IDENTICAL FOR SEVERAL SPLITS ENUMERATED DURING THE SEARCH OF THE BEST SPLIT TO OBTAIN A DETERMINISTIC

BEHAVIOUR DURING FITTING RANDOMSTATE HAS TO BE FIXED

THE DEFAULT VALUE MAXFEATURESAUTO USESNFEATURES RATHER THAN NFEATURES 3 THE LATTER

WAS ORIGINALLY SUGGESTED IN 1 WHEREAS THE FORMER WAS MORE RECENTLY JUSTIFIED EMPIRICALLY IN 2

REFERENCES

RF91CAB2DC4271 RF91CAB2DC4272

510 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES

```
FROM SKLEARNENSEMBLE IMPORT RANDOMFORESTREGRESSOR
FROM SKLEARNDATASETS IMPORT MAKEREGRESSION
X Y MAKEREGRESSIONNFEATURES4 NINFORMATIVE2
RANDOMSTATE0 SHUFFLE FALSE
REGR RANDOMFORESTREGRESSORMAXDEPTH2 RANDOMSTATE0
NESTIMATORS100
REGRFITX Y
RANDOMFORESTREGRESSORBOOTSTRAPTRUE CRITERIONMSE MAXDEPTH2
MAXFEATURESAUTO MAXLEAFNODESNONE
MINIMPURITYDECREASE00 MINIMPURITYSPLITNONE
MINSAMPLESLEAF1 MINSAMPLESSPLIT2
MINWEIGHTFRACTIONLEAF00 NESTIMATORS100 NJOBSNONE
OOBSCOREFALSE RANDOMSTATE0 VERBOSE0 WARMSTARTFALSE
PRINTREGRFEATUREIMPORTANCES
018146984 081473937 000145312 000233767
PRINTREGRPREDICT0 0 0 0
832987858
```

METHODS

```
APPLY SELF X APPLY TREES IN THE FOREST TO X RETURN LEAF INDICES
DECISIONPATH SELF X RETURN THE DECISION PATH IN THE FOREST
FITSELF X Y SAMPLEWEIGHT BUILD A FOREST OF TREES FROM THE TRAINING SET X Y
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR
PREDICT SELF X PREDICT REGRESSION TARGET FOR X
SCORE SELF X Y SAMPLEWEIGHT RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PRE
DICTION
SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR
INIT SELFNESTIMATORS'WARN' CRITERION'MSE' MAXDEPTHNONE MINSAMPLESSPLIT2
MINSAMPLESLEAF1 MINWEIGHTFRACTIONLEAF00 MAXFEATURES'AUTO'
MAXLEAFNODESNONE MINIMPURITYDECREASE00 MINIMPURITYSPLITNONE
BOOTSTRAPTRUE OOBSCOREFALSE NJOBSNONE RANDOMSTATENONE VERBOSE0
WARMSTARTFALSE
APPLYSELF
APPLY TREES IN THE FOREST TO X RETURN LEAF INDICES
PARAMETERS
XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTER
NALLY ITS DTYPE WILL BE CONVERTED TO DTYPENPFLOAT32 IF A SPARSE MATRIX IS PROVIDED
IT WILL BE CONVERTED INTO A SPARSE CSRMATRIX
RETURNS
XLEAVES ARRAYLIKE SHAPE NSAMPLES NESTIMATORS FOR EACH DATAPOINT X IN X AND FOR
EACH TREE IN THE FOREST RETURN THE INDEX OF THE LEAF X ENDS UP IN
DECISIONPATH SELF
RETURN THE DECISION PATH IN THE FOREST
33 MODEL SELECTION AND EVALUATION 511
```

SCIKITLEARN USER GUIDE RELEASE 0213  
NEW IN VERSION 018  
PARAMETERS  
XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTERNALLY ITS DTYPE WILL BE CONVERTED TO DTYPENPFLOAT32 IF A SPARSE MATRIX IS PROVIDED IT WILL BE CONVERTED INTO A SPARSE CSRMATRIX  
RETURNS  
INDICATOR SPARSE CSR ARRAY SHAPE NSAMPLES NNODES RETURN A NODE INDICATOR MATRIX WHERE NON ZERO ELEMENTS INDICATES THAT THE SAMPLES GOES THROUGH THE NODES  
NNODESPTR ARRAY OF SIZE NESTIMATORS 1 THE COLUMNS FROM INDICATORNNODESPTRINNNODESPTRI1 GIVES THE INDICATOR VALUE FOR THE ITH ESTIMATOR  
FEATUREIMPORTANCES  
RETURN THE FEATURE IMPORTANCES THE HIGHER THE MORE IMPORTANT THE FEATURE  
RETURNS  
FEATUREIMPORTANCES ARRAY SHAPE NFEATURES THE VALUES OF THIS ARRAY SUM TO 1 UNLESS ALL TREES ARE SINGLE NODE TREES CONSISTING OF ONLY THE ROOT NODE IN WHICH CASE IT WILL BE AN ARRAY OF ZEROS  
FITSELFXYSAMPLEWEIGHTNONE  
BUILD A FOREST OF TREES FROM THE TRAINING SET X Y  
PARAMETERS  
XARRAYLIKE OR SPARSE MATRIX OF SHAPE NSAMPLES NFEATURES THE TRAINING INPUT SAMPLES INTERNALLY ITS DTYPE WILL BE CONVERTED TO DTYPENPFLOAT32 IF A SPARSE MATRIX IS PROVIDED IT WILL BE CONVERTED INTO A SPARSE CSCMATRIX  
YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS THE TARGET VALUES CLASS LABELS IN CLASSIFICATION REAL NUMBERS IN REGRESSION  
SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OR NONE SAMPLE WEIGHTS IF NONE THEN SAMPLES ARE EQUALLY WEIGHTED SPLITS THAT WOULD CREATE CHILD NODES WITH NET ZERO OR NEGATIVE WEIGHT ARE IGNORED WHILE SEARCHING FOR A SPLIT IN EACH NODE IN THE CASE OF CLASSIFICATION SPLITS ARE ALSO IGNORED IF THEY WOULD RESULT IN ANY SINGLE CLASS CARRYING A NEGATIVE WEIGHT IN EITHER CHILD NODE  
RETURNS  
SELF OBJECT  
GETPARAMS SELFDEEPTREE  
GET PARAMETERS FOR THIS ESTIMATOR  
PARAMETERS  
DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS  
RETURNS  
PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES  
PREDICTSELF  
PREDICT REGRESSION TARGET FOR X  
512 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

THE PREDICTED REGRESSION TARGET OF AN INPUT SAMPLE IS COMPUTED AS THE MEAN PREDICTED REGRESSION TARGETS OF THE TREES IN THE FOREST

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX OF SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTER  
NALLY ITS DTYPE WILL BE CONVERTED TO DTYPENPFLOAT32 IF A SPARSE MATRIX IS PROVIDED  
IT WILL BE CONVERTED INTO A SPARSE CSRMATRIX

RETURNS

YARRAY OF SHAPE NSAMPLES OR NSAMPLES NOUTPUTS THE PREDICTED VALUES

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{U}{V}$  WHERE U IS THE RESIDUAL SUM OF SQUARES  $Y_{TRUE} - Y_{PRED}$   
2SUM AND V IS THE TOTAL SUM OF SQUARES  $Y_{TRUE} - Y_{TRUEMEAN}$  2SUM THE BEST POSSIBLE SCORE  
IS 10 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS  
PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 00

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES FOR SOME ESTIMATORS THIS MAY  
BE A PRECOMPUTED KERNEL MATRIX INSTEAD SHAPE NSAMPLES NSAMPLESFITTED WHERE  
NSAMPLESFITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE VALUES FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT R2 OF SELF PREDICTX WRT Y

NOTES

THE R2 SCORE USED WHEN CALLING SCORE ON A REGRESSOR WILL USE MULTIOUTPUTUNIFORMAVERAGE  
FROM VERSION 023 TO KEEP CONSISTENT WITH METRICSR2SCORE THIS WILL INFLUENCE THE SCORE  
METHOD OF ALL THE MULTIOUTPUT REGRESSORS EXCEPT FOR MULTIOUTPUTMULTIOUTPUTREGRESSOR  
TO SPECIFY THE DEFAULT VALUE MANUALLY AND AVOID THE WARNING PLEASE EITHER CALL METRICSR2SCORE  
DIRECTLY OR MAKE A CUSTOM SCORER WITH METRICSMAKESCORER THE BUILTIN SCORER R2 USES  
MULTIOUTPUTUNIFORMAVERAGE

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE  
PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT  
OF A NESTED OBJECT

RETURNS

SELF

EXAMPLES USING SKLEARNENSEMBLERANDOMFORESTREGRESSOR

- PREDICTION LATENCY
- PLOT INDIVIDUAL AND VOTING REGRESSION PREDICTIONS

33 MODEL SELECTION AND EVALUATION 513

SCIKITLEARN USER GUIDE RELEASE 0213

- COMPARING RANDOM FORESTS AND THE MULTIOUTPUT META ESTIMATOR
- IMPUTING MISSING VALUES BEFORE BUILDING AN ESTIMATOR

SKLEARNENSEMBLE EXTRATREESCLASSIFIER

CLASSSSKLEARNENSEMBLE EXTRATREESCLASSIFIER NESTIMATORS'WARN' CRITE

RION'GINI' MAXDEPTHNONE

MINSAMPLESSPLIT2 MINSAMPLESLEAF1

MINWEIGHTFRACTIONLEAF00

MAXFEATURES'AUTO' MAXLEAFNODESNONE

MINIMPURITYDECREASE00

MINIMPURITYSPLITNONE BOOTSTRAPFALSE

OOBSCOREFALSE NJOBSNONE

RANDOMSTATENONE VERBOSE0

WARMSTARTFALSE CLASSWEIGHTNONE

AN EXTRATREES CLASSIFIER

THIS CLASS IMPLEMENTS A META ESTIMATOR THAT FITS A NUMBER OF RANDOMIZED DECISION TREES AKA EXTRATREES ON VARIOUS SUBSAMPLES OF THE DATASET AND USES AVERAGING TO IMPROVE THE PREDICTIVE ACCURACY AND CONTROL OVERFITTING

READ MORE IN THE USER GUIDE

PARAMETERS

NESTIMATORS INTEGER OPTIONAL DEFAULT10 THE NUMBER OF TREES IN THE FOREST

CHANGED IN VERSION 020 THE DEFAULT VALUE OF NESTIMATORS WILL CHANGE FROM 10 IN

VERSION 020 TO 100 IN VERSION 022

CRITERION STRING OPTIONAL DEFAULT"GINI" THE FUNCTION TO MEASURE THE QUALITY OF A SPLIT SUP

PORTED CRITERIA ARE "GINI" FOR THE GINI IMPURITY AND "ENTROPY" FOR THE INFORMATION GAIN

MAXDEPTH INTEGER OR NONE OPTIONAL DEFAULTNONE THE MAXIMUM DEPTH OF THE TREE IF

NONE THEN NODES ARE EXPANDED UNTIL ALL LEAVES ARE PURE OR UNTIL ALL LEAVES CONTAIN LESS THAN

MINSAMPLESSPLIT SAMPLES

MINSAMPLESSPLIT INT FLOAT OPTIONAL DEFAULT2 THE MINIMUM NUMBER OF SAMPLES REQUIRED

TO SPLIT AN INTERNAL NODE

- IF INT THEN CONSIDER MINSAMPLESSPLIT AS THE MINIMUM NUMBER
- IF FLOAT THEN MINSAMPLESSPLIT IS A FRACTION AND CEILMINSAMPLESSPLIT

NSAMPLES ARE THE MINIMUM NUMBER OF SAMPLES FOR EACH SPLIT

CHANGED IN VERSION 018 ADDED FLOAT VALUES FOR FRACTIONS

MINSAMPLESLEAF INT FLOAT OPTIONAL DEFAULT1 THE MINIMUM NUMBER OF SAMPLES REQUIRED

TO BE AT A LEAF NODE A SPLIT POINT AT ANY DEPTH WILL ONLY BE CONSIDERED IF IT LEAVES AT LEAST

MINSAMPLESLEAF TRAINING SAMPLES IN EACH OF THE LEFT AND RIGHT BRANCHES THIS MAY

HAVE THE EFFECT OF SMOOTHING THE MODEL ESPECIALLY IN REGRESSION

- IF INT THEN CONSIDER MINSAMPLESLEAF AS THE MINIMUM NUMBER
- IF FLOAT THEN MINSAMPLESLEAF IS A FRACTION AND CEILMINSAMPLESLEAF

NSAMPLES ARE THE MINIMUM NUMBER OF SAMPLES FOR EACH NODE

CHANGED IN VERSION 018 ADDED FLOAT VALUES FOR FRACTIONS

MINWEIGHTFRACTIONLEAF FLOAT OPTIONAL DEFAULT0 THE MINIMUM WEIGHTED FRACTION OF THE

SUM TOTAL OF WEIGHTS OF ALL THE INPUT SAMPLES REQUIRED TO BE AT A LEAF NODE SAMPLES HAVE

EQUAL WEIGHT WHEN SAMPLEWEIGHT IS NOT PROVIDED

514 CHAPTER 3 USER GUIDE



SCIKITLEARN USER GUIDE RELEASE 0213

MAXFEATURES INT FLOAT STRING OR NONE OPTIONAL DEFAULT" AUTO" THE NUMBER OF FEATURES TO CONSIDER WHEN LOOKING FOR THE BEST SPLIT

- IF INT THEN CONSIDER MAXFEATURES FEATURES AT EACH SPLIT
- IF FLOAT THEN MAXFEATURES IS A FRACTION AND INTMAXFEATURES NFEATURES FEATURES ARE CONSIDERED AT EACH SPLIT
- IF "AUTO" THEN MAXFEATURESSQRTNFEATURES
- IF "SQRT" THEN MAXFEATURESSQRTNFEATURES
- IF "LOG2" THEN MAXFEATURESLOG2NFEATURES
- IF NONE THEN MAXFEATURESNFEATURES

NOTE THE SEARCH FOR A SPLIT DOES NOT STOP UNTIL AT LEAST ONE VALID PARTITION OF THE NODE SAMPLES IS FOUND EVEN IF IT REQUIRES TO EFFECTIVELY INSPECT MORE THAN MAXFEATURES FEATURES

MAXLEAFNODES INT OR NONE OPTIONAL DEFAULTNONE GROW TREES WITH MAXLEAFNODES IN BESTFIRST FASHION BEST NODES ARE DEFINED AS RELATIVE REDUCTION IN IMPURITY IF NONE THEN UNLIMITED NUMBER OF LEAF NODES

MINIMPURITYDECREASE FLOAT OPTIONAL DEFAULT0 A NODE WILL BE SPLIT IF THIS SPLIT INDUCES A DECREASE OF THE IMPURITY GREATER THAN OR EQUAL TO THIS VALUE

THE WEIGHTED IMPURITY DECREASE EQUATION IS THE FOLLOWING

$$N(T) \cdot N \cdot \text{IMPURITY} - N(T) \cdot \text{RIGHTIMPURITY} - N(TL) \cdot \text{LEFTIMPURITY}$$

WHERE N IS THE TOTAL NUMBER OF SAMPLES N(T) IS THE NUMBER OF SAMPLES AT THE CURRENT NODE N(TL) IS THE NUMBER OF SAMPLES IN THE LEFT CHILD AND N(T) - N(TL) IS THE NUMBER OF SAMPLES IN THE RIGHT CHILD

N(T)N(TL) AND N(TL) ALL REFER TO THE WEIGHTED SUM IF SAMPLEWEIGHT IS PASSED

NEW IN VERSION 019

MINIMPURITYSPLIT FLOAT DEFAULT1E7 THRESHOLD FOR EARLY STOPPING IN TREE GROWTH A NODE WILL SPLIT IF ITS IMPURITY IS ABOVE THE THRESHOLD OTHERWISE IT IS A LEAF

DEPRECATED SINCE VERSION 019 MINIMPURITYSPLIT HAS BEEN DEPRECATED IN FAVOR OF MINIMPURITYDECREASE IN 019 THE DEFAULT VALUE OF MINIMPURITYSPLIT WILL CHANGE FROM 1E7 TO 0 IN 023 AND IT WILL BE REMOVED IN 025 USE MINIMPURITYDECREASE INSTEAD

BOOTSTRAP BOOLEAN OPTIONAL DEFAULTFALSE WHETHER BOOTSTRAP SAMPLES ARE USED WHEN BUILDING TREES IF FALSE THE WHOLE DATASET IS USED TO BUILD EACH TREE

OOBSCORE BOOL OPTIONAL DEFAULTFALSE WHETHER TO USE OUTFBAG SAMPLES TO ESTIMATE THE GENERALIZATION ACCURACY

NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF JOBS TO RUN IN PARALLEL FOR BOTH FIT AND PREDICT NONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT

1MEANS USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RANDOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NRANDOM

VERBOSE INT OPTIONAL DEFAULT0 CONTROLS THE VERBOSITY WHEN FITTING AND PREDICTING

33 MODEL SELECTION AND EVALUATION 515

SCIKITLEARN USER GUIDE RELEASE 0213

WARMSTART BOOL OPTIONAL DEFAULTFALSE WHEN SET TO TRUE REUSE THE SOLUTION OF THE PREVIOUS CALL TO FIT AND ADD MORE ESTIMATORS TO THE ENSEMBLE OTHERWISE JUST FIT A WHOLE NEW FOREST SEE THE GLOSSARY

CLASSWEIGHT DICT LIST OF DICTS “BALANCED” “BALANCEDSUBSAMPLE” OR NONE OPTIONAL DEFAULTNONE WEIGHTS ASSOCIATED WITH CLASSES IN THE FORM CLASSLABEL WEIGHT

IF NOT GIVEN ALL CLASSES ARE SUPPOSED TO HAVE WEIGHT ONE FOR MULTIOUTPUT PROBLEMS A LIST OF DICTS CAN BE PROVIDED IN THE SAME ORDER AS THE COLUMNS OF Y

NOTE THAT FOR MULTIOUTPUT INCLUDING MULTILABEL WEIGHTS SHOULD BE DEFINED FOR EACH CLASS OF EVERY COLUMN IN ITS OWN DICT FOR EXAMPLE FOR FOURCLASS MULTILABEL CLASSIFICATION WEIGHTS SHOULD BE 0 1 1 1 0 1 1 5 0 1 1 1 0 1 1 1 INSTEAD OF 11 25

31 41

THE “BALANCED” MODE USES THE VALUES OF Y TO AUTOMATICALLY ADJUST WEIGHTS INVERSELY PROPORTIONAL TO CLASS FREQUENCIES IN THE INPUT DATA AS NSAMPLES NCLASSES NP

BINCOUNTY

THE “BALANCEDSUBSAMPLE” MODE IS THE SAME AS “BALANCED” EXCEPT THAT WEIGHTS ARE COMPUTED BASED ON THE BOOTSTRAP SAMPLE FOR EVERY TREE GROWN

FOR MULTIOUTPUT THE WEIGHTS OF EACH COLUMN OF Y WILL BE MULTIPLIED

NOTE THAT THESE WEIGHTS WILL BE MULTIPLIED WITH SAMPLEWEIGHT PASSED THROUGH THE FIT METHOD IF SAMPLEWEIGHT IS SPECIFIED

ATTRIBUTES

ESTIMATORS LIST OF DECISIONTREECLASSIFIER THE COLLECTION OF FITTED SUBESTIMATORS

CLASSES ARRAY OF SHAPE NCLASSES OR A LIST OF SUCH ARRAYS THE CLASSES LABELS SINGLE OUTPUT PROBLEM OR A LIST OF ARRAYS OF CLASS LABELS MULTIOUTPUT PROBLEM

NCLASSES INT OR LIST THE NUMBER OF CLASSES SINGLE OUTPUT PROBLEM OR A LIST CONTAINING THE NUMBER OF CLASSES FOR EACH OUTPUT MULTIOUTPUT PROBLEM

FEATUREIMPORTANCES ARRAY OF SHAPE NFEATURES RETURN THE FEATURE IMPORTANCES

THE HIGHER THE MORE IMPORTANT THE FEATURE

NFEATURES INT THE NUMBER OF FEATURES WHEN FIT IS PERFORMED

NOUTPUTS INT THE NUMBER OF OUTPUTS WHEN FIT IS PERFORMED

OOBSCORE FLOAT SCORE OF THE TRAINING DATASET OBTAINED USING AN OUTFBAG ESTIMATE

OOBDECISIONFUNCTION ARRAY OF SHAPE NSAMPLES NCLASSES DECISION FUNCTION COMPUTED WITH OUTFBAG ESTIMATE ON THE TRAINING SET IF NESTIMATORS IS SMALL IT MIGHT

BE POSSIBLE THAT A DATA POINT WAS NEVER LEFT OUT DURING THE BOOTSTRAP IN THIS CASE OOBDECISIONFUNCTION MIGHT CONTAIN NAN

SEE ALSO

SKLEARNTREEEXTRATREECLASSIFIER BASE CLASSIFIER FOR THIS ENSEMBLE

RANDOMFORESTCLASSIFIER ENSEMBLE CLASSIFIER BASED ON TREES WITH OPTIMAL SPLITS

NOTES

THE DEFAULT VALUES FOR THE PARAMETERS CONTROLLING THE SIZE OF THE TREES EG MAXDEPTH

MINSAMPLESLEAF ETC LEAD TO FULLY GROWN AND UNPRUNED TREES WHICH CAN POTENTIALLY BE VERY LARGE ON 516 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

SOME DATA SETS TO REDUCE MEMORY CONSUMPTION THE COMPLEXITY AND SIZE OF THE TREES SHOULD BE CONTROLLED BY SETTING THOSE PARAMETER VALUES

REFERENCES

RC8F28BFAD63F1

METHODS

APPLY SELF X APPLY TREES IN THE FOREST TO X RETURN LEAF INDICES

DECISIONPATH SELF X RETURN THE DECISION PATH IN THE FOREST

FITSELF X Y SAMPLEWEIGHT BUILD A FOREST OF TREES FROM THE TRAINING SET X Y

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PREDICT SELF X PREDICT CLASS FOR X

PREDICTLOGPROBA SELF X PREDICT CLASS LOGPROBABILITIES FOR X

PREDICTPROBA SELF X PREDICT CLASS PROBABILITIES FOR X

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELFNESTIMATORS'WARN' CRITERION'GINI' MAXDEPTHNONE MINSAMPLESSPLIT2

MINSAMPLESLEAF1 MINWEIGHTFRACTIONLEAF00 MAXFEATURES'AUTO'

MAXLEAFNODESNONE MINIMPURITYDECREASE00 MINIMPURITYSPLITNONE BOOT

STRAPFALSE OOBSCOREFALSE NJOBSNONE RANDOMSTATENONE VERBOSE0

WARMSTARTFALSE CLASSWEIGHTNONE

APPLYSELF X

APPLY TREES IN THE FOREST TO X RETURN LEAF INDICES

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTERNALLY ITS DTYPE WILL BE CONVERTED TO DTYPENPFLOAT32 IF A SPARSE MATRIX IS PROVIDED IT WILL BE CONVERTED INTO A SPARSE CSRMATRIX

RETURNS

XLEAVES ARRAYLIKE SHAPE NSAMPLES NESTIMATORS FOR EACH DATAPOINT X IN X AND FOR EACH TREE IN THE FOREST RETURN THE INDEX OF THE LEAF X ENDS UP IN

DECISIONPATH SELF X

RETURN THE DECISION PATH IN THE FOREST

NEW IN VERSION 018

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTERNALLY ITS DTYPE WILL BE CONVERTED TO DTYPENPFLOAT32 IF A SPARSE MATRIX IS PROVIDED IT WILL BE CONVERTED INTO A SPARSE CSRMATRIX

RETURNS

INDICATOR SPARSE CSR ARRAY SHAPE NSAMPLES NNODES RETURN A NODE INDICATOR MATRIX WHERE NON ZERO ELEMENTS INDICATES THAT THE SAMPLES GOES THROUGH THE NODES

33 MODEL SELECTION AND EVALUATION 517

SCIKITLEARN USER GUIDE RELEASE 0213

**NNODESPTR** ARRAY OF SIZE **NESTIMATORS** 1 THE COLUMNS FROM INDICA  
**TORNNODESPTR**IN**NNODESPTR**1 GIVES THE INDICATOR VALUE FOR THE ITH ESTIMATOR

**FEATUREIMPORTANCES**  
RETURN THE FEATURE IMPORTANCES THE HIGHER THE MORE IMPORTANT THE FEATURE  
RETURNS  
FEATUREIMPORTANCES ARRAY SHAPE **NFEATURES** THE VALUES OF THIS ARRAY SUM TO 1 UNLESS  
ALL TREES ARE SINGLE NODE TREES CONSISTING OF ONLY THE ROOT NODE IN WHICH CASE IT WILL BE AN  
ARRAY OF ZEROS

**FITSELFXY**SAMPLEWEIGHT**None**  
BUILD A FOREST OF TREES FROM THE TRAINING SET **X Y**

**PARAMETERS**  
**X**ARRAYLIKE OR SPARSE MATRIX OF SHAPE **NSAMPLES NFEATURES** THE TRAINING INPUT SAM  
PLES INTERNALLY ITS DTYPE WILL BE CONVERTED TO **DTYPENPFLOAT32** IF A SPARSE MATRIX  
IS PROVIDED IT WILL BE CONVERTED INTO A SPARSE **CSCMATRIX**

**Y**ARRAYLIKE SHAPE **NSAMPLES** OR **NSAMPLES** NOUTPUTS THE TARGET VALUES CLASS LABELS  
IN CLASSIFICATION REAL NUMBERS IN REGRESSION

**SAMPLEWEIGHT** ARRAYLIKE SHAPE **NSAMPLES** OR **None** SAMPLE WEIGHTS IF **None** THEN  
SAMPLES ARE EQUALLY WEIGHTED SPLITS THAT WOULD CREATE CHILD NODES WITH NET ZERO OR NEG  
ATIVE WEIGHT ARE IGNORED WHILE SEARCHING FOR A SPLIT IN EACH NODE IN THE CASE OF CLASSIFI  
CATION SPLITS ARE ALSO IGNORED IF THEY WOULD RESULT IN ANY SINGLE CLASS CARRYING A NEGATIVE  
WEIGHT IN EITHER CHILD NODE

RETURNS  
SELF OBJECT

**GETPARAMS** SELF**DEEP****True**  
GET PARAMETERS FOR THIS ESTIMATOR

**PARAMETERS**  
**DEEP** BOOLEAN OPTIONAL IF **True** WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBOBJECTS THAT ARE ESTIMATORS

RETURNS  
PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

**PREDICTSELF****X**  
PREDICT CLASS FOR **X**

THE PREDICTED CLASS OF AN INPUT SAMPLE IS A VOTE BY THE TREES IN THE FOREST WEIGHTED BY THEIR PROBABILITY  
ESTIMATES THAT IS THE PREDICTED CLASS IS THE ONE WITH HIGHEST MEAN PROBABILITY ESTIMATE ACROSS THE TREES

**PARAMETERS**  
**X**ARRAYLIKE OR SPARSE MATRIX OF SHAPE **NSAMPLES NFEATURES** THE INPUT SAMPLES INTER  
NALLY ITS DTYPE WILL BE CONVERTED TO **DTYPENPFLOAT32** IF A SPARSE MATRIX IS PROVIDED  
IT WILL BE CONVERTED INTO A SPARSE **CSRMATRIX**

RETURNS  
**Y**ARRAY OF SHAPE **NSAMPLES** OR **NSAMPLES** NOUTPUTS THE PREDICTED CLASSES

518 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

PREDICTLOGPROBA SELF  
PREDICT CLASS LOGPROBABILITIES FOR X  
THE PREDICTED CLASS LOGPROBABILITIES OF AN INPUT SAMPLE IS COMPUTED AS THE LOG OF THE MEAN PREDICTED CLASS PROBABILITIES OF THE TREES IN THE FOREST

PARAMETERS  
XARRAYLIKE OR SPARSE MATRIX OF SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTERNALLY ITS DTYPE WILL BE CONVERTED TO DYPENPFLOAT32 IF A SPARSE MATRIX IS PROVIDED IT WILL BE CONVERTED INTO A SPARSE CSRMATRIX

RETURNS  
PARRAY OF SHAPE NSAMPLES NCLASSES OR A LIST OF NOUTPUTS SUCH ARRAYS IF NOUTPUTS 1 THE CLASS PROBABILITIES OF THE INPUT SAMPLES THE ORDER OF THE CLASSES CORRESPONDS TO THAT IN THE ATTRIBUTE CLASSES

PREDICTPROBA SELF  
PREDICT CLASS PROBABILITIES FOR X  
THE PREDICTED CLASS PROBABILITIES OF AN INPUT SAMPLE ARE COMPUTED AS THE MEAN PREDICTED CLASS PROBABILITIES OF THE TREES IN THE FOREST THE CLASS PROBABILITY OF A SINGLE TREE IS THE FRACTION OF SAMPLES OF THE SAME CLASS IN A LEAF

PARAMETERS  
XARRAYLIKE OR SPARSE MATRIX OF SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTERNALLY ITS DTYPE WILL BE CONVERTED TO DYPENPFLOAT32 IF A SPARSE MATRIX IS PROVIDED IT WILL BE CONVERTED INTO A SPARSE CSRMATRIX

RETURNS  
PARRAY OF SHAPE NSAMPLES NCLASSES OR A LIST OF NOUTPUTS SUCH ARRAYS IF NOUTPUTS 1 THE CLASS PROBABILITIES OF THE INPUT SAMPLES THE ORDER OF THE CLASSES CORRESPONDS TO THAT IN THE ATTRIBUTE CLASSES

SCORESELFXYSAMPLEWEIGHTNONE  
RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS  
IN MULTILABEL CLASSIFICATION THIS IS THE SUBSET ACCURACY WHICH IS A HARSH METRIC SINCE YOU REQUIRE FOR EACH SAMPLE THAT EACH LABEL SET BE CORRECTLY PREDICTED

PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES  
YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE LABELS FOR X  
SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS  
SCORE FLOAT MEAN ACCURACY OF SELF  
SETPARAMS SELF  
SET THE PARAMETERS OF THIS ESTIMATOR  
THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS  
SELF

33 MODEL SELECTION AND EVALUATION 519

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES USING SKLEARNENSEMBLEEXTRATREESCLASSIFIER

- PIXEL IMPORTANCES WITH A PARALLEL FOREST OF TREES
- FEATURE IMPORTANCES WITH FORESTS OF TREES
- HASHING FEATURE TRANSFORMATION USING TOTALLY RANDOM TREES
- PLOT THE DECISION SURFACES OF ENSEMBLES OF TREES ON THE IRIS DATASET

SKLEARNENSEMBLE EXTRATREESREGRESSOR

CLASS SKLEARNENSEMBLE EXTRATREESREGRESSOR NESTIMATORS'WARN' CRITERION'MSE' MAXDEPTHNONE

MINSAMPLESSPLIT2 MINSAMPLESLEAF1

MINWEIGHTFRACTIONLEAF00

MAXFEATURES'AUTO' MAXLEAFNODESNONE

MINIMPURITYDECREASE00

MINIMPURITYSPLITNONE BOOTSTRAPFALSE

OOBSCOREFALSE NJOBSNONE

RANDOMSTATENONE VERBOSE0

WARMSTARTFALSE

AN EXTRATREES REGRESSOR

THIS CLASS IMPLEMENTS A META ESTIMATOR THAT FITS A NUMBER OF RANDOMIZED DECISION TREES AKA EXTRATREES ON VARIOUS SUBSAMPLES OF THE DATASET AND USES AVERAGING TO IMPROVE THE PREDICTIVE ACCURACY AND CONTROL OVERFITTING

READ MORE IN THE USER GUIDE

PARAMETERS

NESTIMATORS INTEGER OPTIONAL DEFAULT10 THE NUMBER OF TREES IN THE FOREST

CHANGED IN VERSION 020 THE DEFAULT VALUE OF NESTIMATORS WILL CHANGE FROM 10 IN VERSION 020 TO 100 IN VERSION 022

CRITERION STRING OPTIONAL DEFAULT"MSE" THE FUNCTION TO MEASURE THE QUALITY OF A SPLIT SUPPORTED CRITERIA ARE "MSE" FOR THE MEAN SQUARED ERROR WHICH IS EQUAL TO VARIANCE REDUCTION AS FEATURE SELECTION CRITERION AND "MAE" FOR THE MEAN ABSOLUTE ERROR

NEW IN VERSION 018 MEAN ABSOLUTE ERROR MAE CRITERION

MAXDEPTH INTEGER OR NONE OPTIONAL DEFAULTNONE THE MAXIMUM DEPTH OF THE TREE IF NONE THEN NODES ARE EXPANDED UNTIL ALL LEAVES ARE PURE OR UNTIL ALL LEAVES CONTAIN LESS THAN MINSAMPLESSPLIT SAMPLES

MINSAMPLESSPLIT INT FLOAT OPTIONAL DEFAULT2 THE MINIMUM NUMBER OF SAMPLES REQUIRED TO SPLIT AN INTERNAL NODE

- IF INT THEN CONSIDER MINSAMPLESSPLIT AS THE MINIMUM NUMBER
- IF FLOAT THEN MINSAMPLESSPLIT IS A FRACTION AND CEILMINSAMPLESSPLIT

NSAMPLES ARE THE MINIMUM NUMBER OF SAMPLES FOR EACH SPLIT

CHANGED IN VERSION 018 ADDED FLOAT VALUES FOR FRACTIONS

MINSAMPLESLEAF INT FLOAT OPTIONAL DEFAULT1 THE MINIMUM NUMBER OF SAMPLES REQUIRED TO BE AT A LEAF NODE A SPLIT POINT AT ANY DEPTH WILL ONLY BE CONSIDERED IF IT LEAVES AT LEAST MINSAMPLESLEAF TRAINING SAMPLES IN EACH OF THE LEFT AND RIGHT BRANCHES THIS MAY HAVE THE EFFECT OF SMOOTHING THE MODEL ESPECIALLY IN REGRESSION

520 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

- IF INT THEN CONSIDER MINSAMPLESLEAF AS THE MINIMUM NUMBER
  - IF FLOAT THEN MINSAMPLESLEAF IS A FRACTION AND CEILMINSAMPLESLEAF NSAMPLES ARE THE MINIMUM NUMBER OF SAMPLES FOR EACH NODE
- CHANGED IN VERSION 018 ADDED FLOAT VALUES FOR FRACTIONS
- MINWEIGHTFRACTIONLEAF FLOAT OPTIONAL DEFAULT0 THE MINIMUM WEIGHTED FRACTION OF THE SUM TOTAL OF WEIGHTS OF ALL THE INPUT SAMPLES REQUIRED TO BE AT A LEAF NODE SAMPLES HAVE EQUAL WEIGHT WHEN SAMPLEWEIGHT IS NOT PROVIDED
- MAXFEATURES INT FLOAT STRING OR NONE OPTIONAL DEFAULT" AUTO" THE NUMBER OF FEATURES TO CONSIDER WHEN LOOKING FOR THE BEST SPLIT
- IF INT THEN CONSIDER MAXFEATURES FEATURES AT EACH SPLIT
  - IF FLOAT THEN MAXFEATURES IS A FRACTION AND INTMAXFEATURES NFEATURES FEATURES ARE CONSIDERED AT EACH SPLIT
  - IF "AUTO" THEN MAXFEATURESNFEATURES
  - IF "SQRT" THEN MAXFEATURESSQRTNFEATURES
  - IF "LOG2" THEN MAXFEATURESLOG2NFEATURES
  - IF NONE THEN MAXFEATURESNFEATURES

NOTE THE SEARCH FOR A SPLIT DOES NOT STOP UNTIL AT LEAST ONE VALID PARTITION OF THE NODE SAMPLES IS FOUND EVEN IF IT REQUIRES TO EFFECTIVELY INSPECT MORE THAN MAXFEATURES FEATURES

MAXLEAFNODES INT OR NONE OPTIONAL DEFAULTNONE GROW TREES WITH MAXLEAFNODES IN BESTFIRST FASHION BEST NODES ARE DEFINED AS RELATIVE REDUCTION IN IMPURITY IF NONE THEN UNLIMITED NUMBER OF LEAF NODES

MINIMPURITYDECREASE FLOAT OPTIONAL DEFAULT0 A NODE WILL BE SPLIT IF THIS SPLIT INDUCES A DECREASE OF THE IMPURITY GREATER THAN OR EQUAL TO THIS VALUE

THE WEIGHTED IMPURITY DECREASE EQUATION IS THE FOLLOWING

$$NT \cdot N \cdot IMPURITY - NTR \cdot NT \cdot RIGHTIMPURITY - NTL \cdot NT \cdot LEFTIMPURITY$$

WHERE N IS THE TOTAL NUMBER OF SAMPLES NT IS THE NUMBER OF SAMPLES AT THE CURRENT NODE NTL IS THE NUMBER OF SAMPLES IN THE LEFT CHILD AND NTR IS THE NUMBER OF SAMPLES IN THE RIGHT CHILD

NNTNTR ANDNTL ALL REFER TO THE WEIGHTED SUM IF SAMPLEWEIGHT IS PASSED

NEW IN VERSION 019

MINIMPURITYSPLIT FLOAT DEFAULT1E7 THRESHOLD FOR EARLY STOPPING IN TREE GROWTH A NODE WILL SPLIT IF ITS IMPURITY IS ABOVE THE THRESHOLD OTHERWISE IT IS A LEAF

DEPRECATED SINCE VERSION 019 MINIMPURITYSPLIT HAS BEEN DEPRECATED IN FAVOR OF MINIMPURITYDECREASE IN 019 THE DEFAULT VALUE OF MINIMPURITYSPLIT WILL CHANGE FROM 1E7 TO 0 IN 023 AND IT WILL BE REMOVED IN 025 USE MINIMPURITYDECREASE INSTEAD

BOOTSTRAP BOOLEAN OPTIONAL DEFAULTFALSE WHETHER BOOTSTRAP SAMPLES ARE USED WHEN BUILDING TREES IF FALSE THE WHOLE DATSET IS USED TO BUILD EACH TREE

OOBSCORE BOOL OPTIONAL DEFAULTFALSE WHETHER TO USE OUTFBAG SAMPLES TO ESTIMATE THE R2 ON UNSEEN DATA

SCIKITLEARN USER GUIDE RELEASE 0213

NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF JOBS TO RUN IN PARALLEL FOR BOTH  
FIT ANDPREDICT NONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT  
1MEANS USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS  
RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RAN  
DOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE  
RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS  
THE RANDOMSTATE INSTANCE USED BY NPRANDOM  
VERBOSE INT OPTIONAL DEFAULT0 CONTROLS THE VERBOSITY WHEN FITTING AND PREDICTING  
WARMSTART BOOL OPTIONAL DEFAULTFALSE WHEN SET TO TRUE REUSE THE SOLUTION OF THE PRE  
VIOUS CALL TO FIT AND ADD MORE ESTIMATORS TO THE ENSEMBLE OTHERWISE JUST FIT A WHOLE NEW  
FOREST SEE THE GLOSSARY

ATTRIBUTES  
ESTIMATORS LIST OF DECISIONTREEREgressor THE COLLECTION OF FITTED SUBESTIMATORS  
FEATUREIMPORTANCES ARRAY OF SHAPE NFEATURES RETURN THE FEATURE IMPORTANCES  
THE HIGHER THE MORE IMPORTANT THE FEATURE  
NFEATURES INT THE NUMBER OF FEATURES  
NOUTPUTS INT THE NUMBER OF OUTPUTS  
OOBSCORE FLOAT SCORE OF THE TRAINING DATASET OBTAINED USING AN OUTOFBAG ESTIMATE  
OOBPREDICTION ARRAY OF SHAPE NSAMPLES PREDICTION COMPUTED WITH OUTOFBAG ESTIMATE  
ON THE TRAINING SET

SEE ALSO  
SKLEARNNTREEEXTRATREEREgressor BASE ESTIMATOR FOR THIS ENSEMBLE  
RANDOMFORESTRegressor ENSEMBLE REGRESSOR USING TREES WITH OPTIMAL SPLITS  
NOTES  
THE DEFAULT VALUES FOR THE PARAMETERS CONTROLLING THE SIZE OF THE TREES EG MAXDEPTH  
MINSAMPLESLEAF ETC LEAD TO FULLY GROWN AND UNPRUNED TREES WHICH CAN POTENTIALLY BE VERY LARGE ON  
SOME DATA SETS TO REDUCE MEMORY CONSUMPTION THE COMPLEXITY AND SIZE OF THE TREES SHOULD BE CONTROLLED BY  
SETTING THOSE PARAMETER VALUES

REFERENCES  
RA7D0C8995FBC1  
METHODS  
APPLY SELF X APPLY TREES IN THE FOREST TO X RETURN LEAF INDICES  
DECISIONPATH SELF X RETURN THE DECISION PATH IN THE FOREST  
FITSELF X Y SAMPLEWEIGHT BUILD A FOREST OF TREES FROM THE TRAINING SET X Y  
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR  
PREDICT SELF X PREDICT REGRESSION TARGET FOR X

CONTINUED ON NEXT PAGE  
522 CHAPTER 3 USER GUIDE



SCIKITLEARN USER GUIDE RELEASE 0213

TABLE 318 – CONTINUED FROM PREVIOUS PAGE

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELFNESTIMATORS'WARN' CRITERION'MSE' MAXDEPTHNONE MINSAMPLESSPLIT2

MINSAMPLESLEAF1 MINWEIGHTFRACTIONLEAF00 MAXFEATURES'AUTO'

MAXLEAFNODESNONE MINIMPURITYDECREASE00 MINIMPURITYSPLITNONE BOOT

STRAPFALSE OOBSCOREFALSE NJOBSNONE RANDOMSTATENONE VERBOSE0

WARMSTARTFALSE

APPLYSELF X

APPLY TREES IN THE FOREST TO X RETURN LEAF INDICES

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTERNALLY ITS DTYPE WILL BE CONVERTED TO DTYPENPFLOAT32 IF A SPARSE MATRIX IS PROVIDED IT WILL BE CONVERTED INTO A SPARSE CSRMATRIX

RETURNS

XLEAVES ARRAYLIKE SHAPE NSAMPLES NESTIMATORS FOR EACH DATAPOINT X IN X AND FOR EACH TREE IN THE FOREST RETURN THE INDEX OF THE LEAF X ENDS UP IN

DECISIONPATH SELF X

RETURN THE DECISION PATH IN THE FOREST

NEW IN VERSION 018

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTERNALLY ITS DTYPE WILL BE CONVERTED TO DTYPENPFLOAT32 IF A SPARSE MATRIX IS PROVIDED IT WILL BE CONVERTED INTO A SPARSE CSRMATRIX

RETURNS

INDICATOR SPARSE CSR ARRAY SHAPE NSAMPLES NNODES RETURN A NODE INDICATOR MATRIX WHERE NON ZERO ELEMENTS INDICATES THAT THE SAMPLES GOES THROUGH THE NODES

NNODESPTR ARRAY OF SIZE NESTIMATORS 1 THE COLUMNS FROM INDICATOR

TORNNODESPTRINNNODESPTRI1 GIVES THE INDICATOR VALUE FOR THE ITH ESTIMATOR

FEATUREIMPORTANCES

RETURN THE FEATURE IMPORTANCES THE HIGHER THE MORE IMPORTANT THE FEATURE

RETURNS

FEATUREIMPORTANCES ARRAY SHAPE NFEATURES THE VALUES OF THIS ARRAY SUM TO 1 UNLESS ALL TREES ARE SINGLE NODE TREES CONSISTING OF ONLY THE ROOT NODE IN WHICH CASE IT WILL BE AN ARRAY OF ZEROS

FITSELF X Y SAMPLEWEIGHTNONE

BUILD A FOREST OF TREES FROM THE TRAINING SET X Y

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX OF SHAPE NSAMPLES NFEATURES THE TRAINING INPUT SAMPLES INTERNALLY ITS DTYPE WILL BE CONVERTED TO DTYPENPFLOAT32 IF A SPARSE MATRIX IS PROVIDED IT WILL BE CONVERTED INTO A SPARSE CSRMATRIX

33 MODEL SELECTION AND EVALUATION 523

SCIKITLEARN USER GUIDE RELEASE 0213

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS THE TARGET VALUES CLASS LABELS  
IN CLASSIFICATION REAL NUMBERS IN REGRESSION

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OR NONE SAMPLE WEIGHTS IF NONE THEN  
SAMPLES ARE EQUALLY WEIGHTED SPLITS THAT WOULD CREATE CHILD NODES WITH NET ZERO OR NEG  
ATIVE WEIGHT ARE IGNORED WHILE SEARCHING FOR A SPLIT IN EACH NODE IN THE CASE OF CLASSIFI  
CATION SPLITS ARE ALSO IGNORED IF THEY WOULD RESULT IN ANY SINGLE CLASS CARRYING A NEGATIVE  
WEIGHT IN EITHER CHILD NODE

RETURNS  
SELF OBJECT

GETPARAMS SELFDEEPTURE  
GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS  
DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBOBJECTS THAT ARE ESTIMATORS

RETURNS  
PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELF  
PREDICT REGRESSION TARGET FOR X  
THE PREDICTED REGRESSION TARGET OF AN INPUT SAMPLE IS COMPUTED AS THE MEAN PREDICTED REGRESSION TARGETS OF  
THE TREES IN THE FOREST

PARAMETERS  
XARRAYLIKE OR SPARSE MATRIX OF SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTER  
NALLY ITS DTYPE WILL BE CONVERTED TO DTYPENPFLOAT32 IF A SPARSE MATRIX IS PROVIDED  
IT WILL BE CONVERTED INTO A SPARSE CSRMATRIX

RETURNS  
YARRAY OF SHAPE NSAMPLES OR NSAMPLES NOUTPUTS THE PREDICTED VALUES

SCORESELFXYSAMPLEWEIGHTNONE  
RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION  
THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{U}{V}$  WHERE U IS THE RESIDUAL SUM OF SQUARES  $Y_{TRUE} - Y_{PRED}$   
 $2 \sum$  AND V IS THE TOTAL SUM OF SQUARES  $Y_{TRUE} - Y_{TRUEMEAN}$   $2 \sum$  THE BEST POSSIBLE SCORE  
IS 10 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS  
PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 00

PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES FOR SOME ESTIMATORS THIS MAY  
BE A PRECOMPUTED KERNEL MATRIX INSTEAD SHAPE NSAMPLES NSAMPLESFITTED WHERE  
NSAMPLESFITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE VALUES FOR X  
SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS  
SCORE FLOAT R2 OF SELF PREDICTX WRT Y

524 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

NOTES

THE R2 SCORE USED WHEN CALLING SCORE ON A REGRESSOR WILL USE MULTIOUTPUTUNIFORMAVERAGE FROM VERSION 023 TO KEEP CONSISTENT WITH METRICSR2SCORE THIS WILL INFLUENCE THE SCORE METHOD OF ALL THE MULTIOUTPUT REGRESSORS EXCEPT FOR MULTIOUTPUTMULTIOUTPUTREGRESSOR TO SPECIFY THE DEFAULT VALUE MANUALLY AND AVOID THE WARNING PLEASE EITHER CALL METRICSR2SCORE DIRECTLY OR MAKE A CUSTOM SCORER WITH METRICSMAKESCORER THE BUILTIN SCORER R2 USES MULTIOUTPUTUNIFORMAVERAGE  
SETPARAMS SELFPARAMS  
SET THE PARAMETERS OF THIS ESTIMATOR  
THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT  
RETURNS

SELF  
EXAMPLES USING SKLEARNENSEMBLEEXTRATREESREGRESSOR

- FACE COMPLETION WITH A MULTIOUTPUT ESTIMATORS
- IMPUTING MISSING VALUES WITH VARIANTS OF ITERATIVEIMPUTER

SKLEARNENSEMBLE GRADIENTBOOSTINGCLASSIFIER  
CLASSSSKLEARNENSEMBLE GRADIENTBOOSTINGCLASSIFIER LOSS'DEVIANCE' LEARNINGRATE01

NESTIMATORS100 SUBSAM  
PLE10 CRITERION'FRIEDMANMSE'

MINSAMPLESSPLIT2  
MINSAMPLESLEAF1  
MINWEIGHTFRACTIONLEAF00

MAXDEPTH3  
MINIMPURITYDECREASE00  
MINIMPURITYSPLITNONE  
INITNONE RANDOMSTATENONE

MAXFEATURESNONE VER  
BOSE0 MAXLEAFNODESNONE  
WARMSTARTFALSE PRE  
SORT'AUTO' VALIDATIONFRACTION01  
NITERNOCHANGENONE TOL00001

GRADIENT BOOSTING FOR CLASSIFICATION  
GB BUILDS AN ADDITIVE MODEL IN A FORWARD STAGewise FASHION IT ALLOWS FOR THE OPTIMIZATION OF ARBITRARY DIFFEREN  
TIABLE LOSS FUNCTIONS IN EACH STAGE NCLASSES REGRESSION TREES ARE FIT ON THE NEGATIVE GRADIENT OF THE BINOMIAL  
OR MULTINOMIAL DEVIANCE LOSS FUNCTION BINARY CLASSIFICATION IS A SPECIAL CASE WHERE ONLY A SINGLE REGRESSION TREE  
IS INDUCED

READ MORE IN THE USER GUIDE  
PARAMETERS  
33 MODEL SELECTION AND EVALUATION 525

SCIKITLEARN USER GUIDE RELEASE 0213

LOSS ‘DEVIANCE’ ‘EXPONENTIAL’ OPTIONAL DEFAULT ‘DEVIANCE’ LOSS FUNCTION TO BE OPTIMIZED  
‘DEVIANCE’ REFERS TO DEVIANCE LOGISTIC REGRESSION FOR CLASSIFICATION WITH PROBABILISTIC OUT  
PUTS FOR LOSS ‘EXPONENTIAL’ GRADIENT BOOSTING RECOVERS THE ADABOOST ALGORITHM  
LEARNINGRATE FLOAT OPTIONAL DEFAULT 0.1 LEARNING RATE SHRINKS THE CONTRIBUTION OF EACH TREE  
BY LEARNINGRATE THERE IS A TRADEOFF BETWEEN LEARNINGRATE AND NESTIMATORS  
NESTIMATORS INT DEFAULT 100 THE NUMBER OF BOOSTING STAGES TO PERFORM GRADIENT BOOSTING  
IS FAIRLY ROBUST TO OVERFITTING SO A LARGE NUMBER USUALLY RESULTS IN BETTER PERFORMANCE  
SUBSAMPLE FLOAT OPTIONAL DEFAULT 0.10 THE FRACTION OF SAMPLES TO BE USED FOR FITTING THE IN  
DIVIDUAL BASE LEARNERS IF SMALLER THAN 0.10 THIS RESULTS IN STOCHASTIC GRADIENT BOOSTING  
SUBSAMPLE INTERACTS WITH THE PARAMETER NESTIMATORS CHOOSING SUBSAMPLE  
0.10 LEADS TO A REDUCTION OF VARIANCE AND AN INCREASE IN BIAS  
CRITERION STRING OPTIONAL DEFAULT “FRIEDMANMSE” THE FUNCTION TO MEASURE THE QUALITY OF  
A SPLIT SUPPORTED CRITERIA ARE “FRIEDMANMSE” FOR THE MEAN SQUARED ERROR WITH IMPROVE  
MENT SCORE BY FRIEDMAN “MSE” FOR MEAN SQUARED ERROR AND “MAE” FOR THE MEAN ABSOLUTE  
ERROR THE DEFAULT VALUE OF “FRIEDMANMSE” IS GENERALLY THE BEST AS IT CAN PROVIDE A BETTER  
APPROXIMATION IN SOME CASES

NEW IN VERSION 0.18

MINSAMPLESSPLIT INT FLOAT OPTIONAL DEFAULT 2 THE MINIMUM NUMBER OF SAMPLES REQUIRED  
TO SPLIT AN INTERNAL NODE

- IF INT THEN CONSIDER MINSAMPLESSPLIT AS THE MINIMUM NUMBER
- IF FLOAT THEN MINSAMPLESSPLIT IS A FRACTION AND CEIL(MINSAMPLESSPLIT

NSAMPLES ARE THE MINIMUM NUMBER OF SAMPLES FOR EACH SPLIT  
CHANGED IN VERSION 0.18 ADDED FLOAT VALUES FOR FRACTIONS

MINSAMPLESLEAF INT FLOAT OPTIONAL DEFAULT 1 THE MINIMUM NUMBER OF SAMPLES REQUIRED  
TO BE AT A LEAF NODE A SPLIT POINT AT ANY DEPTH WILL ONLY BE CONSIDERED IF IT LEAVES AT LEAST  
MINSAMPLESLEAF TRAINING SAMPLES IN EACH OF THE LEFT AND RIGHT BRANCHES THIS MAY  
HAVE THE EFFECT OF SMOOTHING THE MODEL ESPECIALLY IN REGRESSION

- IF INT THEN CONSIDER MINSAMPLESLEAF AS THE MINIMUM NUMBER
- IF FLOAT THEN MINSAMPLESLEAF IS A FRACTION AND CEIL(MINSAMPLESLEAF

NSAMPLES ARE THE MINIMUM NUMBER OF SAMPLES FOR EACH NODE  
CHANGED IN VERSION 0.18 ADDED FLOAT VALUES FOR FRACTIONS

MINWEIGHTFRACTIONLEAF FLOAT OPTIONAL DEFAULT 0 THE MINIMUM WEIGHTED FRACTION OF THE  
SUM TOTAL OF WEIGHTS OF ALL THE INPUT SAMPLES REQUIRED TO BE AT A LEAF NODE SAMPLES HAVE  
EQUAL WEIGHT WHEN SAMPLEWEIGHT IS NOT PROVIDED

MAXDEPTH INTEGER OPTIONAL DEFAULT 3 MAXIMUM DEPTH OF THE INDIVIDUAL REGRESSION ESTIMA  
TORS THE MAXIMUM DEPTH LIMITS THE NUMBER OF NODES IN THE TREE TUNE THIS PARAMETER FOR  
BEST PERFORMANCE THE BEST VALUE DEPENDS ON THE INTERACTION OF THE INPUT VARIABLES

MINIMPURITYDECREASE FLOAT OPTIONAL DEFAULT 0 A NODE WILL BE SPLIT IF THIS SPLIT INDUCES  
A DECREASE OF THE IMPURITY GREATER THAN OR EQUAL TO THIS VALUE

THE WEIGHTED IMPURITY DECREASE EQUATION IS THE FOLLOWING

$$N \cdot I - N_L \cdot I_L - N_R \cdot I_R$$

NTL NT LEFTIMPURITY

SCIKITLEARN USER GUIDE RELEASE 0213

WHERE N IS THE TOTAL NUMBER OF SAMPLES NTL IS THE NUMBER OF SAMPLES AT THE CURRENT NODE  
NTL IS THE NUMBER OF SAMPLES IN THE LEFT CHILD AND NTR IS THE NUMBER OF SAMPLES IN  
THE RIGHT CHILD

NNTNTR AND NTL ALL REFER TO THE WEIGHTED SUM IF SAMPLEWEIGHT IS PASSED

NEW IN VERSION 019

MINIMPURITYSPLIT FLOAT DEFAULT 1E7 THRESHOLD FOR EARLY STOPPING IN TREE GROWTH A NODE  
WILL SPLIT IF ITS IMPURITY IS ABOVE THE THRESHOLD OTHERWISE IT IS A LEAF

DEPRECATED SINCE VERSION 019 MINIMPURITYSPLIT HAS BEEN DEPRECATED IN FAVOR OF

MINIMPURITYDECREASE IN 019 THE DEFAULT VALUE OF MINIMPURITYSPLIT  
WILL CHANGE FROM 1E7 TO 0 IN 023 AND IT WILL BE REMOVED IN 025 USE

MINIMPURITYDECREASE INSTEAD

INIT ESTIMATOR OR 'ZERO' OPTIONAL DEFAULT NONE AN ESTIMATOR OBJECT THAT IS USED TO COMPUTE  
THE INITIAL PREDICTIONS INIT HAS TO PROVIDE FIT AND PREDICT PROBABILITIES IF 'ZERO' THE  
INITIAL RAW PREDICTIONS ARE SET TO ZERO BY DEFAULT A DUMMY ESTIMATOR PREDICTING THE  
CLASSES PRIORS IS USED

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULT NONE IF INT RANDOM  
NUMBER GENERATOR IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE  
RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS  
THE RANDOMSTATE INSTANCE USED BY NRANDOM

MAXFEATURES INT FLOAT STRING OR NONE OPTIONAL DEFAULT NONE THE NUMBER OF FEATURES TO  
CONSIDER WHEN LOOKING FOR THE BEST SPLIT

- IF INT THEN CONSIDER MAXFEATURES FEATURES AT EACH SPLIT
- IF FLOAT THEN MAXFEATURES IS A FRACTION AND INTMAXFEATURES  
NFEATURES FEATURES ARE CONSIDERED AT EACH SPLIT
- IF "AUTO" THEN MAXFEATURES SQRT NFEATURES
- IF "SQRT" THEN MAXFEATURES SQRT NFEATURES
- IF "LOG2" THEN MAXFEATURES LOG2 NFEATURES
- IF NONE THEN MAXFEATURES NFEATURES

CHOOSING MAXFEATURES NFEATURES LEADS TO A REDUCTION OF VARIANCE AND AN INCREASE  
IN BIAS

NOTE THE SEARCH FOR A SPLIT DOES NOT STOP UNTIL AT LEAST ONE VALID PARTITION OF THE NODE SAMPLES  
IS FOUND EVEN IF IT REQUIRES TO EFFECTIVELY INSPECT MORE THAN MAXFEATURES FEATURES

VERBOSE INT DEFAULT 0 ENABLE VERBOSE OUTPUT IF 1 THEN IT PRINTS PROGRESS AND PERFORMANCE  
ONCE IN A WHILE THE MORE TREES THE LOWER THE FREQUENCY IF GREATER THAN 1 THEN IT PRINTS  
PROGRESS AND PERFORMANCE FOR EVERY TREE

MAXLEAFNODES INT OR NONE OPTIONAL DEFAULT NONE GROW TREES WITH MAXLEAFNODES

IN BESTFIRST FASHION BEST NODES ARE DEFINED AS RELATIVE REDUCTION IN IMPURITY IF NONE THEN  
UNLIMITED NUMBER OF LEAF NODES

WARMSTART BOOL DEFAULT FALSE WHEN SET TO TRUE REUSE THE SOLUTION OF THE PREVIOUS CALL TO  
FIT AND ADD MORE ESTIMATORS TO THE ENSEMBLE OTHERWISE JUST ERASE THE PREVIOUS SOLUTION SEE  
THE GLOSSARY

PRESORT BOOL OR 'AUTO' OPTIONAL DEFAULT 'AUTO' WHETHER TO PRESORT THE DATA TO SPEED UP THE  
FINDING OF BEST SPLITS IN FITTING AUTO MODE BY DEFAULT WILL USE PRESORTING ON DENSE DATA AND

SCIKITLEARN USER GUIDE RELEASE 0213

DEFAULT TO NORMAL SORTING ON SPARSE DATA SETTING PRESORT TO TRUE ON SPARSE DATA WILL RAISE AN ERROR

NEW IN VERSION 017 PRESORT PARAMETER

VALIDATIONFRACTION FLOAT OPTIONAL DEFAULT 01 THE PROPORTION OF TRAINING DATA TO SET ASIDE AS VALIDATION SET FOR EARLY STOPPING MUST BE BETWEEN 0 AND 1 ONLY USED IF NITERNOCHANGE IS SET TO AN INTEGER

NEW IN VERSION 020

NITERNOCHANGE INT DEFAULT NONE NITERNOCHANGE IS USED TO DECIDE IF EARLY STOPPING WILL BE USED TO TERMINATE TRAINING WHEN VALIDATION SCORE IS NOT IMPROVING BY DEFAULT IT IS SET TO NONE TO DISABLE EARLY STOPPING IF SET TO A NUMBER IT WILL SET ASIDE VALIDATIONFRACTION SIZE OF THE TRAINING DATA AS VALIDATION AND TERMINATE TRAINING WHEN VALIDATION SCORE IS NOT IMPROVING IN ALL OF THE PREVIOUS NITERNOCHANGE NUMBERS OF ITERATIONS THE SPLIT IS STRATIFIED

NEW IN VERSION 020

TOLFLOAT OPTIONAL DEFAULT 1E4 TOLERANCE FOR THE EARLY STOPPING WHEN THE LOSS IS NOT IMPROVING BY AT LEAST TOL FOR NITERNOCHANGE ITERATIONS IF SET TO A NUMBER THE TRAINING STOPS

NEW IN VERSION 020

ATTRIBUTES

NESTIMATORS INT THE NUMBER OF ESTIMATORS AS SELECTED BY EARLY STOPPING IF NITERNOCHANGE IS SPECIFIED OTHERWISE IT IS SET TO NESTIMATORS

NEW IN VERSION 020

FEATUREIMPORTANCES ARRAY SHAPE NFEATURES RETURN THE FEATURE IMPORTANCES THE HIGHER THE MORE IMPORTANT THE FEATURE

OOBIMPROVEMENT ARRAY SHAPE NESTIMATORS THE IMPROVEMENT IN LOSS DEVIANCE ON THE OUTFBAG SAMPLES RELATIVE TO THE PREVIOUS ITERATION OOBIMPROVEMENT0 IS THE IMPROVEMENT IN LOSS OF THE FIRST STAGE OVER THE INIT ESTIMATOR

TRAINSORE ARRAY SHAPE NESTIMATORS THE ITH SCORE TRAINSCOREI IS THE DEVIANCE LOSS OF THE MODEL AT ITERATION ION THE INBAG SAMPLE IF SUBSAMPLE 1 THIS IS THE DEVIANCE ON THE TRAINING DATA

LOSS LOSSFUNCTION THE CONCRETE LOSSFUNCTION OBJECT

INIT ESTIMATOR THE ESTIMATOR THAT PROVIDES THE INITIAL PREDICTIONS SET VIA THE INIT ARGUMENT ORLOSSINITESTIMATOR

ESTIMATORS NDARRAY OF DECISIONTREEREGRESSORSHAPE NESTIMATORS LOSSK THE COLLECTION OF FITTED SUBESTIMATORS LOSSK IS 1 FOR BINARY CLASSIFICATION OTHERWISE NCLASSES

SEE ALSO

SKLEARNENSEMBLEHISTGRADIENTBOOSTINGCLASSIFIER

SKLEARNNTREEDECISIONTREECLASSIFIER RANDOMFORESTCLASSIFIER

ADABOOSTCLASSIFIER

528 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

NOTES

THE FEATURES ARE ALWAYS RANDOMLY PERMUTED AT EACH SPLIT THEREFORE THE BEST FOUND SPLIT MAY VARY EVEN WITH THE SAME TRAINING DATA AND MAXFEATURESNFEATURES IF THE IMPROVEMENT OF THE CRITERION IS IDENTICAL FOR SEVERAL SPLITS ENUMERATED DURING THE SEARCH OF THE BEST SPLIT TO OBTAIN A DETERMINISTIC BEHAVIOUR DURING FITTING RANDOMSTATE HAS TO BE FIXED

REFERENCES

J FRIEDMAN GREEDY FUNCTION APPROXIMATION A GRADIENT BOOSTING MACHINE THE ANNALS OF STATISTICS V OL 29 NO 5 2001

10 FRIEDMAN STOCHASTIC GRADIENT BOOSTING 1999

T HASTIE R TIBSHIRANI AND J FRIEDMAN ELEMENTS OF STATISTICAL LEARNING ED 2 SPRINGER 2009

METHODS

APPLY SELF X APPLY TREES IN THE ENSEMBLE TO X RETURN LEAF INDICES

DECISIONFUNCTION SELF X COMPUTE THE DECISION FUNCTION OF X

FITSELF X Y SAMPLEWEIGHT MONITOR FIT THE GRADIENT BOOSTING MODEL

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PREDICT SELF X PREDICT CLASS FOR X

PREDICTLOGPROBA SELF X PREDICT CLASS LOGPROBABILITIES FOR X

PREDICTPROBA SELF X PREDICT CLASS PROBABILITIES FOR X

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

STAGEDDECISIONFUNCTION SELF X COMPUTE DECISION FUNCTION OF XFOR EACH ITERATION

STAGEDPREDICT SELF X PREDICT CLASS AT EACH STAGE FOR X

STAGEDPREDICTPROBA SELF X PREDICT CLASS PROBABILITIES AT EACH STAGE FOR X

INIT SELF LOSS'DEVIANANCE' LEARNINGRATE01 NESTIMATORS100 SUBSAM

PLE10 CRITERION'FRIEDMANMSE' MINSAMPLESSPLIT2 MINSAMPLESLEAF1

MINWEIGHTFRACTIONLEAF00 MAXDEPTH3 MINIMPURITYDECREASE00

MINIMPURITYSPLITNONE INITNONE RANDOMSTATENONE MAXFEATURESNONE

VERBOSE0 MAXLEAFNODESNONE WARMSTARTFALSE PRESORT'AUTO' VALIDA

TIONFRACTION01 NITERNOCHANGENONE TOL00001

APPLYSELF X

APPLY TREES IN THE ENSEMBLE TO X RETURN LEAF INDICES

NEW IN VERSION 017

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTER

NALLY ITS DTYPE WILL BE CONVERTED TO DTYPE NPFLOAT32 IF A SPARSE MATRIX IS PROVIDED

IT WILL BE CONVERTED TO A SPARSE CSR MATRIX

RETURNS

XLEAVES ARRAYLIKE SHAPE NSAMPLES NESTIMATORS NCLASSES FOR EACH DATAPOINT X IN X

AND FOR EACH TREE IN THE ENSEMBLE RETURN THE INDEX OF THE LEAF X ENDS UP IN EACH ESTIMATOR

33 MODEL SELECTION AND EVALUATION 529

SCIKITLEARN USER GUIDE RELEASE 0213

IN THE CASE OF BINARY CLASSIFICATION NCLASSES IS 1

DECISIONFUNCTION SELF X

COMPUTE THE DECISION FUNCTION OF X

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTER

NALLY IT WILL BE CONVERTED TO DTYPENPFLOAT32 AND IF A SPARSE MATRIX IS PROVIDED TO

A SPARSECSRMATRIX

RETURNS

SCORE ARRAY SHAPE NSAMPLES NCLASSES OR NSAMPLES THE DECISION FUNCTION OF THE IN

PUT SAMPLES WHICH CORRESPONDS TO THE RAW VALUES PREDICTED FROM THE TREES OF THE ENSEMBLE

THE ORDER OF THE CLASSES CORRESPONDS TO THAT IN THE ATTRIBUTE CLASSES REGRESSION AND

BINARY CLASSIFICATION PRODUCE AN ARRAY OF SHAPE NSAMPLES

FEATUREIMPORTANCES

RETURN THE FEATURE IMPORTANCES THE HIGHER THE MORE IMPORTANT THE FEATURE

RETURNS

FEATUREIMPORTANCES ARRAY SHAPE NFEATURES THE VALUES OF THIS ARRAY SUM TO 1 UNLESS

ALL TREES ARE SINGLE NODE TREES CONSISTING OF ONLY THE ROOT NODE IN WHICH CASE IT WILL BE AN

ARRAY OF ZEROS

FITSELFXYSAMPLEWEIGHTNONE MONITORNONE

FIT THE GRADIENT BOOSTING MODEL

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTER

NALLY IT WILL BE CONVERTED TO DTYPENPFLOAT32 AND IF A SPARSE MATRIX IS PROVIDED TO

A SPARSECSRMATRIX

YARRAYLIKE SHAPE NSAMPLES TARGET VALUES STRINGS OR INTEGERS IN CLASSIFICATION REAL

NUMBERS IN REGRESSION FOR CLASSIFICATION LABELS MUST CORRESPOND TO CLASSES

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OR NONE SAMPLE WEIGHTS IF NONE THEN

SAMPLES ARE EQUALLY WEIGHTED SPLITS THAT WOULD CREATE CHILD NODES WITH NET ZERO OR NEG

ATIVE WEIGHT ARE IGNORED WHILE SEARCHING FOR A SPLIT IN EACH NODE IN THE CASE OF CLASSIFI

CATION SPLITS ARE ALSO IGNORED IF THEY WOULD RESULT IN ANY SINGLE CLASS CARRYING A NEGATIVE

WEIGHT IN EITHER CHILD NODE

MONITOR CALLABLE OPTIONAL THE MONITOR IS CALLED AFTER EACH ITERATION WITH THE CURRENT ITER

ATION A REFERENCE TO THE ESTIMATOR AND THE LOCAL VARIABLES OF FITSTAGES AS KEYWORD

ARGUMENTSCALLABLEI SELF LOCALS IF THE CALLABLE RETURNS TRUE THE FIT

TING PROCEDURE IS STOPPED THE MONITOR CAN BE USED FOR VARIOUS THINGS SUCH AS COMPUTING

HELDOUT ESTIMATES EARLY STOPPING MODEL INTROSPECT AND SNAPSHOTING

RETURNS

SELF OBJECT

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

530 CHAPTER 3 USER GUIDE



SCIKITLEARN USER GUIDE RELEASE 0213

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELF

PREDICT CLASS FOR X

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTER  
NALLY IT WILL BE CONVERTED TO DTYPENPFLOAT32 AND IF A SPARSE MATRIX IS PROVIDED TO  
A SPARSECSRMATRIX

RETURNS

YARRAY SHAPE NSAMPLES THE PREDICTED VALUES

PREDICTLOGPROBA SELF

PREDICT CLASS LOGPROBABILITIES FOR X

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTER  
NALLY IT WILL BE CONVERTED TO DTYPENPFLOAT32 AND IF A SPARSE MATRIX IS PROVIDED TO  
A SPARSECSRMATRIX

RETURNS

PARRAY SHAPE NSAMPLES NCLASSES THE CLASS LOGPROBABILITIES OF THE INPUT SAMPLES THE  
ORDER OF THE CLASSES CORRESPONDS TO THAT IN THE ATTRIBUTE CLASSES

RAISES

ATTRIBUTEERROR IF THELOSS DOES NOT SUPPORT PROBABILITIES

PREDICTPROBA SELF

PREDICT CLASS PROBABILITIES FOR X

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTER  
NALLY IT WILL BE CONVERTED TO DTYPENPFLOAT32 AND IF A SPARSE MATRIX IS PROVIDED TO  
A SPARSECSRMATRIX

RETURNS

PARRAY SHAPE NSAMPLES NCLASSES THE CLASS PROBABILITIES OF THE INPUT SAMPLES THE  
ORDER OF THE CLASSES CORRESPONDS TO THAT IN THE ATTRIBUTE CLASSES

RAISES

ATTRIBUTEERROR IF THELOSS DOES NOT SUPPORT PROBABILITIES

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS

IN MULTILABEL CLASSIFICATION THIS IS THE SUBSET ACCURACY WHICH IS A HARSH METRIC SINCE YOU REQUIRE FOR EACH  
SAMPLE THAT EACH LABEL SET BE CORRECTLY PREDICTED

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE LABELS FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

33 MODEL SELECTION AND EVALUATION 531

SCIKITLEARN USER GUIDE RELEASE 0213

RETURNS

SCORE FLOAT MEAN ACCURACY OF SELF PREDICT X WRT Y

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENT PARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

STAGED DECISION FUNCTION SELF X

COMPUTE DECISION FUNCTION OF X FOR EACH ITERATION

THIS METHOD ALLOWS MONITORING IE DETERMINE ERROR ON TESTING SET AFTER EACH STAGE  
PARAMETERS

X ARRAY LIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTER  
NALLY IT WILL BE CONVERTED TO DTYPENP FLOAT32 AND IF A SPARSE MATRIX IS PROVIDED TO  
A SPARSE CSR MATRIX

RETURNS

SCORE GENERATOR OF ARRAY SHAPE NSAMPLES K THE DECISION FUNCTION OF THE INPUT SAM  
PLES WHICH CORRESPONDS TO THE RAW VALUES PREDICTED FROM THE TREES OF THE ENSEMBLE THE  
CLASSES CORRESPONDS TO THAT IN THE ATTRIBUTE CLASSES REGRESSION AND BINARY CLASSIFICATION  
ARE SPECIAL CASES WITH K = 1 OTHERWISE KNCLASSES

STAGED PREDICT SELF X

PREDICT CLASS AT EACH STAGE FOR X

THIS METHOD ALLOWS MONITORING IE DETERMINE ERROR ON TESTING SET AFTER EACH STAGE  
PARAMETERS

X ARRAY LIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTER  
NALLY IT WILL BE CONVERTED TO DTYPENP FLOAT32 AND IF A SPARSE MATRIX IS PROVIDED TO  
A SPARSE CSR MATRIX

RETURNS

Y GENERATOR OF ARRAY OF SHAPE NSAMPLES THE PREDICTED VALUE OF THE INPUT SAMPLES

STAGED PREDICT PROBA SELF X

PREDICT CLASS PROBABILITIES AT EACH STAGE FOR X

THIS METHOD ALLOWS MONITORING IE DETERMINE ERROR ON TESTING SET AFTER EACH STAGE  
PARAMETERS

X ARRAY LIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTER  
NALLY IT WILL BE CONVERTED TO DTYPENP FLOAT32 AND IF A SPARSE MATRIX IS PROVIDED TO  
A SPARSE CSR MATRIX

RETURNS

Y GENERATOR OF ARRAY OF SHAPE NSAMPLES THE PREDICTED VALUE OF THE INPUT SAMPLES

532 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES USING SKLEARNENSEMBLEGRADIENTBOOSTINGCLASSIFIER

- GRADIENT BOOSTING REGULARIZATION
- EARLY STOPPING OF GRADIENT BOOSTING
- FEATURE TRANSFORMATIONS WITH ENSEMBLES OF TREES
- GRADIENT BOOSTING OUTFBAG ESTIMATES
- FEATURE DISCRETIZATION

SKLEARNENSEMBLE GRADIENTBOOSTINGREGRESSOR

CLASSSSKLEARNENSEMBLE GRADIENTBOOSTINGREGRESSOR LOSS'LS' LEARNINGRATE01

NESTIMATORS100 SUBSAM

PLE10 CRITERION'FRIEDMANMSE'

MINSAMPLESSPLIT2

MINSAMPLESLEAF1

MINWEIGHTFRACTIONLEAF00

MAXDEPTH3

MINIMPURITYDECREASE00

MINIMPURITYSPLITNONE

INITNONE RANDOMSTATENONE

MAXFEATURESNONE ALPHA09

VERBOSE0 MAXLEAFNODESNONE

WARMSTARTFALSE PRE

SORT'AUTO' VALIDATIONFRACTION01

NITERNOCHANGENONE TOL00001

GRADIENT BOOSTING FOR REGRESSION

GB BUILDS AN ADDITIVE MODEL IN A FORWARD STAGewise FASHION IT ALLOWS FOR THE OPTIMIZATION OF ARBITRARY DIFFER

ENTIABLE LOSS FUNCTIONS IN EACH STAGE A REGRESSION TREE IS FIT ON THE NEGATIVE GRADIENT OF THE GIVEN LOSS FUNCTION

READ MORE IN THE USER GUIDE

PARAMETERS

LOSS 'LS' 'LAD' 'HUBER' 'QUANTILE' OPTIONAL DEFAULT'LS' LOSS FUNCTION TO BE OPTIMIZED 'LS'

REFERS TO LEAST SQUARES REGRESSION 'LAD' LEAST ABSOLUTE DEVIATION IS A HIGHLY ROBUST LOSS

FUNCTION SOLELY BASED ON ORDER INFORMATION OF THE INPUT VARIABLES 'HUBER' IS A COMBINATION

OF THE TWO 'QUANTILE' ALLOWS QUANTILE REGRESSION USE ALPHA TO SPECIFY THE QUANTILE

LEARNINGRATE FLOAT OPTIONAL DEFAULT01 LEARNING RATE SHRINKS THE CONTRIBUTION OF EACH TREE

BYLEARNINGRATE THERE IS A TRADEOFF BETWEEN LEARNINGRATE AND NESTIMATORS

NESTIMATORS INT DEFAULT100 THE NUMBER OF BOOSTING STAGES TO PERFORM GRADIENT BOOSTING

IS FAIRLY ROBUST TO OVERFITTING SO A LARGE NUMBER USUALLY RESULTS IN BETTER PERFORMANCE

SUBSAMPLE FLOAT OPTIONAL DEFAULT10 THE FRACTION OF SAMPLES TO BE USED FOR FITTING THE IN

DIVIDUAL BASE LEARNERS IF SMALLER THAN 10 THIS RESULTS IN STOCHASTIC GRADIENT BOOSTING

SUBSAMPLE INTERACTS WITH THE PARAMETER NESTIMATORS CHOOSING SUBSAMPLE

10 LEADS TO A REDUCTION OF VARIANCE AND AN INCREASE IN BIAS

CRITERION STRING OPTIONAL DEFAULT"FRIEDMANMSE" THE FUNCTION TO MEASURE THE QUALITY OF

A SPLIT SUPPORTED CRITERIA ARE "FRIEDMANMSE" FOR THE MEAN SQUARED ERROR WITH IMPROVE

MENT SCORE BY FRIEDMAN "MSE" FOR MEAN SQUARED ERROR AND "MAE" FOR THE MEAN ABSOLUTE

ERROR THE DEFAULT VALUE OF "FRIEDMANMSE" IS GENERALLY THE BEST AS IT CAN PROVIDE A BETTER

APPROXIMATION IN SOME CASES

33 MODEL SELECTION AND EVALUATION 533

SCIKITLEARN USER GUIDE RELEASE 0213

NEW IN VERSION 018

MINSAMPLESSPLIT INT FLOAT OPTIONAL DEFAULT2 THE MINIMUM NUMBER OF SAMPLES REQUIRED TO SPLIT AN INTERNAL NODE

- IF INT THEN CONSIDER MINSAMPLESSPLIT AS THE MINIMUM NUMBER
- IF FLOAT THEN MINSAMPLESSPLIT IS A FRACTION AND CEILMINSAMPLESSPLIT

NSAMPLES ARE THE MINIMUM NUMBER OF SAMPLES FOR EACH SPLIT

CHANGED IN VERSION 018 ADDED FLOAT VALUES FOR FRACTIONS

MINSAMPLESLEAF INT FLOAT OPTIONAL DEFAULT1 THE MINIMUM NUMBER OF SAMPLES REQUIRED

TO BE AT A LEAF NODE A SPLIT POINT AT ANY DEPTH WILL ONLY BE CONSIDERED IF IT LEAVES AT LEAST MINSAMPLESLEAF TRAINING SAMPLES IN EACH OF THE LEFT AND RIGHT BRANCHES THIS MAY

HAVE THE EFFECT OF SMOOTHING THE MODEL ESPECIALLY IN REGRESSION

- IF INT THEN CONSIDER MINSAMPLESLEAF AS THE MINIMUM NUMBER
- IF FLOAT THEN MINSAMPLESLEAF IS A FRACTION AND CEILMINSAMPLESLEAF

NSAMPLES ARE THE MINIMUM NUMBER OF SAMPLES FOR EACH NODE

CHANGED IN VERSION 018 ADDED FLOAT VALUES FOR FRACTIONS

MINWEIGHTFRACTIONLEAF FLOAT OPTIONAL DEFAULT0 THE MINIMUM WEIGHTED FRACTION OF THE SUM TOTAL OF WEIGHTS OF ALL THE INPUT SAMPLES REQUIRED TO BE AT A LEAF NODE SAMPLES HAVE EQUAL WEIGHT WHEN SAMPLEWEIGHT IS NOT PROVIDED

MAXDEPTH INTEGER OPTIONAL DEFAULT3 MAXIMUM DEPTH OF THE INDIVIDUAL REGRESSION ESTIMATORS THE MAXIMUM DEPTH LIMITS THE NUMBER OF NODES IN THE TREE TUNE THIS PARAMETER FOR

BEST PERFORMANCE THE BEST VALUE DEPENDS ON THE INTERACTION OF THE INPUT VARIABLES

MINIMPURITYDECREASE FLOAT OPTIONAL DEFAULT0 A NODE WILL BE SPLIT IF THIS SPLIT INDUCES

A DECREASE OF THE IMPURITY GREATER THAN OR EQUAL TO THIS VALUE

THE WEIGHTED IMPURITY DECREASE EQUATION IS THE FOLLOWING

$$NT \cdot N \cdot \text{IMPURITY} - NTR \cdot NT \cdot \text{RIGHTIMPURITY}$$

$$NTL \cdot NT \cdot \text{LEFTIMPURITY}$$

WHERE N IS THE TOTAL NUMBER OF SAMPLES NT IS THE NUMBER OF SAMPLES AT THE CURRENT NODE

NTL IS THE NUMBER OF SAMPLES IN THE LEFT CHILD AND NTR IS THE NUMBER OF SAMPLES IN

THE RIGHT CHILD

NNTNTR ANDNTL ALL REFER TO THE WEIGHTED SUM IF SAMPLEWEIGHT IS PASSED

NEW IN VERSION 019

MINIMPURITYSPLIT FLOAT DEFAULT1E7 THRESHOLD FOR EARLY STOPPING IN TREE GROWTH A NODE

WILL SPLIT IF ITS IMPURITY IS ABOVE THE THRESHOLD OTHERWISE IT IS A LEAF

DEPRECATED SINCE VERSION 019 MINIMPURITYSPLIT HAS BEEN DEPRECATED IN FAVOR OF

MINIMPURITYDECREASE IN 019 THE DEFAULT VALUE OF MINIMPURITYSPLIT

WILL CHANGE FROM 1E7 TO 0 IN 023 AND IT WILL BE REMOVED IN 025 USE

MINIMPURITYDECREASE INSTEAD

INIT ESTIMATOR OR 'ZERO' OPTIONAL DEFAULTNONE AN ESTIMATOR OBJECT THAT IS USED TO COMPUTE

THE INITIAL PREDICTIONS INIT HAS TO PROVIDE FIT ANDPREDICT IF 'ZERO' THE INITIAL RAW

PREDICTIONS ARE SET TO ZERO BY DEFAULT A DUMMYESTIMATOR IS USED PREDICTING EITHER THE

AVERAGE TARGET VALUE FOR LOSS'LS' OR A QUANTILE FOR THE OTHER LOSSES

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RAN

DOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE

SCIKITLEARN USER GUIDE RELEASE 0213

RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NPRANDOM

MAXFEATURES INT FLOAT STRING OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF FEATURES TO CONSIDER WHEN LOOKING FOR THE BEST SPLIT

- IF INT THEN CONSIDER MAXFEATURES FEATURES AT EACH SPLIT
- IF FLOAT THEN MAXFEATURES IS A FRACTION AND INTMAXFEATURES NFEATURES FEATURES ARE CONSIDERED AT EACH SPLIT
- IF "AUTO" THEN MAXFEATURESNFEATURES
- IF "SQRT" THEN MAXFEATURESSQRTNFEATURES
- IF "LOG2" THEN MAXFEATURESLOG2NFEATURES
- IF NONE THEN MAXFEATURESNFEATURES

CHOOSINGMAXFEATURES NFEATURES LEADS TO A REDUCTION OF VARIANCE AND AN INCREASE IN BIAS

NOTE THE SEARCH FOR A SPLIT DOES NOT STOP UNTIL AT LEAST ONE VALID PARTITION OF THE NODE SAMPLES IS FOUND EVEN IF IT REQUIRES TO EFFECTIVELY INSPECT MORE THAN MAXFEATURES FEATURES

ALPHA FLOAT DEFAULT0.9 THE ALPHAQUANTILE OF THE HUBER LOSS FUNCTION AND THE QUANTILE LOSS FUNCTION ONLY IF LOSSHUBER ORLOSSQUANTILE

VERBOSE INT DEFAULT 0 ENABLE VERBOSE OUTPUT IF 1 THEN IT PRINTS PROGRESS AND PERFORMANCE ONCE IN A WHILE THE MORE TREES THE LOWER THE FREQUENCY IF GREATER THAN 1 THEN IT PRINTS PROGRESS AND PERFORMANCE FOR EVERY TREE

MAXLEAFNODES INT OR NONE OPTIONAL DEFAULTNONE GROW TREES WITH MAXLEAFNODES

IN BESTFIRST FASHION BEST NODES ARE DEFINED AS RELATIVE REDUCTION IN IMPURITY IF NONE THEN UNLIMITED NUMBER OF LEAF NODES

WARMSTART BOOL DEFAULT FALSE WHEN SET TO TRUE REUSE THE SOLUTION OF THE PREVIOUS CALL TO FIT AND ADD MORE ESTIMATORS TO THE ENSEMBLE OTHERWISE JUST ERASE THE PREVIOUS SOLUTION SEE THE GLOSSARY

PRESORT BOOL OR 'AUTO' OPTIONAL DEFAULT'AUTO' WHETHER TO PRESORT THE DATA TO SPEED UP THE FINDING OF BEST SPLITS IN FITTING AUTO MODE BY DEFAULT WILL USE PRESORTING ON DENSE DATA AND DEFAULT TO NORMAL SORTING ON SPARSE DATA SETTING PRESORT TO TRUE ON SPARSE DATA WILL RAISE AN ERROR

NEW IN VERSION 0.17 OPTIONAL PARAMETER PRESORT

VALIDATIONFRACTION FLOAT OPTIONAL DEFAULT 0.1 THE PROPORTION OF TRAINING DATA TO SET ASIDE AS VALIDATION SET FOR EARLY STOPPING MUST BE BETWEEN 0 AND 1 ONLY USED IF NITERNOCHANGE IS SET TO AN INTEGER

NEW IN VERSION 0.20

NITERNOCHANGE INT DEFAULT NONE NITERNOCHANGE IS USED TO DECIDE IF EARLY STOPPING WILL BE USED TO TERMINATE TRAINING WHEN VALIDATION SCORE IS NOT IMPROVING BY DEFAULT IT IS SET TO NONE TO DISABLE EARLY STOPPING IF SET TO A NUMBER IT WILL SET ASIDE VALIDATIONFRACTION SIZE OF THE TRAINING DATA AS VALIDATION AND TERMINATE TRAINING WHEN VALIDATION SCORE IS NOT IMPROVING IN ALL OF THE PREVIOUS NITERNOCHANGE NUMBERS OF ITERATIONS

NEW IN VERSION 0.20

33 MODEL SELECTION AND EVALUATION 535

SCIKITLEARN USER GUIDE RELEASE 0213  
TOLFLOAT OPTIONAL DEFAULT 1E4 TOLERANCE FOR THE EARLY STOPPING WHEN THE LOSS IS NOT IMPROVING BY AT LEAST TOL FOR NITERNOCHANGE ITERATIONS IF SET TO A NUMBER THE TRAINING STOPS  
NEW IN VERSION 020  
ATTRIBUTES  
FEATUREIMPORTANCES ARRAY SHAPE NFEATURES RETURN THE FEATURE IMPORTANCES THE HIGHER THE MORE IMPORTANT THE FEATURE  
OOBIMPROVEMENT ARRAY SHAPE NESTIMATORS THE IMPROVEMENT IN LOSS DEVIANCE ON THE OUTFBAG SAMPLES RELATIVE TO THE PREVIOUS ITERATION OOBIMPROVEMENT0 IS THE IMPROVEMENT IN LOSS OF THE FIRST STAGE OVER THE INIT ESTIMATOR  
TRAINSCORE ARRAY SHAPE NESTIMATORS THE ITH SCORE TRAINSCOREI IS THE DEVIANCE LOSS OF THE MODEL AT ITERATION ION THE INBAG SAMPLE IF SUBSAMPLE 1 THIS IS THE DEVIANCE ON THE TRAINING DATA  
LOSS LOSSFUNCTION THE CONCRETE LOSSFUNCTION OBJECT  
INIT ESTIMATOR THE ESTIMATOR THAT PROVIDES THE INITIAL PREDICTIONS SET VIA THE INIT ARGUMENT  
ORLOSSINITESTIMATOR  
ESTIMATORS ARRAY OF DECISIONTREEREgressor SHAPE NESTIMATORS 1 THE COLLECTION OF FITTED SUBESTIMATORS  
SEE ALSO  
SKLEARNENSEMBLEHISTGRADIENTBOOSTINGREGRESSOR  
SKLEARNTREEDECISIONTREEREgressor RANDOMFORESTREGRESSOR  
NOTES  
THE FEATURES ARE ALWAYS RANDOMLY PERMUTED AT EACH SPLIT THEREFORE THE BEST FOUND SPLIT MAY VARY EVEN WITH THE SAME TRAINING DATA AND MAXFEATURESNFEATURES IF THE IMPROVEMENT OF THE CRITERION IS IDENTICAL FOR SEVERAL SPLITS ENUMERATED DURING THE SEARCH OF THE BEST SPLIT TO OBTAIN A DETERMINISTIC BEHAVIOUR DURING FITTING  
RANDOMSTATE HAS TO BE FIXED  
REFERENCES  
J FRIEDMAN GREEDY FUNCTION APPROXIMATION A GRADIENT BOOSTING MACHINE THE ANNALS OF STATISTICS V OL 29 NO 5 2001  
10 FRIEDMAN STOCHASTIC GRADIENT BOOSTING 1999  
T HASTIE R TIBSHIRANI AND J FRIEDMAN ELEMENTS OF STATISTICAL LEARNING ED 2 SPRINGER 2009  
METHODS  
APPLY SELF X APPLY TREES IN THE ENSEMBLE TO X RETURN LEAF INDICES  
FITSELF X Y SAMPLEWEIGHT MONITOR FIT THE GRADIENT BOOSTING MODEL  
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR  
PREDICT SELF X PREDICT REGRESSION TARGET FOR X  
CONTINUED ON NEXT PAGE  
536 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

TABLE 320 – CONTINUED FROM PREVIOUS PAGE

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

STAGEDPREDICT SELF X PREDICT REGRESSION TARGET AT EACH STAGE FOR X

INIT SELF LOSS'LS' LEARNINGRATE01 NESTIMATORS100 SUBSAMPLE10

CRITERION'FRIEDMANMSE' MINSAMPLESSPLIT2 MINSAMPLESLEAF1

MINWEIGHTFRACTIONLEAF00 MAXDEPTH3 MINIMPURITYDECREASE00

MINIMPURITYSPLITNONE INITNONE RANDOMSTATENONE MAXFEATURESNONE AL

PHA09 VERBOSE0 MAXLEAFNODESNONE WARMSTARTFALSE PRESORT'AUTO' VALIDATIONFRACTION01 NITERNOCHANGENONE TOL00001

APPLYSELF X

APPLY TREES IN THE ENSEMBLE TO X RETURN LEAF INDICES

NEW IN VERSION 017

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTERNALLY ITS DTYPE WILL BE CONVERTED TO DYPENPFLOAT32 IF A SPARSE MATRIX IS PROVIDED IT WILL BE CONVERTED TO A SPARSE CSRMATRIX

RETURNS

XLEAVES ARRAYLIKE SHAPE NSAMPLES NESTIMATORS FOR EACH DATAPOINT X IN X AND FOR EACH TREE IN THE ENSEMBLE RETURN THE INDEX OF THE LEAF X ENDS UP IN EACH ESTIMATOR

FEATUREIMPORTANCES

RETURN THE FEATURE IMPORTANCES THE HIGHER THE MORE IMPORTANT THE FEATURE

RETURNS

FEATUREIMPORTANCES ARRAY SHAPE NFEATURES THE VALUES OF THIS ARRAY SUM TO 1 UNLESS ALL TREES ARE SINGLE NODE TREES CONSISTING OF ONLY THE ROOT NODE IN WHICH CASE IT WILL BE AN ARRAY OF ZEROS

FITSELF X Y SAMPLEWEIGHTNONE MONITORNONE

FIT THE GRADIENT BOOSTING MODEL

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTERNALLY IT WILL BE CONVERTED TO DYPENPFLOAT32 AND IF A SPARSE MATRIX IS PROVIDED TO A SPARSECSRMATRIX

YARRAYLIKE SHAPE NSAMPLES TARGET VALUES STRINGS OR INTEGERS IN CLASSIFICATION REAL NUMBERS IN REGRESSION FOR CLASSIFICATION LABELS MUST CORRESPOND TO CLASSES

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OR NONE SAMPLE WEIGHTS IF NONE THEN SAMPLES ARE EQUALLY WEIGHTED SPLITS THAT WOULD CREATE CHILD NODES WITH NET ZERO OR NEGATIVE WEIGHT ARE IGNORED WHILE SEARCHING FOR A SPLIT IN EACH NODE IN THE CASE OF CLASSIFICATION SPLITS ARE ALSO IGNORED IF THEY WOULD RESULT IN ANY SINGLE CLASS CARRYING A NEGATIVE WEIGHT IN EITHER CHILD NODE

MONITOR CALLABLE OPTIONAL THE MONITOR IS CALLED AFTER EACH ITERATION WITH THE CURRENT ITERATION A REFERENCE TO THE ESTIMATOR AND THE LOCAL VARIABLES OF FITSTAGES AS KEYWORD

33 MODEL SELECTION AND EVALUATION 537

SCIKITLEARN USER GUIDE RELEASE 0213

ARGUMENTSCALLABLEI SELF LOCALS IF THE CALLABLE RETURNS TRUE THE FITTING PROCEDURE IS STOPPED THE MONITOR CAN BE USED FOR VARIOUS THINGS SUCH AS COMPUTING HELDOUT ESTIMATES EARLY STOPPING MODEL INTROSPECT AND SNAPSHOTING

RETURNS

SELF OBJECT

GETPARAMS SELFDEEPTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELF

PREDICT REGRESSION TARGET FOR X

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTERNALLY IT WILL BE CONVERTED TO DTYPENPFLOAT32 AND IF A SPARSE MATRIX IS PROVIDED TO A SPARSECSRMATRIX

RETURNS

YARRAY SHAPE NSAMPLES THE PREDICTED VALUES

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{U}{V}$  WHERE U IS THE RESIDUAL SUM OF SQUARES  $Y_{TRUE} - Y_{PRED}$   $2SUM$  AND V IS THE TOTAL SUM OF SQUARES  $Y_{TRUE} - Y_{TRUEMEAN}$   $2SUM$  THE BEST POSSIBLE SCORE IS 10 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 00

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES FOR SOME ESTIMATORS THIS MAY BE A PRECOMPUTED KERNEL MATRIX INSTEAD SHAPE NSAMPLES NSAMPLESFITTED WHERE NSAMPLESFITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE VALUES FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT R2 OF SELF PREDICTX WRT Y

NOTES

THE R2 SCORE USED WHEN CALLING SCORE ON A REGRESSOR WILL USE MULTIOUTPUTUNIFORMAVERAGE FROM VERSION 023 TO KEEP CONSISTENT WITH METRICSR2SCORE THIS WILL INFLUENCE THE SCORE METHOD OF ALL THE MULTIOUTPUT REGRESSORS EXCEPT FOR MULTIOUTPUTMULTIOUTPUTREGRESSOR TO SPECIFY THE DEFAULT VALUE MANUALLY AND AVOID THE WARNING PLEASE EITHER CALL METRICSR2SCORE DIRECTLY OR MAKE A CUSTOM SCORER WITH METRICSMAKESCORER THE BUILTIN SCORER R2 USES MULTIOUTPUTUNIFORMAVERAGE

538 CHAPTER 3 USER GUIDE



SCIKITLEARN USER GUIDE RELEASE 0213

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

STAGEDPREDICT SELF X

PREDICT REGRESSION TARGET AT EACH STAGE FOR X

THIS METHOD ALLOWS MONITORING IE DETERMINE ERROR ON TESTING SET AFTER EACH STAGE PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTERNALLY IT WILL BE CONVERTED TO DTYPENPFLOAT32 AND IF A SPARSE MATRIX IS PROVIDED TO A SPARSECSRMATRIX

RETURNS

YGENERATOR OF ARRAY OF SHAPE NSAMPLES THE PREDICTED VALUE OF THE INPUT SAMPLES

EXAMPLES USING SKLEARNENSEMBLEGRADIENTBOOSTINGREGRESSOR

- MODEL COMPLEXITY INFLUENCE
- PLOT INDIVIDUAL AND VOTING REGRESSION PREDICTIONS
- PREDICTION INTERVALS FOR GRADIENT BOOSTING REGRESSION
- GRADIENT BOOSTING REGRESSION
- PARTIAL DEPENDENCE PLOTS

333 MODEL EVALUATION QUANTIFYING THE QUALITY OF PREDICTIONS

THERE ARE 3 DIFFERENT APIS FOR EVALUATING THE QUALITY OF A MODEL'S PREDICTIONS

•ESTIMATOR SCORE METHOD ESTIMATORS HAVE A SCORE METHOD PROVIDING A DEFAULT EVALUATION CRITERION FOR THE PROBLEM THEY ARE DESIGNED TO SOLVE THIS IS NOT DISCUSSED ON THIS PAGE BUT IN EACH ESTIMATOR'S DOCUMENTATION

•SCORING PARAMETER MODELEVALUATION TOOLS USING CROSSVALIDATION SUCH ASMODELSELECTION CROSSVALSCORE ANDMODELSELECTIONGRIDSEARCHCV RELY ON AN INTERNAL SCORING STRATEGY THIS

IS DISCUSSED IN THE SECTION THE SCORING PARAMETER DEFINING MODEL EVALUATION RULES

•METRIC FUNCTIONS THEMETRICS MODULE IMPLEMENTS FUNCTIONS ASSESSING PREDICTION ERROR FOR SPECIFIC PURPOSES THESE METRICS ARE DETAILED IN SECTIONS ON CLASSIFICATION METRICS MULTILABEL RANKING METRICS REGRESSION METRICS ANDCLUSTERING METRICS

FINALLY DUMMY ESTIMATORS ARE USEFUL TO GET A BASELINE VALUE OF THOSE METRICS FOR RANDOM PREDICTIONS

SEE ALSO

FOR "PAIRWISE" METRICS BETWEEN SAMPLES AND NOT ESTIMATORS OR PREDICTIONS SEE THE PAIRWISE METRICS AFFINITIES AND KERNELS SECTION

33 MODEL SELECTION AND EVALUATION 539

SCIKITLEARN USER GUIDE RELEASE 0213

THESCORING PARAMETER DEFINING MODEL EVALUATION RULES

MODEL SELECTION AND EVALUATION USING TOOLS SUCH AS MODELSELECTIONGRIDSEARCHCV AND MODELSELECTIONCROSSVALSCORE TAKE ASCORING PARAMETER THAT CONTROLS WHAT METRIC THEY APPLY TO THE ESTIMATORS EVALUATED

COMMON CASES PREDEFINED VALUES

FOR THE MOST COMMON USE CASES YOU CAN DESIGNATE A SCORER OBJECT WITH THE SCORING PARAMETER THE TABLE BELOW SHOWS ALL POSSIBLE VALUES ALL SCORER OBJECTS FOLLOW THE CONVENTION THAT HIGHER RETURN VALUES ARE BETTER THAN LOWER RETURN VALUES THUS METRICS WHICH MEASURE THE DISTANCE BETWEEN THE MODEL AND THE DATA LIKE METRICS MEANSQUAREDERROR ARE AVAILABLE AS NEGMEANSQUAREDERROR WHICH RETURN THE NEGATED VALUE OF THE METRIC

SCORING FUNCTION COMMENT

CLASSIFICATION

‘ACCURACY’ METRICSACCURACYScore

‘BALANCEDACCURACY’ METRICSBALANCEDACCURACYScore

‘AVERAGEPRECISION’ METRICSAVERAGEPRECISIONScore

‘BRIERSCORELOSS’ METRICSBRIERSCORELOSS

‘F1’ METRICSF1Score FOR BINARY TARGETS

‘F1MICRO’ METRICSF1Score MICROAVERAGED

‘F1MACRO’ METRICSF1Score MACROAVERAGED

‘F1WEIGHTED’ METRICSF1Score WEIGHTED AVERAGE

‘F1SAMPLES’ METRICSF1Score BY MULTILABEL SAMPLE

‘NEGLOGLOSS’ METRICSLGLOSS REQUIRESPREDICTPROBA SUPPORT

‘PRECISION’ ETC METRICSPRECISIONScore SUFFIXES APPLY AS WITH ‘F1’

‘RECALL’ ETC METRICSPRECALLScore SUFFIXES APPLY AS WITH ‘F1’

‘JACCARD’ ETC METRICSJACCARDScore SUFFIXES APPLY AS WITH ‘F1’

‘ROCAUC’ METRICSRCAUCScore

CLUSTERING

‘ADJUSTEDMUTUALINFOScore’ METRICSDJUSTEDMUTUALINFOScore

‘ADJUSTEDDRANDScore’ METRICSDJUSTEDDRANDScore

‘COMPLETENESSScore’ METRICSCOMPLETENESSScore

‘FOWLKESMALLOWSScore’ METRICSFOWLKESMALLOWSScore

‘HOMOGENEITYScore’ METRICSHOMOGENEITYScore

‘MUTUALINFOScore’ METRICSMUTUALINFOScore

‘NORMALIZEDMUTUALINFOScore’ METRICSNORMALIZEDMUTUALINFOScore

‘VMEASUREScore’ METRICSVMEASUREScore

REGRESSION

‘EXPLAINEDVARIANCE’ METRICSEXPLAINEDVARIANCEScore

‘MAXERROR’ METRICSMAXERROR

‘NEGMEANABSOLUTEERROR’ METRICSMEANABSOLUTEERROR

‘NEGMEANSQUAREDERROR’ METRICSMEANSQUAREDERROR

‘NEGMEANSQUAREDLOGERROR’ METRICSMEANSQUAREDLOGERROR

‘NEGMEDIANABSOLUTEERROR’ METRICSMEDIANABSOLUTEERROR

‘R2’ METRICSR2Score

USAGE EXAMPLES

FROM SKLEARN IMPORT SVM DATASETS

FROM SKLEARNMODELSELECTION IMPORT CROSSVALSCORE

IRIS DATASETSLOADIRIS

540 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

X Y IRISDATA IRISTARGET

CLF SVMSCVCGAMMASCALE RANDOMSTATE0

CROSSVALSCORECLF X Y SCORINGRECALLMACRO

CV5

ARRAY096 096 096 093 1

MODEL SVMSCV

CROSSVALSCOREMODEL X Y CV5 SCORINGWRONGCHOICE

TRACEBACK MOST RECENT CALL LAST

VALUEERROR WRONGCHOICE IS NOT A VALID SCORING VALUE USE SORTEDSKLEARNMETRICS

→SCORERSKEYS TO GET VALID OPTIONS

NOTE THE VALUES LISTED BY THE VALUEERROR EXCEPTION CORRESPOND TO THE FUNCTIONS MEASURING PREDICTION ACCURACY DESCRIBED IN THE FOLLOWING SECTIONS THE SCORER OBJECTS FOR THOSE FUNCTIONS ARE STORED IN THE DICTIONARY SKLEARN METRICSSCORERS

DEFINING YOUR SCORING STRATEGY FROM METRIC FUNCTIONS

THE MODULE SKLEARNMETRICS ALSO EXPOSES A SET OF SIMPLE FUNCTIONS MEASURING A PREDICTION ERROR GIVEN GROUND TRUTH AND PREDICTION

- FUNCTIONS ENDING WITH SCORE RETURN A VALUE TO MAXIMIZE THE HIGHER THE BETTER
- FUNCTIONS ENDING WITH ERROR ORLOSS RETURN A VALUE TO MINIMIZE THE LOWER THE BETTER WHEN CONVERTING INTO A SCORER OBJECT USING MAKESCORER SET THEGREATERISBETTER PARAMETER TO FALSE TRUE BY DEFAULT SEE THE PARAMETER DESCRIPTION BELOW

METRICS AVAILABLE FOR VARIOUS MACHINE LEARNING TASKS ARE DETAILED IN SECTIONS BELOW

MANY METRICS ARE NOT GIVEN NAMES TO BE USED AS SCORING VALUES SOMETIMES BECAUSE THEY REQUIRE ADDITIONAL PARAMETERS SUCH AS FBETAScore IN SUCH CASES YOU NEED TO GENERATE AN APPROPRIATE SCORING OBJECT THE SIMPLEST WAY TO GENERATE A CALLABLE OBJECT FOR SCORING IS BY USING MAKESCORER THAT FUNCTION CONVERTS METRICS INTO CALLABLES THAT CAN BE USED FOR MODEL EVALUATION

ONE TYPICAL USE CASE IS TO WRAP AN EXISTING METRIC FUNCTION FROM THE LIBRARY WITH NONDEFAULT VALUES FOR ITS PARAMETERS SUCH AS THE BETA PARAMETER FOR THE FBETAScore FUNCTION

```
FROM SKLEARNMETRICS IMPORT FBETAScore MAKESCORER
FTWOSCORER MAKESCORERFBETAScore BETA2
FROM SKLEARNMODELSELECTION IMPORT GRIDSEARCHCV
FROM SKLEARNsvm IMPORT LINEARSVC
GRID GRIDSEARCHCVLINEARSVC PARAMGRIDC 1 10
SCORINGFTWOSCORER CV5
```

THE SECOND USE CASE IS TO BUILD A COMPLETELY CUSTOM SCORER OBJECT FROM A SIMPLE PYTHON FUNCTION USING MAKESCORER WHICH CAN TAKE SEVERAL PARAMETERS

- THE PYTHON FUNCTION YOU WANT TO USE MYCUSTOMLOSSFUNC IN THE EXAMPLE BELOW
- WHETHER THE PYTHON FUNCTION RETURNS A SCORE GREATERISBETTERTRUE THE DEFAULT OR A LOSS GREATERISBETTERFALSE IF A LOSS THE OUTPUT OF THE PYTHON FUNCTION IS NEGATED BY THE SCORER OBJECT CONFORMING TO THE CROSS VALIDATION CONVENTION THAT SCORERS RETURN HIGHER VALUES FOR BETTER MODELS
- FOR CLASSIFICATION METRICS ONLY WHETHER THE PYTHON FUNCTION YOU PROVIDED REQUIRES CONTINUOUS DECISION CERTAIN TIES NEEDSTHRESHOLDTRUE THE DEFAULT VALUE IS FALSE
- ANY ADDITIONAL PARAMETERS SUCH AS BETA ORLABELS INF1SCORE

33 MODEL SELECTION AND EVALUATION 541

SCIKITLEARN USER GUIDE RELEASE 0213

HERE IS AN EXAMPLE OF BUILDING CUSTOM SCORERS AND OF USING THE GREATERISBETTER PARAMETER

```
import numpy as np
def mycustomlossfunc(ytrue, ypred):
    diff = np.abs(ytrue - ypred)
    return np.log(1 + diff)
```

```
score = mycustomlossfunc(x, y)
which will be np.log(2.0693) given the values for x
and y defined below
score = makescorer(mycustomlossfunc, greaterisbetter=False)
x = 1
y = 0
```

```
from sklearn.dummy import DummyClassifier
clf = DummyClassifier(strategy='most_frequent', random_state=0)
clf.fit(x, y)
mycustomlossfunc(clf.predict(x), y)
0.69
score = clf.score(x, y)
0.69
```

IMPLEMENTING YOUR OWN SCORING OBJECT

YOU CAN GENERATE EVEN MORE FLEXIBLE MODEL SCORERS BY CONSTRUCTING YOUR OWN SCORING OBJECT FROM SCRATCH WITHOUT USING THE MAKESCORER FACTORY FOR A CALLABLE TO BE A SCORER IT NEEDS TO MEET THE PROTOCOL SPECIFIED BY THE FOLLOWING TWO RULES

- IT CAN BE CALLED WITH PARAMETERS ESTIMATOR, X, Y, WHERE ESTIMATOR IS THE MODEL THAT SHOULD BE EVALUATED, X IS VALIDATION DATA AND Y IS THE GROUND TRUTH TARGET FOR X IN THE SUPERVISED CASE OR NONE IN THE UNSUPERVISED CASE
  - IT RETURNS A FLOATING POINT NUMBER THAT QUANTIFIES THE ESTIMATOR PREDICTION QUALITY ON X WITH REFERENCE TO Y. AGAIN BY CONVENTION HIGHER NUMBERS ARE BETTER SO IF YOUR SCORER RETURNS LOSS THAT VALUE SHOULD BE NEGATED.
- NOTE: USING CUSTOM SCORERS IN FUNCTIONS WHERE `n_jobs = 1` WHILE DEFINING THE CUSTOM SCORING FUNCTION ALONGSIDE THE CALLING FUNCTION SHOULD WORK OUT OF THE BOX WITH THE DEFAULT JOBLIB BACKEND. LOHY IMPORTING IT FROM ANOTHER MODULE WILL BE A MORE ROBUST APPROACH AND WORK INDEPENDENTLY OF THE JOBLIB BACKEND.

FOR EXAMPLE TO USE `n_jobs` GREATER THAN 1 IN THE EXAMPLE BELOW, `CUSTOMSCORINGFUNCTION` FUNCTION IS SAVED IN A USER-CREATED MODULE `CUSTOMSCORERMODULE.PY` AND IMPORTED

```
from customscorermodule import customscoringfunction
crossval_score_model
```

```
x_train
y_train
scoring = makescorer(customscoringfunction, greaterisbetter=False)
cv5
n_jobs=1
```

SCIKITLEARN USER GUIDE RELEASE 0213

USING MULTIPLE METRIC EVALUATION

SCIKITLEARN ALSO PERMITS EVALUATION OF MULTIPLE METRICS IN GRIDSEARCHCV RANDOMIZEDSEARCHCV AND CROSSVALIDATE

THERE ARE TWO WAYS TO SPECIFY MULTIPLE SCORING METRICS FOR THE SCORING PARAMETER

- AS AN ITERABLE OF STRING METRICS

SCORING ACCURACY PRECISION

- AS ADICT MAPPING THE SCORER NAME TO THE SCORING FUNCTION

FROM SKLEARNMETRICS IMPORT ACCURACYScore

FROM SKLEARNMETRICS IMPORT MAKESCORER

SCORING ACCURACY MAKESCORERACCURACYScore

PREC PRECISION

NOTE THAT THE DICT VALUES CAN EITHER BE SCORER FUNCTIONS OR ONE OF THE PREDEFINED METRIC STRINGS

CURRENTLY ONLY THOSE SCORER FUNCTIONS THAT RETURN A SINGLE SCORE CAN BE PASSED INSIDE THE DICT SCORER FUNCTIONS THAT RETURN MULTIPLE VALUES ARE NOT PERMITTED AND WILL REQUIRE A WRAPPER TO RETURN A SINGLE METRIC

FROM SKLEARNMODELSELECTION IMPORT CROSSVALIDATE

FROM SKLEARNMETRICS IMPORT CONFUSIONMATRIX

A SAMPLE TOY BINARY CLASSIFICATION DATASET

X Y DATASETSMAKECLASSIFICATIONNNCLASSES2 RANDOMSTATE0

SVM LINEARSVCRANDOMSTATE0

DEF TNYTRUE YPRED RETURNCONFUSIONMATRIXYTRUE YPRED0 0

DEF FPYTRUE YPRED RETURNCONFUSIONMATRIXYTRUE YPRED0 1

DEF FNYTRUE YPRED RETURNCONFUSIONMATRIXYTRUE YPRED1 0

DEF TPYTRUE YPRED RETURNCONFUSIONMATRIXYTRUE YPRED1 1

SCORING TP MAKESCORERTP TN MAKESCORERTN

FP MAKESCORERFP FN MAKESCORERFN

CVRESULTS CROSSVALIDATESVMFITX Y X Y

SCORINGSCORING CV5

GETTING THE TEST SET TRUE POSITIVE SCORES

PRINTCVRESULTSTESTTP

10 9 8 7 8

GETTING THE TEST SET FALSE NEGATIVE SCORES

PRINTCVRESULTSTESTFN

0 1 2 3 2

CLASSIFICATION METRICS

THESKLEARNMETRICS MODULE IMPLEMENTS SEVERAL LOSS SCORE AND UTILITY FUNCTIONS TO MEASURE CLASSIFICATION PERFORMANCE SOME METRICS MIGHT REQUIRE PROBABILITY ESTIMATES OF THE POSITIVE CLASS CONFIDENCE VALUES OR BINARY DECISIONS VALUES MOST IMPLEMENTATIONS ALLOW EACH SAMPLE TO PROVIDE A WEIGHTED CONTRIBUTION TO THE OVERALL SCORE THROUGH THE SAMPLEWEIGHT PARAMETER

SOME OF THESE ARE RESTRICTED TO THE BINARY CLASSIFICATION CASE

PRECISIONRECALLCURVE YTRUE PROBASPRED COMPUTE PRECISIONRECALL PAIRS FOR DIFFERENT PROBABILITY THRESHOLDS

ROCCURVE YTRUE YSCORE POSLABEL COMPUTE RECEIVER OPERATING CHARACTERISTIC ROC

BALANCEDACCURACYScore YTRUE YPRED COMPUTE THE BALANCED ACCURACY

33 MODEL SELECTION AND EVALUATION 543

SCIKITLEARN USER GUIDE RELEASE 0213

OTHERS ALSO WORK IN THE MULTICLASS CASE

COHENKAPPASCORE Y1 Y2 LABELS WEIGHTS    COHEN’S KAPPA A STATISTIC THAT MEASURES INTERANNOTATOR AGREEMENT

CONFUSIONMATRIX YTRUE YPRED LABELS    COMPUTE CONFUSION MATRIX TO EVALUATE THE ACCURACY OF A CLASSIFICATION

HINGELOSS YTRUE PREDDECISION LABELS    AVERAGE HINGE LOSS NONREGULARIZED

MATTHEWSCORRCOEY YTRUE YPRED    COMPUTE THE MATTHEWS CORRELATION COEFFICIENT MCC

SOME ALSO WORK IN THE MULTILABEL CASE

ACCURACYSYSCORE YTRUE YPRED NORMALIZE    ACCURACY CLASSIFICATION SCORE

CLASSIFICATIONREPORT YTRUE YPRED    BUILD A TEXT REPORT SHOWING THE MAIN CLASSIFICATION METRICS

F1SCORE YTRUE YPRED LABELS    COMPUTE THE F1 SCORE ALSO KNOWN AS BALANCED FSCORE OR FMEASURE

FBETASCORE YTRUE YPRED BETA LABELS    COMPUTE THE FBETA SCORE

HAMMINGLOSS YTRUE YPRED LABELS    COMPUTE THE AVERAGE HAMMING LOSS

JACCARDSYSCORE YTRUE YPRED LABELS    JACCARD SIMILARITY COEFFICIENT SCORE

LOGLOSS YTRUE YPRED EPS NORMALIZE    LOG LOSS AKA LOGISTIC LOSS OR CROSSENTROPY LOSS

MULTILABELCONFUSIONMATRIX YTRUE YPRED COMPUTE A CONFUSION MATRIX FOR EACH CLASS OR SAMPLE

PRECISIONRECALLFSCORESUPPORT YTRUE

YPREDCOMPUTE PRECISION RECALL FMEASURE AND SUPPORT FOR EACH CLASS

PRECISIONSCORE YTRUE YPRED LABELS    COMPUTE THE PRECISION

RECALLSCORE YTRUE YPRED LABELS    COMPUTE THE RECALL

ZEROONELOSS YTRUE YPRED NORMALIZE    ZEROONE CLASSIFICATION LOSS

AND SOME WORK WITH BINARY AND MULTILABEL BUT NOT MULTICLASS PROBLEMS

AVERAGEPRECISIONSCORE YTRUE YSCORE    COMPUTE AVERAGE PRECISION AP FROM PREDICTION SCORES

ROCAUCSCORE YTRUE YSCORE AVERAGE    COMPUTE AREA UNDER THE RECEIVER OPERATING CHARACTERIS TIC CURVE ROC AUC FROM PREDICTION SCORES

IN THE FOLLOWING SUBSECTIONS WE WILL DESCRIBE EACH OF THOSE FUNCTIONS PRECEDED BY SOME NOTES ON COMMON API AND METRIC DEFINITION

FROM BINARY TO MULTICLASS AND MULTILABEL

SOME METRICS ARE ESSENTIALLY DEFINED FOR BINARY CLASSIFICATION TASKS EG F1SCORE ROCAUCSCORE IN THESE CASES BY DEFAULT ONLY THE POSITIVE LABEL IS EVALUATED ASSUMING BY DEFAULT THAT THE POSITIVE CLASS IS LABELLED 1THOUGH THIS MAY BE CONFIGURABLE THROUGH THE POSLABEL PARAMETER

IN EXTENDING A BINARY METRIC TO MULTICLASS OR MULTILABEL PROBLEMS THE DATA IS TREATED AS A COLLECTION OF BINARY PROBLEMS ONE FOR EACH CLASS THERE ARE THEN A NUMBER OF WAYS TO AVERAGE BINARY METRIC CALCULATIONS ACROSS THE SET OF CLASSES EACH OF WHICH MAY BE USEFUL IN SOME SCENARIO WHERE AVAILABLE YOU SHOULD SELECT AMONG THESE USING THE AVERAGE PARAMETER

- MACRO SIMPLY CALCULATES THE MEAN OF THE BINARY METRICS GIVING EQUAL WEIGHT TO EACH CLASS IN PROBLEMS WHERE INFREQUENT CLASSES ARE NONETHELESS IMPORTANT MACROAVERAGING MAY BE A MEANS OF HIGHLIGHTING THEIR PERFORMANCE ON THE OTHER HAND THE ASSUMPTION THAT ALL CLASSES ARE EQUALLY IMPORTANT IS OFTEN UNTRUE SUCH THAT MACROAVERAGING WILL OVEREMPHASIZE THE TYPICALLY LOW PERFORMANCE ON AN INFREQUENT CLASS

544 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

- WEIGHTED ACCOUNTS FOR CLASS IMBALANCE BY COMPUTING THE AVERAGE OF BINARY METRICS IN WHICH EACH CLASS’S SCORE IS WEIGHTED BY ITS PRESENCE IN THE TRUE DATA SAMPLE
  - MICRO GIVES EACH SAMPLECLASS PAIR AN EQUAL CONTRIBUTION TO THE OVERALL METRIC EXCEPT AS A RESULT OF SAMPLE WEIGHT RATHER THAN SUMMING THE METRIC PER CLASS THIS SUMS THE DIVIDENDS AND DIVISORS THAT MAKE UP THE PERCLASS METRICS TO CALCULATE AN OVERALL QUOTIENT MICROAVERAGING MAY BE PREFERRED IN MULTILABEL SETTINGS INCLUDING MULTICLASS CLASSIFICATION WHERE A MAJORITY CLASS IS TO BE IGNORED
  - SAMPLES APPLIES ONLY TO MULTILABEL PROBLEMS IT DOES NOT CALCULATE A PERCLASS MEASURE INSTEAD CALCULATING THE METRIC OVER THE TRUE AND PREDICTED CLASSES FOR EACH SAMPLE IN THE EVALUATION DATA AND RETURNING THEIR SAMPLEWEIGHT WEIGHTED AVERAGE
  - SELECTING AVERAGENONE WILL RETURN AN ARRAY WITH THE SCORE FOR EACH CLASS
- WHILE MULTICLASS DATA IS PROVIDED TO THE METRIC LIKE BINARY TARGETS AS AN ARRAY OF CLASS LABELS MULTILABEL DATA IS SPECIFIED AS AN INDICATOR MATRIX IN WHICH CELL I J HAS VALUE 1 IF SAMPLE I HAS LABELJ AND VALUE 0 OTHERWISE

ACCURACY SCORE  
THE ACCURACY SCORE FUNCTION COMPUTES THE ACCURACY EITHER THE FRACTION DEFAULT OR THE COUNT NORMALIZE FALSE OF CORRECT PREDICTIONS  
IN MULTILABEL CLASSIFICATION THE FUNCTION RETURNS THE SUBSET ACCURACY IF THE ENTIRE SET OF PREDICTED LABELS FOR A SAMPLE STRICTLY MATCH WITH THE TRUE SET OF LABELS THEN THE SUBSET ACCURACY IS 1.0 OTHERWISE IT IS 0.0  
IF  $\hat{y}_i$  IS THE PREDICTED VALUE OF THE  $i$ TH SAMPLE AND  $y_i$  IS THE CORRESPONDING TRUE VALUE THEN THE FRACTION OF CORRECT PREDICTIONS OVER  $N$  SAMPLES IS DEFINED AS

$$\text{ACCURACY} = \frac{1}{N} \sum_{i=1}^N \mathbb{1}_{\hat{y}_i = y_i}$$

```
WHERE  $\mathbb{1}$  IS THE INDICATOR FUNCTION
import numpy as np
from sklearn.metrics import accuracy_score
YPRED = [0, 2, 1, 3]
YTRUE = [0, 1, 2, 3]
accuracy_score(YTRUE, YPRE
```

0.5  
accuracy\_score(YTRUE, YPRE, normalize=False)  
2  
IN THE MULTILABEL CASE WITH BINARY LABEL INDICATORS  
accuracy\_score(np.array([0, 1, 1, 1]), np.ones(2, 2))

0.5  
EXAMPLE  
• SEE TEST WITH PERMUTATIONS THE SIGNIFICANCE OF A CLASSIFICATION SCORE FOR AN EXAMPLE OF ACCURACY SCORE USAGE USING PERMUTATIONS OF THE DATASET  
33 MODEL SELECTION AND EVALUATION 545

BALANCED ACCURACY SCORE

THEBALANCEDACCURACYSCORE FUNCTION COMPUTES THE BALANCED ACCURACY WHICH AVOIDS INFLATED PERFORMANCE ESTIMATES ON IMBALANCED DATASETS IT IS THE MACROAVERAGE OF RECALL SCORES PER CLASS OR EQUIVALENTLY RAW ACCURACY WHERE EACH SAMPLE IS WEIGHTED ACCORDING TO THE INVERSE PREVALENCE OF ITS TRUE CLASS THUS FOR BALANCED DATASETS THE SCORE IS EQUAL TO ACCURACY

IN THE BINARY CASE BALANCED ACCURACY IS EQUAL TO THE ARITHMETIC MEAN OF SENSITIVITY TRUE POSITIVE RATE AND SPECIFICITY TRUE NEGATIVE RATE OR THE AREA UNDER THE ROC CURVE WITH BINARY PREDICTIONS RATHER THAN SCORES

IF THE CLASSIFIER PERFORMS EQUALLY WELL ON EITHER CLASS THIS TERM REDUCES TO THE CONVENTIONAL ACCURACY IE THE NUMBER OF CORRECT PREDICTIONS DIVIDED BY THE TOTAL NUMBER OF PREDICTIONS

IN CONTRAST IF THE CONVENTIONAL ACCURACY IS ABOVE CHANCE ONLY BECAUSE THE CLASSIFIER TAKES ADVANTAGE OF AN IMBALANCED TEST SET THEN THE BALANCED ACCURACY AS APPROPRIATE WILL DROP TO1

00000000

THE SCORE RANGES FROM 0 TO 1 OR WHEN ADJUSTEDTRUE IS USED IT RESCALED TO THE RANGE1

1-00000000TO 1 INCLUSIVE

WITH PERFORMANCE AT RANDOM SCORING 0

IF00IS THE TRUE VALUE OF THE 0TH SAMPLE AND 00IS THE CORRESPONDING SAMPLE WEIGHT THEN WE ADJUST THE SAMPLE WEIGHT TO

0000Σ

01000000

WHERE 10IS THE INDICATOR FUNCTION GIVEN PREDICTED 00FOR SAMPLE 0 BALANCED ACCURACY IS DEFINED AS

BALANCEDACCURACY 00 0 1Σ00Σ

010000 00

WITHADJUSTEDTRUE BALANCED ACCURACY REPORTS THE RELATIVE INCREASE FROM BALANCEDACCURACY 00 0 1

00000000 IN THE BINARY CASE THIS IS ALSO KNOWN AS YODEN’S J STATISTIC OR INFORMEDNESS

NOTE THE MULTICLASS DEFINITION HERE SEEMS THE MOST REASONABLE EXTENSION OF THE METRIC USED IN BINARY CLASSIFICATION THOUGH THERE IS NO CERTAIN CONSENSUS IN THE LITERATURE

- OUR DEFINITION MOSLEY2013 KELLEHER2015 ANDGUYON2015 WHERE GUYON2015 ADOPT THE ADJUSTED VER SION TO ENSURE THAT RANDOM PREDICTIONS HAVE A SCORE OF 0AND PERFECT PREDICTIONS HAVE A SCORE OF 1
- CLASS BALANCED ACCURACY AS DESCRIBED IN MOSLEY2013 THE MINIMUM BETWEEN THE PRECISION AND THE RECALL FOR EACH CLASS IS COMPUTED THOSE VALUES ARE THEN AVERAGED OVER THE TOTAL NUMBER OF CLASSES TO GET THE BALANCED ACCURACY
- BALANCED ACCURACY AS DESCRIBED IN URBANOWICZ2015 THE AVERAGE OF SENSITIVITY AND SPECIFICITY IS COMPUTED FOR EACH CLASS AND THEN AVERAGED OVER TOTAL NUMBER OF CLASSES

REFERENCES

COHEN’S KAPPA

THE FUNCTION COHENKAPPASCORE COMPUTES COHEN’S KAPPA STATISTIC THIS MEASURE IS INTENDED TO COMPARE LABEL INGS BY DIFFERENT HUMAN ANNOTATORS NOT A CLASSIFIER VERSUS A GROUND TRUTH

THE KAPPA SCORE SEE DOCSTRING IS A NUMBER BETWEEN 1 AND 1 SCORES ABOVE 8 ARE GENERALLY CONSIDERED GOOD AGREE MENT ZERO OR LOWER MEANS NO AGREEMENT PRACTICALLY RANDOM LABELS



SCIKITLEARN USER GUIDE RELEASE 0213

KAPPA SCORES CAN BE COMPUTED FOR BINARY OR MULTICLASS PROBLEMS BUT NOT FOR MULTILABEL PROBLEMS EXCEPT BY MANUALLY COMPUTING A PERLABEL SCORE AND NOT FOR MORE THAN TWO ANNOTATORS

```
FROM SKLEARNMETRICS IMPORT COHENKAPPASCORE
YTRUE 2 0 2 2 0 1
YPRED 0 0 2 2 0 2
COHENKAPPASCOREYTRUE YPRED
04285714285714286
```

CONFUSION MATRIX

THECONFUSIONMATRIX FUNCTION EVALUATES CLASSIFICATION ACCURACY BY COMPUTING THE CONFUSION MATRIX WITH EACH ROW CORRESPONDING TO THE TRUE CLASS [HTTPS://ENWIKIPEDIA.ORG/WIKI/CONFUSIONMATRIX](https://en.wikipedia.org/wiki/Confusion_matrix) WIKIPEDIA AND OTHER REFERENCES MAY USE DIFFERENT CONVENTION FOR AXES

BY DEFINITION ENTRY  $\begin{bmatrix} a & b \\ c & d \end{bmatrix}$  IN A CONFUSION MATRIX IS THE NUMBER OF OBSERVATIONS ACTUALLY IN GROUP  $\begin{bmatrix} 1 \\ 0 \end{bmatrix}$  BUT PREDICTED TO BE IN GROUP  $\begin{bmatrix} 0 \\ 1 \end{bmatrix}$  HERE IS AN EXAMPLE

```
FROM SKLEARNMETRICS IMPORT CONFUSIONMATRIX
YTRUE 2 0 2 2 0 1
YPRED 0 0 2 2 0 2
CONFUSIONMATRIXYTRUE YPRED
ARRAY2 0 0
0 0 1
1 0 2
```

HERE IS A VISUAL REPRESENTATION OF SUCH A CONFUSION MATRIX THIS FIGURE COMES FROM THE CONFUSION MATRIX EXAMPLE FOR BINARY PROBLEMS WE CAN

GET COUNTS OF TRUE NEGATIVES FALSE POSITIVES FALSE NEGATIVES AND TRUE POSITIVES AS FOLLOWS

```
YTRUE 0 0 0 1 1 1 1 1
YPRED 0 1 0 1 0 1 0 1
TN FP FN TP CONFUSIONMATRIXYTRUE YPREDRAVEL
```

33 MODEL SELECTION AND EVALUATION 547

SCIKITLEARN USER GUIDE RELEASE 0213

TN FP FN TP  
2 1 2 3

EXAMPLE

- SEE CONFUSION MATRIX FOR AN EXAMPLE OF USING A CONFUSION MATRIX TO EVALUATE CLASSIFIER OUTPUT QUALITY
- SEE RECOGNIZING HANDWRITTEN DIGITS FOR AN EXAMPLE OF USING A CONFUSION MATRIX TO CLASSIFY HANDWRITTEN DIGITS

- SEE CLASSIFICATION OF TEXT DOCUMENTS USING SPARSE FEATURES FOR AN EXAMPLE OF USING A CONFUSION MATRIX TO CLASSIFY TEXT DOCUMENTS

CLASSIFICATION REPORT

THECLASSIFICATIONREPORT FUNCTION BUILDS A TEXT REPORT SHOWING THE MAIN CLASSIFICATION METRICS HERE IS A SMALL EXAMPLE WITH CUSTOM TARGETNAMES AND INFERRED LABELS

FROM SKLEARNMETRICS IMPORT CLASSIFICATIONREPORT

YTRUE 0 1 2 2 0

YPRED 0 0 2 1 0

TARGETNAMES CLASS 0 CLASS 1 CLASS 2

PRINTCLASSIFICATIONREPORTYTRUE YPRED TARGETNAMESTARGETNAMES

PRECISION RECALL F1SCORE SUPPORT

CLASS 0 067 100 080 2

CLASS 1 000 000 000 1

CLASS 2 100 050 067 2

ACCURACY 060 5

MACRO AVG 056 050 049 5

WEIGHTED AVG 067 060 059 5

EXAMPLE

- SEE RECOGNIZING HANDWRITTEN DIGITS FOR AN EXAMPLE OF CLASSIFICATION REPORT USAGE FOR HANDWRITTEN DIGITS
- SEE CLASSIFICATION OF TEXT DOCUMENTS USING SPARSE FEATURES FOR AN EXAMPLE OF CLASSIFICATION REPORT USAGE FOR TEXT DOCUMENTS
- SEE PARAMETER ESTIMATION USING GRID SEARCH WITH CROSSVALIDATION FOR AN EXAMPLE OF CLASSIFICATION REPORT USAGE FOR GRID SEARCH WITH NESTED CROSSVALIDATION

HAMMING LOSS

THEHAMMINGLOSS COMPUTES THE AVERAGE HAMMING LOSS OR HAMMING DISTANCE BETWEEN TWO SETS OF SAMPLES

IF  $\hat{y}_i$  IS THE PREDICTED VALUE FOR THE  $i$ TH LABEL OF A GIVEN SAMPLE  $y_i$  IS THE CORRESPONDING TRUE VALUE AND  $N$  LABELS IS THE NUMBER OF CLASSES OR LABELS THEN THE HAMMING LOSS  $\frac{1}{N} \sum_{i=1}^N \frac{|\hat{y}_i - y_i|}{N}$  BETWEEN TWO SAMPLES IS DEFINED AS

$$\frac{1}{N} \sum_{i=1}^N \frac{|\hat{y}_i - y_i|}{N}$$

$$\frac{1}{N} \sum_{i=1}^N \frac{|\hat{y}_i - y_i|}{N}$$

010000

SCIKITLEARN USER GUIDE RELEASE 0213  
WHERE 1 IS THE INDICATOR FUNCTION  
FROM SKLEARNMETRICS IMPORT HAMMINGLOSS

YPRED 1 2 3 4  
YTRUE 2 2 3 4  
HAMMINGLOSSYTRUE YPRED  
025

IN THE MULTILABEL CASE WITH BINARY LABEL INDICATORS  
HAMMINGLOSSNPARRAY0 1 1 1 NPZEROS2 2  
075

NOTE IN MULTICLASS CLASSIFICATION THE HAMMING LOSS CORRESPONDS TO THE HAMMING DISTANCE BETWEEN YTRUE AND YPRED WHICH IS SIMILAR TO THE ZERO ONE LOSS FUNCTION HOWEVER WHILE ZEROONE LOSS PENALIZES PREDICTION SETS THAT DO NOT STRICTLY MATCH TRUE SETS THE HAMMING LOSS PENALIZES INDIVIDUAL LABELS THUS THE HAMMING LOSS UPPER BOUNDED BY THE ZEROONE LOSS IS ALWAYS BETWEEN ZERO AND ONE INCLUSIVE AND PREDICTING A PROPER SUBSET OR SUPERSSET OF THE TRUE LABELS WILL GIVE A HAMMING LOSS BETWEEN ZERO AND ONE EXCLUSIVE  
PRECISION RECALL AND FMEASURES

INTUITIVELY PRECISION IS THE ABILITY OF THE CLASSIFIER NOT TO LABEL AS POSITIVE A SAMPLE THAT IS NEGATIVE AND RECALL IS THE ABILITY OF THE CLASSIFIER TO FIND ALL THE POSITIVE SAMPLES

THE FMEASURE AND FMEASURES CAN BE INTERPRETED AS A WEIGHTED HARMONIC MEAN OF THE PRECISION AND RECALL A RECALL MEASURE REACHES ITS BEST VALUE AT 1 AND ITS WORST SCORE AT 0 WITH PRECISION AND RECALL ARE EQUIVALENT AND THE RECALL AND THE PRECISION ARE EQUALLY IMPORTANT

THEPRECISIONRECALLCURVE COMPUTES A PRECISIONRECALL CURVE FROM THE GROUND TRUTH LABEL AND A SCORE GIVEN BY THE CLASSIFIER BY VARYING A DECISION THRESHOLD

THEAVERAGEPRECISIONSCORE FUNCTION COMPUTES THE AVERAGE PRECISION AP FROM PREDICTION SCORES THE VALUE IS BETWEEN 0 AND 1 AND HIGHER IS BETTER AP IS DEFINED AS

$AP = \frac{1}{N} \sum_{n=1}^N \frac{P_n + R_n}{2}$   
WHERE P AND R ARE THE PRECISION AND RECALL AT THE NTH THRESHOLD WITH RANDOM PREDICTIONS THE AP IS THE FRACTION OF POSITIVE SAMPLES

REFERENCES MANNING2008 AND EVERINGHAM2010 PRESENT ALTERNATIVE VARIANTS OF AP THAT INTERPOLATE THE PRECISION RECALL CURVE CURRENTLY AVERAGEPRECISIONSCORE DOES NOT IMPLEMENT ANY INTERPOLATED VARIANT REFERENCES DAVIS2006 AND FLACH2015 DESCRIBE WHY A LINEAR INTERPOLATION OF POINTS ON THE PRECISIONRECALL CURVE PROVIDES AN OVERLY OPTIMISTIC MEASURE OF CLASSIFIER PERFORMANCE THIS LINEAR INTERPOLATION IS USED WHEN COMPUTING AREA UNDER THE CURVE WITH THE TRAPEZOIDAL RULE IN AUC

SEVERAL FUNCTIONS ALLOW YOU TO ANALYZE THE PRECISION RECALL AND FMEASURES SCORE  
AVERAGEPRECISIONSCORE YTRUE YSCORE COMPUTE AVERAGE PRECISION AP FROM PREDICTION SCORES  
F1SCORE YTRUE YPRED LABELS COMPUTE THE F1 SCORE ALSO KNOWN AS BALANCED FSCORE OR  
FMEASURE  
FBETASCORE YTRUE YPRED BETA LABELS COMPUTE THE FBETA SCORE  
PRECISIONRECALLCURVE YTRUE PROBASPRED COMPUTE PRECISIONRECALL PAIRS FOR DIFFERENT PROBABILITY THRESHOLDS

CONTINUED ON NEXT PAGE  
33 MODEL SELECTION AND EVALUATION 549

PRECISIONRECALLFSCORESUPPORT YTRUE

YPREDCOMPUTE PRECISION RECALL FMEASURE AND SUPPORT FOR EACH CLASS

PRECISIONSCORE YTRUE YPRED LABELS COMPUTE THE PRECISION

RECALLSCORE YTRUE YPRED LABELS COMPUTE THE RECALL

NOTE THAT THE PRECISIONRECALLCURVE FUNCTION IS RESTRICTED TO THE BINARY CASE THE

AVERAGEPRECISIONSCORE FUNCTION WORKS ONLY IN BINARY CLASSIFICATION AND MULTILABEL INDICATOR FORMAT

EXAMPLES

- SEE CLASSIFICATION OF TEXT DOCUMENTS USING SPARSE FEATURES FOR AN EXAMPLE OF F1SCORE USAGE TO CLASSIFY TEXT DOCUMENTS

- SEE PARAMETER ESTIMATION USING GRID SEARCH WITH CROSSVALIDATION FOR AN EXAMPLE OF PRECISIONSCORE ANDRECALLSCORE USAGE TO ESTIMATE PARAMETERS USING GRID SEARCH WITH NESTED CROSSVALIDATION

- SEE PRECISIONRECALL FOR AN EXAMPLE OF PRECISIONRECALLCURVE USAGE TO EVALUATE CLASSIFIER OUTPUT QUALITY

REFERENCES

BINARY CLASSIFICATION

IN A BINARY CLASSIFICATION TASK THE TERMS “POSITIVE” AND “NEGATIVE” REFER TO THE CLASSIFIER’S PREDICTION AND THE TERMS “TRUE” AND “FALSE” REFER TO WHETHER THAT PREDICTION CORRESPONDS TO THE EXTERNAL JUDGMENT SOMETIMES KNOWN AS THE “OBSERVATION” GIVEN THESE DEFINITIONS WE CAN FORMULATE THE FOLLOWING TABLE

ACTUAL CLASS OBSERVATION

PREDICTED CLASS EXPECTATION TP TRUE POSITIVE CORRECT RESULT FP FALSE POSITIVE UNEXPECTED RESULT

FN FALSE NEGATIVE MISSING RESULT TN TRUE NEGATIVE CORRECT ABSENCE OF RESULT

IN THIS CONTEXT WE CAN DEFINE THE NOTIONS OF PRECISION RECALL AND FMEASURE

PRECISION  $\frac{TP}{TP+FP}$

$\frac{TP}{TP+FP}$

RECALL  $\frac{TP}{TP+FN}$

$\frac{TP}{TP+FN}$

$\frac{1}{\frac{1}{2PRECISION} + \frac{1}{2RECALL}}$

$\frac{2PRECISION RECALL}{PRECISION + RECALL}$

HERE ARE SOME SMALL EXAMPLES IN BINARY CLASSIFICATION

FROM SKLEARN IMPORT METRICS

YPRED 0 1 0 0

YTRUE 0 1 0 1

METRICSPRECISIONSCOREYTRUE YPRED

10

METRICSCOREYTRUE YPRED

05

550 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

METRICSF1SCOREYTRUE YPRED

066

METRICSFBETASCOREYTRUE YPRED BETA05

083

METRICSFBETASCOREYTRUE YPRED BETA1

066

METRICSFBETASCOREYTRUE YPRED BETA2

055

METRICSPRECISIONRECALLFSCORESUPPORTYTRUE YPRED BETA05

ARRAY066 1 ARRAY1 05 ARRAY071 083 ARRAY2

↩→2

IMPORT NUMPY AS NP

FROM SKLEARNMETRICS IMPORT PRECISIONRECALLCURVE

FROM SKLEARNMETRICS IMPORT AVERAGEPRECISIONSCORE

YTRUE NPARRAY0 0 1 1

YSCORES NPARRAY01 04 035 08

PRECISION RECALL THRESHOLD PRECISIONRECALLCURVEYTRUE YSCORES

PRECISION

ARRAY066 05 1 1

RECALL

ARRAY1 05 05 0

THRESHOLD

ARRAY035 04 08

AVERAGEPRECISIONSCOREYTRUE YSCORES

083

MULTICLASS AND MULTILABEL CLASSIFICATION

IN MULTICLASS AND MULTILABEL CLASSIFICATION TASK THE NOTIONS OF PRECISION RECALL AND FMEASURES CAN BE APPLIED TO EACH LABEL INDEPENDENTLY THERE ARE A FEW WAYS TO COMBINE RESULTS ACROSS LABELS SPECIFIED BY THE AVERAGE ARGUMENT TO THE AVERAGEPRECISIONSCORE MULTILABEL ONLY F1SCORE FBETASCORE PRECISIONRECALLFSCORESUPPORT PRECISIONSCORE ANDRECALLSCORE FUNCTIONS AS DESCRIBED

ABOVE NOTE THAT IF ALL LABELS ARE INCLUDED “MICRO” AVERAGING IN A MULTICLASS SETTING WILL PRODUCE PRECISION RECALL AND THAT ARE ALL IDENTICAL TO ACCURACY ALSO NOTE THAT “WEIGHTED” AVERAGING MAY PRODUCE AN FSCORE THAT IS NOT BETWEEN PRECISION AND RECALL

TO MAKE THIS MORE EXPLICIT CONSIDER THE FOLLOWING NOTATION

- $\hat{Y}$  THE SET OF PREDICTED  $\{y_i\}$   $\{y_i\}$  PAIRS
- $Y$  THE SET OF TRUE  $\{y_i\}$   $\{y_i\}$  PAIRS
- $L$  THE SET OF LABELS
- $S$  THE SET OF SAMPLES
- $S_i$  THE SUBSET OF  $S$  WITH SAMPLE  $i$  IE  $\{s \in S \mid s[i] = i\}$
- $S_l$  THE SUBSET OF  $S$  WITH LABEL  $l$
- SIMILARLY  $S_l$  AND  $S_i$  ARE SUBSETS OF  $S$
- $S_i \cap S_l \neq \emptyset$

$S_i$  FOR SOME SETS  $S_i$  AND  $S_l$

- $S_i \cap S_l \neq \emptyset$

CONVENTIONS VARY ON HANDLING  $\emptyset$  THIS IMPLEMENTATION USES  $\emptyset$   $\neq 0$  AND SIMILAR FOR

33 MODEL SELECTION AND EVALUATION 551

SCIKITLEARN USER GUIDE RELEASE 0213

• $\frac{1}{n} \sum_{i=1}^n \frac{1}{\frac{1}{n} \sum_{j=1}^n x_{ij}}$

THEN THE METRICS ARE DEFINED AS

AVERAGE PRECISION RECALL FBETA

MICRO  $\frac{1}{n} \sum_{i=1}^n \frac{1}{\frac{1}{n} \sum_{j=1}^n x_{ij}}$

SAMPLES1

$\sum_{i=1}^n \frac{1}{\frac{1}{n} \sum_{j=1}^n x_{ij}}$

$\sum_{i=1}^n \frac{1}{\frac{1}{n} \sum_{j=1}^n x_{ij}}$

$\sum_{i=1}^n \frac{1}{\frac{1}{n} \sum_{j=1}^n x_{ij}}$

MACRO1

$\sum_{i=1}^n \frac{1}{\frac{1}{n} \sum_{j=1}^n x_{ij}}$

$\sum_{i=1}^n \frac{1}{\frac{1}{n} \sum_{j=1}^n x_{ij}}$

$\sum_{i=1}^n \frac{1}{\frac{1}{n} \sum_{j=1}^n x_{ij}}$

WEIGHTED1 $\sum_{i=1}^n \frac{1}{\frac{1}{n} \sum_{j=1}^n x_{ij}}$

$\sum_{i=1}^n \frac{1}{\frac{1}{n} \sum_{j=1}^n x_{ij}}$

$\sum_{i=1}^n \frac{1}{\frac{1}{n} \sum_{j=1}^n x_{ij}}$

$\sum_{i=1}^n \frac{1}{\frac{1}{n} \sum_{j=1}^n x_{ij}}$

$\sum_{i=1}^n \frac{1}{\frac{1}{n} \sum_{j=1}^n x_{ij}}$

$\sum_{i=1}^n \frac{1}{\frac{1}{n} \sum_{j=1}^n x_{ij}}$

NONE( $\frac{1}{n} \sum_{j=1}^n x_{ij}$ )( $\frac{1}{n} \sum_{j=1}^n x_{ij}$ )( $\frac{1}{n} \sum_{j=1}^n x_{ij}$ )

FROM SKLEARN IMPORT METRICS

YTRUE 0 1 2 0 1 2

YPRED 0 2 1 0 0 1

METRICSPRECISIONSCOREYTRUE YPRED AVERAGEMACRO

022

METRICSPRECISIONSCOREYTRUE YPRED AVERAGEMICRO

033

METRICSF1SCOREYTRUE YPRED AVERAGEWEIGHTED

026

METRICSF1SCOREYTRUE YPRED AVERAGEMACRO BETA05

023

METRICSPRECISIONRECALLFSCORESUPPORTYTRUE YPRED BETA05 AVERAGE NONE

ARRAY066 0 0 ARRAY1 0 0 ARRAY071 0

↔ 0 ARRAY2 2 2

FOR MULTICLASS CLASSIFICATION WITH A “NEGATIVE CLASS” IT IS POSSIBLE TO EXCLUDE SOME LABELS

METRICSPRECISIONSCOREYTRUE YPRED LABELS1 2 AVERAGEMICRO

EXCLUDING 0 NO LABELS WERE CORRECTLY RECALLED

00

SIMILARLY LABELS NOT PRESENT IN THE DATA SAMPLE MAY BE ACCOUNTED FOR IN MACROAVERAGING

METRICSPRECISIONSCOREYTRUE YPRED LABELS0 1 2 3 AVERAGEMACRO

0166

JACCARD SIMILARITY COEFFICIENT SCORE

THEJACCARDScore FUNCTION COMPUTES THE AVERAGE OF JACCARD SIMILARITY COEFFICIENTS ALSO CALLED THE JACCARD INDEX

BETWEEN PAIRS OF LABEL SETS

THE JACCARD SIMILARITY COEFFICIENT OF THE  $i$ TH SAMPLES WITH A GROUND TRUTH LABEL SET  $G$  AND PREDICTED LABEL SET  $P$  IS

DEFINED AS

$$\frac{|G \cap P|}{|G \cup P|}$$

JACCARDScore WORKS LIKEPRECISIONRECALLFSCORESUPPORT AS A NAIVELY SETWISE MEASURE APPLYING

NATIVELY TO BINARY TARGETS AND EXTENDED TO APPLY TO MULTILABEL AND MULTICLASS THROUGH THE USE OF FROM BINARY TO

MULTICLASS AND MULTILABEL SEEABOVE

IN THE BINARY CASE

552 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

```
import numpy as np
from sklearn.metrics import jaccard_score
y_true = np.array([1, 1, 0])
y_pred = np.array([1, 1, 1])
jaccard_score(y_true, y_pred)
0.6666
```

IN THE MULTILABEL CASE WITH BINARY LABEL INDICATORS

```
jaccard_score(y_true, y_pred, average='samples')
0.5833
jaccard_score(y_true, y_pred, average='macro')
0.6666
jaccard_score(y_true, y_pred, average='none')
array([0.5, 0.5, 1.])
```

MULTICLASS PROBLEMS ARE BINARIZED AND TREATED LIKE THE CORRESPONDING MULTILABEL PROBLEM

```
y_pred = [0, 2, 1, 2]
y_true = [0, 1, 2, 2]
jaccard_score(y_true, y_pred, average='none')
array([1. , 0. , 0.33])
```

JACCARDSCOREYTRUE YPRED AVERAGEMACRO 044

JACCARDSCOREYTRUE YPRED AVERAGEMICRO 033

HINGE LOSS

THEHINGELOSS FUNCTION COMPUTES THE AVERAGE DISTANCE BETWEEN THE MODEL AND THE DATA USING HINGE LOSS A ONE SIDED METRIC THAT CONSIDERS ONLY PREDICTION ERRORS HINGE LOSS IS USED IN MAXIMAL MARGIN CLASSIFIERS SUCH AS SUPPORT VECTOR MACHINES

IF THE LABELS ARE ENCODED WITH 1 AND -1  $y$  IS THE TRUE VALUE AND  $\hat{y}$  IS THE PREDICTED DECISIONS AS OUTPUT BY DECISIONFUNCTION THEN THE HINGE LOSS IS DEFINED AS

$$L(y, \hat{y}) = \max(0, 1 - y\hat{y})$$

IF THERE ARE MORE THAN TWO LABELS HINGELOSS USES A MULTICLASS VARIANT DUE TO CRAMMER SINGER HERE IS THE PAPER DESCRIBING IT

IF  $\hat{y}$  IS THE PREDICTED DECISION FOR TRUE LABEL AND  $\max_{k \neq y} \hat{y}_k$  IS THE MAXIMUM OF THE PREDICTED DECISIONS FOR ALL OTHER LABELS WHERE PREDICTED DECISIONS ARE OUTPUT BY DECISION FUNCTION THEN MULTICLASS HINGE LOSS IS DEFINED BY

$$L(y, \hat{y}) = \max(0, 1 - \hat{y}_y + \max_{k \neq y} \hat{y}_k)$$

HERE A SMALL EXAMPLE DEMONSTRATING THE USE OF THE HINGELOSS FUNCTION WITH A SVM CLASSIFIER IN A BINARY CLASS PROBLEM

```
from sklearn import svm
from sklearn.metrics import hinge_loss
X = [0, 1]
Y = [1, 1]
```

33 MODEL SELECTION AND EVALUATION 553

SCIKITLEARN USER GUIDE RELEASE 0213  
EST SVMLINEARSVCRANDOMSTATE0  
ESTFITX Y  
LINEARSVCC10 CLASSWEIGHTNONE DUALTRUE FITINTERCEPTTRUE  
INTERCEPTSCALING1 LOSSSQUAREDHINGE MAXITER1000  
MULTICLASSOVR PENALTYL2 RANDOMSTATE0 TOL00001  
VERBOSE0  
PREDDCISION ESTDECISIONFUNCTION2 3 05  
PREDDCISION  
ARRAY218 236 009  
HINGELOSS1 1 1 PREDDCISION  
03  
HERE IS AN EXAMPLE DEMONSTRATING THE USE OF THE HINGELOSS FUNCTION WITH A SVM CLASSIFIER IN A MULTICLASS PROBLEM  
X NPARRAY0 1 2 3  
Y NPARRAY0 1 2 3  
LABELS NPARRAY0 1 2 3  
EST SVMLINEARSVC  
ESTFITX Y  
LINEARSVCC10 CLASSWEIGHTNONE DUALTRUE FITINTERCEPTTRUE  
INTERCEPTSCALING1 LOSSSQUAREDHINGE MAXITER1000  
MULTICLASSOVR PENALTYL2 RANDOMSTATENONE TOL00001  
VERBOSE0  
PREDDCISION ESTDECISIONFUNCTION1 2 3  
YTRUE 0 2 3  
HINGELOSSYTRUE PREDDCISION LABELS  
056  
LOG LOSS  
LOG LOSS ALSO CALLED LOGISTIC REGRESSION LOSS OR CROSSENTROPY LOSS IS DEFINED ON PROBABILITY ESTIMATES IT IS COMMONLY  
USED IN MULTINOMIAL LOGISTIC REGRESSION AND NEURAL NETWORKS AS WELL AS IN SOME VARIANTS OF EXPECTATIONMAXIMIZATION  
AND CAN BE USED TO EVALUATE THE PROBABILITY OUTPUTS PREDICTPROBA OF A CLASSIFIER INSTEAD OF ITS DISCRETE PREDIC  
TIONS  
FOR BINARY CLASSIFICATION WITH A TRUE LABEL  $y \in \{0,1\}$  AND A PROBABILITY ESTIMATE  $p$   $PR$   $p$  1 THE LOG LOSS PER SAMPLE  
IS THE NEGATIVE LOGLIKELIHOOD OF THE CLASSIFIER GIVEN THE TRUE LABEL  
$$-\log p$$
  
THIS EXTENDS TO THE MULTICLASS CASE AS FOLLOWS LET THE TRUE LABELS FOR A SET OF SAMPLES BE ENCODED AS A 10FK BINARY  
INDICATOR MATRIX  $Y$  IE  $Y_{ij} = 1$  IF SAMPLE  $i$  HAS LABEL  $j$  TAKEN FROM A SET OF  $K$  LABELS LET  $P$  BE A MATRIX OF PROBABILITY  
ESTIMATES WITH  $P_{ij}$   $PR$   $P_{ij}$  1 THEN THE LOG LOSS OF THE WHOLE SET IS  
$$-\log p$$
  
$$-\log p$$
  
$$-\log p$$
  
TO SEE HOW THIS GENERALIZES THE BINARY LOG LOSS GIVEN ABOVE NOTE THAT IN THE BINARY CASE  $y \in \{0,1\}$  AND  $p \in [0,1]$   
1  $-p$  1 SO EXPANDING THE INNER SUM OVER  $y \in \{0,1\}$  GIVES THE BINARY LOG LOSS  
THE LOGLOSS FUNCTION COMPUTES LOG LOSS GIVEN A LIST OF GROUNDTRUTH LABELS AND A PROBABILITY MATRIX AS RETURNED BY  
AN ESTIMATOR'S PREDICTPROBA METHOD  
FROM SKLEARNMETRICS IMPORT LOGLOSS  
YTRUE 0 0 1 1  
554 CHAPTER 3 USER GUIDE



YPRED 9 1 8 2 3 7 01 99

LOGLOSSYTRUE YPRED

01738

THE FIRST9 1 INYPRED DENOTES 90 PROBABILITY THAT THE FIRST SAMPLE HAS LABEL 0 THE LOG LOSS IS NONNEGATIVE MATTHEWS CORRELATION COEFFICIENT

THEMATTHEWSCORRCOEF FUNCTION COMPUTES THE MATTHEW’S CORRELATION COEFFICIENT MCC FOR BINARY CLASSES QUOTING WIKIPEDIA

“THE MATTHEWS CORRELATION COEFFICIENT IS USED IN MACHINE LEARNING AS A MEASURE OF THE QUALITY OF BINARY TWOCLASS CLASSIFICATIONS IT TAKES INTO ACCOUNT TRUE AND FALSE POSITIVES AND NEGATIVES AND IS GENERALLY REGARDED AS A BALANCED MEASURE WHICH CAN BE USED EVEN IF THE CLASSES ARE OF VERY DIFFERENT SIZES THE MCC IS IN ESSENCE A CORRELATION COEFFICIENT VALUE BETWEEN -1 AND 1 A COEFFICIENT OF 1 REPRESENTS A PERFECT PREDICTION 0 AN AVERAGE RANDOM PREDICTION AND -1 AN INVERSE PREDICTION THE STATISTIC IS ALSO KNOWN AS THE PHI COEFFICIENT”

IN THE BINARY TWOCLASS CASE TP AND FN ARE RESPECTIVELY THE NUMBER OF TRUE POSITIVES TRUE NEGATIVES FALSE POSITIVES AND FALSE NEGATIVES THE MCC IS DEFINED AS

$$\frac{TP \times TN - FP \times FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}}$$

IN THE MULTICLASS CASE THE MATTHEWS CORRELATION COEFFICIENT CAN BE DEFINED IN TERMS OF A CONFUSIONMATRIX FOR CLASSES TO SIMPLIFY THE DEFINITION CONSIDER THE FOLLOWING INTERMEDIATE VARIABLES

• TP

THE NUMBER OF TIMES CLASS 0 TRULY OCCURRED

• PP

THE NUMBER OF TIMES CLASS 0 WAS PREDICTED

• CP

THE TOTAL NUMBER OF SAMPLES CORRECTLY PREDICTED

• N

THE TOTAL NUMBER OF SAMPLES

THEN THE MULTICLASS MCC IS DEFINED AS

$$\frac{CP - \sum_i PP_i^2}{\sqrt{N^2 - \sum_i PP_i^2}}$$

$$\frac{CP^2 - \sum_i PP_i^2}{N^2 - \sum_i PP_i^2}$$

$$\frac{CP^2 - \sum_i PP_i^2}{N^2 - \sum_i PP_i^2}$$

WHEN THERE ARE MORE THAN TWO LABELS THE VALUE OF THE MCC WILL NO LONGER RANGE BETWEEN -1 AND 1 INSTEAD THE MINIMUM VALUE WILL BE SOMEWHERE BETWEEN -1 AND 0 DEPENDING ON THE NUMBER AND DISTRIBUTION OF GROUND TRUE LABELS THE MAXIMUM VALUE IS ALWAYS 1

HERE IS A SMALL EXAMPLE ILLUSTRATING THE USAGE OF THE MATTHEWSCORRCOEF FUNCTION

FROM SKLEARNMETRICS IMPORT MATTHEWSCORRCOEF

YTRUE 1 1 1 1

YPRED 1 1 1 1

MATTHEWSCORRCOEFYTRUE YPRED

0.33

MULTILABEL CONFUSION MATRIX

THEMULTILABELCONFUSIONMATRIX FUNCTION COMPUTES CLASSWISE DEFAULT OR SAMPLEWISE SAMPLE

WISETRUE MULTILABEL CONFUSION MATRIX TO EVALUATE THE ACCURACY OF A CLASSIFICATION MULTILABELCONFUSIONMATRIX

33 MODEL SELECTION AND EVALUATION 555

SCIKITLEARN USER GUIDE RELEASE 0213

ALSO TREATS MULTICLASS DATA AS IF IT WERE MULTILABEL AS THIS IS A TRANSFORMATION COMMONLY APPLIED TO EVALUATE MULTICLASS PROBLEMS WITH BINARY CLASSIFICATION METRICS SUCH AS PRECISION RECALL ETC

WHEN CALCULATING CLASSWISE MULTILABEL CONFUSION MATRIX  $\hat{y}$  THE COUNT OF TRUE NEGATIVES FOR CLASS  $i$  IS  $\sum_j (1 - y_{ij}) (1 - \hat{y}_{ij})$  FALSE NEGATIVES IS  $\sum_j (1 - y_{ij}) y_{ij}$  TRUE POSITIVES IS  $\sum_j y_{ij} \hat{y}_{ij}$  AND FALSE POSITIVES IS  $\sum_j \hat{y}_{ij} (1 - y_{ij})$

HERE IS AN EXAMPLE DEMONSTRATING THE USE OF THE MULTILABELCONFUSIONMATRIX FUNCTION WITH MULTILABEL INDICATOR MATRIX INPUT

```
import numpy as np
from sklearn.metrics import multilabel_confusion_matrix

y_true = np.array([0, 1, 0, 1, 0, 1, 0, 1, 0, 1])
y_pred = np.array([0, 0, 1, 0, 0, 1, 1, 0, 1, 0])

multilabel_confusion_matrix(y_true, y_pred)
```

OR A CONFUSION MATRIX CAN BE CONSTRUCTED FOR EACH SAMPLE'S LABELS

```
multilabel_confusion_matrix(y_true, y_pred, samplewise=True)
```

HERE IS AN EXAMPLE DEMONSTRATING THE USE OF THE MULTILABELCONFUSIONMATRIX FUNCTION WITH MULTICLASS INPUT

```
y_true = ['CAT', 'ANT', 'CAT', 'CAT', 'ANT', 'BIRD']
y_pred = ['ANT', 'ANT', 'CAT', 'CAT', 'ANT', 'CAT']

multilabel_confusion_matrix(y_true, y_pred, labels=['ANT', 'BIRD', 'CAT'])
```

HERE ARE SOME EXAMPLES DEMONSTRATING THE USE OF THE MULTILABELCONFUSIONMATRIX FUNCTION TO CALCULATE RECALL OR SENSITIVITY SPECIFICITY FALL OUT AND MISS RATE FOR EACH CLASS IN A PROBLEM WITH MULTILABEL INDICATOR MATRIX INPUT

CALCULATING RECALL ALSO CALLED THE TRUE POSITIVE RATE OR THE SENSITIVITY FOR EACH CLASS

```
y_true = np.array([0, 1, 0, 1, 1, 0])
y_pred = np.array([0, 1, 0, 1, 1, 0])
```

556 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

YPRED NPARRAY0 1 0

0 0 1

1 1 0

MCM MULTILABELCONFUSIONMATRIXTRUE YPRED

TN MCM 0 0

TP MCM 1 1

FN MCM 1 0

FP MCM 0 1

TP TP FN

ARRAY1 05 0

CALCULATING SPECIFICITY ALSO CALLED THE TRUE NEGATIVE RATE FOR EACH CLASS

TN TN FP

ARRAY1 0 05

CALCULATING FALL OUT ALSO CALLED THE FALSE POSITIVE RATE FOR EACH CLASS

FP FP TN

ARRAY0 1 05

CALCULATING MISS RATE ALSO CALLED THE FALSE NEGATIVE RATE FOR EACH CLASS

FN FN TP

ARRAY0 05 1

RECEIVER OPERATING CHARACTERISTIC ROC

THE FUNCTION ROCCURVE COMPUTES THE RECEIVER OPERATING CHARACTERISTIC CURVE OR ROC CURVE QUOTING WIKIPEDIA

“A RECEIVER OPERATING CHARACTERISTIC ROC OR SIMPLY ROC CURVE IS A GRAPHICAL PLOT WHICH ILLUSTRATES THE PERFORMANCE OF A BINARY CLASSIFIER SYSTEM AS ITS DISCRIMINATION THRESHOLD IS VARIED IT IS CREATED BY PLOTTING THE FRACTION OF TRUE POSITIVES OUT OF THE POSITIVES TPR TRUE POSITIVE RATE VS THE FRACTION OF FALSE POSITIVES OUT OF THE NEGATIVES FPR FALSE POSITIVE RATE AT VARIOUS THRESHOLD SETTINGS TPR IS ALSO KNOWN AS SENSITIVITY AND FPR IS ONE MINUS THE SPECIFICITY OR TRUE NEGATIVE RATE”

THIS FUNCTION REQUIRES THE TRUE BINARY VALUE AND THE TARGET SCORES WHICH CAN EITHER BE PROBABILITY ESTIMATES OF THE POSITIVE CLASS CONFIDENCE VALUES OR BINARY DECISIONS HERE IS A SMALL EXAMPLE OF HOW TO USE THE ROCCURVE FUNCTION

IMPORT NUMPY AS NP

FROM SKLEARNMETRICS IMPORT ROCCURVE

Y NPARRAY1 1 2 2

SCORES NPARRAY01 04 035 08

FPR TPR THRESHOLDS ROCCURVEY SCORES POSLABEL2

FPR

ARRAY0 0 05 05 1

TPR

ARRAY0 05 05 1 1

THRESHOLDS

ARRAY18 08 04 035 01

33 MODEL SELECTION AND EVALUATION 557

SCIKITLEARN USER GUIDE RELEASE 0213

THIS FIGURE SHOWS AN EXAMPLE OF SUCH AN ROC CURVE

THEROCAUCSCORE FUNCTION COMPUTES THE AREA UNDER THE RECEIVER OPERATING CHARACTERISTIC ROC CURVE WHICH IS ALSO DENOTED BY AUC OR AUROC BY COMPUTING THE AREA UNDER THE ROC CURVE THE CURVE INFORMATION IS SUMMARIZED IN ONE NUMBER FOR MORE INFORMATION SEE THE WIKIPEDIA ARTICLE ON AUC

```
import numpy as np
from sklearn.metrics import roc_auc_score
y_true = np.array([0, 1, 1])
y_scores = np.array([0.04, 0.35, 0.8])
roc_auc_score(y_true, y_scores)
```

0.75

IN MULTILABEL CLASSIFICATION THE ROCAUCSCORE FUNCTION IS EXTENDED BY AVERAGING OVER THE LABELS AS ABOVE COMPARED TO METRICS SUCH AS THE SUBSET ACCURACY THE HAMMING LOSS OR THE F1 SCORE ROC DOESN'T REQUIRE OPTIMIZING A THRESHOLD FOR EACH LABEL THE ROCAUCSCORE FUNCTION CAN ALSO BE USED IN MULTICLASS CLASSIFICATION IF THE PREDICTED OUTPUTS HAVE BEEN BINARIZED

IN APPLICATIONS WHERE A HIGH FALSE POSITIVE RATE IS NOT TOLERABLE THE PARAMETER MAXFPR OF ROCAUCSCORE CAN BE USED TO SUMMARIZE THE ROC CURVE UP TO THE GIVEN LIMIT

558 CHAPTER 3 USER GUIDE

EXAMPLES

- SEE RECEIVER OPERATING CHARACTERISTIC ROC FOR AN EXAMPLE OF USING ROC TO EVALUATE THE QUALITY OF THE OUTPUT OF A CLASSIFIER
- SEE RECEIVER OPERATING CHARACTERISTIC ROC WITH CROSS VALIDATION FOR AN EXAMPLE OF USING ROC TO EVALUATE CLASSIFIER OUTPUT QUALITY USING CROSSVALIDATION
- SEE SPECIES DISTRIBUTION MODELING FOR AN EXAMPLE OF USING ROC TO MODEL SPECIES DISTRIBUTION

ZERO ONE LOSS

THEZEROONELOSS FUNCTION COMPUTES THE SUM OR THE AVERAGE OF THE 01 CLASSIFICATION LOSS  $\{0-1\}$  OVER  $n$ SAMPLES BY DEFAULT THE FUNCTION NORMALIZES OVER THE SAMPLE TO GET THE SUM OF THE  $\{0-1\}$  SETNORMALIZE TOFALSE IN MULTILABEL CLASSIFICATION THE ZEROONELOSS SCORES A SUBSET AS ONE IF ITS LABELS STRICTLY MATCH THE PREDICTIONS AND AS A ZERO IF THERE ARE ANY ERRORS BY DEFAULT THE FUNCTION RETURNS THE PERCENTAGE OF IMPERFECTLY PREDICTED SUBSETS TO GET THE COUNT OF SUCH SUBSETS INSTEAD SET NORMALIZE TOFALSE IF  $\hat{y}_i$ IS THE PREDICTED VALUE OF THE  $i$ TH SAMPLE AND  $y_i$ IS THE CORRESPONDING TRUE VALUE THEN THE 01 LOSS  $\{0-1\}$ IS DEFINED AS

$$\{0-1\} = \frac{1}{n} \sum_{i=1}^n 1_{\hat{y}_i \neq y_i}$$

WHERE  $1_{\cdot}$ IS THE INDICATOR FUNCTION

```
FROM SKLEARNMETRICS IMPORT ZEROONELOSS
```

```
YPRED 1 2 3 4
```

```
YTRUE 2 2 3 4
```

```
ZEROONELOSSYTRUE YPRED
```

025

SCIKITLEARN USER GUIDE RELEASE 0213  
ZEROONELOSSYTRUE YPRED NORMALIZE FALSE

1  
IN THE MULTILABEL CASE WITH BINARY LABEL INDICATORS WHERE THE FIRST LABEL SET 01 HAS AN ERROR  
ZEROONELOSSNPARRAY0 1 1 1 NPONES2 2  
05

ZEROONELOSSNPARRAY0 1 1 1 NPONES2 2 NORMALIZE FALSE  
1

EXAMPLE  
• SEE RECURSIVE FEATURE ELIMINATION WITH CROSSVALIDATION FOR AN EXAMPLE OF ZERO ONE LOSS USAGE TO PERFORM  
RECURSIVE FEATURE ELIMINATION WITH CROSSVALIDATION

BRIER SCORE LOSS  
THEBRIERSCORELOSS FUNCTION COMPUTES THE BRIER SCORE FOR BINARY CLASSES QUOTING WIKIPEDIA  
“THE BRIER SCORE IS A PROPER SCORE FUNCTION THAT MEASURES THE ACCURACY OF PROBABILISTIC PREDICTIONS IT IS  
APPLICABLE TO TASKS IN WHICH PREDICTIONS MUST ASSIGN PROBABILITIES TO A SET OF MUTUALLY EXCLUSIVE DISCRETE  
OUTCOMES”  
THIS FUNCTION RETURNS A SCORE OF THE MEAN SQUARE DIFFERENCE BETWEEN THE ACTUAL OUTCOME AND THE PREDICTED PROBABILITY  
OF THE POSSIBLE OUTCOME THE ACTUAL OUTCOME HAS TO BE 1 OR 0 TRUE OR FALSE WHILE THE PREDICTED PROBABILITY OF THE  
ACTUAL OUTCOME CAN BE A VALUE BETWEEN 0 AND 1  
THE BRIER SCORE LOSS IS ALSO BETWEEN 0 TO 1 AND THE LOWER THE SCORE THE MEAN SQUARE DIFFERENCE IS SMALLER THE MORE  
ACCURATE THE PREDICTION IS IT CAN BE THOUGHT OF AS A MEASURE OF THE “CALIBRATION” OF A SET OF PROBABILISTIC PREDICTIONS

$$\frac{1}{n} \sum_{i=1}^n (p_i - y_i)^2$$
  
WHERE  $n$  IS THE TOTAL NUMBER OF PREDICTIONS  $p_i$  IS THE PREDICTED PROBABILITY OF THE ACTUAL OUTCOME  $y_i$   
HERE IS A SMALL EXAMPLE OF USAGE OF THIS FUNCTION

```
import numpy as np
from sklearn.metrics import brier_score_loss

y_true = np.array([1, 1, 0])
y_true_categorical = np.array(['spam', 'ham', 'ham', 'spam'])
y_prob = np.array([0.1, 0.9, 0.8, 0.4])
y_pred = np.array([1, 1, 0])

brier_score_loss(y_true, y_prob)
0.055
brier_score_loss(y_true, y_prob, pos_label=0)
0.055
brier_score_loss(y_true_categorical, y_prob, pos_label='ham')
0.055
brier_score_loss(y_true, y_prob, 0.5)
0.0
```

560 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLE

- SEE PROBABILITY CALIBRATION OF CLASSIFIERS FOR AN EXAMPLE OF BRIER SCORE LOSS USAGE TO PERFORM PROBABILITY CALIBRATION OF CLASSIFIERS

REFERENCES

- G BRIER VERIFICATION OF FORECASTS EXPRESSED IN TERMS OF PROBABILITY MONTHLY WEATHER REVIEW 781 1950

MULTILABEL RANKING METRICS

IN MULTILABEL LEARNING EACH SAMPLE CAN HAVE ANY NUMBER OF GROUND TRUTH LABELS ASSOCIATED WITH IT THE GOAL IS TO GIVE HIGH SCORES AND BETTER RANK TO THE GROUND TRUTH LABELS

COVERAGE ERROR

THECOVERAGEERROR FUNCTION COMPUTES THE AVERAGE NUMBER OF LABELS THAT HAVE TO BE INCLUDED IN THE FINAL PREDICTION SUCH THAT ALL TRUE LABELS ARE PREDICTED THIS IS USEFUL IF YOU WANT TO KNOW HOW MANY TOPSCOREDLABELS YOU HAVE TO PREDICT IN AVERAGE WITHOUT MISSING ANY TRUE ONE THE BEST VALUE OF THIS METRICS IS 1 THE AVERAGE NUMBER OF TRUE LABELS

NOTE OUR IMPLEMENTATION’S SCORE IS 1 GREATER THAN THE ONE GIVEN IN TSOU MAKAS ET AL 2010 THIS EXTENDS IT TO HANDLE THE DEGENERATE CASE IN WHICH AN INSTANCE HAS 0 TRUE LABELS

FORMALLY GIVEN A BINARY INDICATOR MATRIX OF THE GROUND TRUTH LABELS  $y \in \{0,1\}^{SAMPLES \times LABELS}$  AND THE SCORE ASSOCIATED WITH EACH LABEL  $j \in \{1, \dots, LABELS\}$  THE COVERAGE IS DEFINED AS

$$coverage_j = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{y_{ij} = 1}$$
$$coverage = \frac{1}{n} \sum_{j=1}^L coverage_j$$
$$coverage_{max} = \max_j coverage_j$$

WITH RANK  $r_j$  GIVEN THE RANK DEFINITION TIES IN YSCORES ARE BROKEN BY GIVING THE MAXIMAL RANK THAT WOULD HAVE BEEN ASSIGNED TO ALL TIED VALUES

HERE IS A SMALL EXAMPLE OF USAGE OF THIS FUNCTION

```
import numpy as np
from sklearn.metrics import coverage_error
y_true = np.array([0, 0, 0, 0, 1])
y_score = np.array([0.75, 0.5, 1, 1, 0.2])
coverage_error(y_true, y_score)
```

25 LABEL RANKING AVERAGE PRECISION

THELABELRANKINGAVERAGEPRECISIONSCORE FUNCTION IMPLEMENTS LABEL RANKING AVERAGE PRECISION LRAP THIS METRIC IS LINKED TO THE AVERAGEPRECISIONSCORE FUNCTION BUT IS BASED ON THE NOTION OF LABEL RANKING INSTEAD OF PRECISION AND RECALL

SCIKITLEARN USER GUIDE RELEASE 0213

LABEL RANKING AVERAGE PRECISION LRAP AVERAGES OVER THE SAMPLES THE ANSWER TO THE FOLLOWING QUESTION FOR EACH GROUND TRUTH LABEL WHAT FRACTION OF HIGHERRANKED LABELS WERE TRUE LABELS THIS PERFORMANCE MEASURE WILL BE HIGHER IF YOU ARE ABLE TO GIVE BETTER RANK TO THE LABELS ASSOCIATED WITH EACH SAMPLE THE OBTAINED SCORE IS ALWAYS STRICTLY GREATER THAN 0 AND THE BEST VALUE IS 1 IF THERE IS EXACTLY ONE RELEVANT LABEL PER SAMPLE LABEL RANKING AVERAGE PRECISION IS EQUIVALENT TO THE MEAN RECIPROCAL RANK

FORMALLY GIVEN A BINARY INDICATOR MATRIX OF THE GROUND TRUTH LABELS  $\mathbf{Y} \in \{0,1\}^{SAMPLES \times LABELS}$  AND THE SCORE ASSOCIATED WITH EACH LABEL  $\mathbf{Y} \in \mathbb{R}^{SAMPLES \times LABELS}$  THE AVERAGE PRECISION IS DEFINED AS

$$\frac{1}{SAMPLES} \sum_{i=1}^{SAMPLES} \frac{1}{RANK(Y[i,:])}$$

WHERE  $RANK(Y[i,:])$  IS THE RANK OF THE  $i$ TH ROW OF  $Y$  IN DESCENDING ORDER OF THE SCORES ASSOCIATED WITH EACH LABEL

FORMALLY GIVEN A BINARY INDICATOR MATRIX OF THE GROUND TRUTH LABELS  $\mathbf{Y} \in \{0,1\}^{SAMPLES \times LABELS}$  AND THE SCORE ASSOCIATED WITH EACH LABEL  $\mathbf{Y} \in \mathbb{R}^{SAMPLES \times LABELS}$  THE RANKING LOSS IS DEFINED AS

HERE IS A SMALL EXAMPLE OF USAGE OF THIS FUNCTION

```
import numpy as np
from sklearn.metrics import label_ranking_average_precision_score
y_true = np.array([0, 0, 0, 0, 1])
y_score = np.array([0.75, 0.5, 1, 1, 0.2])
label_ranking_average_precision_score(y_true, y_score)
```

RANKING LOSS

THE LABEL RANKING LOSS FUNCTION COMPUTES THE RANKING LOSS WHICH AVERAGES OVER THE SAMPLES THE NUMBER OF LABEL PAIRS THAT ARE INCORRECTLY ORDERED IE TRUE LABELS HAVE A LOWER SCORE THAN FALSE LABELS WEIGHTED BY THE INVERSE OF THE NUMBER OF ORDERED PAIRS OF FALSE AND TRUE LABELS THE LOWEST ACHIEVABLE RANKING LOSS IS ZERO

FORMALLY GIVEN A BINARY INDICATOR MATRIX OF THE GROUND TRUTH LABELS  $\mathbf{Y} \in \{0,1\}^{SAMPLES \times LABELS}$  AND THE SCORE ASSOCIATED WITH EACH LABEL  $\mathbf{Y} \in \mathbb{R}^{SAMPLES \times LABELS}$  THE RANKING LOSS IS DEFINED AS

$$\frac{1}{SAMPLES} \sum_{i=1}^{SAMPLES} \frac{1}{RANK(Y[i,:])}$$

WHERE  $RANK(Y[i,:])$  IS THE RANK OF THE  $i$ TH ROW OF  $Y$  IN DESCENDING ORDER OF THE SCORES ASSOCIATED WITH EACH LABEL

FORMALLY GIVEN A BINARY INDICATOR MATRIX OF THE GROUND TRUTH LABELS  $\mathbf{Y} \in \{0,1\}^{SAMPLES \times LABELS}$  AND THE SCORE ASSOCIATED WITH EACH LABEL  $\mathbf{Y} \in \mathbb{R}^{SAMPLES \times LABELS}$  THE RANKING LOSS IS DEFINED AS

HERE IS A SMALL EXAMPLE OF USAGE OF THIS FUNCTION

```
import numpy as np
from sklearn.metrics import label_ranking_loss
y_true = np.array([0, 0, 0, 0, 1])
y_score = np.array([0.75, 0.5, 1, 1, 0.2])
label_ranking_loss(y_true, y_score)
```

WITH THE FOLLOWING PREDICTION WE HAVE PERFECT AND MINIMAL LOSS

```
y_true = np.array([0, 1, 0, 2, 1, 0, 2, 0, 9])
y_score = np.array([0, 1, 0, 2, 1, 0, 2, 0, 9])
label_ranking_loss(y_true, y_score)
```



SCIKITLEARN USER GUIDE RELEASE 0213

REFERENCES

• TSOU MAKAS G KATAKIS I VLAHAVAS I 2010 MINING MULTILABEL DATA IN DATA MINING AND KNOWLEDGE

DISCOVERY HANDBOOK PP 667685 SPRINGER US

REGRESSION METRICS

THE SKLEARN METRICS MODULE IMPLEMENTS SEVERAL LOSS SCORE AND UTILITY FUNCTIONS TO MEASURE REGRESSION

PERFORMANCE SOME OF THOSE HAVE BEEN ENHANCED TO HANDLE THE MULTIOUTPUT CASE MEANS SQUARED ERROR

MEAN ABSOLUTE ERROR EXPLAINED VARIANCE SCORE AND R<sup>2</sup> SCORE

THESE FUNCTIONS HAVE AN MULTIOUTPUT KEYWORD ARGUMENT WHICH SPECIFIES THE WAY THE SCORES OR LOSSES FOR EACH

INDIVIDUAL TARGET SHOULD BE AVERAGED THE DEFAULT IS UNIFORM AVERAGE WHICH SPECIFIES A UNIFORMLY WEIGHTED

MEAN OVER OUTPUTS IF AN NDARRAY OF SHAPEN OUTPUTS IS PASSED THEN ITS ENTRIES ARE INTERPRETED AS WEIGHTS

AND AN ACCORDING WEIGHTED AVERAGE IS RETURNED IF MULTIOUTPUT IS RAW VALUES IS SPECIFIED THEN ALL UNALTERED

INDIVIDUAL SCORES OR LOSSES WILL BE RETURNED IN AN ARRAY OF SHAPE N OUTPUTS

THE R<sup>2</sup> SCORE AND EXPLAINED VARIANCE SCORE ACCEPT AN ADDITIONAL VALUE VARIANCE WEIGHTED FOR

THE MULTIOUTPUT PARAMETER THIS OPTION LEADS TO A WEIGHTING OF EACH INDIVIDUAL SCORE BY THE VARIANCE OF THE

CORRESPONDING TARGET VARIABLE THIS SETTING QUANTIFIES THE GLOBALLY CAPTURED UNSCALED VARIANCE IF THE TARGET VARI

ABLES ARE OF DIFFERENT SCALE THEN THIS SCORE PUTS MORE IMPORTANCE ON WELL EXPLAINING THE HIGHER VARIANCE VARIABLES

MULTIOUTPUT VARIANCE WEIGHTED IS THE DEFAULT VALUE FOR R<sup>2</sup> SCORE FOR BACKWARD COMPATIBILITY THIS

WILL BE CHANGED TO UNIFORM AVERAGE IN THE FUTURE

EXPLAINED VARIANCE SCORE

THE EXPLAINED VARIANCE SCORE COMPUTES THE EXPLAINED VARIANCE REGRESSION SCORE

IF  $\hat{y}$  IS THE ESTIMATED TARGET OUTPUT  $y$  THE CORRESPONDING CORRECT TARGET OUTPUT AND  $\sigma^2$  IS VARIANCE THE SQUARE OF THE

STANDARD DEVIATION THEN THE EXPLAINED VARIANCE IS ESTIMATED AS FOLLOW

$$\frac{(\hat{y} - y)^2}{\sigma^2} \quad 1 - \frac{(\hat{y} - y)^2}{\sigma^2}$$

$\sigma^2$

THE BEST POSSIBLE SCORE IS 10 LOWER VALUES ARE WORSE

HERE IS A SMALL EXAMPLE OF USAGE OF THE EXPLAINED VARIANCE SCORE FUNCTION

FROM SKLEARN METRICS IMPORT EXPLAINED VARIANCE SCORE

Y TRUE 3 05 2 7

Y PRED 25 00 2 8

EXPLAINED VARIANCE SCORE Y TRUE Y PRED

0957

Y TRUE 05 1 1 1 7 6

Y PRED 0 2 1 2 8 5

EXPLAINED VARIANCE SCORE Y TRUE Y PRED MULTIOUTPUT RAW VALUES

ARRAY 0967 1

EXPLAINED VARIANCE SCORE Y TRUE Y PRED MULTIOUTPUT 03 07

0990

33 MODEL SELECTION AND EVALUATION 563

MAX ERROR  
THEMAXERROR FUNCTION COMPUTES THE MAXIMUM RESIDUAL ERROR A METRIC THAT CAPTURES THE WORST CASE ERROR BETWEEN THE PREDICTED VALUE AND THE TRUE VALUE IN A PERFECTLY FITTED SINGLE OUTPUT REGRESSION MODEL MAXERROR WOULD BE0 ON THE TRAINING SET AND THOUGH THIS WOULD BE HIGHLY UNLIKELY IN THE REAL WORLD THIS METRIC SHOWS THE EXTENT OF ERROR THAT THE MODEL HAD WHEN IT WAS FITTED  
IF $\hat{y}_i$ IS THE PREDICTED VALUE OF THE  $i$ TH SAMPLE AND  $y_i$ IS THE CORRESPONDING TRUE VALUE THEN THE MAX ERROR IS DEFINED AS  
MAX ERROR  $\max_i |\hat{y}_i - y_i|$   
HERE IS A SMALL EXAMPLE OF USAGE OF THE MAXERROR FUNCTION  
FROM SKLEARNMETRICS IMPORT MAXERROR  
YTRUE 3 2 7 1  
YPRED 9 2 7 1  
MAXERRORYTRUE YPRED  
6

THEMAXERROR DOES NOT SUPPORT MULTIOUTPUT  
MEAN ABSOLUTE ERROR  
THEMEANABSOLUTEERROR FUNCTION COMPUTES MEAN ABSOLUTE ERROR A RISK METRIC CORRESPONDING TO THE EXPECTED VALUE OF THE ABSOLUTE ERROR LOSS OR  $L_1$ NORM LOSS  
IF $\hat{y}_i$ IS THE PREDICTED VALUE OF THE  $i$ TH SAMPLE AND  $y_i$ IS THE CORRESPONDING TRUE VALUE THEN THE MEAN ABSOLUTE ERROR MAE ESTIMATED OVER  $n$ SAMPLES IS DEFINED AS  
MAE  $\frac{1}{n} \sum_{i=1}^n |\hat{y}_i - y_i|$

HERE IS A SMALL EXAMPLE OF USAGE OF THE MEANABSOLUTEERROR FUNCTION  
FROM SKLEARNMETRICS IMPORT MEANABSOLUTEERROR  
YTRUE 3 05 2 7  
YPRED 25 00 2 8  
MEANABSOLUTEERRORYTRUE YPRED  
05  
YTRUE 05 1 1 1 7 6  
YPRED 0 2 1 2 8 5  
MEANABSOLUTEERRORYTRUE YPRED  
075  
MEANABSOLUTEERRORYTRUE YPRED MULTIOUTPUTRAWVALUES  
ARRAY05 1  
MEANABSOLUTEERRORYTRUE YPRED MULTIOUTPUT03 07

085  
MEAN SQUARED ERROR  
THEMEANSQUAREDERROR FUNCTION COMPUTES MEAN SQUARE ERROR A RISK METRIC CORRESPONDING TO THE EXPECTED VALUE OF THE SQUARED QUADRATIC ERROR OR LOSS  
564 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

IF  $\hat{y}_i$  IS THE PREDICTED VALUE OF THE  $i$ TH SAMPLE AND  $y_i$  IS THE CORRESPONDING TRUE VALUE THEN THE MEAN SQUARED ERROR MSE ESTIMATED OVER  $n$  SAMPLES IS DEFINED AS

$$MSE = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2$$

HERE IS A SMALL EXAMPLE OF USAGE OF THE MEANSQUAREDERROR FUNCTION

```
from sklearn.metrics import meansquarederror
```

```
y_true = [3, 0.5, 2, 7]
y_pred = [2.5, 0, 2, 8]
meansquarederror(y_true, y_pred)
0.375
```

```
y_true = [0.5, 1, 1, 1, 7, 6]
y_pred = [0, 2, 1, 2, 8, 5]
meansquarederror(y_true, y_pred)
0.7083
```

EXAMPLES

- SEE GRADIENT BOOSTING REGRESSION FOR AN EXAMPLE OF MEAN SQUARED ERROR USAGE TO EVALUATE GRADIENT BOOSTING REGRESSION

MEAN SQUARED LOGARITHMIC ERROR

THE MEANSQUAREDLOGERROR FUNCTION COMPUTES A RISK METRIC CORRESPONDING TO THE EXPECTED VALUE OF THE SQUARED LOGARITHMIC QUADRATIC ERROR OR LOSS

IF  $\hat{y}_i$  IS THE PREDICTED VALUE OF THE  $i$ TH SAMPLE AND  $y_i$  IS THE CORRESPONDING TRUE VALUE THEN THE MEAN SQUARED LOGARITHMIC ERROR MSLE ESTIMATED OVER  $n$  SAMPLES IS DEFINED AS

$$MSLE = \frac{1}{n} \sum_{i=1}^n (\log(\hat{y}_i + 1) - \log(y_i + 1))^2$$

WHERE  $\log$  MEANS THE NATURAL LOGARITHM OF  $x$ . THIS METRIC IS BEST TO USE WHEN TARGETS HAVING EXPONENTIAL GROWTH SUCH AS POPULATION COUNTS AVERAGE SALES OF A COMMODITY OVER A SPAN OF YEARS ETC. NOTE THAT THIS METRIC PENALIZES AN UNDERPREDICTED ESTIMATE GREATER THAN AN OVERPREDICTED ESTIMATE.

HERE IS A SMALL EXAMPLE OF USAGE OF THE MEANSQUAREDLOGERROR FUNCTION

```
from sklearn.metrics import meansquaredlogerror
```

```
y_true = [3, 5, 2.5, 7]
y_pred = [2.5, 5, 4, 8]
meansquaredlogerror(y_true, y_pred)
0.039
```

```
y_true = [0.5, 1, 1, 2, 7, 6]
y_pred = [0.5, 2, 1, 2.5, 8, 8]
meansquaredlogerror(y_true, y_pred)
0.044
```

MEDIAN ABSOLUTE ERROR

THE MEDIANABSOLUTEERROR IS PARTICULARLY INTERESTING BECAUSE IT IS ROBUST TO OUTLIERS. THE LOSS IS CALCULATED BY TAKING THE MEDIAN OF ALL ABSOLUTE DIFFERENCES BETWEEN THE TARGET AND THE PREDICTION.

SCIKITLEARN USER GUIDE RELEASE 0213

IF  $\hat{y}_i$  IS THE PREDICTED VALUE OF THE  $i$ TH SAMPLE AND  $y_i$  IS THE CORRESPONDING TRUE VALUE THEN THE MEDIAN ABSOLUTE ERROR MEDAE ESTIMATED OVER  $n$ SAMPLES IS DEFINED AS

MEDAE  $= \text{MEDIAN } |y_1 - \hat{y}_1| \dots |y_n - \hat{y}_n|$

THE MEDIANABSOLUTEERROR DOES NOT SUPPORT MULTIOUTPUT

HERE IS A SMALL EXAMPLE OF USAGE OF THE MEDIANABSOLUTEERROR FUNCTION

```
FROM SKLEARNMETRICS IMPORT MEDIANABSOLUTEERROR
YTRUE 3 05 2 7
YPRED 25 00 2 8
MEDIANABSOLUTEERROR YTRUE YPRED
05
```

R2SCORE THE COEFFICIENT OF DETERMINATION

THE R2SCORE FUNCTION COMPUTES THE COEFFICIENT OF DETERMINATION USUALLY DENOTED AS  $R^2$

IT REPRESENTS THE PROPORTION OF VARIANCE OF  $Y$  THAT HAS BEEN EXPLAINED BY THE INDEPENDENT VARIABLES IN THE MODEL IT PROVIDES AN INDICATION OF GOODNESS OF FIT AND THEREFORE A MEASURE OF HOW WELL UNSEEN SAMPLES ARE LIKELY TO BE PREDICTED BY THE MODEL THROUGH THE PROPORTION OF EXPLAINED VARIANCE

AS SUCH VARIANCE IS DATASET DEPENDENT  $R^2$  MAY NOT BE MEANINGFULLY COMPARABLE ACROSS DIFFERENT DATASETS BEST POSSIBLE SCORE IS 1.0 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS PREDICTS THE EXPECTED VALUE OF  $Y$  DISREGARDING THE INPUT FEATURES WOULD GET A  $R^2$  SCORE OF 0.0

IF  $\hat{y}_i$  IS THE PREDICTED VALUE OF THE  $i$ TH SAMPLE AND  $y_i$  IS THE CORRESPONDING TRUE VALUE FOR TOTAL  $n$ SAMPLES THE ESTIMATED  $R^2$  IS DEFINED AS

$$R^2 = \frac{\sum_{i=1}^n (y_i - \bar{y})(\hat{y}_i - \bar{\hat{y}})}{\sqrt{\sum_{i=1}^n (y_i - \bar{y})^2 \sum_{i=1}^n (\hat{y}_i - \bar{\hat{y}})^2}}$$

WHERE  $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$

$\bar{\hat{y}} = \frac{1}{n} \sum_{i=1}^n \hat{y}_i$

AND  $\sum_{i=1}^n (\hat{y}_i - \bar{\hat{y}})^2 = \sum_{i=1}^n \hat{y}_i^2 - n \bar{\hat{y}}^2$

NOTE THAT  $R^2$  SCORE CALCULATES UNADJUSTED  $R^2$  WITHOUT CORRECTING FOR BIAS IN SAMPLE VARIANCE OF  $Y$

HERE IS A SMALL EXAMPLE OF USAGE OF THE  $R^2$  SCORE FUNCTION

```
FROM SKLEARNMETRICS IMPORT R2SCORE
YTRUE 3 05 2 7
YPRED 25 00 2 8
R2SCORE YTRUE YPRED
0.948
YTRUE 05 1 1 1 7 6
YPRED 0 2 1 2 8 5
R2SCORE YTRUE YPRED MULTIOUTPUT VARIANCEWEIGHTED
0.938
YTRUE 05 1 1 1 7 6
YPRED 0 2 1 2 8 5
R2SCORE YTRUE YPRED MULTIOUTPUT UNIFORM AVERAGE
0.936
R2SCORE YTRUE YPRED MULTIOUTPUT RAWVALUES
0.936
```

ARRAY 0965 0908

R2SCORE YTRUE YPRED MULTIOUTPUT 03 07

566 CHAPTER 3 USER GUIDE

0925

EXAMPLE

- SEE LASSO AND ELASTIC NET FOR SPARSE SIGNALS FOR AN EXAMPLE OF R2SCORE USAGE TO EVALUATE LASSO AND ELASTIC NET ON SPARSE SIGNALS

CLUSTERING METRICS

THE SKLEARNMETRICS MODULE IMPLEMENTS SEVERAL LOSS SCORE AND UTILITY FUNCTIONS FOR MORE INFORMATION SEE THE CLUSTERING PERFORMANCE EVALUATION SECTION FOR INSTANCE CLUSTERING AND BICLUSTERING EVALUATION FOR BICLUSTERING

DUMMY ESTIMATORS

WHEN DOING SUPERVISED LEARNING A SIMPLE SANITY CHECK CONSISTS OF COMPARING ONE'S ESTIMATOR AGAINST SIMPLE RULES OF THUMB DUMMYCLASSIFIER IMPLEMENTS SEVERAL SUCH SIMPLE STRATEGIES FOR CLASSIFICATION

- STRATIFIED GENERATES RANDOM PREDICTIONS BY RESPECTING THE TRAINING SET CLASS DISTRIBUTION

- MOSTFREQUENT ALWAYS PREDICTS THE MOST FREQUENT LABEL IN THE TRAINING SET

- PRIOR ALWAYS PREDICTS THE CLASS THAT MAXIMIZES THE CLASS PRIOR LIKE MOSTFREQUENT AND

PREDICTPROBA RETURNS THE CLASS PRIOR

- UNIFORM GENERATES PREDICTIONS UNIFORMLY AT RANDOM

- CONSTANT ALWAYS PREDICTS A CONSTANT LABEL THAT IS PROVIDED BY THE USER A MAJOR MOTIVATION OF THIS

METHOD IS F1SCORING WHEN THE POSITIVE CLASS IS IN THE MINORITY

NOTE THAT WITH ALL THESE STRATEGIES THE PREDICT METHOD COMPLETELY IGNORES THE INPUT DATA

TO ILLUSTRATE DUMMYCLASSIFIER FIRST LET'S CREATE AN IMBALANCED DATASET

FROM SKLEARN DATASETS IMPORT LOADIRIS

FROM SKLEARN MODELSELECTION IMPORT TRAINTESTSPLIT

IRIS = LOADIRIS

X, Y = IRIS.data, IRIS.target

Y = Y[1:]

X\_train, X\_test, Y\_train, Y\_test = train\_test\_split(X, Y, random\_state=0)

NEXT LET'S COMPARE THE ACCURACY OF SVC AND MOSTFREQUENT

FROM SKLEARN DUMMY IMPORT DUMMYCLASSIFIER

FROM SKLEARN SVM IMPORT SVC

clf = SVC(kernel='linear', C=1, fit\_x\_train=Y\_train)

clf.score(X\_test, Y\_test)

0.63

clf = DUMMYCLASSIFIER(STRATEGY='mostfrequent', random\_state=0)

clf.fit(X\_train, Y\_train)

DUMMYCLASSIFIER.CONSTANT\_NONE, random\_state=0, strategy='mostfrequent')

clf.score(X\_test, Y\_test)

0.57

WE SEE THAT SVC DOESN'T DO MUCH BETTER THAN A DUMMY CLASSIFIER NOW LET'S CHANGE THE KERNEL

33 MODEL SELECTION AND EVALUATION 567

SCIKITLEARN USER GUIDE RELEASE 0213  
CLF SVCGAMMASCALE KERNELRBF C1FITXTRAIN YTRAIN  
CLFSCOREXTEST YTEST  
094

WE SEE THAT THE ACCURACY WAS BOOSTED TO ALMOST 100 A CROSS VALIDATION STRATEGY IS RECOMMENDED FOR A BETTER ESTIMATE OF THE ACCURACY IF IT IS NOT TOO CPU COSTLY FOR MORE INFORMATION SEE THE CROSSVALIDATION EVALUATING ESTIMATOR PERFORMANCE SECTION MOREOVER IF YOU WANT TO OPTIMIZE OVER THE PARAMETER SPACE IT IS HIGHLY RECOMMENDED TO USE AN APPROPRIATE METHODOLOGY SEE THE TUNING THE HYPERPARAMETERS OF AN ESTIMATOR SECTION FOR DETAILS MORE GENERALLY WHEN THE ACCURACY OF A CLASSIFIER IS TOO CLOSE TO RANDOM IT PROBABLY MEANS THAT SOMETHING WENT WRONG FEATURES ARE NOT HELPFUL A HYPERPARAMETER IS NOT CORRECTLY TUNED THE CLASSIFIER IS SUFFERING FROM CLASS IMBALANCE ETC DUMMYREGRESSOR ALSO IMPLEMENTS FOUR SIMPLE RULES OF THUMB FOR REGRESSION

- MEAN ALWAYS PREDICTS THE MEAN OF THE TRAINING TARGETS
  - MEDIAN ALWAYS PREDICTS THE MEDIAN OF THE TRAINING TARGETS
  - QUANTILE ALWAYS PREDICTS A USER PROVIDED QUANTILE OF THE TRAINING TARGETS
  - CONSTANT ALWAYS PREDICTS A CONSTANT VALUE THAT IS PROVIDED BY THE USER
- IN ALL THESE STRATEGIES THE PREDICT METHOD COMPLETELY IGNORES THE INPUT DATA

334 MODEL PERSISTENCE

AFTER TRAINING A SCIKITLEARN MODEL IT IS DESIRABLE TO HAVE A WAY TO PERSIST THE MODEL FOR FUTURE USE WITHOUT HAVING TO RETRAIN THE FOLLOWING SECTION GIVES YOU AN EXAMPLE OF HOW TO PERSIST A MODEL WITH PICKLE WE’LL ALSO REVIEW A FEW SECURITY AND MAINTAINABILITY ISSUES WHEN WORKING WITH PICKLE SERIALIZATION AN ALTERNATIVE TO PICKLING IS TO EXPORT THE MODEL TO ANOTHER FORMAT USING ONE OF THE MODEL EXPORT TOOLS LISTED UNDER RELATED PROJECTS UNLIKE PICKLING ONCE EXPORTED YOU CANNOT RECOVER THE FULL SCIKITLEARN ESTIMATOR OBJECT BUT YOU CAN DEPLOY THE MODEL FOR PREDICTION USUALLY BY USING TOOLS SUPPORTING OPEN MODEL INTERCHANGE FORMATS SUCH AS ‘ONNX’ OR‘PMML’

PERSISTENCE EXAMPLE  
IT IS POSSIBLE TO SAVE A MODEL IN SCIKITLEARN BY USING PYTHON’S BUILTIN PERSISTENCE MODEL NAMELY PICKLE

```
FROM SKLEARN IMPORT SVM
FROM SKLEARN IMPORT DATASETS
CLF SVMSCGAMMASCALE
IRIS DATASETSLOADIRIS
X Y IRISDATA IRISTARGET
CLFFITX Y
SVCC10 CACHESIZE200 CLASSWEIGHTNONE COEF000
DECISIONFUNCTIONSHAPEOVR DEGREE3 GAMMASCALE KERNELRBF
MAXITER1 PROBABILITYFALSE RANDOMSTATENONE SHRINKINGTRUE
TOL0001 VERBOSEFALSE
IMPORT PICKLE
S PICKLEDUMPSCLF
CLF2 PICKLELOADSS
CLF2PREDICTX01
ARRAY0
Y0
0
```

SCIKITLEARN USER GUIDE RELEASE 0213

IN THE SPECIFIC CASE OF SCIKITLEARN IT MAY BE BETTER TO USE JOBLIB’S REPLACEMENT OF PICKLE DUMP LOAD WHICH IS MORE EFFICIENT ON OBJECTS THAT CARRY LARGE NUMPY ARRAYS INTERNALLY AS IS OFTEN THE CASE FOR FITTED SCIKITLEARN ESTIMATORS BUT CAN ONLY PICKLE TO THE DISK AND NOT TO A STRING

```
FROM JOBLIB IMPORT DUMP LOAD
DUMPCLF FILENAMEJOBLIB
```

LATER YOU CAN LOAD BACK THE PICKLED MODEL POSSIBLY IN ANOTHER PYTHON PROCESS WITH

```
CLF LOADFILENAMEJOBLIB
```

NOTEDUMP ANDLOAD FUNCTIONS ALSO ACCEPT FILELIKE OBJECT INSTEAD OF FILENAMES MORE INFORMATION ON DATA PERSISTENCE WITH JOBLIB IS AVAILABLE HERE

SECURITY MAINTAINABILITY LIMITATIONS

PICKLE AND JOBLIB BY EXTENSION HAS SOME ISSUES REGARDING MAINTAINABILITY AND SECURITY BECAUSE OF THIS

- NEVER UNPICKLE UNTRUSTED DATA AS IT COULD LEAD TO MALICIOUS CODE BEING EXECUTED UPON LOADING
- WHILE MODELS SAVED USING ONE VERSION OF SCIKITLEARN MIGHT LOAD IN OTHER VERSIONS THIS IS ENTIRELY UNSUPPORTED AND INADVISABLE IT SHOULD ALSO BE KEPT IN MIND THAT OPERATIONS PERFORMED ON SUCH DATA COULD GIVE DIFFERENT AND UNEXPECTED RESULTS

IN ORDER TO REBUILD A SIMILAR MODEL WITH FUTURE VERSIONS OF SCIKITLEARN ADDITIONAL METADATA SHOULD BE SAVED ALONG THE PICKLED MODEL

- THE TRAINING DATA EG A REFERENCE TO AN IMMUTABLE SNAPSHOT
- THE PYTHON SOURCE CODE USED TO GENERATE THE MODEL
- THE VERSIONS OF SCIKITLEARN AND ITS DEPENDENCIES
- THE CROSS VALIDATION SCORE OBTAINED ON THE TRAINING DATA

THIS SHOULD MAKE IT POSSIBLE TO CHECK THAT THE CROSSVALIDATION SCORE IS IN THE SAME RANGE AS BEFORE

SINCE A MODEL INTERNAL REPRESENTATION MAY BE DIFFERENT ON TWO DIFFERENT ARCHITECTURES DUMPING A MODEL ON ONE ARCHITECTURE AND LOADING IT ON ANOTHER ARCHITECTURE IS NOT SUPPORTED

IF YOU WANT TO KNOW MORE ABOUT THESE ISSUES AND EXPLORE OTHER POSSIBLE SERIALIZATION METHODS PLEASE REFER TO THIS TALK BY ALEX GAYNOR

335 VALIDATION CURVES PLOTTING SCORES TO EVALUATE MODELS

EVERY ESTIMATOR HAS ITS ADVANTAGES AND DRAWBACKS ITS GENERALIZATION ERROR CAN BE DECOMPOSED IN TERMS OF BIAS VARIANCE AND NOISE THE BIAS OF AN ESTIMATOR IS ITS AVERAGE ERROR FOR DIFFERENT TRAINING SETS THE VARIANCE OF AN ESTIMATOR INDICATES HOW SENSITIVE IT IS TO VARYING TRAINING SETS NOISE IS A PROPERTY OF THE DATA

IN THE FOLLOWING PLOT WE SEE A FUNCTION  $\cos^3$

2AND SOME NOISY SAMPLES FROM THAT FUNCTION WE USE THREE

DIFFERENT ESTIMATORS TO FIT THE FUNCTION LINEAR REGRESSION WITH POLYNOMIAL FEATURES OF DEGREE 1 4 AND 15 WE SEE THAT THE FIRST ESTIMATOR CAN AT BEST PROVIDE ONLY A POOR FIT TO THE SAMPLES AND THE TRUE FUNCTION BECAUSE IT IS TOO SIMPLE HIGH BIAS THE SECOND ESTIMATOR APPROXIMATES IT ALMOST PERFECTLY AND THE LAST ESTIMATOR APPROXIMATES THE TRAINING DATA PERFECTLY BUT DOES NOT FIT THE TRUE FUNCTION VERY WELL IE IT IS VERY SENSITIVE TO VARYING TRAINING DATA HIGH VARIANCE BIAS AND VARIANCE ARE INHERENT PROPERTIES OF ESTIMATORS AND WE USUALLY HAVE TO SELECT LEARNING ALGORITHMS AND HYPER PARAMETERS SO THAT BOTH BIAS AND VARIANCE ARE AS LOW AS POSSIBLE SEE BIASVARIANCE DILEMMA ANOTHER WAY TO REDUCE

33 MODEL SELECTION AND EVALUATION 569

SCIKITLEARN USER GUIDE RELEASE 0213

THE VARIANCE OF A MODEL IS TO USE MORE TRAINING DATA HOWEVER YOU SHOULD ONLY COLLECT MORE TRAINING DATA IF THE TRUE FUNCTION IS TOO COMPLEX TO BE APPROXIMATED BY AN ESTIMATOR WITH A LOWER VARIANCE

IN THE SIMPLE ONEDIMENSIONAL PROBLEM THAT WE HAVE SEEN IN THE EXAMPLE IT IS EASY TO SEE WHETHER THE ESTIMATOR SUFFERS FROM BIAS OR VARIANCE HOWEVER IN HIGHDIMENSIONAL SPACES MODELS CAN BECOME VERY DIFFICULT TO VISUALIZE FOR THIS REASON IT IS OFTEN HELPFUL TO USE THE TOOLS DESCRIBED BELOW

EXAMPLES

- UNDERFITTING VS OVERFITTING
- PLOTING VALIDATION CURVES
- PLOTING LEARNING CURVES

VALIDATION CURVE

TO VALIDATE A MODEL WE NEED A SCORING FUNCTION SEE MODEL EVALUATION QUANTIFYING THE QUALITY OF PREDICTIONS FOR EXAMPLE ACCURACY FOR CLASSIFIERS THE PROPER WAY OF CHOOSING MULTIPLE HYPERPARAMETERS OF AN ESTIMATOR ARE OF COURSE GRID SEARCH OR SIMILAR METHODS SEE TUNING THE HYPERPARAMETERS OF AN ESTIMATOR THAT SELECT THE HYPERPARAMETER WITH THE MAXIMUM SCORE ON A VALIDATION SET OR MULTIPLE VALIDATION SETS NOTE THAT IF WE OPTIMIZED THE HYPERPARAMETERS BASED ON A VALIDATION SCORE THE VALIDATION SCORE IS BIASED AND NOT A GOOD ESTIMATE OF THE GENERALIZATION ANY LONGER TO GET A PROPER ESTIMATE OF THE GENERALIZATION WE HAVE TO COMPUTE THE SCORE ON ANOTHER TEST SET HOWEVER IT IS SOMETIMES HELPFUL TO PLOT THE INFLUENCE OF A SINGLE HYPERPARAMETER ON THE TRAINING SCORE AND THE VALIDATION SCORE TO FIND OUT WHETHER THE ESTIMATOR IS OVERFITTING OR UNDERFITTING FOR SOME HYPERPARAMETER VALUES

THE FUNCTION VALIDATIONCURVE CAN HELP IN THIS CASE

```
import numpy as np
from sklearn.model_selection import validation_curve
from sklearn.datasets import load_iris
from sklearn.linear_model import Ridge
np.random.seed(0)
iris = load_iris()
X, y = iris.data, iris.target
indices = np.arange(y.shape[0])
np.random.shuffle(indices)
570 CHAPTER 3 USER GUIDE
```



```
SCIKITLEARN USER GUIDE RELEASE 0213
X Y XINDICES YINDICES
TRAINSCORES VALIDSCORES VALIDATIONCURVERIDGE X Y ALPHA
NPLOGSPACE7 3 3
CV5
TRAINSCORES
ARRAY093 094 092 091 092
093 094 092 091 092
051 052 049 047 049
VALIDSCORES
ARRAY090 084 094 096 093
090 084 094 096 093
046 025 050 049 052
```

IF THE TRAINING SCORE AND THE VALIDATION SCORE ARE BOTH LOW THE ESTIMATOR WILL BE UNDERFITTING IF THE TRAINING SCORE IS HIGH AND THE VALIDATION SCORE IS LOW THE ESTIMATOR IS OVERFITTING AND OTHERWISE IT IS WORKING VERY WELL A LOW TRAINING SCORE AND A HIGH VALIDATION SCORE IS USUALLY NOT POSSIBLE ALL THREE CASES CAN BE FOUND IN THE PLOT BELOW WHERE WE VARY THE PARAMETER  $\gamma$  OF AN SVM ON THE DIGITS DATASET

LEARNING CURVE  
A LEARNING CURVE SHOWS THE VALIDATION AND TRAINING SCORE OF AN ESTIMATOR FOR VARYING NUMBERS OF TRAINING SAMPLES IT IS A TOOL TO FIND OUT HOW MUCH WE BENEFIT FROM ADDING MORE TRAINING DATA AND WHETHER THE ESTIMATOR SUFFERS MORE FROM A VARIANCE ERROR OR A BIAS ERROR IF BOTH THE VALIDATION SCORE AND THE TRAINING SCORE CONVERGE TO A VALUE THAT IS TOO LOW WITH INCREASING SIZE OF THE TRAINING SET WE WILL NOT BENEFIT MUCH FROM MORE TRAINING DATA IN THE FOLLOWING PLOT YOU CAN SEE AN EXAMPLE NAIVE BAYES ROUGHLY CONVERGES TO A LOW SCORE  
WE WILL PROBABLY HAVE TO USE AN ESTIMATOR OR A PARAMETRIZATION OF THE CURRENT ESTIMATOR THAT CAN LEARN MORE COMPLEX CONCEPTS IE HAS A LOWER BIAS IF THE TRAINING SCORE IS MUCH GREATER THAN THE VALIDATION SCORE FOR THE MAXIMUM NUMBER OF TRAINING SAMPLES ADDING MORE TRAINING SAMPLES WILL MOST LIKELY INCREASE GENERALIZATION IN THE FOLLOWING PLOT YOU CAN SEE THAT THE SVM COULD BENEFIT FROM MORE TRAINING EXAMPLES  
WE CAN USE THE FUNCTION LEARNINGCURVE TO GENERATE THE VALUES THAT ARE REQUIRED TO PLOT SUCH A LEARNING CURVE  
NUMBER OF SAMPLES THAT HAVE BEEN USED THE AVERAGE SCORES ON THE TRAINING SETS AND THE AVERAGE SCORES ON THE VALIDATION SETS

```
FROM SKLEARNMODELSELECTION IMPORT LEARNINGCURVE
FROM SKLEARN SVM IMPORT SVC
33 MODEL SELECTION AND EVALUATION 571
```



SCIKITLEARN USER GUIDE RELEASE 0213  
TRAINSIZES TRAINSCORES VALIDSCORES LEARNINGCURVE  
SVCKERNELLINEAR X Y TRAINSIZES50 80 110 CV5

TRAINSIZES  
ARRAY 50 80 110  
TRAINSCORES  
ARRAY098 098 098 098 098  
098 1 098 098 098  
098 1 098 098 099

VALIDSCORES  
ARRAY1 093 1 1 096  
1 096 1 1 096  
1 096 1 1 096

34 INSPECTION  
341 PARTIAL DEPENDENCE PLOTS  
PARTIAL DEPENDENCE PLOTS PDP SHOW THE DEPENDENCE BETWEEN THE TARGET RESPONSE1AND A SET OF 'TARGET' FEATURES  
MARGINALIZING OVER THE VALUES OF ALL OTHER FEATURES THE 'COMPLEMENT' FEATURES INTUITIVELY WE CAN INTERPRET THE PARTIAL  
DEPENDENCE AS THE EXPECTED TARGET RESPONSE AS A FUNCTION OF THE 'TARGET' FEATURES  
DUE TO THE LIMITS OF HUMAN PERCEPTION THE SIZE OF THE TARGET FEATURE SET MUST BE SMALL USUALLY ONE OR TWO THUS THE  
TARGET FEATURES ARE USUALLY CHOSEN AMONG THE MOST IMPORTANT FEATURES  
THE FIGURE BELOW SHOWS FOUR ONEWAY AND ONE TWOWAY PARTIAL DEPENDENCE PLOTS FOR THE CALIFORNIA HOUSING DATASET  
WITH AGRADIENTBOOSTINGREGRESSOR  
ONEWAY PDPS TELL US ABOUT THE INTERACTION BETWEEN THE TARGET RESPONSE AND THE TARGET FEATURE EG LINEAR NONLINEAR  
THE UPPER LEFT PLOT IN THE ABOVE FIGURE SHOWS THE EFFECT OF THE MEDIAN INCOME IN A DISTRICT ON THE MEDIAN HOUSE PRICE  
1FOR CLASSIFICATION THE TARGET RESPONSE MAY BE THE PROBABILITY OF A CLASS THE POSITIVE CLASS FOR BINARY CLASSIFICATION OR THE D  
34 INSPECTION 573

SCIKITLEARN USER GUIDE RELEASE 0213

WE CAN CLEARLY SEE A LINEAR RELATIONSHIP AMONG THEM NOTE THAT PDPS ASSUME THAT THE TARGET FEATURES ARE INDEPENDENT FROM THE COMPLEMENT FEATURES AND THIS ASSUMPTION IS OFTEN VIOLATED IN PRACTICE PDPS WITH TWO TARGET FEATURES SHOW THE INTERACTIONS AMONG THE TWO FEATURES FOR EXAMPLE THE TWOVARIABLE PDP IN THE ABOVE FIGURE SHOWS THE DEPENDENCE OF MEDIAN HOUSE PRICE ON JOINT VALUES OF HOUSE AGE AND AVERAGE OCCUPANTS PER HOUSEHOLD WE CAN CLEARLY SEE AN INTERACTION BETWEEN THE TWO FEATURES FOR AN AVERAGE OCCUPANCY GREATER THAN TWO THE HOUSE PRICE IS NEARLY INDEPENDENT OF THE HOUSE AGE WHEREAS FOR VALUES LESS THAN 2 THERE IS A STRONG DEPENDENCE ON AGE THE SKLEARNINSPECTION MODULE PROVIDES A CONVENIENCE FUNCTION PLOTPARTIALDEPENDENCE TO CREATE ONEWAY AND TWOWAY PARTIAL DEPENDENCE PLOTS IN THE BELOW EXAMPLE WE SHOW HOW TO CREATE A GRID OF PARTIAL DEPENDENCE PLOTS TWO ONEWAY PDPS FOR THE FEATURES 0AND1AND A TWOWAY PDP BETWEEN THE TWO FEATURES

```
FROM SKLEARNDATASETS IMPORT MAKEHASTIE102
FROM SKLEARNENSEMBLE IMPORT GRADIENTBOOSTINGCLASSIFIER
FROM SKLEARNINSPECTION IMPORT PLOTPARTIALDEPENDENCE
X Y MAKEHASTIE102RANDOMSTATE0
CLF GRADIENTBOOSTINGCLASSIFIERNESTIMATORS100 LEARNINGRATE10
MAXDEPTH1 RANDOMSTATE0FITX Y
FEATURES 0 1 0 1
PLOTPARTIALDEPENDENCECLF X FEATURES
```

YOU CAN ACCESS THE NEWLY CREATED FIGURE AND AXES OBJECTS USING PLTGCF ANDPLTGCA FOR MULTICLASS CLASSIFICATION YOU NEED TO SET THE CLASS LABEL FOR WHICH THE PDPS SHOULD BE CREATED VIA THE TARGET ARGUMENT

```
FROM SKLEARNDATASETS IMPORT LOADIRIS
IRIS LOADIRIS
MCCCLF GRADIENTBOOSTINGCLASSIFIERNESTIMATORS10
MAXDEPTH1FITIRISDATA IRISTARGET
FEATURES 3 2 3 2
PLOTPARTIALDEPENDENCEMCCCLF X FEATURES TARGET0
```

THE SAME PARAMETER TARGET IS USED TO SPECIFY THE TARGET IN MULTIOUTPUT REGRESSION SETTINGS IF YOU NEED THE RAW VALUES OF THE PARTIAL DEPENDENCE FUNCTION RATHER THAN THE PLOTS YOU CAN USE THE SKLEARN INSPECTIONPARTIALDEPENDENCE FUNCTION

```
FROM SKLEARNINSPECTION IMPORT PARTIALDEPENDENCE
PDP AXES PARTIALDEPENDENCECLF X 0
PDP
ARRAY 2466 2466
AXES
ARRAY1624 1592
```

THE VALUES AT WHICH THE PARTIAL DEPENDENCE SHOULD BE EVALUATED ARE DIRECTLY GENERATED FROM X FOR 2WAY PARTIAL DEPENDENCE A 2DGRID OF VALUES IS GENERATED THE VALUES FIELD RETURNED BY SKLEARNINSPECTION PARTIALDEPENDENCE GIVES THE ACTUAL VALUES USED IN THE GRID FOR EACH TARGET FEATURE THEY ALSO CORRESPOND TO THE AXIS OF THE PLOTS FOR EACH VALUE OF THE 'TARGET' FEATURES IN THE GRID THE PARTIAL DEPENDENCE FUNCTION NEEDS TO MARGINALIZE THE PREDICTIONS OF THE ESTIMATOR OVER ALL POSSIBLE VALUES OF THE 'COMPLEMENT' FEATURES WITH THE BRUTE METHOD THIS IS DONE BY REPLACING EVERY TARGET FEATURE VALUE OF XBY THOSE IN THE GRID AND COMPUTING THE AVERAGE PREDICTION IN DECISION TREES THIS CAN BE EVALUATED EFFICIENTLY WITHOUT REFERENCE TO THE TRAINING DATA RECURSION METHOD FOR EACH GRID POINT A WEIGHTED TREE TRAVERSAL IS PERFORMED IF A SPLIT NODE INVOLVES A 'TARGET' FEATURE THE CORRESPONDING LEFT OR RIGHT BRANCH IS FOLLOWED OTHERWISE BOTH BRANCHES ARE FOLLOWED EACH BRANCH IS WEIGHTED BY THE FRACTION OF

574 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

TRAINING SAMPLES THAT ENTERED THAT BRANCH FINALLY THE PARTIAL DEPENDENCE IS GIVEN BY A WEIGHTED AVERAGE OF ALL VISITED LEAVES NOTE THAT WITH THE RECURSION METHODXIS ONLY USED TO GENERATE THE GRID NOT TO COMPUTE THE AVERAGED PREDICTIONS THE AVERAGED PREDICTIONS WILL ALWAYS BE COMPUTED ON THE DATA WITH WHICH THE TREES WERE TRAINED

EXAMPLES

•PARTIAL DEPENDENCE PLOTS

REFERENCES

35 DATASET TRANSFORMATIONS

SCIKITLEARN PROVIDES A LIBRARY OF TRANSFORMERS WHICH MAY CLEAN SEE PREPROCESSING DATA REDUCE SEE UNSUPERVISED DIMENSIONALITY REDUCTION EXPAND SEE KERNEL APPROXIMATION OR GENERATE SEE FEATURE EXTRACTION FEATURE REPRESENTATIONS

LIKE OTHER ESTIMATORS THESE ARE REPRESENTED BY CLASSES WITH A FIT METHOD WHICH LEARNS MODEL PARAMETERS EG MEAN AND STANDARD DEVIATION FOR NORMALIZATION FROM A TRAINING SET AND A TRANSFORM METHOD WHICH APPLIES THIS TRANSFORMATION MODEL TO UNSEEN DATA FITTRANSFORM MAY BE MORE CONVENIENT AND EFFICIENT FOR MODELLING AND TRANSFORMING THE TRAINING DATA SIMULTANEOUSLY

COMBINING SUCH TRANSFORMERS EITHER IN PARALLEL OR SERIES IS COVERED IN PIPELINES AND COMPOSITE ESTIMATORS PAIR WISE METRICS AFFINITIES AND KERNELS COVERS TRANSFORMING FEATURE SPACES INTO AFFINITY MATRICES WHILE TRANSFORMING THE PREDICTION TARGET Y CONSIDERS TRANSFORMATIONS OF THE TARGET SPACE EG CATEGORICAL LABELS FOR USE IN SCIKITLEARN

351 PIPELINES AND COMPOSITE ESTIMATORS

TRANSFORMERS ARE USUALLY COMBINED WITH CLASSIFIERS REGRESSORS OR OTHER ESTIMATORS TO BUILD A COMPOSITE ESTIMATOR THE MOST COMMON TOOL IS A PIPELINE PIPELINE IS OFTEN USED IN COMBINATION WITH FEATUREUNION WHICH CONCATENATES THE OUTPUT OF TRANSFORMERS INTO A COMPOSITE FEATURE SPACE TRANSFORMEDTARGETREGRESSOR DEALS WITH TRANSFORMING THE TARGET IE LOGTRANSFORM Y IN CONTRAST PIPELINES ONLY TRANSFORM THE OBSERVED DATA X

PIPELINE CHAINING ESTIMATORS

PIPELINE CAN BE USED TO CHAIN MULTIPLE ESTIMATORS INTO ONE THIS IS USEFUL AS THERE IS OFTEN A FIXED SEQUENCE OF STEPS IN PROCESSING THE DATA FOR EXAMPLE FEATURE SELECTION NORMALIZATION AND CLASSIFICATION PIPELINE SERVES MULTIPLE PURPOSES HERE

CONVENIENCE AND ENCAPSULATION YOU ONLY HAVE TO CALL FIT ANDPREDICT ONCE ON YOUR DATA TO FIT A WHOLE SEQUENCE OF ESTIMATORS

JOINT PARAMETER SELECTION YOU CAN GRID SEARCH OVER PARAMETERS OF ALL ESTIMATORS IN THE PIPELINE AT ONCE SAFETY PIPELINES HELP AVOID LEAKING STATISTICS FROM YOUR TEST DATA INTO THE TRAINED MODEL IN CROSSVALIDATION BY ENSURING THAT THE SAME SAMPLES ARE USED TO TRAIN THE TRANSFORMERS AND PREDICTORS

ALL ESTIMATORS IN A PIPELINE EXCEPT THE LAST ONE MUST BE TRANSFORMERS IE MUST HAVE A TRANSFORM METHOD THE LAST ESTIMATOR MAY BE ANY TYPE TRANSFORMER CLASSIFIER ETC

35 DATASET TRANSFORMATIONS 575

SCIKITLEARN USER GUIDE RELEASE 0213

USAGE

CONSTRUCTION

THE PIPELINE IS BUILT USING A LIST OF KEY VALUE PAIRS WHERE THE KEY IS A STRING CONTAINING THE NAME YOU WANT TO GIVE THIS STEP AND VALUE IS AN ESTIMATOR OBJECT

```
FROM SKLEARNPIPELINE IMPORT PIPELINE
FROM SKLEARN SVM IMPORT SVC
FROM SKLEARN DECOMPOSITION IMPORT PCA
ESTIMATORS REDUCEDIM PCA CLF SVC
PIPE PIPELINE ESTIMATORS
PIPE
PIPELINE MEMORY NONE
STEPS REDUCEDIM PCA COPY TRUE
CLF SV CC10 VERBOSE FALSE
```

THE UTILITY FUNCTION MAKE PIPELINE IS A SHORTHAND FOR CONSTRUCTING PIPELINES IT TAKES A VARIABLE NUMBER OF ESTIMATORS AND RETURNS A PIPELINE FILLING IN THE NAMES AUTOMATICALLY

```
FROM SKLEARN PIPELINE IMPORT MAKE PIPELINE
FROM SKLEARN NNAIVE BAYES IMPORT MULTINOMIAL NB
FROM SKLEARN PREPROCESSING IMPORT BINARIZER
MAKE PIPELINE BINARIZER MULTINOMIAL NB
PIPELINE MEMORY NONE
STEPS BINARIZER BINARIZER COPY TRUE THRESHOLD 00
MULTINOMIAL NB MULTINOMIAL NB ALPHA 10
CLASS PRIOR NONE
FIT PRIOR TRUE
VERBOSE FALSE
ACCESSING STEPS
```

THE ESTIMATORS OF A PIPELINE ARE STORED AS A LIST IN THE STEPS ATTRIBUTE BUT CAN BE ACCESSED BY INDEX OR NAME BY INDEXING WITH IDX THE PIPELINE

```
PIPE STEPS 0
REDUCEDIM PCA COPY TRUE ITERATED POWER AUTO NCOMPONENTS NONE
RANDOM STATE NONE SVDS SOLVER AUTO TOL 00
WHITEN FALSE
PIPE 0
PCA COPY TRUE ITERATED POWER AUTO NCOMPONENTS NONE RANDOM STATE NONE
SVDS SOLVER AUTO TOL 00 WHITEN FALSE
PIPE REDUCEDIM
PCA COPY TRUE
PIPELINE'S NAMED STEPS ATTRIBUTE ALLOWS ACCESSING STEPS BY NAME WITH TAB COMPLETION IN INTERACTIVE ENVIRONMENTS
PIPE NAMED STEPS REDUCEDIM IS PIPE REDUCEDIM
TRUE
```

A SUB PIPELINE CAN ALSO BE EXTRACTED USING THE SLICING NOTATION COMMONLY USED FOR PYTHON SEQUENCES SUCH AS LISTS OR STRINGS ALTHOUGH ONLY A STEP OF 1 IS PERMITTED THIS IS CONVENIENT FOR PERFORMING ONLY SOME OF THE TRANSFORMATIONS OR THEIR INVERSE

576 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

PIPE1

PIPELINEMEMORYNONE STEPSREDUCEDIM PCACOPYTRUE

PIPE1

PIPELINEMEMORYNONE STEPSCLF SVCC10

NESTED PARAMETERS

PARAMETERS OF THE ESTIMATORS IN THE PIPELINE CAN BE ACCESSED USING THE ESTIMATORPARAMETER SYNTAX

PIPESETPARAMSCLFC10

PIPELINEMEMORYNONE

STEPSREDUCEDIM PCACOPYTRUE ITERATEDPOWERAUTO

CLF SVCC10 CACHESIZE200 CLASSWEIGHTNONE

VERBOSEFALSE

THIS IS PARTICULARLY IMPORTANT FOR DOING GRID SEARCHES

FROM SKLEARNMODELSELECTION IMPORT GRIDSEARCHCV

PARAMGRID DICTREDUCEDIMNCOMPONENTS2 5 10

CLFC01 10 100

GRIDSEARCH GRIDSEARCHCVPIPE PARAMGRIDPARAMGRID

INDIVIDUAL STEPS MAY ALSO BE REPLACED AS PARAMETERS AND NONFINAL STEPS MAY BE IGNORED BY SETTING THEM TO PASSTHROUGH

FROM SKLEARNLINEARMODEL IMPORT LOGISTICREGRESSION

PARAMGRID DICTREDUCEDIMPASSTHROUGH PCA5 PCA10

CLFSVC LOGISTICREGRESSION

CLFC01 10 100

GRIDSEARCH GRIDSEARCHCVPIPE PARAMGRIDPARAMGRID

THE ESTIMATORS OF THE PIPELINE CAN BE RETRIEVED BY INDEX

PIPE0

PCACOPYTRUE

OR BY NAME

PIPEREDUCEDIM

PCACOPYTRUE

EXAMPLES

- PIPELINE ANOVA SVM
- SAMPLE PIPELINE FOR TEXT FEATURE EXTRACTION AND EVALUATION
- PIPELINING CHAINING A PCA AND A LOGISTIC REGRESSION
- EXPLICIT FEATURE MAP APPROXIMATION FOR RBF KERNELS
- SVMANOVA SVM WITH UNIVARIATE FEATURE SELECTION
- SELECTING DIMENSIONALITY REDUCTION WITH PIPELINE AND GRIDSEARCHCV

35 DATASET TRANSFORMATIONS 577

SCIKITLEARN USER GUIDE RELEASE 0213

SEE ALSO

•TUNING THE HYPERPARAMETERS OF AN ESTIMATOR

NOTES

CALLINGFIT ON THE PIPELINE IS THE SAME AS CALLING FIT ON EACH ESTIMATOR IN TURN TRANSFORM THE INPUT AND PASS IT ON TO THE NEXT STEP THE PIPELINE HAS ALL THE METHODS THAT THE LAST ESTIMATOR IN THE PIPELINE HAS IE IF THE LAST ESTIMATOR IS A CLASSIFIER THE PIPELINE CAN BE USED AS A CLASSIFIER IF THE LAST ESTIMATOR IS A TRANSFORMER AGAIN SO IS THE PIPELINE CACHING TRANSFORMERS AVOID REPEATED COMPUTATION

FITTING TRANSFORMERS MAY BE COMPUTATIONALLY EXPENSIVE WITH ITS MEMORY PARAMETER SET PIPELINE WILL CACHE EACH TRANSFORMER AFTER CALLING FIT THIS FEATURE IS USED TO AVOID COMPUTING THE FIT TRANSFORMERS WITHIN A PIPELINE IF THE PARAMETERS AND INPUT DATA ARE IDENTICAL A TYPICAL EXAMPLE IS THE CASE OF A GRID SEARCH IN WHICH THE TRANSFORMERS CAN BE FITTED ONLY ONCE AND REUSED FOR EACH CONFIGURATION

THE PARAMETER MEMORY IS NEEDED IN ORDER TO CACHE THE TRANSFORMERS MEMORY CAN BE EITHER A STRING CONTAINING THE DIRECTORY WHERE TO CACHE THE TRANSFORMERS OR A JOBLIBMEMORY OBJECT

```
FROM TEMPFILE IMPORT MKDTEMP
FROM SHUTIL IMPORT RMTREE
FROM SKLEARNDECOMPOSITION IMPORT PCA
FROM SKLEARN SVM IMPORT SVC
FROM SKLEARN PIPELINE IMPORT PIPELINE
ESTIMATORS REDUCEDIM PCA CLF SVC
CACHEDIR MKDTEMP
PIPE PIPELINEESTIMATORS MEMORYCACHEDIR
```

```
PIPE
PIPELINE
STEPSREDUCEDIM PCACOPYTRUE
CLF SVCC10 VERBOSEFALSE
CLEAR THE CACHE DIRECTORY WHEN YOU DONT NEED IT ANYMORE
RMTREECACHEDIR
```

WARNING SIDE EFFECT OF CACHING TRANSFORMERS  
USING APIPELINE WITHOUT CACHE ENABLED IT IS POSSIBLE TO INSPECT THE ORIGINAL INSTANCE SUCH AS

```
FROM SKLEARN DATASETS IMPORT LOADDIGITS
DIGITS LOADDIGITS
PCA1 PCA
SVM1 SVCGAMMASCALE
PIPE PIPELINEREDUCEDIM PCA1 CLF SVM1
PIPEFITDIGITSDATA DIGITSTARGET
```

```
PIPELINEMEMORYNONE
STEPSREDUCEDIM PCA CLF SVC
VERBOSEFALSE
THE PCA INSTANCE CAN BE INSPECTED DIRECTLY
PRINTPCA1COMPONENTS
177484909E19 407058917E18
578 CHAPTER 3 USER GUIDE
```



SCIKITLEARN USER GUIDE RELEASE 0213

ENABLING CACHING TRIGGERS A CLONE OF THE TRANSFORMERS BEFORE FITTING THEREFORE THE TRANSFORMER INSTANCE GIVEN TO THE PIPELINE CANNOT BE INSPECTED DIRECTLY IN FOLLOWING EXAMPLE ACCESSING THE PCA INSTANCEPCA2 WILL RAISE AN ATTRIBUTEERROR SINCEPCA2 WILL BE AN UNFITTED TRANSFORMER INSTEAD USE THE ATTRIBUTE NAMEDSTEPS TO INSPECT ESTIMATORS WITHIN THE PIPELINE

```
CACHEDIR MKDTEMP
PCA2 PCA
SVM2 SVCGAMMASCALE
CACHEDPIPE PIPELINEREDUCEDIM PCA2 CLF SVM2
MEMORYCACHEDIR
CACHEDPIPEFITDIGITS DATA DIGITSTARGET
```

PIPELINEMEMORY

STEPSREDUCEDIM PCA CLF SVC

VERBOSEFALSE

PRINTCACHEDPIPENAMEDSTEPSREDUCEDIMCOMPONENTS

177484909E19 407058917E18

REMOVE THE CACHE DIRECTORY

RMTREECACHEDIR

EXAMPLES

- SELECTING DIMENSIONALITY REDUCTION WITH PIPELINE AND GRIDSEARCHCV

TRANSFORMING TARGET IN REGRESSION

TRANSFORMEDTARGETREGRESSOR TRANSFORMS THE TARGETS YBEFORE FITTING A REGRESSION MODEL THE PREDICTIONS ARE MAPPED BACK TO THE ORIGINAL SPACE VIA AN INVERSE TRANSFORM IT TAKES AS AN ARGUMENT THE REGRESSOR THAT WILL BE USED FOR PREDICTION AND THE TRANSFORMER THAT WILL BE APPLIED TO THE TARGET VARIABLE

```
IMPORT NUMPY AS NP
FROM SKLEARNDATASETS IMPORT LOADBOSTON
FROM SKLEARNCOMPOSE IMPORT TRANSFORMEDTARGETREGRESSOR
FROM SKLEARNPREPROCESSING IMPORT QUANTILETRANSFORMER
FROM SKLEARNLINEARMODEL IMPORT LINEARREGRESSION
FROM SKLEARNMODELSELECTION IMPORT TRAINTESTSPLIT
BOSTON LOADBOSTON
X BOSTONDATA
Y BOSTONTARGET
TRANSFORMER QUANTILETRANSFORMEROUTPUTDISTRIBUTIONNORMAL
REGRESSOR LINEARREGRESSION
REGR TRANSFORMEDTARGETREGRESSORREGRESSORREGRESSOR
TRANSFORMERTRANSFORMER
XTRAIN XTEST YTRAIN YTEST TRAINTESTSPLITX Y RANDOMSTATE0
REGRFITXTRAIN YTRAIN
TRANSFORMEDTARGETREGRESSOR
PRINTR2 SCORE 02FFORMATREGRSCOREXTEST YTEST
R2 SCORE 067
RAWTARGETREGR LINEARREGRESSIONFITXTRAIN YTRAIN
PRINTR2 SCORE 02FFORMATRAWTARGETREGRSCOREXTEST YTEST
R2 SCORE 064
35 DATASET TRANSFORMATIONS 579
```

SCIKITLEARN USER GUIDE RELEASE 0213

FOR SIMPLE TRANSFORMATIONS INSTEAD OF A TRANSFORMER OBJECT A PAIR OF FUNCTIONS CAN BE PASSED DEFINING THE TRANSFORMATION AND ITS INVERSE MAPPING

```
def funcX
    return nplogX
def inversefuncX
    return npexpX
```

SUBSEQUENTLY THE OBJECT IS CREATED AS

```
regr = TransformedTargetRegressor(regressor=regressor,
    func=func,
    inversefunc=inversefunc,
    reg_fit=train, y_train=train,
    transformed_target=regressor,
    print_r2_score=0.2, format_regr_score_test=y_test,
    r2_score=0.65)
```

BY DEFAULT THE PROVIDED FUNCTIONS ARE CHECKED AT EACH FIT TO BE THE INVERSE OF EACH OTHER HOWEVER IT IS POSSIBLE TO BYPASS THIS CHECKING BY SETTING CHECKINVERSE TO FALSE

```
def inversefuncX
    return X
```

```
regr = TransformedTargetRegressor(regressor=regressor,
    func=func,
    inversefunc=inversefunc,
    check_inverse=False,
    reg_fit=train, y_train=train,
    transformed_target=regressor,
    print_r2_score=0.2, format_regr_score_test=y_test,
    r2_score=0.45)
```

NOTE THE TRANSFORMATION CAN BE TRIGGERED BY SETTING EITHER TRANSFORMER OR THE PAIR OF FUNCTIONS FUNC AND INVERSEFUNC HOWEVER SETTING BOTH OPTIONS WILL RAISE AN ERROR

EXAMPLES

- EFFECT OF TRANSFORMING THE TARGETS IN REGRESSION MODEL

FEATUREUNION COMPOSITE FEATURE SPACES

FEATUREUNION COMBINES SEVERAL TRANSFORMER OBJECTS INTO A NEW TRANSFORMER THAT COMBINES THEIR OUTPUT A

FEATUREUNION TAKES A LIST OF TRANSFORMER OBJECTS DURING FITTING EACH OF THESE IS FIT TO THE DATA INDEPENDENTLY

THE TRANSFORMERS ARE APPLIED IN PARALLEL AND THE FEATURE MATRICES THEY OUTPUT ARE CONCATENATED SIDE BY SIDE INTO A LARGER MATRIX

WHEN YOU WANT TO APPLY DIFFERENT TRANSFORMATIONS TO EACH FIELD OF THE DATA SEE THE RELATED CLASS SKLEARNCOMPOSE

COLUMNTRANSFORMER SEE USER GUIDE

FEATUREUNION SERVES THE SAME PURPOSES AS PIPELINE CONVENIENCE AND JOINT PARAMETER ESTIMATION AND VALIDATION

FEATUREUNION AND PIPELINE CAN BE COMBINED TO CREATE COMPLEX MODELS

580 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

AFEATUREUNION HAS NO WAY OF CHECKING WHETHER TWO TRANSFORMERS MIGHT PRODUCE IDENTICAL FEATURES IT ONLY PRODUCES A UNION WHEN THE FEATURE SETS ARE DISJOINT AND MAKING SURE THEY ARE THE CALLER’S RESPONSIBILITY

USAGE

AFEATUREUNION IS BUILT USING A LIST OF KEY VALUE PAIRS WHERE THE KEY IS THE NAME YOU WANT TO GIVE TO A GIVEN TRANSFORMATION AN ARBITRARY STRING IT ONLY SERVES AS AN IDENTIFIER AND VALUE IS AN ESTIMATOR OBJECT

```
FROM SKLEARNPIPELINE IMPORT FEATUREUNION
FROM SKLEARNDECOMPOSITION IMPORT PCA
FROM SKLEARNDECOMPOSITION IMPORT KERNELPCA
ESTIMATORS LINEARPCA PCA KERNELPCA KERNELPCA
COMBINED FEATUREUNIONESTIMATORS
COMBINED
FEATUREUNIONNJOBSNONE
TRANSFORMERLISTLINEARPCA PCACOPYTRUE
KERNELPCA KERNELPCAALPHA10
TRANSFORMERWEIGHTSNONE VERBOSEFALSE
LIKE PIPELINES FEATURE UNIONS HAVE A SHORTHAND CONSTRUCTOR CALLED MAKEUNION THAT DOES NOT REQUIRE EXPLICIT NAMING
OF THE COMPONENTS
LIKEPIPELINE INDIVIDUAL STEPS MAY BE REPLACED USING SETPARAMS AND IGNORED BY SETTING TO DROP
COMBINEDSETPARAMSKERNELPCADROP
```

FEATUREUNIONNJOBSNONE

TRANSFORMERLISTLINEARPCA PCACOPYTRUE

KERNELPCA DROP

TRANSFORMERWEIGHTSNONE VERBOSEFALSE

EXAMPLES

- CONCATENATING MULTIPLE FEATURE EXTRACTION METHODS

COLUMNTRANSFORMER FOR HETEROGENEOUS DATA

WARNING THECOMPOSECOLUMNTRANSFORMER CLASS IS EXPERIMENTAL AND THE API IS SUBJECT TO CHANGE

MANY DATASETS CONTAIN FEATURES OF DIFFERENT TYPES SAY TEXT FLOATS AND DATES WHERE EACH TYPE OF FEATURE REQUIRES SEPARATE PREPROCESSING OR FEATURE EXTRACTION STEPS OFTEN IT IS EASIEST TO PREPROCESS DATA BEFORE APPLYING SCIKITLEARN METHODS FOR EXAMPLE USING PANDAS PROCESSING YOUR DATA BEFORE PASSING IT TO SCIKITLEARN MIGHT BE PROBLEMATIC FOR ONE OF THE FOLLOWING REASONS

1 INCORPORATING STATISTICS FROM TEST DATA INTO THE PREPROCESSORS MAKES CROSSVALIDATION SCORES UNRELIABLE KNOWN AS DATA LEAKAGE FOR EXAMPLE IN THE CASE OF SCALERS OR IMPUTING MISSING VALUES

2 YOU MAY WANT TO INCLUDE THE PARAMETERS OF THE PREPROCESSORS IN A PARAMETER SEARCH

THECOLUMNTRANSFORMER HELPS PERFORMING DIFFERENT TRANSFORMATIONS FOR DIFFERENT COLUMNS OF THE DATA WITHIN A PIPELINE THAT IS SAFE FROM DATA LEAKAGE AND THAT CAN BE PARAMETRIZED COLUMNTRANSFORMER WORKS ON ARRAYS SPARSE MATRICES AND PANDAS DATAFRAMES

35 DATASET TRANSFORMATIONS 581

SCIKITLEARN USER GUIDE RELEASE 0213  
TO EACH COLUMN A DIFFERENT TRANSFORMATION CAN BE APPLIED SUCH AS PREPROCESSING OR A SPECIFIC FEATURE EXTRACTION METHOD

```
import pandas as pd
X = pd.DataFrame({
    'city': ['LONDON', 'LONDON', 'PARIS', 'SALLISAW'],
    'title': ['HIS LAST BOW', 'HOW WATSON LEARNED THE TRICK',
             'A MOVEABLE FEAST', 'THE GRAPES OF WRATH'],
    'expertrating': [5, 3, 4, 5],
    'userrating': [4, 5, 4, 3]
})
```

FOR THIS DATA WE MIGHT WANT TO ENCODE THE CITY COLUMN AS A CATEGORICAL VARIABLE USING PREPROCESSING ONEHOTENCODER BUT APPLY A FEATUREEXTRACTIONTEXTCOUNTVECTORIZER TO THE TITLE COLUMN AS WE MIGHT USE MULTIPLE FEATURE EXTRACTION METHODS ON THE SAME COLUMN WE GIVE EACH TRANSFORMER A UNIQUE NAME SAY CITYCATEGORY AND TITLEBOW. BY DEFAULT THE REMAINING RATING COLUMNS ARE IGNORED

```
from sklearn.compose import ColumnTransformer
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.preprocessing import OneHotEncoder

column_trans = ColumnTransformer(
    [('citycategory', OneHotEncoder(dtype=int), ['city']),
     ('titlebow', CountVectorizer(), ['title'])],
    remainder='drop',
    column_transformers=[('columntransformer', ColumnTransformer(
        [('titlebow', CountVectorizer(), ['title'])],
        remainder='drop',
        sparse_threshold=0.3,
        transformer_weights=None,
        transformers=None,
        verbose=0)],
        ['title'],
        ['titlebow'])],
    verbose=0)
```

```
citycategoryX0LONDON citycategoryX0PARIS citycategoryX0SALLISAW
titlebowbow titlebowfeast titlebowgrapes titlebowhis
titlebowhow titlebowlast titlebowlearned titlebowmoveable
titlebowof titlebowthe titlebowtrick titlebowwatson
titlebowwrath
columntranstransformXtoarray
```

```
array([[1, 0, 0, 1, 0, 0, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0],
       [1, 0, 0, 0, 0, 0, 1, 0, 1, 0, 0, 1, 1, 1, 0, 0],
       [0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0],
       [0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 1, 1, 0, 0, 1]])
```

IN THE ABOVE EXAMPLE THE COUNTVECTORIZER EXPECTS A 1D ARRAY AS INPUT AND THEREFORE THE COLUMNS WERE SPECIFIED AS A STRING 'title' HOWEVER PREPROCESSINGONEHOTENCODER AS MOST OF OTHER TRANSFORMERS EXPECTS 2D DATA THEREFORE IN THAT CASE YOU NEED TO SPECIFY THE COLUMN AS A LIST OF STRINGS ['city'] APART FROM A SCALAR OR A SINGLE ITEM LIST THE COLUMN SELECTION CAN BE SPECIFIED AS A LIST OF MULTIPLE ITEMS AN INTEGER ARRAY A SLICE OR A BOOLEAN MASK STRINGS CAN REFERENCE COLUMNS IF THE INPUT IS A DATAFRAME INTEGERS ARE ALWAYS INTERPRETED AS THE POSITIONAL COLUMNS

WE CAN KEEP THE REMAINING RATING COLUMNS BY SETTING REMAINDERPASSTHROUGH. THE VALUES ARE APPENDED TO THE END OF THE TRANSFORMATION

```
column_trans = ColumnTransformer(
    [('citycategory', OneHotEncoder(dtype=int), ['city']),
     ('titlebow', CountVectorizer(), ['title'])],
    remainder='passthrough',
    verbose=0)
```

582 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213  
TITLEBOW COUNTVECTORIZER TITLE  
REMAINDERPASSTHROUGH  
COLUMNTRANSFITTRANSFORMX

ARRAY1 0 0 1 0 0 1 0 1 0 0 0 0 0 0 5 4  
1 0 0 0 0 0 0 1 0 1 0 0 1 1 1 0 3 5  
0 1 0 0 1 0 0 0 0 0 1 0 0 0 0 0 4 4  
0 0 1 0 0 1 0 0 0 0 0 1 1 0 0 1 5 3

THEREMAINDER PARAMETER CAN BE SET TO AN ESTIMATOR TO TRANSFORM THE REMAINING RATING COLUMNS THE TRANSFORMED  
VALUES ARE APPENDED TO THE END OF THE TRANSFORMATION  
FROM SKLEARNPREPROCESSING IMPORT MINMAXSCALER  
COLUMNTRANS COLUMNTRANSFORMER  
CITYCATEGORY ONEHOTENCODER CITY  
TITLEBOW COUNTVECTORIZER TITLE  
REMAINDERMINMAXSCALER  
COLUMNTRANSFITTRANSFORMX 2

ARRAY1 05  
0 1  
05 05  
1 0

THEMAKECOLUMNTRANSFORMER FUNCTION IS AVAILABLE TO MORE EASILY CREATE A COLUMNTRANSFORMER OBJECT  
SPECIFICALLY THE NAMES WILL BE GIVEN AUTOMATICALLY THE EQUIVALENT FOR THE ABOVE EXAMPLE WOULD BE

FROM SKLEARNCOMPOSE IMPORT MAKECOLUMNTRANSFORMER  
COLUMNTRANS MAKECOLUMNTRANSFORMER  
ONEHOTENCODER CITY  
COUNTVECTORIZER TITLE  
REMAINDERMINMAXSCALER  
COLUMNTRANS  
COLUMNTRANSFORMERNJOBSNONE REMAINDERMINMAXSCALERCOPYTRUE  
SPARSETHRESHOLD03  
TRANSFORMERWEIGHTSNONE  
TRANSFORMERSONEHOTENCODER

EXAMPLES  
•COLUMN TRANSFORMER WITH HETEROGENEOUS DATA SOURCES  
•COLUMN TRANSFORMER WITH MIXED TYPES

352 FEATURE EXTRACTION  
THESKLEARNFEATUREEXTRACTION MODULE CAN BE USED TO EXTRACT FEATURES IN A FORMAT SUPPORTED BY MACHINE  
LEARNING ALGORITHMS FROM DATASETS CONSISTING OF FORMATS SUCH AS TEXT AND IMAGE  
NOTE FEATURE EXTRACTION IS VERY DIFFERENT FROM FEATURE SELECTION THE FORMER CONSISTS IN TRANSFORMING ARBITRARY DATA  
SUCH AS TEXT OR IMAGES INTO NUMERICAL FEATURES USABLE FOR MACHINE LEARNING THE LATTER IS A MACHINE LEARNING TECHNIQUE  
35 DATASET TRANSFORMATIONS 583

SCIKITLEARN USER GUIDE RELEASE 0213  
APPLIED ON THESE FEATURES  
LOADING FEATURES FROM DICTS  
THE CLASSDICTVECTORIZER CAN BE USED TO CONVERT FEATURE ARRAYS REPRESENTED AS LISTS OF STANDARD PYTHON DICT  
OBJECTS TO THE NUMPYSCIPY REPRESENTATION USED BY SCIKITLEARN ESTIMATORS  
WHILE NOT PARTICULARLY FAST TO PROCESS PYTHON’S DICT HAS THE ADVANTAGES OF BEING CONVENIENT TO USE BEING SPARSE  
ABSENT FEATURES NEED NOT BE STORED AND STORING FEATURE NAMES IN ADDITION TO VALUES  
DICTVECTORIZER IMPLEMENTS WHAT IS CALLED ONEOFK OR “ONEHOT” CODING FOR CATEGORICAL AKA NOMINAL DISCRETE  
FEATURES CATEGORICAL FEATURES ARE “ATTRIBUTEVALUE” PAIRS WHERE THE VALUE IS RESTRICTED TO A LIST OF DISCRETE OF POSSIBILITIES  
WITHOUT ORDERING EG TOPIC IDENTIFIERS TYPES OF OBJECTS TAGS NAMES  
IN THE FOLLOWING “CITY” IS A CATEGORICAL ATTRIBUTE WHILE “TEMPERATURE” IS A TRADITIONAL NUMERICAL FEATURE  
MEASUREMENTS  
CITY DUBAI TEMPERATURE 33  
CITY LONDON TEMPERATURE 12  
CITY SAN FRANCISCO TEMPERATURE 18

```
FROM SKLEARNFEATUREEXTRACTION IMPORT DICTVECTORIZER
VEC DICTVECTORIZER
VECFITTRANSFORMMEASUREMENTSTOARRAY
ARRAY 1 0 0 33
    0 1 0 12
    0 0 1 18
VECGETFEATURENAMES
CITYDUBAI CITYLONDON CITYSAN FRANCISCO TEMPERATURE
DICTVECTORIZER IS ALSO A USEFUL REPRESENTATION TRANSFORMATION FOR TRAINING SEQUENCE CLASSIFIERS IN NATURAL LAN
GUAGE PROCESSING MODELS THAT TYPICALLY WORK BY EXTRACTING FEATURE WINDOWS AROUND A PARTICULAR WORD OF INTEREST
FOR EXAMPLE SUPPOSE THAT WE HAVE A FIRST ALGORITHM THAT EXTRACTS PART OF SPEECH POS TAGS THAT WE WANT TO USE AS
COMPLEMENTARY TAGS FOR TRAINING A SEQUENCE CLASSIFIER EG A CHUNKER THE FOLLOWING DICT COULD BE SUCH A WINDOW OF
FEATURES EXTRACTED AROUND THE WORD ‘SAT’ IN THE SENTENCE ‘THE CAT SAT ON THE MAT’
POSWINDOW
```

```
WORD2 THE
POS2 DT
WORD1 CAT
POS1 NN
WORD1 ON
POS1 PP
```

IN A REAL APPLICATION ONE WOULD EXTRACT MANY SUCH DICTIONARIES

THIS DESCRIPTION CAN BE VECTORIZED INTO A SPARSE TWODIMENSIONAL MATRIX SUITABLE FOR FEEDING INTO A CLASSIFIER MAYBE  
AFTER BEING PIPED INTO A TEXTTFIDFTRANSFORMER FOR NORMALIZATION  
VEC DICTVECTORIZER  
POSVECTORIZED VECFITTRANSFORMPOSWINDOW  
584 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

POSVECTORIZED

1X6 SPARSE MATRIX OF TYPE NUMPYFLOAT64

WITH 6 STORED ELEMENTS IN COMPRESSED SPARSE FORMAT

POSVECTORIZEDTOARRAY

ARRAY1 1 1 1 1 1

VECGETFEATURENAMES

POS1PP POS1NN POS2DT WORD1ON WORD1CAT WORD2THE

AS YOU CAN IMAGINE IF ONE EXTRACTS SUCH A CONTEXT AROUND EACH INDIVIDUAL WORD OF A CORPUS OF DOCUMENTS THE RESULTING MATRIX WILL BE VERY WIDE MANY ONEHOTFEATURES WITH MOST OF THEM BEING VALUED TO ZERO MOST OF THE TIME SO AS TO MAKE THE RESULTING DATA STRUCTURE ABLE TO FIT IN MEMORY THE DICTVECTORIZER CLASS USES A SCIPYSPARSE MATRIX BY DEFAULT INSTEAD OF A NUMPYNDARRAY

FEATURE HASHING

THE CLASSFEATUREHASHER IS A HIGHSPEED LOWMEMORY VECTORIZER THAT USES A TECHNIQUE KNOWN AS FEATURE HASHING OR THE “HASHING TRICK” INSTEAD OF BUILDING A HASH TABLE OF THE FEATURES ENCOUNTERED IN TRAINING AS THE VECTORIZERS DO INSTANCES OF FEATUREHASHER APPLY A HASH FUNCTION TO THE FEATURES TO DETERMINE THEIR COLUMN INDEX IN SAMPLE MATRICES DIRECTLY THE RESULT IS INCREASED SPEED AND REDUCED MEMORY USAGE AT THE EXPENSE OF INSPECTABILITY THE HASHER DOES NOT REMEMBER WHAT THE INPUT FEATURES LOOKED LIKE AND HAS NO INVERSETRANSFORM METHOD

SINCE THE HASH FUNCTION MIGHT CAUSE COLLISIONS BETWEEN UNRELATED FEATURES A SIGNED HASH FUNCTION IS USED AND THE SIGN OF THE HASH VALUE DETERMINES THE SIGN OF THE VALUE STORED IN THE OUTPUT MATRIX FOR A FEATURE THIS WAY COLLISIONS ARE LIKELY TO CANCEL OUT RATHER THAN ACCUMULATE ERROR AND THE EXPECTED MEAN OF ANY OUTPUT FEATURE’S VALUE IS ZERO

THIS MECHANISM IS ENABLED BY DEFAULT WITH ALTERNATESIGNTRUE AND IS PARTICULARLY USEFUL FOR SMALL HASH TABLE SIZES NFEATURES 10000 FOR LARGE HASH TABLE SIZES IT CAN BE DISABLED TO ALLOW THE OUTPUT TO BE PASSED TO ESTIMATORS LIKE SKLEARNNAIVEBAYESMULTINOMIALNB ORSKLEARNFEATURESELECTIONCHI2

FEATURE SELECTORS THAT EXPECT NONNEGATIVE INPUTS

FEATUREHASHER ACCEPTS EITHER MAPPINGS LIKE PYTHON’S DICT AND ITS VARIANTS IN THE COLLECTIONS MODULE

FEATURE VALUE PAIRS OR STRINGS DEPENDING ON THE CONSTRUCTOR PARAMETER INPUTTYPE MAPPING ARE TREATED AS LISTS OFFEATURE VALUE PAIRS WHILE SINGLE STRINGS HAVE AN IMPLICIT VALUE OF 1 SO FEAT1 FEAT2 FEAT3 IS INTERPRETED AS FEAT1 1 FEAT2 1 FEAT3 1 IF A SINGLE FEATURE OCCURS MULTIPLE TIMES IN A SAMPLE THE ASSOCIATED VALUES WILL BE SUMMED SO FEAT 2 ANDFEAT 35 BECOME FEAT 55 THE OUTPUT FROM FEATUREHASHER IS ALWAYS A SCIPYSPARSE MATRIX IN THE CSR FORMAT

FEATURE HASHING CAN BE EMPLOYED IN DOCUMENT CLASSIFICATION BUT UNLIKE TEXTCOUNTVECTORIZER

FEATUREHASHER DOES NOT DO WORD SPLITTING OR ANY OTHER PREPROCESSING EXCEPT UNICODETOUTF8 ENCODING SEE VECTORIZING A LARGE TEXT CORPUS WITH THE HASHING TRICK BELOW FOR A COMBINED TOKENIZERHASHER

AS AN EXAMPLE CONSIDER A WORDLEVEL NATURAL LANGUAGE PROCESSING TASK THAT NEEDS FEATURES EXTRACTED FROM TOKEN PARTOFSPEECH PAIRS ONE COULD USE A PYTHON GENERATOR FUNCTION TO EXTRACT FEATURES

DEFTOKENFEATURESTOKEN PARTOFSPEECH

IFTOKENISDIGIT

YIELDNUMERIC

ELSE

YIELDTOKENFORMATTOKENLOWER

YIELDTOKENPOS FORMATTOKEN PARTOFSPEECH

IFTOKENOISUPPER

YIELDUPPERCASEINITIAL

IFTOKENISUPPER

YIELDALLUPPERCASE

YIELDPOSFORMATPARTOFSPEECH

THEN THERAWX TO BE FED TO FEATUREHASHERTRANSFORM CAN BE CONSTRUCTED USING

35 DATASET TRANSFORMATIONS 585

SCIKITLEARN USER GUIDE RELEASE 0213

RAWX TOKENFEATURESTOK POSTAGGERTOK FORTOKINCORPUS  
AND FED TO A HASHER WITH  
HASHER FEATUREHASHERINPUTTYPESTRING  
X HASHERTRANSFORMRAWX  
TO GET ASCIPYSPARSE MATRIXX

NOTE THE USE OF A GENERATOR COMPREHENSION WHICH INTRODUCES LAZINESS INTO THE FEATURE EXTRACTION TOKENS ARE ONLY  
PROCESSED ON DEMAND FROM THE HASHER

IMPLEMENTATION DETAILS

FEATUREHASHER USES THE SIGNED 32BIT VARIANT OF MURMURHASH3 AS A RESULT AND BECAUSE OF LIMITATIONS IN SCIPY  
SPARSE THE MAXIMUM NUMBER OF FEATURES SUPPORTED IS CURRENTLY 231−1

THE ORIGINAL FORMULATION OF THE HASHING TRICK BY WEINBERGER ET AL USED TWO SEPARATE HASH FUNCTIONS *h*AND $\pi$ TO DETER  
MINE THE COLUMN INDEX AND SIGN OF A FEATURE RESPECTIVELY THE PRESENT IMPLEMENTATION WORKS UNDER THE ASSUMPTION  
THAT THE SIGN BIT OF MURMURHASH3 IS INDEPENDENT OF ITS OTHER BITS

SINCE A SIMPLE MODULO IS USED TO TRANSFORM THE HASH FUNCTION TO A COLUMN INDEX IT IS ADVISABLE TO USE A POWER OF TWO  
AS THENFEATURES PARAMETER OTHERWISE THE FEATURES WILL NOT BE MAPPED EVENLY TO THE COLUMNS

REFERENCES

- KILIAN WEINBERGER ANIRBAN DASGUPTA JOHN LANGFORD ALEX SMOLA AND JOSH ATTENBERG 2009 FEATURE HASH  
ING FOR LARGE SCALE MULTITASK LEARNING PROC ICML
- MURMURHASH3

TEXT FEATURE EXTRACTION

THE BAG OF WORDS REPRESENTATION

TEXT ANALYSIS IS A MAJOR APPLICATION FIELD FOR MACHINE LEARNING ALGORITHMS HOWEVER THE RAW DATA A SEQUENCE OF  
SYMBOLS CANNOT BE FED DIRECTLY TO THE ALGORITHMS THEMSELVES AS MOST OF THEM EXPECT NUMERICAL FEATURE VECTORS WITH A  
FIXED SIZE RATHER THAN THE RAW TEXT DOCUMENTS WITH VARIABLE LENGTH

IN ORDER TO ADDRESS THIS SCIKITLEARN PROVIDES UTILITIES FOR THE MOST COMMON WAYS TO EXTRACT NUMERICAL FEATURES FROM  
TEXT CONTENT NAMELY

- TOKENIZING STRINGS AND GIVING AN INTEGER ID FOR EACH POSSIBLE TOKEN FOR INSTANCE BY USING WHITESPACES AND  
PUNCTUATION AS TOKEN SEPARATORS
- COUNTING THE OCCURRENCES OF TOKENS IN EACH DOCUMENT
- NORMALIZING AND WEIGHTING WITH DIMINISHING IMPORTANCE TOKENS THAT OCCUR IN THE MAJORITY OF SAMPLES DOCU  
MENTS

IN THIS SCHEME FEATURES AND SAMPLES ARE DEFINED AS FOLLOWS

- EACH INDIVIDUAL TOKEN OCCURRENCE FREQUENCY NORMALIZED OR NOT IS TREATED AS A FEATURE
- THE VECTOR OF ALL THE TOKEN FREQUENCIES FOR A GIVEN DOCUMENT IS CONSIDERED A MULTIVARIATE SAMPLE

586 CHAPTER 3 USER GUIDE



SCIKITLEARN USER GUIDE RELEASE 0213

A CORPUS OF DOCUMENTS CAN THUS BE REPRESENTED BY A MATRIX WITH ONE ROW PER DOCUMENT AND ONE COLUMN PER TOKEN  
EG WORD OCCURRING IN THE CORPUS

WE CALL VECTORIZATION THE GENERAL PROCESS OF TURNING A COLLECTION OF TEXT DOCUMENTS INTO NUMERICAL FEATURE VECTORS THIS  
SPECIFIC STRATEGY TOKENIZATION COUNTING AND NORMALIZATION IS CALLED THE BAG OF WORDS OR “BAG OF NGRAMS” REPRESEN  
TATION DOCUMENTS ARE DESCRIBED BY WORD OCCURRENCES WHILE COMPLETELY IGNORING THE RELATIVE POSITION INFORMATION OF  
THE WORDS IN THE DOCUMENT

SPARSITY

AS MOST DOCUMENTS WILL TYPICALLY USE A VERY SMALL SUBSET OF THE WORDS USED IN THE CORPUS THE RESULTING MATRIX WILL  
HAVE MANY FEATURE VALUES THAT ARE ZEROS TYPICALLY MORE THAN 99 OF THEM

FOR INSTANCE A COLLECTION OF 10000 SHORT TEXT DOCUMENTS SUCH AS EMAILS WILL USE A VOCABULARY WITH A SIZE IN THE ORDER  
OF 100000 UNIQUE WORDS IN TOTAL WHILE EACH DOCUMENT WILL USE 100 TO 1000 UNIQUE WORDS INDIVIDUALLY

IN ORDER TO BE ABLE TO STORE SUCH A MATRIX IN MEMORY BUT ALSO TO SPEED UP ALGEBRAIC OPERATIONS MATRIX VECTOR IMPL  
EMENTATIONS WILL TYPICALLY USE A SPARSE REPRESENTATION SUCH AS THE IMPLEMENTATIONS AVAILABLE IN THE SCIPYSPARSE  
PACKAGE

COMMON VECTORIZER USAGE

COUNTVECTORIZER IMPLEMENTS BOTH TOKENIZATION AND OCCURRENCE COUNTING IN A SINGLE CLASS

FROM SKLEARNFEATUREEXTRACTIONTEXT IMPORT COUNTVECTORIZER

THIS MODEL HAS MANY PARAMETERS HOWEVER THE DEFAULT VALUES ARE QUITE REASONABLE PLEASE SEE THE REFERENCE DOCUMEN  
TATION FOR THE DETAILS

VECTORIZER COUNTVECTORIZER

VECTORIZER

COUNTVECTORIZERANALYZERWORD BINARYFALSE DECODEERRORSTRICT

DTYPE NUMPYINT64 ENCODINGUTF8 INPUTCONTENT

LOWERCASETRUE MAXDF10 MAXFEATURESNONE MINDF1

NGRAMRANGE1 1 PREPROCESSORNONE STOPWORDSNONE

STRIPACCENTSNONE TOKENPATTERNUBWWB

TOKENIZERNONE VOCABULARYNONE

LET’S USE IT TO TOKENIZE AND COUNT THE WORD OCCURRENCES OF A MINIMALISTIC CORPUS OF TEXT DOCUMENTS

CORPUS

THIS IS THE FIRST DOCUMENT

THIS IS THE SECOND SECOND DOCUMENT

AND THE THIRD ONE

IS THIS THE FIRST DOCUMENT

X VECTORIZERFITTRANSFORMCORPUS

X

4X9 SPARSE MATRIX OF TYPE NUMPYINT64

WITH 19 STORED ELEMENTS IN COMPRESSED SPARSE FORMAT

THE DEFAULT CONFIGURATION TOKENIZES THE STRING BY EXTRACTING WORDS OF AT LEAST 2 LETTERS THE SPECIFIC FUNCTION THAT DOES  
THIS STEP CAN BE REQUESTED EXPLICITLY

35 DATASET TRANSFORMATIONS 587

SCIKITLEARN USER GUIDE RELEASE 0213  
ANALYZE VECTORIZERBUILDANALYZER  
ANALYZETHIS IS A TEXT DOCUMENT TO ANALYZE  
THIS IS TEXT DOCUMENT TO ANALYZE  
TRUE  
EACH TERM FOUND BY THE ANALYZER DURING THE FIT IS ASSIGNED A UNIQUE INTEGER INDEX CORRESPONDING TO A COLUMN IN THE  
RESULTING MATRIX THIS INTERPRETATION OF THE COLUMNS CAN BE RETRIEVED AS FOLLOWS  
VECTORIZERGETFEATURENAMES  
AND DOCUMENT FIRST IS ONE  
SECOND THE THIRD THIS  
TRUE  
XTOARRAY  
ARRAY0 1 1 1 0 0 1 0 1  
0 1 0 1 0 2 1 0 1  
1 0 0 0 1 0 1 1 0  
0 1 1 1 0 0 1 0 1  
THE CONVERSE MAPPING FROM FEATURE NAME TO COLUMN INDEX IS STORED IN THE VOCABULARY ATTRIBUTE OF THE VECTORIZER  
VECTORIZERVOCABULARYGETDOCUMENT  
1  
HENCE WORDS THAT WERE NOT SEEN IN THE TRAINING CORPUS WILL BE COMPLETELY IGNORED IN FUTURE CALLS TO THE TRANSFORM  
METHOD  
VECTORIZERTRANSFORMSOMETHING COMPLETELY NEWTOARRAY  
  
ARRAY0 0 0 0 0 0 0 0 0  
NOTE THAT IN THE PREVIOUS CORPUS THE FIRST AND THE LAST DOCUMENTS HAVE EXACTLY THE SAME WORDS HENCE ARE ENCODED IN  
EQUAL VECTORS IN PARTICULAR WE LOSE THE INFORMATION THAT THE LAST DOCUMENT IS AN INTERROGATIVE FORM TO PRESERVE SOME  
OF THE LOCAL ORDERING INFORMATION WE CAN EXTRACT 2GRAMS OF WORDS IN ADDITION TO THE 1GRAMS INDIVIDUAL WORDS  
BIGRAMVECTORIZER COUNTVECTORIZERNGRAMRANGE1 2  
TOKENPATTERNRBWB MINDF1  
ANALYZE BIGRAMVECTORIZERBUILDANALYZER  
ANALYZEBIGRAMS ARE COOL  
BI GRAMS ARE COOL BI GRAMS GRAMS ARE ARE COOL  
TRUE  
THE VOCABULARY EXTRACTED BY THIS VECTORIZER IS HENCE MUCH BIGGER AND CAN NOW RESOLVE AMBIGUITIES ENCODED IN LOCAL  
POSITIONING PATTERNS  
X2 BIGRAMVECTORIZERFITTRANSFORMCORPUSTOARRAY  
X2  
  
ARRAY0 0 1 1 1 1 1 0 0 0 0 0 1 1 0 0 0 0 1 1 0  
0 0 1 0 0 1 1 0 0 2 1 1 1 0 1 0 0 0 1 1 0  
1 1 0 0 0 0 0 0 1 0 0 0 1 0 0 1 1 1 0 0 0  
0 0 1 1 1 1 0 1 0 0 0 0 1 1 0 0 0 0 1 0 1  
IN PARTICULAR THE INTERROGATIVE FORM “IS THIS” IS ONLY PRESENT IN THE LAST DOCUMENT  
FEATUREINDEX BIGRAMVECTORIZERVOCABULARYGETIS THIS  
X2 FEATUREINDEX  
588 CHAPTER 3 USER GUIDE

ARRAY0 0 0 1

USING STOP WORDS

STOP WORDS ARE WORDS LIKE “AND” “THE” “HIM” WHICH ARE PRESUMED TO BE UNINFORMATIVE IN REPRESENTING THE CONTENT OF A TEXT AND WHICH MAY BE REMOVED TO AVOID THEM BEING CONSTRUED AS SIGNAL FOR PREDICTION SOMETIMES HOWEVER SIMILAR WORDS ARE USEFUL FOR PREDICTION SUCH AS IN CLASSIFYING WRITING STYLE OR PERSONALITY

THERE ARE SEVERAL KNOWN ISSUES IN OUR PROVIDED ‘ENGLISH’ STOP WORD LIST SEE NQY18

PLEASE TAKE CARE IN CHOOSING A STOP WORD LIST POPULAR STOP WORD LISTS MAY INCLUDE WORDS THAT ARE HIGHLY INFORMATIVE TO SOME TASKS SUCH AS COMPUTER

YOU SHOULD ALSO MAKE SURE THAT THE STOP WORD LIST HAS HAD THE SAME PREPROCESSING AND TOKENIZATION APPLIED AS THE ONE USED IN THE VECTORIZER THE WORD WE’VE IS SPLIT INTO WEANDVEBY COUNTVECTORIZER’S DEFAULT TOKENIZER SO IF WE’VE IS IN STOPWORDS BUT VEIS NOT VEWILL BE RETAINED FROM WE’VE IN TRANSFORMED TEXT OUR VECTORIZERS WILL TRY TO IDENTIFY AND WARN ABOUT SOME KINDS OF INCONSISTENCIES

REFERENCES

TF-IDF TERM WEIGHTING

IN A LARGE TEXT CORPUS SOME WORDS WILL BE VERY PRESENT EG “THE” “A” “IS” IN ENGLISH HENCE CARRYING VERY LITTLE MEANINGFUL INFORMATION ABOUT THE ACTUAL CONTENTS OF THE DOCUMENT IF WE WERE TO FEED THE DIRECT COUNT DATA DIRECTLY TO A CLASSIFIER THOSE VERY FREQUENT TERMS WOULD SHADOW THE FREQUENCIES OF RARER YET MORE INTERESTING TERMS IN ORDER TO REWEIGHT THE COUNT FEATURES INTO FLOATING POINT VALUES SUITABLE FOR USAGE BY A CLASSIFIER IT IS VERY COMMON TO USE THE TF-IDF TRANSFORM

TF MEANS TERMFREQUENCY WHILE TF-IDF MEANS TERMFREQUENCY TIMES INVERSE DOCUMENTFREQUENCY  $TFIDF_{TD}$

$TFIDF_{TD} \times IDF_T$

USING THE `TFIDFTRANSFORMER` ’S DEFAULT SETTINGS `TFIDFTRANSFORMERNORML2 USEIDFTRUE`

`SMOOTHIDFTRUE` `SUBLINEARTFFFALSE` THE TERM FREQUENCY THE NUMBER OF TIMES A TERM OCCURS IN A GIVEN DOCUMENT IS MULTIPLIED WITH IDF COMPONENT WHICH IS COMPUTED AS

$IDF_t = \log \frac{1}{p_t}$

$1DF_t = \frac{1}{p_t}$

WHERE  $p_t$  IS THE TOTAL NUMBER OF DOCUMENTS IN THE DOCUMENT SET AND  $DF_t$  IS THE NUMBER OF DOCUMENTS IN THE DOCUMENT SET THAT CONTAIN TERM  $t$  THE RESULTING  $TFIDF$  VECTORS ARE THEN NORMALIZED BY THE EUCLIDEAN NORM

$\frac{1}{\sqrt{\sum_{t=1}^n (TFIDF_t)^2}}$

$\frac{1}{\sqrt{2}}$

$\frac{1}{\sqrt{2}} \frac{1}{\sqrt{2}} \dots \frac{1}{\sqrt{2}}$

THIS WAS ORIGINALLY A TERM WEIGHTING SCHEME DEVELOPED FOR INFORMATION RETRIEVAL AS A RANKING FUNCTION FOR SEARCH ENGINES RESULTS THAT HAS ALSO FOUND GOOD USE IN DOCUMENT CLASSIFICATION AND CLUSTERING

THE FOLLOWING SECTIONS CONTAIN FURTHER EXPLANATIONS AND EXAMPLES THAT ILLUSTRATE HOW THE  $TFIDF$ S ARE COMPUTED EXACTLY AND HOW THE  $TFIDF$ S COMPUTED IN SCIKITLEARN’S `TFIDFTRANSFORMER` AND `TFIDFVECTORIZER` DIFFER SLIGHTLY FROM THE STANDARD TEXTBOOK NOTATION THAT DEFINES THE IDF AS

$IDF_t = \log \frac{1}{p_t}$

$1DF_t = \frac{1}{p_t}$

IN THE `TFIDFTRANSFORMER` AND `TFIDFVECTORIZER` WITH `SMOOTHIDFFALSE` THE “1” COUNT IS ADDED TO THE IDF INSTEAD OF THE IDF’S DENOMINATOR

$IDF_t = \log \frac{1}{p_t + 1}$

$DF_t = \frac{1}{p_t + 1}$

SCIKITLEARN USER GUIDE RELEASE 0213

THIS NORMALIZATION IS IMPLEMENTED BY THE TFIDFTRANSFORMER CLASS

```
FROM SKLEARNFEATUREEXTRACTIONTEXT IMPORT TFIDFTRANSFORMER
TRANSFORMER TFIDFTRANSFORMERSMOOTHIDF FALSE
TRANSFORMER
TFIDFTRANSFORMERNORML2 SMOOTHIDFFALSE SUBLINEARTFFALSE
USEIDFTRUE
```

AGAIN PLEASE SEE THE REFERENCE DOCUMENTATION FOR THE DETAILS ON ALL THE PARAMETERS

LET’S TAKE AN EXAMPLE WITH THE FOLLOWING COUNTS THE FIRST TERM IS PRESENT 100 OF THE TIME HENCE NOT VERY INTERESTING

THE TWO OTHER FEATURES ONLY IN LESS THAN 50 OF THE TIME HENCE PROBABLY MORE REPRESENTATIVE OF THE CONTENT OF THE DOCUMENTS

COUNTS 3 0 1

2 0 0

3 0 0

4 0 0

3 2 0

3 0 2

TFIDF TRANSFORMERFITTRANSFORMCOUNTS

TFIDF

6X3 SPARSE MATRIX OF TYPE NUMPYFLOAT64

WITH 9 STORED ELEMENTS IN COMPRESSED SPARSE FORMAT

TFIDFTOARRAY

ARRAY081940995 0 057320793

1 0 0

1 0 0

1 0 0

047330339 088089948 0

058149261 0 081355169

EACH ROW IS NORMALIZED TO HAVE UNIT EUCLIDEAN NORM

$\sqrt{\frac{1^2+0^2+0^2}{1^2+0^2+0^2}}$

$\sqrt{\frac{1^2+0^2+0^2}{1^2+0^2+0^2}}$

FOR EXAMPLE WE CAN COMPUTE THE TFIDF OF THE FIRST TERM IN THE FIRST DOCUMENT IN THE COUNTS ARRAY AS FOLLOWS

$\frac{1}{6}$

$DF_{TERM1} = 6$

$IDF_{TERM1} = \log \frac{1}{DF_{TERM1}}$

$DF_{TERM1} = 1 \log \frac{1}{6} = 1$

$TFIDF_{TERM1} = TF \times IDF = 3 \times 1 = 3$

NOW IF WE REPEAT THIS COMPUTATION FOR THE REMAINING 2 TERMS IN THE DOCUMENT WE GET

$TFIDF_{TERM2} = 0 \times \log \frac{1}{6} = 1 \times 0 = 0$

$TFIDF_{TERM3} = 1 \times \log \frac{1}{2} = 1 \times 0.6931 \approx 0.6931$

AND THE VECTOR OF RAW TFIDFS

TFIDF RAW 3020986

THEN APPLYING THE EUCLIDEAN L2 NORM WE OBTAIN THE FOLLOWING TFIDFS FOR DOCUMENT 1

$\sqrt{3^2+0^2+0.6931^2}$

3202209862 081900573

590 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

FURTHERMORE THE DEFAULT PARAMETER SMOOTHIDFTRUE ADDS “1” TO THE NUMERATOR AND DENOMINATOR AS IF AN EXTRA DOCUMENT WAS SEEN CONTAINING EVERY TERM IN THE COLLECTION EXACTLY ONCE WHICH PREVENTS ZERO DIVISIONS

$IDF = \log(1 + \frac{1}{\text{term frequency}})$

USING THIS MODIFICATION THE TFIDF OF THE THIRD TERM IN DOCUMENT 1 CHANGES TO 18473

$TFIDF \text{ TERM3} = 1 \times \log(73) \approx 18473$

AND THE L2NORMALIZED TFIDF CHANGES TO

$\frac{3018473}{\sqrt{3202184732 + 08515005243}}$

TRANSFORMER TFIDFTRANSFORMER

TRANSFORMERFITTRANSFORMCOUNTSTOARRAY

ARRAY085151335 0 052433293

1 0 0

1 0 0

1 0 0

055422893 083236428 0

063035731 0 077630514

THE WEIGHTS OF EACH FEATURE COMPUTED BY THE FIT METHOD CALL ARE STORED IN A MODEL ATTRIBUTE

TRANSFORMERIDF

ARRAY1 225 184

AS TF-IDF IS VERY OFTEN USED FOR TEXT FEATURES THERE IS ALSO ANOTHER CLASS CALLED TFIDFVECTORIZER THAT COMBINES ALL THE OPTIONS OF COUNTVECTORIZER ANDTFIDFTRANSFORMER IN A SINGLE MODEL

FROM SKLEARNFEATUREEXTRACTIONTEXT IMPORT TFIDFVECTORIZER

VECTORIZER TFIDFVECTORIZER

VECTORIZERFITTRANSFORMCORPUS

4X9 SPARSE MATRIX OF TYPE NUMPYFLOAT64

WITH 19 STORED ELEMENTS IN COMPRESSED SPARSE FORMAT

WHILE THE TF-IDF NORMALIZATION IS OFTEN VERY USEFUL THERE MIGHT BE CASES WHERE THE BINARY OCCURRENCE MARKERS MIGHT OFFER BETTER FEATURES THIS CAN BE ACHIEVED BY USING THE BINARY PARAMETER OF COUNTVECTORIZER IN PARTICULAR

SOME ESTIMATORS SUCH AS BERNOULLI NAIVE BAYES EXPLICITLY MODEL DISCRETE BOOLEAN RANDOM VARIABLES ALSO VERY SHORT TEXTS ARE LIKELY TO HAVE NOISY TF-IDF VALUES WHILE THE BINARY OCCURRENCE INFO IS MORE STABLE

AS USUAL THE BEST WAY TO ADJUST THE FEATURE EXTRACTION PARAMETERS IS TO USE A CROSSVALIDATED GRID SEARCH FOR INSTANCE BY PIPELINING THE FEATURE EXTRACTOR WITH A CLASSIFIER

•SAMPLE PIPELINE FOR TEXT FEATURE EXTRACTION AND EVALUATION

DECODING TEXT FILES

TEXT IS MADE OF CHARACTERS BUT FILES ARE MADE OF BYTES THESE BYTES REPRESENT CHARACTERS ACCORDING TO SOME ENCODING TO WORK WITH TEXT FILES IN PYTHON THEIR BYTES MUST BE DECODED TO A CHARACTER SET CALLED UNICODE COMMON ENCODINGS ARE ASCII LATIN1 WESTERN EUROPE KOI8R RUSSIAN AND THE UNIVERSAL ENCODINGS UTF8 AND UTF16 MANY OTHERS EXIST

NOTE AN ENCODING CAN ALSO BE CALLED A ‘CHARACTER SET’ BUT THIS TERM IS LESS ACCURATE SEVERAL ENCODINGS CAN EXIST FOR A SINGLE CHARACTER SET

35 DATASET TRANSFORMATIONS 591

SCIKITLEARN USER GUIDE RELEASE 0213

THE TEXT FEATURE EXTRACTORS IN SCIKITLEARN KNOW HOW TO DECODE TEXT FILES BUT ONLY IF YOU TELL THEM WHAT ENCODING THE FILES ARE IN THE COUNTVECTORIZER TAKES ANENCODING PARAMETER FOR THIS PURPOSE FOR MODERN TEXT FILES THE CORRECT ENCODING IS PROBABLY UTF8 WHICH IS THEREFORE THE DEFAULT ENCODINGUTF8 IF THE TEXT YOU ARE LOADING IS NOT ACTUALLY ENCODED WITH UTF8 HOWEVER YOU WILL GET A UNICODEDECODEERROR THE VECTORIZERS CAN BE TOLD TO BE SILENT ABOUT DECODING ERRORS BY SETTING THE DECODEERROR PARAMETER TO EITHER IGNORE ORREPLACE SEE THE DOCUMENTATION FOR THE PYTHON FUNCTION BYTESDECODE FOR MORE DETAILS TYPE HELPBYTESDECODE AT THE PYTHON PROMPT

IF YOU ARE HAVING TROUBLE DECODING TEXT HERE ARE SOME THINGS TO TRY

- FIND OUT WHAT THE ACTUAL ENCODING OF THE TEXT IS THE FILE MIGHT COME WITH A HEADER OR README THAT TELLS YOU THE ENCODING OR THERE MIGHT BE SOME STANDARD ENCODING YOU CAN ASSUME BASED ON WHERE THE TEXT COMES FROM
- YOU MAY BE ABLE TO FIND OUT WHAT KIND OF ENCODING IT IS IN GENERAL USING THE UNIX COMMAND FILE THE PYTHON CHARDET MODULE COMES WITH A SCRIPT CALLED CHARDTECTPY THAT WILL GUESS THE SPECIFIC ENCODING THOUGH YOU CANNOT RELY ON ITS GUESS BEING CORRECT
- YOU COULD TRY UTF8 AND DISREGARD THE ERRORS YOU CAN DECODE BYTE STRINGS WITH BYTES DECODEERRORSREPLACE TO REPLACE ALL DECODING ERRORS WITH A MEANINGLESS CHARACTER OR SET DECODEERRORREPLACE IN THE VECTORIZER THIS MAY DAMAGE THE USEFULNESS OF YOUR FEATURES
- REAL TEXT MAY COME FROM A VARIETY OF SOURCES THAT MAY HAVE USED DIFFERENT ENCODINGS OR EVEN BE SLOPPILY DECODED IN A DIFFERENT ENCODING THAN THE ONE IT WAS ENCODED WITH THIS IS COMMON IN TEXT RETRIEVED FROM THE WEB THE PYTHON PACKAGE FTFY CAN AUTOMATICALLY SORT OUT SOME CLASSES OF DECODING ERRORS SO YOU COULD TRY DECODING THE UNKNOWN TEXT AS LATIN1 AND THEN USING FTFY TO FIX ERRORS
- IF THE TEXT IS IN A MISHMASH OF ENCODINGS THAT IS SIMPLY TOO HARD TO SORT OUT WHICH IS THE CASE FOR THE 20 NEWSGROUPS DATASET YOU CAN FALL BACK ON A SIMPLE SINGLEBYTE ENCODING SUCH AS LATIN1 SOME TEXT MAY DISPLAY INCORRECTLY BUT AT LEAST THE SAME SEQUENCE OF BYTES WILL ALWAYS REPRESENT THE SAME FEATURE FOR EXAMPLE THE FOLLOWING SNIPPET USES CHARDET NOT SHIPPED WITH SCIKITLEARN MUST BE INSTALLED SEPARATELY TO FIGURE OUT THE ENCODING OF THREE TEXTS IT THEN VECTORIZES THE TEXTS AND PRINTS THE LEARNED VOCABULARY THE OUTPUT IS NOT SHOWN HERE

```
import chardet
TEXT1 = b"SEI MIR GEGR XC3XBCXC3X9F T MEIN SAUERKRAUT"
TEXT2 = b"HOLDSELIG SIND DEINE GER XFCCH"
TEXT3 = b" B XFFXFE AX00UX00FX00 X00 FX00LX00XFCX00 GX00EX00LX00NX00
↪X00DX00EX00SX00 X00 GX00EX00SX00AX00NX00GX00EX00SX00X00
↪X00HX00EX00RX00ZX00LX00IX00EX00BX00CX00HX00EX00NX00X00
↪X00TX00RX00AX00GX00 X00 IX00CX00HX00 X00 DX00IX00CX00HX00
↪X00FX00OX00RX00TX00"
DECODED = [chardet.detect(x).encoding
            for x in [TEXT1, TEXT2, TEXT3]]
V = CountVectorizer(fit_decode_vocab)
for term in V.printv:
    depending on the version of chardet it might get the first one wrong
    for an introduction to unicode and character encodings in general see joel spolsky's absolute minimum every
    software developer must know about unicode
    applications and examples
    the bag of words representation is quite simplistic but surprisingly useful in practice
    in particular in a supervised setting it can be successfully combined with fast and scalable linear models to train
    document classifiers for instance
592 chapter 3 user guide
```

SCIKITLEARN USER GUIDE RELEASE 0213

- CLASSIFICATION OF TEXT DOCUMENTS USING SPARSE FEATURES  
IN AN UNSUPERVISED SETTING IT CAN BE USED TO GROUP SIMILAR DOCUMENTS TOGETHER BY APPLYING CLUSTERING ALGORITHMS SUCH ASKMEANS
  - CLUSTERING TEXT DOCUMENTS USING KMEANS
- FINALLY IT IS POSSIBLE TO DISCOVER THE MAIN TOPICS OF A CORPUS BY RELAXING THE HARD ASSIGNMENT CONSTRAINT OF CLUSTERING FOR INSTANCE BY USING NONNEGATIVE MATRIX FACTORIZATION NMF OR NNMF
- TOPIC EXTRACTION WITH NONNEGATIVE MATRIX FACTORIZATION AND LATENT DIRICHLET ALLOCATION

LIMITATIONS OF THE BAG OF WORDS REPRESENTATION

A COLLECTION OF UNIGRAMS WHAT BAG OF WORDS IS CANNOT CAPTURE PHRASES AND MULTIWORD EXPRESSIONS EFFECTIVELY DISREGARDING ANY WORD ORDER DEPENDENCE ADDITIONALLY THE BAG OF WORDS MODEL DOESN'T ACCOUNT FOR POTENTIAL MISSPELLINGS OR WORD DERIVATIONS

NGRAMS TO THE RESCUE INSTEAD OF BUILDING A SIMPLE COLLECTION OF UNIGRAMS N1 ONE MIGHT PREFER A COLLECTION OF BIGRAMS N2 WHERE OCCURRENCES OF PAIRS OF CONSECUTIVE WORDS ARE COUNTED  
ONE MIGHT ALTERNATIVELY CONSIDER A COLLECTION OF CHARACTER NGRAMS A REPRESENTATION RESILIENT AGAINST MISSPELLINGS AND DERIVATIONS

FOR EXAMPLE LET'S SAY WE'RE DEALING WITH A CORPUS OF TWO DOCUMENTS WORDS WPRDS THE SECOND DOCUMENT CONTAINS A MISSPELLING OF THE WORD 'WORDS' A SIMPLE BAG OF WORDS REPRESENTATION WOULD CONSIDER THESE TWO AS VERY DISTINCT DOCUMENTS DIFFERING IN BOTH OF THE TWO POSSIBLE FEATURES A CHARACTER 2GRAM REPRESENTATION HOWEVER WOULD FIND THE DOCUMENTS MATCHING IN 4 OUT OF 8 FEATURES WHICH MAY HELP THE PREFERRED CLASSIFIER DECIDE BETTER

NGRAMVECTORIZER COUNTVECTORIZERANALYZERCHARWB NGRAMRANGE2 2  
COUNTS NGRAMVECTORIZERFITTRANSFORMWORDS WPRDS  
NGRAMVECTORIZERGETFEATURENAMES  
 W DS OR PR RD S WO WP  
TRUE  
COUNTSTOARRAYASTYPEINT  
ARRAY1 1 1 0 1 1 1 0  
1 1 0 1 1 1 0 1  
IN THE ABOVE EXAMPLE CHARWB ANALYZER IS USED WHICH CREATES NGRAMS ONLY FROM CHARACTERS INSIDE WORD BOUNDARIES PADDED WITH SPACE ON EACH SIDE THE CHAR ANALYZER ALTERNATIVELY CREATES NGRAMS THAT SPAN ACROSS WORDS  
NGRAMVECTORIZER COUNTVECTORIZERANALYZERCHARWB NGRAMRANGE5 5  
NGRAMVECTORIZERFITTRANSFORMJUMPY FOX

1X4 SPARSE MATRIX OF TYPE NUMPYINT64  
WITH 4 STORED ELEMENTS IN COMPRESSED SPARSE FORMAT  
NGRAMVECTORIZERGETFEATURENAMES  
 FOX JUMP JUMPY UMPY  
TRUE  
NGRAMVECTORIZER COUNTVECTORIZERANALYZERCHAR NGRAMRANGE5 5  
NGRAMVECTORIZERFITTRANSFORMJUMPY FOX

1X5 SPARSE MATRIX OF TYPE NUMPYINT64  
WITH 5 STORED ELEMENTS IN COMPRESSED SPARSE FORMAT  
NGRAMVECTORIZERGETFEATURENAMES  
 JUMPY MPY F PY FO UMPY Y FOX  
TRUE  
35 DATASET TRANSFORMATIONS 593

SCIKITLEARN USER GUIDE RELEASE 0213

THE WORD BOUNDARIESAWARE VARIANT CHARWB IS ESPECIALLY INTERESTING FOR LANGUAGES THAT USE WHITESPACES FOR WORD SEPARATION AS IT GENERATES SIGNIFICANTLY LESS NOISY FEATURES THAN THE RAW CHAR VARIANT IN THAT CASE FOR SUCH LANGUAGES IT CAN INCREASE BOTH THE PREDICTIVE ACCURACY AND CONVERGENCE SPEED OF CLASSIFIERS TRAINED USING SUCH FEATURES WHILE RETAINING THE ROBUSTNESS WITH REGARDS TO MISSPELLINGS AND WORD DERIVATIONS WHILE SOME LOCAL POSITIONING INFORMATION CAN BE PRESERVED BY EXTRACTING NGRAMS INSTEAD OF INDIVIDUAL WORDS BAG OF WORDS AND BAG OF NGRAMS DESTROY MOST OF THE INNER STRUCTURE OF THE DOCUMENT AND HENCE MOST OF THE MEANING CARRIED BY THAT INTERNAL STRUCTURE

IN ORDER TO ADDRESS THE WIDER TASK OF NATURAL LANGUAGE UNDERSTANDING THE LOCAL STRUCTURE OF SENTENCES AND PARAGRAPHS SHOULD THUS BE TAKEN INTO ACCOUNT MANY SUCH MODELS WILL THUS BE CASTED AS “STRUCTURED OUTPUT” PROBLEMS WHICH ARE CURRENTLY OUTSIDE OF THE SCOPE OF SCIKITLEARN

VECTORIZING A LARGE TEXT CORPUS WITH THE HASHING TRICK

THE ABOVE VECTORIZATION SCHEME IS SIMPLE BUT THE FACT THAT IT HOLDS AN IN MEMORY MAPPING FROM THE STRING TOKENS TO THE INTEGER FEATURE INDICES THE VOCABULARY ATTRIBUTE CAUSES SEVERAL PROBLEMS WHEN DEALING WITH LARGE DATASETS

- THE LARGER THE CORPUS THE LARGER THE VOCABULARY WILL GROW AND HENCE THE MEMORY USE TOO
- FITTING REQUIRES THE ALLOCATION OF INTERMEDIATE DATA STRUCTURES OF SIZE PROPORTIONAL TO THAT OF THE ORIGINAL DATASET
- BUILDING THE WORD MAPPING REQUIRES A FULL PASS OVER THE DATASET HENCE IT IS NOT POSSIBLE TO FIT TEXT CLASSIFIERS IN A STRICTLY ONLINE MANNER
- PICKLING AND UNPICKLING VECTORIZERS WITH A LARGE VOCABULARY CAN BE VERY SLOW TYPICALLY MUCH SLOWER THAN PICKLING UNPICKLING FLAT DATA STRUCTURES SUCH AS A NUMPY ARRAY OF THE SAME SIZE
- IT IS NOT EASILY POSSIBLE TO SPLIT THE VECTORIZATION WORK INTO CONCURRENT SUB TASKS AS THE VOCABULARY ATTRIBUTE WOULD HAVE TO BE A SHARED STATE WITH A FINE GRAINED SYNCHRONIZATION BARRIER THE MAPPING FROM TOKEN STRING TO FEATURE INDEX IS DEPENDENT ON ORDERING OF THE FIRST OCCURRENCE OF EACH TOKEN HENCE WOULD HAVE TO BE SHARED POTENTIALLY HARMING THE CONCURRENT WORKERS’ PERFORMANCE TO THE POINT OF MAKING THEM SLOWER THAN THE SEQUENTIAL VARIANT

IT IS POSSIBLE TO OVERCOME THOSE LIMITATIONS BY COMBINING THE “HASHING TRICK” FEATURE HASHING IMPLEMENTED BY THE SKLEARNFEATUREEXTRACTIONFEATUREHASHER CLASS AND THE TEXT PREPROCESSING AND TOKENIZATION FEATURES OF THE COUNTVECTORIZER

THIS COMBINATION IS IMPLEMENTING IN HASHINGVECTORIZER A TRANSFORMER CLASS THAT IS MOSTLY API COMPATIBLE WITH COUNTVECTORIZER HASHINGVECTORIZER IS STATELESS MEANING THAT YOU DON’T HAVE TO CALL FIT ON IT

FROM SKLEARNFEATUREEXTRACTIONTEXT IMPORT HASHINGVECTORIZER

HV HASHINGVECTORIZER(NFEATURES=10

HVTRANSFORMCORPUS

4X10 SPARSE MATRIX OF TYPE NUMPYFLOAT64

WITH 16 STORED ELEMENTS IN COMPRESSED SPARSE FORMAT

YOU CAN SEE THAT 16 NONZERO FEATURE TOKENS WERE EXTRACTED IN THE VECTOR OUTPUT THIS IS LESS THAN THE 19 NONZEROS EXTRACTED PREVIOUSLY BY THE COUNTVECTORIZER ON THE SAME TOY CORPUS THE DISCREPANCY COMES FROM HASH FUNCTION COLLISIONS BECAUSE OF THE LOW VALUE OF THE NFEATURES PARAMETER

IN A REAL WORLD SETTING THE NFEATURES PARAMETER CAN BE LEFT TO ITS DEFAULT VALUE OF 220 ROUGHLY ONE MILLION POSSIBLE FEATURES IF MEMORY OR DOWNSTREAM MODELS SIZE IS AN ISSUE SELECTING A LOWER VALUE SUCH AS 218 MIGHT HELP WITHOUT INTRODUCING TOO MANY ADDITIONAL COLLISIONS ON TYPICAL TEXT CLASSIFICATION TASKS



SCIKITLEARN USER GUIDE RELEASE 0213  
NOTE THAT THE DIMENSIONALITY DOES NOT AFFECT THE CPU TRAINING TIME OF ALGORITHMS WHICH OPERATE ON CSR MATRICES  
LINEARSVCDUALTRUE PERCEPTRON SGDClassifier PASSIVEAGGRESSIVE BUT IT DOES FOR ALGO  
RITHMS THAT WORK WITH CSC MATRICES LINEARSVCDUALFALSE LASSO ETC  
LET’S TRY AGAIN WITH THE DEFAULT SETTING  
HV HASHINGVECTORIZER  
HVTRANSFORMCORPUS

4X1048576 SPARSE MATRIX OF TYPE NUMPYFLOAT64  
WITH 19 STORED ELEMENTS IN COMPRESSED SPARSE FORMAT  
WE NO LONGER GET THE COLLISIONS BUT THIS COMES AT THE EXPENSE OF A MUCH LARGER DIMENSIONALITY OF THE OUTPUT SPACE OF  
COURSE OTHER TERMS THAN THE 19 USED HERE MIGHT STILL COLLIDE WITH EACH OTHER  
THEHASHINGVECTORIZER ALSO COMES WITH THE FOLLOWING LIMITATIONS

- IT IS NOT POSSIBLE TO INVERT THE MODEL NO INVERSETRANSFORM METHOD NOR TO ACCESS THE ORIGINAL STRING REPRESENTATION OF THE FEATURES BECAUSE OF THE ONEWAY NATURE OF THE HASH FUNCTION THAT PERFORMS THE MAPPING
- IT DOES NOT PROVIDE IDF WEIGHTING AS THAT WOULD INTRODUCE STATEFULNESS IN THE MODEL A TFIDFTRANSFORMER CAN BE APPENDED TO IT IN A PIPELINE IF REQUIRED

PERFORMING OUTFCORE SCALING WITH HASHINGVECTORIZER  
AN INTERESTING DEVELOPMENT OF USING A HASHINGVECTORIZER IS THE ABILITY TO PERFORM OUTFCORE SCALING THIS MEANS THAT WE CAN LEARN FROM DATA THAT DOES NOT FIT INTO THE COMPUTER’S MAIN MEMORY  
A STRATEGY TO IMPLEMENT OUTFCORE SCALING IS TO STREAM DATA TO THE ESTIMATOR IN MINIBATCHES EACH MINIBATCH IS VECTORIZED USING HASHINGVECTORIZER SO AS TO GUARANTEE THAT THE INPUT SPACE OF THE ESTIMATOR HAS ALWAYS THE SAME DIMENSIONALITY THE AMOUNT OF MEMORY USED AT ANY TIME IS THUS BOUNDED BY THE SIZE OF A MINIBATCH ALTHOUGH THERE IS NO LIMIT TO THE AMOUNT OF DATA THAT CAN BE INGESTED USING SUCH AN APPROACH FROM A PRACTICAL POINT OF VIEW THE LEARNING TIME IS OFTEN LIMITED BY THE CPU TIME ONE WANTS TO SPEND ON THE TASK  
FOR A FULLFLEDGED EXAMPLE OF OUTFCORE SCALING IN A TEXT CLASSIFICATION TASK SEE OUTFCORE CLASSIFICATION OF TEXT DOCUMENTS

CUSTOMIZING THE VECTORIZER CLASSES  
IT IS POSSIBLE TO CUSTOMIZE THE BEHAVIOR BY PASSING A CALLABLE TO THE VECTORIZER CONSTRUCTOR

```
def mytokenizers  
    return ssplit
```

```
Vectorizer CountVectorizerTokenizerMyTokenizer  
VectorizerBuildAnalyzerSome Punctuation  
Some Punctuation  
True  
In particular we name
```

- PREPROCESSOR A CALLABLE THAT TAKES AN ENTIRE DOCUMENT AS INPUT AS A SINGLE STRING AND RETURNS A POSSIBLY TRANSFORMED VERSION OF THE DOCUMENT STILL AS AN ENTIRE STRING THIS CAN BE USED TO REMOVE HTML TAGS LOWERCASE THE ENTIRE DOCUMENT ETC
- TOKENIZER A CALLABLE THAT TAKES THE OUTPUT FROM THE PREPROCESSOR AND SPLITS IT INTO TOKENS THEN RETURNS A LIST OF THESE

35 DATASET TRANSFORMATIONS 595

SCIKITLEARN USER GUIDE RELEASE 0213

•ANALYZER A CALLABLE THAT REPLACES THE PREPROCESSOR AND TOKENIZER THE DEFAULT ANALYZERS ALL CALL THE PREPRO  
CESSOR AND TOKENIZER BUT CUSTOM ANALYZERS WILL SKIP THIS NGRAM EXTRACTION AND STOP WORD FILTERING TAKE PLACE  
AT THE ANALYZER LEVEL SO A CUSTOM ANALYZER MAY HAVE TO REPRODUCE THESE STEPS  
LUCENE USERS MIGHT RECOGNIZE THESE NAMES BUT BE AWARE THAT SCIKITLEARN CONCEPTS MAY NOT MAP ONETOONE ONTO  
LUCENE CONCEPTS  
TO MAKE THE PREPROCESSOR TOKENIZER AND ANALYZERS AWARE OF THE MODEL PARAMETERS IT IS POSSIBLE TO DERIVE FROM THE CLASS  
AND OVERRIDE THE BUILDPREPROCESSOR BUILDTOKENIZER ANDBUILDANALYZER FACTORY METHODS INSTEAD  
OF PASSING CUSTOM FUNCTIONS

SOME TIPS AND TRICKS

• IF DOCUMENTS ARE PRETOKENIZED BY AN EXTERNAL PACKAGE THEN STORE THEM IN FILES OR STRINGS WITH THE TOKENS  
SEPARATED BY WHITESPACE AND PASS ANALYZERSTRSPLIT  
• FANCY TOKENLEVEL ANALYSIS SUCH AS STEMMING LEMMATIZING COMPOUND SPLITTING FILTERING BASED ON PARTOF  
SPEECH ETC ARE NOT INCLUDED IN THE SCIKITLEARN CODEBASE BUT CAN BE ADDED BY CUSTOMIZING EITHER THE TOKENIZER  
OR THE ANALYZER HERE'S A COUNTVECTORIZER WITH A TOKENIZER AND LEMMATIZER USING NLTK  
FROM NLTK IMPORT WORDTOKENIZE  
FROM NLTKSTEM IMPORT WORDNETLEMMATIZER  
CLASS LEMMATOKENIZER OBJECT  
DEF INITSELF  
SELFVNL WORDNETLEMMATIZER  
DEF CALLSELF DOC  
RETURN SELFVNLLEMMATIZET FORTINWORDTOKENIZEDOC

VECT COUNTVECTORIZERTOKENIZERLEMMATOKENIZER  
NOTE THAT THIS WILL NOT FILTER OUT PUNCTUATION  
THE FOLLOWING EXAMPLE WILL FOR INSTANCE TRANSFORM SOME BRITISH SPELLING TO AMERICAN SPELLING  
IMPORT RE  
DEF TOBRITISHTOKENS  
FOR TINTOKENS  
T RESUBROUR R1OR T  
T RESUBRBTRE R1ER T  
T RESUBRIYEINGATION R1Z2 T  
T RESUBROGUE OG T  
YIELD T

CLASS CUSTOMVECTORIZER COUNTVECTORIZER  
DEF BUILDTOKENIZERSELF  
TOKENIZE SUPERBUILDTOKENIZER  
RETURN LAMBDA DOC LISTTOBRITISHTOKENIZEDOC

PRINTCUSTOMVECTORIZERBUILDANALYZERUCOLOR COLOUR  
COLOR COLOR  
FOR OTHER STYLES OF PREPROCESSING EXAMPLES INCLUDE STEMMING LEMMATIZATION OR NORMALIZING NUMERICAL TOKENS  
WITH THE LATTER ILLUSTRATED IN  
-BICLUSTERING DOCUMENTS WITH THE SPECTRAL COCLUSTERING ALGORITHM  
CUSTOMIZING THE VECTORIZER CAN ALSO BE USEFUL WHEN HANDLING ASIAN LANGUAGES THAT DO NOT USE AN EXPLICIT WORD SEPARATOR  
SUCH AS WHITESPACE  
596 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

IMAGE FEATURE EXTRACTION

PATCH EXTRACTION

THEEXTRACTPATCHES2D FUNCTION EXTRACTS PATCHES FROM AN IMAGE STORED AS A TWODIMENSIONAL ARRAY OR THREEDIMENSIONAL WITH COLOR INFORMATION ALONG THE THIRD AXIS FOR REBUILDING AN IMAGE FROM ALL ITS PATCHES USE RECONSTRUCTFROMPATCHES2D FOR EXAMPLE LET USE GENERATE A 4X4 PIXEL PICTURE WITH 3 COLOR CHANNELS EG IN RGB FORMAT

```
import numpy as np
from sklearn.feature_extraction import image
oneimage = np.arange(4 * 3).reshape(4, 3)
oneimage = 0 # R CHANNEL OF A FAKE RGB PICTURE
array([[0, 3, 6, 9],
       [12, 15, 18, 21],
       [24, 27, 30, 33],
       [36, 39, 42, 45]])
patches = image.extract_patches_2d(oneimage, (2, 2), max_patches=2,
                                   random_state=0)
patches.shape
(2, 2, 2, 3)
patches[0]
array([[0, 3],
       [12, 15],
       [15, 18],
       [27, 30]])
patches = image.extract_patches_2d(oneimage, (2, 2))
patches.shape
(9, 2, 2, 3)
patches[4][0]
array([15, 18],
      [27, 30])
```

LET US NOW TRY TO RECONSTRUCT THE ORIGINAL IMAGE FROM THE PATCHES BY AVERAGING ON OVERLAPPING AREAS

```
reconstructed = image.reconstruct_from_patches_2d(patches, (4, 4, 3))
np.testing.assert_array_equal(oneimage, reconstructed)
```

THEPATCHEXTRACTOR CLASS WORKS IN THE SAME WAY AS EXTRACTPATCHES2D ONLY IT SUPPORTS MULTIPLE IMAGES AS INPUT IT IS IMPLEMENTED AS AN ESTIMATOR SO IT CAN BE USED IN PIPELINES SEE

```
fiveimages = np.arange(5 * 4 * 3).reshape(5, 4, 3)
patches = image.patch_extractor(2, 2).transform(fiveimages)
patches.shape
(45, 2, 2, 3)
```

CONNECTIVITY GRAPH OF AN IMAGE

SEVERAL ESTIMATORS IN THE SCIKITLEARN CAN USE CONNECTIVITY INFORMATION BETWEEN FEATURES OR SAMPLES FOR INSTANCE WARD CLUSTERING HIERARCHICAL CLUSTERING CAN CLUSTER TOGETHER ONLY NEIGHBORING PIXELS OF AN IMAGE THUS FORMING CONTIGUOUS PATCHES

FOR THIS PURPOSE THE ESTIMATORS USE A 'CONNECTIVITY' MATRIX GIVING WHICH SAMPLES ARE CONNECTED

35 DATASET TRANSFORMATIONS 597

SCIKITLEARN USER GUIDE RELEASE 0213

THE FUNCTION IMGTOGRAPH RETURNS SUCH A MATRIX FROM A 2D OR 3D IMAGE SIMILARLY GRIDTOGRAPH BUILD A CONNECTIVITY MATRIX FOR IMAGES GIVEN THE SHAPE OF THESE IMAGE  
THESE MATRICES CAN BE USED TO IMPOSE CONNECTIVITY IN ESTIMATORS THAT USE CONNECTIVITY INFORMATION SUCH AS WARD CLUSTERING HIERARCHICAL CLUSTERING BUT ALSO TO BUILD PRECOMPUTED KERNELS OR SIMILARITY MATRICES  
NOTE EXAMPLES

- A DEMO OF STRUCTURED WARD HIERARCHICAL CLUSTERING ON AN IMAGE OF COINS
- SPECTRAL CLUSTERING FOR IMAGE SEGMENTATION
- FEATURE AGGLOMERATION VS UNIVARIATE SELECTION

353 PREPROCESSING DATA

THESKLEARNPREPROCESSING PACKAGE PROVIDES SEVERAL COMMON UTILITY FUNCTIONS AND TRANSFORMER CLASSES TO CHANGE RAW FEATURE VECTORS INTO A REPRESENTATION THAT IS MORE SUITABLE FOR THE DOWNSTREAM ESTIMATORS  
IN GENERAL LEARNING ALGORITHMS BENEFIT FROM STANDARDIZATION OF THE DATA SET IF SOME OUTLIERS ARE PRESENT IN THE SET ROBUST SCALERS OR TRANSFORMERS ARE MORE APPROPRIATE THE BEHAVIORS OF THE DIFFERENT SCALERS TRANSFORMERS AND NORMALIZERS ON A DATASET CONTAINING MARGINAL OUTLIERS IS HIGHLIGHTED IN COMPARE THE EFFECT OF DIFFERENT SCALERS ON DATA WITH OUTLIERS  
STANDARDIZATION OR MEAN REMOVAL AND VARIANCE SCALING  
STANDARDIZATION OF DATASETS IS A COMMON REQUIREMENT FOR MANY MACHINE LEARNING ESTIMATORS IMPLEMENTED IN SCIKITLEARN THEY MIGHT BEHAVE BADLY IF THE INDIVIDUAL FEATURES DO NOT MORE OR LESS LOOK LIKE STANDARD NORMALLY DISTRIBUTED DATA GAUSSIAN WITH ZERO MEAN AND UNIT VARIANCE  
IN PRACTICE WE OFTEN IGNORE THE SHAPE OF THE DISTRIBUTION AND JUST TRANSFORM THE DATA TO CENTER IT BY REMOVING THE MEAN VALUE OF EACH FEATURE THEN SCALE IT BY DIVIDING NONCONSTANT FEATURES BY THEIR STANDARD DEVIATION  
FOR INSTANCE MANY ELEMENTS USED IN THE OBJECTIVE FUNCTION OF A LEARNING ALGORITHM SUCH AS THE RBF KERNEL OF SUPPORT VECTOR MACHINES OR THE L1 AND L2 REGULARIZERS OF LINEAR MODELS ASSUME THAT ALL FEATURES ARE CENTERED AROUND ZERO AND HAVE VARIANCE IN THE SAME ORDER IF A FEATURE HAS A VARIANCE THAT IS ORDERS OF MAGNITUDE LARGER THAN OTHERS IT MIGHT DOMINATE THE OBJECTIVE FUNCTION AND MAKE THE ESTIMATOR UNABLE TO LEARN FROM OTHER FEATURES CORRECTLY AS EXPECTED  
THE FUNCTION SCALE PROVIDES A QUICK AND EASY WAY TO PERFORM THIS OPERATION ON A SINGLE ARRAYLIKE DATASET  
FROM SKLEARN IMPORT PREPROCESSING  
IMPORT NUMPY AS NP

598 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

XTRAIN NPARRAY 1 1 2  
2 0 0  
0 1 1

XSCALED PREPROCESSINGSCALEXTRAIN  
XSCALED  
ARRAY 0 122 133  
122 0 026  
122 122 106

SCALED DATA HAS ZERO MEAN AND UNIT VARIANCE  
XSCALEDMEANAXIS0  
ARRAY0 0 0  
XSCALEDSTDAXIS0  
ARRAY1 1 1

THEPREPROCESSING MODULE FURTHER PROVIDES A UTILITY CLASS STANDARDSCALER THAT IMPLEMENTS THE  
TRANSFORMER API TO COMPUTE THE MEAN AND STANDARD DEVIATION ON A TRAINING SET SO AS TO BE ABLE TO LATER REAPPLY  
THE SAME TRANSFORMATION ON THE TESTING SET THIS CLASS IS HENCE SUITABLE FOR USE IN THE EARLY STEPS OF A SKLEARN  
PIPELINEPIPELINE

SCALER PREPROCESSINGSTANDARDSCALERFITXTRAIN  
SCALER  
STANDARDSCALERCOPYTRUE WITHMEANTRUE WITHSTDTRUE  
SCALERMEAN  
ARRAY1 0 033  
SCALERSCALE  
ARRAY081 081 124  
SCALERTRANSFORMXTRAIN  
ARRAY 0 122 133  
122 0 026  
122 122 106

THE SCALER INSTANCE CAN THEN BE USED ON NEW DATA TO TRANSFORM IT THE SAME WAY IT DID ON THE TRAINING SET  
XTEST 1 1 0  
SCALERTRANSFORMXTEST  
ARRAY244 122 026

IT IS POSSIBLE TO DISABLE EITHER CENTERING OR SCALING BY EITHER PASSING WITHMEANFALSE ORWITHSTDFALSE TO  
THE CONSTRUCTOR OF STANDARDSCALER  
SCALING FEATURES TO A RANGE

AN ALTERNATIVE STANDARDIZATION IS SCALING FEATURES TO LIE BETWEEN A GIVEN MINIMUM AND MAXIMUM VALUE OFTEN BETWEEN  
ZERO AND ONE OR SO THAT THE MAXIMUM ABSOLUTE VALUE OF EACH FEATURE IS SCALED TO UNIT SIZE THIS CAN BE ACHIEVED USING  
MINMAXSCALER ORMAXABSSCALER RESPECTIVELY

THE MOTIVATION TO USE THIS SCALING INCLUDE ROBUSTNESS TO VERY SMALL STANDARD DEVIATIONS OF FEATURES AND PRESERVING ZERO  
ENTRIES IN SPARSE DATA

HERE IS AN EXAMPLE TO SCALE A TOY DATA MATRIX TO THE 0 1 RANGE

35 DATASET TRANSFORMATIONS 599

SCIKITLEARN USER GUIDE RELEASE 0213

XTRAIN NPARRAY 1 1 2  
2 0 0  
0 1 1

MINMAXSCALER PREPROCESSINGMINMAXSCALER  
XTRAINMINMAX MINMAXSCALERFITTRANSFORMXTRAIN  
XTRAINMINMAX  
ARRAY05 0 1  
1 05 033333333  
0 1 0

THE SAME INSTANCE OF THE TRANSFORMER CAN THEN BE APPLIED TO SOME NEW TEST DATA UNSEEN DURING THE FIT CALL THE SAME  
SCALING AND SHIFTING OPERATIONS WILL BE APPLIED TO BE CONSISTENT WITH THE TRANSFORMATION PERFORMED ON THE TRAIN DATA

XTEST NPARRAY3 1 4  
XTESTMINMAX MINMAXSCALERTRANSFORMXTEST  
XTESTMINMAX  
ARRAY15 0 166666667

IT IS POSSIBLE TO INTROSPECT THE SCALER ATTRIBUTES TO FIND ABOUT THE EXACT NATURE OF THE TRANSFORMATION LEARNED ON THE  
TRAINING DATA

MINMAXSCALERSCALE  
ARRAY05 05 033  
MINMAXSCALERMIN  
ARRAY0 05 033

IFMINMAXSCALER IS GIVEN AN EXPLICIT FEATURERANGEMIN MAX THE FULL FORMULA IS

$$XSTD \times \frac{X - XMINAXISO}{XMAXAXISO - XMINAXISO}$$
  
XSCALED XSTD MAX MIN MIN

MAXABSSCALER WORKS IN A VERY SIMILAR FASHION BUT SCALES IN A WAY THAT THE TRAINING DATA LIES WITHIN THE RANGE 1  
1BY DIVIDING THROUGH THE LARGEST MAXIMUM VALUE IN EACH FEATURE IT IS MEANT FOR DATA THAT IS ALREADY CENTERED AT ZERO  
OR SPARSE DATA

HERE IS HOW TO USE THE TOY DATA FROM THE PREVIOUS EXAMPLE WITH THIS SCALER

XTRAIN NPARRAY 1 1 2  
2 0 0  
0 1 1

MAXABSSCALER PREPROCESSINGMAXABSSCALER  
XTRAINMAXABS MAXABSSCALERFITTRANSFORMXTRAIN  
XTRAINMAXABS DOCTEST NORMALIZEWHITESPACE  
ARRAY 05 1 1  
1 0 0  
0 1 05

XTEST NPARRAY 3 1 4  
XTESTMAXABS MAXABSSCALERTRANSFORMXTEST  
XTESTMAXABS  
ARRAY15 1 2  
MAXABSSCALERSCALE  
ARRAY2 1 2

SCIKITLEARN USER GUIDE RELEASE 0213

AS WITHSCALE THE MODULE FURTHER PROVIDES CONVENIENCE FUNCTIONS MINMAXSCALE ANDMAXABSSCALE IF YOU DON'T WANT TO CREATE AN OBJECT

SCALING SPARSE DATA

CENTERING SPARSE DATA WOULD DESTROY THE SPARSENESS STRUCTURE IN THE DATA AND THUS RARELY IS A SENSIBLE THING TO DO HOWEVER IT CAN MAKE SENSE TO SCALE SPARSE INPUTS ESPECIALLY IF FEATURES ARE ON DIFFERENT SCALES

MAXABSSCALER ANDMAXABSSCALE WERE SPECIFICALLY DESIGNED FOR SCALING SPARSE DATA AND ARE THE RECOMMENDED WAY TO GO ABOUT THIS HOWEVER SCALE ANDSTANDARDSCALER CAN ACCEPT SCIPYSPARSE MATRICES AS INPUT AS LONG ASWITHMEANFALSE IS EXPLICITLY PASSED TO THE CONSTRUCTOR OTHERWISE A VALUEERROR WILL BE RAISED AS SILENTLY CENTERING WOULD BREAK THE SPARSITY AND WOULD OFTEN CRASH THE EXECUTION BY ALLOCATING EXCESSIVE AMOUNTS OF MEMORY UNINTENTIONALLY ROBUSTSCALER CANNOT BE FITTED TO SPARSE INPUTS BUT YOU CAN USE THE TRANSFORM METHOD ON SPARSE INPUTS

NOTE THAT THE SCALERS ACCEPT BOTH COMPRESSED SPARSE ROWS AND COMPRESSED SPARSE COLUMNS FORMAT SEE SCIPY SPARSECSRMATRIX ANDSCIPYSPARSESCSMATRIX ANY OTHER SPARSE INPUT WILL BE CONVERTED TO THE COMPRESSED SPARSE ROWS REPRESENTATION TO AVOID UNNECESSARY MEMORY COPIES IT IS RECOMMENDED TO CHOOSE THE CSR OR CSC REPRESENTATION UPSTREAM

FINALLY IF THE CENTERED DATA IS EXPECTED TO BE SMALL ENOUGH EXPLICITLY CONVERTING THE INPUT TO AN ARRAY USING THE TOARRAY METHOD OF SPARSE MATRICES IS ANOTHER OPTION

SCALING DATA WITH OUTLIERS

IF YOUR DATA CONTAINS MANY OUTLIERS SCALING USING THE MEAN AND VARIANCE OF THE DATA IS LIKELY TO NOT WORK VERY WELL IN THESE CASES YOU CAN USE ROBUSTSCALE ANDROBUSTSCALER AS DROPIN REPLACEMENTS INSTEAD THEY USE MORE ROBUST ESTIMATES FOR THE CENTER AND RANGE OF YOUR DATA

REFERENCES

FURTHER DISCUSSION ON THE IMPORTANCE OF CENTERING AND SCALING DATA IS AVAILABLE ON THIS FAQ SHOULD I NORMAL IZESTANDARDIZERESCALE THE DATA

SCALING VS WHITENING

IT IS SOMETIMES NOT ENOUGH TO CENTER AND SCALE THE FEATURES INDEPENDENTLY SINCE A DOWNSTREAM MODEL CAN FURTHER MAKE SOME ASSUMPTION ON THE LINEAR INDEPENDENCE OF THE FEATURES

TO ADDRESS THIS ISSUE YOU CAN USE SKLEARNDECOMPOSITIONPCA WITHWHITENTRUE TO FURTHER REMOVE THE LINEAR CORRELATION ACROSS FEATURES

SCALING A 1D ARRAY

ALL ABOVE FUNCTIONS IE SCALE MINMAXSCALE MAXABSSCALE ANDROBUSTSCALE ACCEPT 1D ARRAY WHICH CAN BE USEFUL IN SOME SPECIFIC CASE

35 DATASET TRANSFORMATIONS 601

SCIKITLEARN USER GUIDE RELEASE 0213

CENTERING KERNEL MATRICES

IF YOU HAVE A KERNEL MATRIX OF A KERNEL  $\phi$  THAT COMPUTES A DOT PRODUCT IN A FEATURE SPACE DEFINED BY FUNCTION  $\phi(h)$  A KERNELCENTERER CAN TRANSFORM THE KERNEL MATRIX SO THAT IT CONTAINS INNER PRODUCTS IN THE FEATURE SPACE DEFINED BY  $\phi(h)$  FOLLOWED BY REMOVAL OF THE MEAN IN THAT SPACE

NONLINEAR TRANSFORMATION

TWO TYPES OF TRANSFORMATIONS ARE AVAILABLE QUANTILE TRANSFORMS AND POWER TRANSFORMS BOTH QUANTILE AND POWER TRANSFORMS ARE BASED ON MONOTONIC TRANSFORMATIONS OF THE FEATURES AND THUS PRESERVE THE RANK OF THE VALUES ALONG EACH FEATURE

QUANTILE TRANSFORMS PUT ALL FEATURES INTO THE SAME DESIRED DISTRIBUTION BASED ON THE FORMULA  $\phi^{-1}(\phi(x))$  WHERE  $\phi$  IS THE CUMULATIVE DISTRIBUTION FUNCTION OF THE FEATURE AND  $\phi^{-1}$  THE QUANTILE FUNCTION OF THE DESIRED OUTPUT DISTRIBUTION  $\phi$  THIS FORMULA IS USING THE TWO FOLLOWING FACTS I IF  $\phi$  IS A RANDOM VARIABLE WITH A CONTINUOUS CUMULATIVE DISTRIBUTION FUNCTION  $\phi$  THEN  $\phi(\phi(x))$  IS UNIFORMLY DISTRIBUTED ON  $[0, 1]$  II IF  $\phi$  IS A RANDOM VARIABLE WITH UNIFORM DISTRIBUTION ON  $[0, 1]$  THEN  $\phi^{-1}(\phi(x))$  HAS DISTRIBUTION  $\phi$  BY PERFORMING A RANK TRANSFORMATION A QUANTILE TRANSFORM SMOOTHS OUT UNUSUAL DISTRIBUTIONS AND IS LESS INFLUENCED BY OUTLIERS THAN SCALING METHODS IT DOES HOWEVER DISTORT CORRELATIONS AND DISTANCES WITHIN AND ACROSS FEATURES

POWER TRANSFORMS ARE A FAMILY OF PARAMETRIC TRANSFORMATIONS THAT AIM TO MAP DATA FROM ANY DISTRIBUTION TO AS CLOSE TO A GAUSSIAN DISTRIBUTION

MAPPING TO A UNIFORM DISTRIBUTION

QUANTILETRANSFORMER ANDQUANTILETRANSFORM PROVIDE A NONPARAMETRIC TRANSFORMATION TO MAP THE DATA TO A UNIFORM DISTRIBUTION WITH VALUES BETWEEN 0 AND 1

```
FROM SKLEARNDATASETS IMPORT LOADIRIS
FROM SKLEARNMODELSELECTION IMPORT TRAINTESTSPLIT
IRIS = LOADIRIS
X, Y = IRIS.DATA, IRIS.TARGET
XTRAIN, XTEST, YTRAIN, YTEST = TRAINTESTSPLIT(X, Y, RANDOMSTATE=0)
QUANTILETRANSFORMER = PREPROCESSINGQUANTILETRANSFORMER(RANDOMSTATE=0)
XTRAINTRANS = QUANTILETRANSFORMER.FIT_TRANSFORM(XTRAIN)
XTESTTRANS = QUANTILETRANSFORMER.TRANSFORM(XTEST)
NPPERCENTILEXTRAIN = [0, 0, 25, 50, 75, 100]
ARRAY [43, 51, 58, 65, 79]
```

THIS FEATURE CORRESPONDS TO THE SEPAL LENGTH IN CM ONCE THE QUANTILE TRANSFORMATION APPLIED THOSE LANDMARKS APPROACH CLOSELY THE PERCENTILES PREVIOUSLY DEFINED

```
NPPERCENTILEXTRAINTRANS = [0, 0, 25, 50, 75, 100]
```

```
ARRAY [0.00, 0.024, 0.049, 0.073, 0.099]
THIS CAN BE CONFIRMED ON A INDEPENDENT TESTING SET WITH SIMILAR REMARKS
NPPERCENTILEXTEST = [0, 0, 25, 50, 75, 100]
```

```
ARRAY [44, 51.25, 57.5, 61.75, 73]
NPPERCENTILEXTESTTRANS = [0, 0, 25, 50, 75, 100]
```

```
ARRAY [0.01, 0.025, 0.046, 0.060, 0.094]
602 CHAPTER 3 USER GUIDE
```



SCIKITLEARN USER GUIDE RELEASE 0213

MAPPING TO A GAUSSIAN DISTRIBUTION

IN MANY MODELING SCENARIOS NORMALITY OF THE FEATURES IN A DATASET IS DESIRABLE POWER TRANSFORMS ARE A FAMILY OF PARAMETRIC MONOTONIC TRANSFORMATIONS THAT AIM TO MAP DATA FROM ANY DISTRIBUTION TO AS CLOSE TO A GAUSSIAN DISTRIBUTION AS POSSIBLE IN ORDER TO STABILIZE VARIANCE AND MINIMIZE SKEWNESS

POWERTRANSFORMER CURRENTLY PROVIDES TWO SUCH POWER TRANSFORMATIONS THE YEOJOHNSON TRANSFORM AND THE BOX COX TRANSFORM

THE YEOJOHNSON TRANSFORM IS GIVEN BY

$$\lambda = \begin{cases} \frac{1}{\sqrt{2\pi}} \left( \frac{1}{\sigma} - \frac{1}{\sigma^2} \right) & \text{if } \sigma \geq 0 \\ \frac{1}{\sqrt{2\pi}} \left( \frac{1}{\sigma} - \frac{1}{\sigma^2} \right) & \text{if } \sigma < 0 \end{cases}$$
$$\lambda = \begin{cases} \frac{1}{\sqrt{2\pi}} \left( \frac{1}{\sigma} - \frac{1}{\sigma^2} \right) & \text{if } \sigma \geq 0 \\ \frac{1}{\sqrt{2\pi}} \left( \frac{1}{\sigma} - \frac{1}{\sigma^2} \right) & \text{if } \sigma < 0 \end{cases}$$

WHILE THE BOXCOX TRANSFORM IS GIVEN BY

$$\lambda = \begin{cases} \frac{1}{\sqrt{2\pi}} \left( \frac{1}{\sigma} - \frac{1}{\sigma^2} \right) & \text{if } \sigma \geq 0 \\ \frac{1}{\sqrt{2\pi}} \left( \frac{1}{\sigma} - \frac{1}{\sigma^2} \right) & \text{if } \sigma < 0 \end{cases}$$
$$\lambda = \begin{cases} \frac{1}{\sqrt{2\pi}} \left( \frac{1}{\sigma} - \frac{1}{\sigma^2} \right) & \text{if } \sigma \geq 0 \\ \frac{1}{\sqrt{2\pi}} \left( \frac{1}{\sigma} - \frac{1}{\sigma^2} \right) & \text{if } \sigma < 0 \end{cases}$$

BOXCOX CAN ONLY BE APPLIED TO STRICTLY POSITIVE DATA IN BOTH METHODS THE TRANSFORMATION IS PARAMETERIZED BY  $\lambda$  WHICH IS DETERMINED THROUGH MAXIMUM LIKELIHOOD ESTIMATION HERE IS AN EXAMPLE OF USING BOXCOX TO MAP SAMPLES DRAWN FROM A LOGNORMAL DISTRIBUTION TO A NORMAL DISTRIBUTION

PT PREPROCESSINGPOWERTRANSFORMERMETHODBOXCOX STANDARDIZE FALSE

XLOGNORMAL NPRANDOMRANDOMSTATE616LOGNORMALSIZE3 3

XLOGNORMAL

ARRAY128 118 084

094 160 038

135 021 109

PTFITTRANSFORMXLOGNORMAL

ARRAY 049 017 015

005 058 057

069 084 010

WHILE THE ABOVE EXAMPLE SETS THE STANDARDIZE OPTION TOFALSE POWERTRANSFORMER WILL APPLY ZERO MEAN

UNIT VARIANCE NORMALIZATION TO THE TRANSFORMED OUTPUT BY DEFAULT

BELOW ARE EXAMPLES OF BOXCOX AND YEOJOHNSON APPLIED TO VARIOUS PROBABILITY DISTRIBUTIONS NOTE THAT WHEN APPLIED TO CERTAIN DISTRIBUTIONS THE POWER TRANSFORMS ACHIEVE VERY GAUSSIANLIKE RESULTS BUT WITH OTHERS THEY ARE INEFFECTIVE THIS HIGHLIGHTS THE IMPORTANCE OF VISUALIZING THE DATA BEFORE AND AFTER TRANSFORMATION

IT IS ALSO POSSIBLE TO MAP DATA TO A NORMAL DISTRIBUTION USING QUANTILETRANSFORMER BY SETTING

OUTPUTDISTRIBUTIONNORMAL USING THE EARLIER EXAMPLE WITH THE IRIS DATASET

QUANTILETRANSFORMER PREPROCESSINGQUANTILETRANSFORMER

OUTPUTDISTRIBUTIONNORMAL RANDOMSTATE0

XTRANS QUANTILETRANSFORMERFITTRANSFORMX

QUANTILETRANSFORMERQUANTILES

ARRAY43 2 1 01

44 22 11 01

44 22 12 01



77 41 67 25  
77 42 67 25  
79 44 69 25

THUS THE MEDIAN OF THE INPUT BECOMES THE MEAN OF THE OUTPUT CENTERED AT 0 THE NORMAL OUTPUT IS CLIPPED SO THAT THE INPUT’S MINIMUM AND MAXIMUM — CORRESPONDING TO THE 1E7 AND 1 1E7 QUANTILES RESPECTIVELY — DO NOT BECOME INFINITE UNDER THE TRANSFORMATION

NORMALIZATION

NORMALIZATION IS THE PROCESS OF SCALING INDIVIDUAL SAMPLES TO HAVE UNIT NORM THIS PROCESS CAN BE USEFUL IF YOU PLAN TO USE A QUADRATIC FORM SUCH AS THE DOTPRODUCT OR ANY OTHER KERNEL TO QUANTIFY THE SIMILARITY OF ANY PAIR OF SAMPLES THIS ASSUMPTION IS THE BASE OF THE VECTOR SPACE MODEL OFTEN USED IN TEXT CLASSIFICATION AND CLUSTERING CONTEXTS THE FUNCTION NORMALIZE PROVIDES A QUICK AND EASY WAY TO PERFORM THIS OPERATION ON A SINGLE ARRAYLIKE DATASET EITHER USING THE L1ORL2NORMS

X 1 1 2  
2 0 0  
0 1 1  
XNORMALIZED PREPROCESSINGNORMALIZEX NORML2  
XNORMALIZED  
ARRAY 040 040 081  
1 0 0  
0 070 070

THEPREPROCESSING MODULE FURTHER PROVIDES A UTILITY CLASS NORMALIZER THAT IMPLEMENTS THE SAME OPERATION USING THETRANSFORMER API EVEN THOUGH THE FIT METHOD IS USELESS IN THIS CASE THE CLASS IS STATELESS AS THIS OPERATION TREATS SAMPLES INDEPENDENTLY THIS CLASS IS HENCE SUITABLE FOR USE IN THE EARLY STEPS OF A SKLEARNPIPELINEPIPELINE NORMALIZER PREPROCESSINGNORMALIZERFITX FIT DOES NOTHING

NORMALIZER  
NORMALIZERCOPYTRUE NORML2  
THE NORMALIZER INSTANCE CAN THEN BE USED ON SAMPLE VECTORS AS ANY TRANSFORMER  
NORMALIZERTRANSFORMX  
ARRAY 040 040 081  
1 0 0  
0 070 070  
NORMALIZERTRANSFORM1 1 0  
ARRAY070 070 0  
SPARSE INPUT  
NORMALIZE ANDNORMALIZER ACCEPT BOTH DENSE ARRAYLIKE AND SPARSE MATRICES FROM SCIPYSPARSE AS INPUT FOR SPARSE INPUT THE DATA IS CONVERTED TO THE COMPRESSED SPARSE ROWS REPRESENTATION SEESCIPIYSPARSE CSRMATRIX BEFORE BEING FED TO EFFICIENT CYTHON ROUTINES TO AVOID UNNECESSARY MEMORY COPIES IT IS RECOMMENDED TO CHOOSE THE CSR REPRESENTATION UPSTREAM  
35 DATASET TRANSFORMATIONS 605

SCIKITLEARN USER GUIDE RELEASE 0213

ENCODING CATEGORICAL FEATURES

OFTEN FEATURES ARE NOT GIVEN AS CONTINUOUS VALUES BUT CATEGORICAL FOR EXAMPLE A PERSON COULD HAVE FEATURES MALE FEMALE FROM EUROPE FROM US FROM ASIA USES FIREFOX USES CHROME USES SAFARI USES INTERNET EXPLORER SUCH FEATURES CAN BE EFFICIENTLY CODED AS INTEGERS FOR INSTANCE MALE FROM US USES INTERNET EXPLORER COULD BE EXPRESSED AS 0 1 3 WHILE FEMALE FROM ASIA USES CHROME WOULD BE 1 2 1

TO CONVERT CATEGORICAL FEATURES TO SUCH INTEGER CODES WE CAN USE THE ORDINAL ENCODER THIS ESTIMATOR TRANSFORMS EACH CATEGORICAL FEATURE TO ONE NEW FEATURE OF INTEGERS 0 TO NCATEGORIES - 1

ENC PREPROCESSING ORDINAL ENCODER

X MALE FROM US USES SAFARI FEMALE FROM EUROPE USES FIREFOX

↳

ENC FITX

ORDINAL ENCODER CATEGORIES AUTO DTYPE NUMPY FLOAT64

ENC TRANSFORM FEMALE FROM US USES SAFARI

ARRAY 0 1 1

SUCH INTEGER REPRESENTATION CAN HOWEVER NOT BE USED DIRECTLY WITH ALL SCIKITLEARN ESTIMATORS AS THESE EXPECT CONTINUOUS INPUT AND WOULD INTERPRET THE CATEGORIES AS BEING ORDERED WHICH IS OFTEN NOT DESIRED IE THE SET OF BROWSERS WAS ORDERED ARBITRARILY

ANOTHER POSSIBILITY TO CONVERT CATEGORICAL FEATURES TO FEATURES THAT CAN BE USED WITH SCIKITLEARN ESTIMATORS IS TO USE A ONE OF K ALSO KNOWN AS ONE HOT OR DUMMY ENCODING THIS TYPE OF ENCODING CAN BE OBTAINED WITH THE ONE HOT ENCODER WHICH TRANSFORMS EACH CATEGORICAL FEATURE WITH NCATEGORIES POSSIBLE VALUES INTO NCATEGORIES BINARY FEATURES WITH ONE OF THEM 1 AND ALL OTHERS 0

CONTINUING THE EXAMPLE ABOVE

ENC PREPROCESSING ONE HOT ENCODER

X MALE FROM US USES SAFARI FEMALE FROM EUROPE USES FIREFOX

↳

ENC FITX

ONE HOT ENCODER CATEGORICAL FEATURES NONE CATEGORIES NONE DROP NONE

DTYPE NUMPY FLOAT64 HANDLE UNKNOWN ERROR

NVALUES NONE SPARSE TRUE

ENC TRANSFORM FEMALE FROM US USES SAFARI

MALE FROM EUROPE USES SAFARI TO ARRAY

ARRAY 1 0 0 1 0 1

0 1 1 0 0 1

BY DEFAULT THE VALUES EACH FEATURE CAN TAKE IS INFERRED AUTOMATICALLY FROM THE DATASET AND CAN BE FOUND IN THE CATEGORIES ATTRIBUTE

ENC CATEGORIES

ARRAY FEMALE MALE DTYPE OBJECT ARRAY FROM EUROPE FROM US

↳ DTYPE OBJECT ARRAY USES FIREFOX USES SAFARI DTYPE OBJECT

IT IS POSSIBLE TO SPECIFY THIS EXPLICITLY USING THE PARAMETER CATEGORIES THERE ARE TWO GENDERS FOUR POSSIBLE CONTINENTS AND FOUR WEB BROWSERS IN OUR DATASET

GENDERS FEMALE MALE

LOCATIONS FROM AFRICA FROM ASIA FROM EUROPE FROM US

BROWSERS USES CHROME USES FIREFOX USES IE USES SAFARI

ENC PREPROCESSING ONE HOT ENCODER CATEGORIES GENDERS LOCATIONS BROWSERS

NOTE THAT FOR THERE ARE MISSING CATEGORICAL VALUES FOR THE 2ND AND 3RD FEATURE

606 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

X MALE FROM US USES SAFARI FEMALE FROM EUROPE USES FIREFOX

↔

ENC FITX

ONEHOT ENCODER CATEGORICAL FEATURES NONE

CATEGORIES DROP NONE

DTYPES NUMPY FLOAT64 HANDLE UNKNOWN ERROR

NVALUES NONE SPARSE TRUE

ENC TRANSFORM FEMALE FROM ASIA USES CHROME TO ARRAY

ARRAY 1 0 0 1 0 0 1 0 0 0

IF THERE IS A POSSIBILITY THAT THE TRAINING DATA MIGHT HAVE MISSING CATEGORICAL FEATURES IT CAN OFTEN BE BETTER TO SPECIFY HANDLE UNKNOWN IGNORE INSTEAD OF SETTING THE CATEGORIES MANUALLY AS ABOVE WHEN HANDLE UNKNOWN IGNORE IS SPECIFIED AND UNKNOWN CATEGORIES ARE ENCOUNTERED DURING TRANSFORM NO ERROR WILL BE RAISED BUT THE RESULTING ONEHOT ENCODED COLUMNS FOR THIS FEATURE WILL BE ALL ZEROS HANDLE UNKNOWN IGNORE IS ONLY SUPPORTED FOR ONEHOT ENCODING

ENC PREPROCESSING ONEHOT ENCODER HANDLE UNKNOWN IGNORE

X MALE FROM US USES SAFARI FEMALE FROM EUROPE USES FIREFOX

↔

ENC FITX

ONEHOT ENCODER CATEGORICAL FEATURES NONE CATEGORIES NONE DROP NONE

DTYPES NUMPY FLOAT64 HANDLE UNKNOWN IGNORE

NVALUES NONE SPARSE TRUE

ENC TRANSFORM FEMALE FROM ASIA USES CHROME TO ARRAY

ARRAY 1 0 0 0 0 0

IT IS ALSO POSSIBLE TO ENCODE EACH COLUMN INTO NCATEGORIES 1 COLUMNS INSTEAD OF NCATEGORIES COLUMNS BY USING THE DROP PARAMETER THIS PARAMETER ALLOWS THE USER TO SPECIFY A CATEGORY FOR EACH FEATURE TO BE DROPPED THIS IS USEFUL TO AVOID COLINEARITY IN THE INPUT MATRIX IN SOME CLASSIFIERS SUCH FUNCTIONALITY IS USEFUL FOR EXAMPLE WHEN USING NONREGULARIZED REGRESSION LINEAR REGRESSION SINCE COLINEARITY WOULD CAUSE THE COVARIANCE MATRIX TO BE NONINVERTIBLE WHEN THIS PARAMETER IS NOT NONE HANDLE UNKNOWN MUST BE SET TO ERROR

X MALE FROM US USES SAFARI FEMALE FROM EUROPE USES FIREFOX

↔

DROP ENC PREPROCESSING ONEHOT ENCODER DROP FIRST FITX

DROP ENC CATEGORIES

ARRAY FEMALE MALE DTYPES OBJECT ARRAY FROM EUROPE FROM US

↔ DTYPES OBJECT ARRAY USES FIREFOX USES SAFARI DTYPES OBJECT

DROP ENC TRANSFORM XT TO ARRAY

ARRAY 1 1 1

0 0 0

SEE LOADING FEATURES FROM DICTS FOR CATEGORICAL FEATURES THAT ARE REPRESENTED AS A DICT NOT AS SCALARS

DISCRETIZATION

DISCRETIZATION OTHERWISE KNOWN AS QUANTIZATION OR BINNING PROVIDES A WAY TO PARTITION CONTINUOUS FEATURES INTO DISCRETE VALUES CERTAIN DATASETS WITH CONTINUOUS FEATURES MAY BENEFIT FROM DISCRETIZATION BECAUSE DISCRETIZATION CAN TRANSFORM THE DATASET OF CONTINUOUS ATTRIBUTES TO ONE WITH ONLY NOMINAL ATTRIBUTES

ONEHOT ENCODED DISCRETIZED FEATURES CAN MAKE A MODEL MORE EXPRESSIVE WHILE MAINTAINING INTERPRETABILITY FOR INSTANCE PREPROCESSING WITH A DISCRETIZER CAN INTRODUCE NONLINEARITY TO LINEAR MODELS

35 DATASET TRANSFORMATIONS 607

KBINS DISCRETIZATION

KBINSDISCRETIZER DISCRETIZES FEATURES INTO KBINS

X NPARRAY 3 5 15

0 6 14

6 3 11

EST PREPROCESSINGKBINSDISCRETIZERNBINS3 2 2 ENCODEORDINALFITX

BY DEFAULT THE OUTPUT IS ONEHOT ENCODED INTO A SPARSE MATRIX SEE ENCODING CATEGORICAL FEATURES AND THIS CAN BE CONFIGURED WITH THE ENCODE PARAMETER FOR EACH FEATURE THE BIN EDGES ARE COMPUTED DURING FIT AND TOGETHER WITH THE NUMBER OF BINS THEY WILL DEFINE THE INTERVALS THEREFORE FOR THE CURRENT EXAMPLE THESE INTERVALS ARE DEFINED AS

- FEATURE 1  $-\infty-1-122\infty$
- FEATURE 2  $-\infty55\infty$
- FEATURE 3  $-\infty1414\infty$

BASED ON THESE BIN INTERVALS XIS TRANSFORMED AS FOLLOWS

ESTTRANSFORMX

ARRAY 0 1 1

1 1 1

2 0 0

THE RESULTING DATASET CONTAINS ORDINAL ATTRIBUTES WHICH CAN BE FURTHER USED IN A SKLEARNPIPELINEPIPELINE DISCRETIZATION IS SIMILAR TO CONSTRUCTING HISTOGRAMS FOR CONTINUOUS DATA HOWEVER HISTOGRAMS FOCUS ON COUNTING FEATURES WHICH FALL INTO PARTICULAR BINS WHEREAS DISCRETIZATION FOCUSES ON ASSIGNING FEATURE VALUES TO THESE BINS KBINSDISCRETIZER IMPLEMENTS DIFFERENT BINNING STRATEGIES WHICH CAN BE SELECTED WITH THE STRATEGY PARAMETER THE 'UNIFORM' STRATEGY USES CONSTANTWIDTH BINS THE 'QUANTILE' STRATEGY USES THE QUANTILES VALUES TO HAVE EQUALLY POPULATED BINS IN EACH FEATURE THE 'KMEANS' STRATEGY DEFINES BINS BASED ON A KMEANS CLUSTERING PROCEDURE PERFORMED ON EACH FEATURE INDEPENDENTLY

EXAMPLES

- USING KBINSDISCRETIZER TO DISCRETIZE CONTINUOUS FEATURES
- FEATURE DISCRETIZATION
- DEMONSTRATING THE DIFFERENT STRATEGIES OF KBINSDISCRETIZER

FEATURE BINARIZATION

FEATURE BINARIZATION IS THE PROCESS OF THRESHOLDING NUMERICAL FEATURES TO GET BOOLEAN VALUES THIS CAN BE USEFUL FOR DOWNSTREAM PROBABILISTIC ESTIMATORS THAT MAKE ASSUMPTION THAT THE INPUT DATA IS DISTRIBUTED ACCORDING TO A MULTIVARIATE BERNOULLI DISTRIBUTION FOR INSTANCE THIS IS THE CASE FOR THE SKLEARNNEURALNETWORKBERNOULLIRBM

IT IS ALSO COMMON AMONG THE TEXT PROCESSING COMMUNITY TO USE BINARY FEATURE VALUES PROBABLY TO SIMPLIFY THE PROBABILISTIC REASONING EVEN IF NORMALIZED COUNTS AKA TERM FREQUENCIES OR TFIDF VALUED FEATURES OFTEN PERFORM SLIGHTLY BETTER IN PRACTICE

AS FOR THE NORMALIZER THE UTILITY CLASS BINARIZER IS MEANT TO BE USED IN THE EARLY STAGES OF SKLEARN PIPELINEPIPELINE THEFIT METHOD DOES NOTHING AS EACH SAMPLE IS TREATED INDEPENDENTLY OF OTHERS

SCIKITLEARN USER GUIDE RELEASE 0213

```
X 1 1 2
  2 0 0
  0 1 1
BINARIZER PREPROCESSINGBINARIZERFITX FIT DOES NOTHING
BINARIZER
BINARIZERCOPYTRUE THRESHOLD00
BINARIZERTRANSFORMX
ARRAY1 0 1
1 0 0
0 1 0
```

IT IS POSSIBLE TO ADJUST THE THRESHOLD OF THE BINARIZER

```
BINARIZER PREPROCESSINGBINARIZERTHRESHOLD11
BINARIZERTRANSFORMX
ARRAY0 0 1
1 0 0
0 0 0
```

AS FOR THESTANDARDSCALER ANDNORMALIZER CLASSES THE PREPROCESSING MODULE PROVIDES A COMPANION FUNCTION BINARIZE TO BE USED WHEN THE TRANSFORMER API IS NOT NECESSARY

NOTE THAT THE BINARIZER IS SIMILAR TO THE KBINSDISCRETIZER WHENK 2 AND WHEN THE BIN EDGE IS AT THE VALUETHRESHOLD  
SPARSE INPUT

BINARIZE ANDBINARIZER ACCEPT BOTH DENSE ARRAYLIKE AND SPARSE MATRICES FROM SCIPYSPARSE AS INPUT  
FOR SPARSE INPUT THE DATA IS CONVERTED TO THE COMPRESSED SPARSE ROWS REPRESENTATION SEESCIPIYSPARSE CSRMATRIX TO AVOID UNNECESSARY MEMORY COPIES IT IS RECOMMENDED TO CHOOSE THE CSR REPRESENTATION UP  
STREAM

IMPUTATION OF MISSING VALUES  
TOOLS FOR IMPUTING MISSING VALUES ARE DISCUSSED AT IMPUTATION OF MISSING VALUES  
GENERATING POLYNOMIAL FEATURES

OFTEN IT'S USEFUL TO ADD COMPLEXITY TO THE MODEL BY CONSIDERING NONLINEAR FEATURES OF THE INPUT DATA A SIMPLE AND COM  
MON METHOD TO USE IS POLYNOMIAL FEATURES WHICH CAN GET FEATURES' HIGHORDER AND INTERACTION TERMS IT IS IMPLEMENTED  
INPOLYNOMIALFEATURES

```
IMPORT NUMPY AS NP
FROM SKLEARNPREPROCESSING IMPORT POLYNOMIALFEATURES
X NPARANGE6RESHAPE3 2
X
```

```
ARRAY0 1
2 3
4 5
POLY POLYNOMIALFEATURES2
POLYFITTRANSFORMX
35 DATASET TRANSFORMATIONS 609
```

SCIKITLEARN USER GUIDE RELEASE 0213

```
ARRAY 1 0 1 0 0 1
1 2 3 4 6 9
1 4 5 16 20 25
THE FEATURES OF X HAVE BEEN TRANSFORMED FROM 1 2 TO 1 2 2
1 1 2 2
2
```

IN SOME CASES ONLY INTERACTION TERMS AMONG FEATURES ARE REQUIRED AND IT CAN BE GOTTEN WITH THE SETTING  
INTERACTIONONLYTRUE

```
X NPARANGE9 RESHAPE3 3
X
ARRAY0 1 2
3 4 5
6 7 8
POLY POLYNOMIALFEATURESDEGREE3 INTERACTIONONLY TRUE
POLYFITTRANSFORMX
ARRAY 1 0 1 2 0 0 2 0
1 3 4 5 12 15 20 60
1 6 7 8 42 48 56 336
```

THE FEATURES OF X HAVE BEEN TRANSFORMED FROM 1 2 3 TO 1 2 3 1 2 3 1 3 2 3 1 2 3  
NOTE THAT POLYNOMIAL FEATURES ARE USED IMPLICITLY IN KERNEL METHODS EG SKLEARN SVM SVC SKLEARN  
DECOMPOSITION KERNEL PCA WHEN USING POLYNOMIAL KERNEL FUNCTIONS  
SEE POLYNOMIAL INTERPOLATION FOR RIDGE REGRESSION USING CREATED POLYNOMIAL FEATURES  
CUSTOM TRANSFORMERS

OFTEN YOU WILL WANT TO CONVERT AN EXISTING PYTHON FUNCTION INTO A TRANSFORMER TO ASSIST IN DATA CLEANING OR PROCESSING  
YOU CAN IMPLEMENT A TRANSFORMER FROM AN ARBITRARY FUNCTION WITH FUNCTIONTRANSFORMER FOR EXAMPLE TO BUILD  
A TRANSFORMER THAT APPLIES A LOG TRANSFORMATION IN A PIPELINE DO

```
IMPORT NUMPY AS NP
FROM SKLEARN PREPROCESSING IMPORT FUNCTIONTRANSFORMER
TRANSFORMER FUNCTIONTRANSFORMER NP LOG1P VALIDATE TRUE
X NPARRAY0 1 2 3
TRANSFORMER TRANSFORMX
ARRAY0 069314718
109861229 138629436
```

YOU CAN ENSURE THAT FUNC AND INVERSE FUNC ARE THE INVERSE OF EACH OTHER BY SETTING CHECKINVERSE TRUE  
AND CALLING FIT BEFORE TRANSFORM PLEASE NOTE THAT A WARNING IS RAISED AND CAN BE TURNED INTO AN ERROR WITH A  
FILTER WARNINGS

```
IMPORT WARNINGS
WARNINGS FILTER WARNINGS ERROR MESSAGE CHECKINVERSE
CATEGORY USER WARNING APPEND FALSE
```

FOR A FULL CODE EXAMPLE THAT DEMONSTRATES USING A FUNCTIONTRANSFORMER TO DO CUSTOM FEATURE SELECTION SEE  
USING FUNCTIONTRANSFORMER TO SELECT COLUMNS

354 IMPUTATION OF MISSING VALUES

FOR VARIOUS REASONS MANY REAL WORLD DATASETS CONTAIN MISSING VALUES OFTEN ENCODED AS BLANKS NANS OR OTHER PLACE  
HOLDERS SUCH DATASETS HOWEVER ARE INCOMPATIBLE WITH SCIKITLEARN ESTIMATORS WHICH ASSUME THAT ALL VALUES IN AN ARRAY  
610 CHAPTER 3 USER GUIDE



SCIKITLEARN USER GUIDE RELEASE 0213

ARE NUMERICAL AND THAT ALL HAVE AND HOLD MEANING A BASIC STRATEGY TO USE INCOMPLETE DATASETS IS TO DISCARD ENTIRE ROWS ANDOR COLUMNS CONTAINING MISSING VALUES HOWEVER THIS COMES AT THE PRICE OF LOSING DATA WHICH MAY BE VALUABLE EVEN THOUGH INCOMPLETE A BETTER STRATEGY IS TO IMPUTE THE MISSING VALUES IE TO INFER THEM FROM THE KNOWN PART OF THE DATA SEE THE GLOSSARY OF COMMON TERMS AND API ELEMENTS ENTRY ON IMPUTATION

UNIVARIATE VS MULTIVARIATE IMPUTATION

ONE TYPE OF IMPUTATION ALGORITHM IS UNIVARIATE WHICH IMPUTES VALUES IN THE ITH FEATURE DIMENSION USING ONLY NON MISSING VALUES IN THAT FEATURE DIMENSION EG IMPUTESIMPLEIMPUTER BY CONTRAST MULTIVARIATE IMPUTATION ALGORITHMS USE THE ENTIRE SET OF AVAILABLE FEATURE DIMENSIONS TO ESTIMATE THE MISSING VALUES EG IMPUTE ITERATIVEIMPUTER

UNIVARIATE FEATURE IMPUTATION

THESIMPLEIMPUTER CLASS PROVIDES BASIC STRATEGIES FOR IMPUTING MISSING VALUES MISSING VALUES CAN BE IMPUTED WITH A PROVIDED CONSTANT VALUE OR USING THE STATISTICS MEAN MEDIAN OR MOST FREQUENT OF EACH COLUMN IN WHICH THE MISSING VALUES ARE LOCATED THIS CLASS ALSO ALLOWS FOR DIFFERENT MISSING VALUES ENCODINGS

THE FOLLOWING SNIPPET DEMONSTRATES HOW TO REPLACE MISSING VALUES ENCODED AS NPNAN USING THE MEAN VALUE OF THE COLUMNS AXIS 0 THAT CONTAIN THE MISSING VALUES

```
import numpy as np
from sklearn.impute import SimpleImputer
imp = SimpleImputer(missing_values=np.nan, strategy='mean')
imp.fit(1 2 npnan 3 7 6
simpleimputer.add_indicator(False, copy=True, fill_value=None)
missing_values=np.nan, strategy='mean', verbose=0
X npnan 2 6 npnan 7 6
print(imp.transform(X))
4 2
6 3666
7 6
```

THESIMPLEIMPUTER CLASS ALSO SUPPORTS SPARSE MATRICES

```
import scipy.sparse as sp
X = sp.cscmatrix(1 2 0 1 8 4
imp = SimpleImputer(missing_values=1, strategy='mean')
imp.fit(X)
simpleimputer.add_indicator(False, copy=True, fill_value=None)
missing_values=1, strategy='mean', verbose=0
X_test = sp.cscmatrix(1 2 6 1 7 6
print(imp.transform(X_test).toarray())
3 2
6 3
7 6
```

NOTE THAT THIS FORMAT IS NOT MEANT TO BE USED TO IMPLICITLY STORE MISSING VALUES IN THE MATRIX BECAUSE IT WOULD DENSIFY IT AT TRANSFORM TIME MISSING VALUES ENCODED BY 0 MUST BE USED WITH DENSE INPUT

THESIMPLEIMPUTER CLASS ALSO SUPPORTS CATEGORICAL DATA REPRESENTED AS STRING VALUES OR PANDAS CATEGORICALS WHEN USING THEMOSTFREQUENT ORCONSTANT STRATEGY

```
import pandas as pd
df = pd.DataFrame(X
npnan = Y
35 dataset transformations 611
```

SCIKITLEARN USER GUIDE RELEASE 0213

A NPNAN  
B Y DTYPECATEGORY

IMP SIMPLEIMPUTERSTRATEGYMOSTFREQUENT  
PRINTIMPFITTRANSFORMDF

A X  
A Y  
A Y  
B Y

MULTIVARIATE FEATURE IMPUTATION

A MORE SOPHISTICATED APPROACH IS TO USE THE ITERATIVEIMPUTER CLASS WHICH MODELS EACH FEATURE WITH MISSING VALUES AS A FUNCTION OF OTHER FEATURES AND USES THAT ESTIMATE FOR IMPUTATION IT DOES SO IN AN ITERATED ROUNDROBIN FASHION AT EACH STEP A FEATURE COLUMN IS DESIGNATED AS OUTPUT YAND THE OTHER FEATURE COLUMNS ARE TREATED AS INPUTS X A REGRESSOR IS FIT ON X Y FOR KNOWN Y THEN THE REGRESSOR IS USED TO PREDICT THE MISSING VALUES OF Y THIS IS DONE FOR EACH FEATURE IN AN ITERATIVE FASHION AND THEN IS REPEATED FOR MAXITER IMPUTATION ROUNDS THE RESULTS OF THE FINAL IMPUTATION ROUND ARE RETURNED

NOTE THIS ESTIMATOR IS STILL EXPERIMENTAL FOR NOW THE PREDICTIONS AND THE API MIGHT CHANGE WITHOUT ANY DEPRECATION CYCLE TO USE IT YOU NEED TO EXPLICITLY IMPORT ENABLEITERATIVEIMPUTER

```
import numpy as np
from sklearn.experimental import enable_iterative_imputer
from sklearn.impute import IterativeImputer
```

```
imp = IterativeImputer(max_iter=10, random_state=0,
                        add_indicator=False, estimator=None,
                        imputation_order='ascending', initial_strategy='mean',
                        max_iter=10, max_value=None, min_value=None,
                        missing_values=np.nan, nearest_features=None,
                        random_state=0, sample_posterior=False, tol=0.001,
                        verbose=0)
```

```
x_test = np.nan * 2 * 6 * np.nan * np.nan * 6
# the model learns that the second feature is double the first
print(np.round(imp.transform(x_test)
```

```
1 2
6 12
3 6
```

BOTHSIMPLEIMPUTER ANDITERATIVEIMPUTER CAN BE USED IN A PIPELINE AS A WAY TO BUILD A COMPOSITE ESTIMATOR THAT SUPPORTS IMPUTATION SEE IMPUTING MISSING VALUES BEFORE BUILDING AN ESTIMATOR FLEXIBILITY OF ITERATIVEIMPUTER

THERE ARE MANY WELLESTABLISHED IMPUTATION PACKAGES IN THE R DATA SCIENCE ECOSYSTEM AMELIA MI MICE MISS FOREST ETC MISSFOREST IS POPULAR AND TURNS OUT TO BE A PARTICULAR INSTANCE OF DIFFERENT SEQUENTIAL IMPUTATION ALGORITHMS THAT CAN ALL BE IMPLEMENTED WITH ITERATIVEIMPUTER BY PASSING IN DIFFERENT REGRESSORS TO BE USED FOR PREDICTING MISSING FEATURE VALUES IN THE CASE OF MISSFOREST THIS REGRESSOR IS A RANDOM FOREST SEE SPHXLRAUTOEXAMPLESPLOTITERATIVEIMPUTERVARIANTSCOMPARISONPY

SCIKITLEARN USER GUIDE RELEASE 0213

MULTIPLE VS SINGLE IMPUTATION

IN THE STATISTICS COMMUNITY IT IS COMMON PRACTICE TO PERFORM MULTIPLE IMPUTATIONS GENERATING FOR EXAMPLE MSEPARETE IMPUTATIONS FOR A SINGLE FEATURE MATRIX EACH OF THESE MIMPUTATIONS IS THEN PUT THROUGH THE SUBSEQUENT ANALYSIS PIPELINE EG FEATURE ENGINEERING CLUSTERING REGRESSION CLASSIFICATION THE MFINAL ANALYSIS RESULTS EG HELDOUT VALIDATION ERRORS ALLOW THE DATA SCIENTIST TO OBTAIN UNDERSTANDING OF HOW ANALYTIC RESULTS MAY DIFFER AS A CONSEQUENCE OF THE INHERENT UNCERTAINTY CAUSED BY THE MISSING VALUES THE ABOVE PRACTICE IS CALLED MULTIPLE IMPUTATION OUR IMPLEMENTATION OF ITERATIVEIMPUTER WAS INSPIRED BY THE R MICE PACKAGE MULTIVARIATE IMPUTATION BY CHAINED EQUATIONS<sup>1</sup> BUT DIFFERS FROM IT BY RETURNING A SINGLE IMPUTATION INSTEAD OF MULTIPLE IMPUTATIONS HOWEVER ITERATIVEIMPUTER CAN ALSO BE USED FOR MULTIPLE IMPUTATIONS BY APPLYING IT REPEATEDLY TO THE SAME DATASET WITH DIFFERENT RANDOM SEEDS WHEN SAMPLEPOSTERIORTRUE SEE<sup>2</sup> CHAPTER 4 FOR MORE DISCUSSION ON MULTIPLE VS SINGLE IMPUTATIONS

IT IS STILL AN OPEN PROBLEM AS TO HOW USEFUL SINGLE VS MULTIPLE IMPUTATION IS IN THE CONTEXT OF PREDICTION AND CLASSIFICA TION WHEN THE USER IS NOT INTERESTED IN MEASURING UNCERTAINTY DUE TO MISSING VALUES NOTE THAT A CALL TO THE TRANSFORM METHOD OFITERATIVEIMPUTER IS NOT ALLOWED TO CHANGE THE NUMBER OF SAMPLES THEREFORE MULTIPLE IMPUTATIONS CANNOT BE ACHIEVED BY A SINGLE CALL TO TRANSFORM

REFERENCES

MARKING IMPUTED VALUES

THEMISSINGINDICATOR TRANSFORMER IS USEFUL TO TRANSFORM A DATASET INTO CORRESPONDING BINARY MATRIX INDICATING THE PRESENCE OF MISSING VALUES IN THE DATASET THIS TRANSFORMATION IS USEFUL IN CONJUNCTION WITH IMPUTATION WHEN USING IMPUTATION PRESERVING THE INFORMATION ABOUT WHICH VALUES HAD BEEN MISSING CAN BE INFORMATIVE NAN IS USUALLY USED AS THE PLACEHOLDER FOR MISSING VALUES HOWEVER IT ENFORCES THE DATA TYPE TO BE FLOAT THE PARAMETER MISSINGVALUES ALLOWS TO SPECIFY OTHER PLACEHOLDER SUCH AS INTEGER IN THE FOLLOWING EXAMPLE WE WILL USE 1AS MISSING VALUES

```
FROM SKLEARNIMPUTE IMPORT MISSINGINDICATOR
X NPARRAY1 1 1 3
4 1 0 1
8 1 1 0
INDICATOR MISSINGINDICATORMISSINGVALUES1
MASKMISSINGVALUESONLY INDICATORFITTRANSFORMX
MASKMISSINGVALUESONLY
ARRAY TRUE TRUE FALSE
FALSE TRUE TRUE
FALSE TRUE FALSE
```

THEFEATURES PARAMETER IS USED TO CHOOSE THE FEATURES FOR WHICH THE MASK IS CONSTRUCTED BY DEFAULT IT IS MISSINGONLY WHICH RETURNS THE IMPUTER MASK OF THE FEATURES CONTAINING MISSING VALUES AT FIT TIME INDICATORFEATURES

```
ARRAY0 1 3
THEFEATURES PARAMETER CAN BE SET TO ALL TO RETURNED ALL FEATURES WHETHER OR NOT THEY CONTAIN MISSING VALUES
INDICATOR MISSINGINDICATORMISSINGVALUES1 FEATURESALL
MASKALL INDICATORFITTRANSFORMX
1STEF VAN BUUREN KARIN GROOTHUISOUDSHOORN 2011 “MICE MULTIVARIATE IMPUTATION BY CHAINED EQUATIONS IN R” JOURNAL OF STATISTI
SOFTWARE 45 167
2RODERICK J A LITTLE AND DONALD B RUBIN 1986 “STATISTICAL ANALYSIS WITH MISSING DATA” JOHN WILEY SONS INC NEW YORK NY USA
35 DATASET TRANSFORMATIONS 613
```

SCIKITLEARN USER GUIDE RELEASE 0213

MASKALL

ARRAY TRUE TRUE FALSE FALSE

FALSE TRUE FALSE TRUE

FALSE TRUE FALSE FALSE

INDICATORFEATURES

ARRAY0 1 2 3

WHEN USING THE MISSINGINDICATOR IN APIPELINE BE SURE TO USE THE FEATUREUNION OR

COLUMNTRANSFORMER TO ADD THE INDICATOR FEATURES TO THE REGULAR FEATURES FIRST WE OBTAIN THE IRIS DATASET

AND ADD SOME MISSING VALUES TO IT

FROM SKLEARNDATASETS IMPORT LOADIRIS

FROM SKLEARNIMPUTE IMPORT SIMPLEIMPUTER MISSINGINDICATOR

FROM SKLEARNMODELSELECTION IMPORT TRAINTESTSPLIT

FROM SKLEARNPIPELINE IMPORT FEATUREUNION MAKEPIPELINE

FROM SKLEARNTREE IMPORT DECISIONTREECLASSIFIER

X Y LOADIRISRETURNXY TRUE

MASK NPRANDMRANDINT0 2 SIZEXSHAPEASTYPENPBOOL

XMASK NPNAN

XTRAIN XTEST YTRAIN TRAINTESTSPLITX Y TESTSIZE100

RANDOMSTATE0

NOW WE CREATE A FEATUREUNION ALL FEATURES WILL BE IMPUTED USING SIMPLEIMPUTER IN ORDER TO ENABLE CLASSI

FIERS TO WORK WITH THIS DATA ADDITIONALLY IT ADDS THE THE INDICATOR VARIABLES FROM MISSINGINDICATOR

TRANSFORMER FEATUREUNION

TRANSFORMERLIST

FEATURES SIMPLEIMPUTERSTRATEGYMEAN

INDICATORS MISSINGINDICATOR

TRANSFORMER TRANSFORMERFITXTRAIN YTRAIN

RESULTS TRANSFORMERTRANSFORMXTEST

RESULTSSHAPE

100 8

OF COURSE WE CANNOT USE THE TRANSFORMER TO MAKE ANY PREDICTIONS WE SHOULD WRAP THIS IN A PIPELINE WITH A

CLASSIFIER EG A DECISIONTREECLASSIFIER TO BE ABLE TO MAKE PREDICTIONS

CLF MAKEPIPELINETRANSFORMER DECISIONTREECLASSIFIER

CLF CLFFITXTRAIN YTRAIN

RESULTS CLFPREDICTXTEST

RESULTSSHAPE

100

355 UNSUPERVISED DIMENSIONALITY REDUCTION

IF YOUR NUMBER OF FEATURES IS HIGH IT MAY BE USEFUL TO REDUCE IT WITH AN UNSUPERVISED STEP PRIOR TO SUPERVISED STEPS

MANY OF THE UNSUPERVISED LEARNING METHODS IMPLEMENT A TRANSFORM METHOD THAT CAN BE USED TO REDUCE THE DIMEN

SIONALITY BELOW WE DISCUSS TWO SPECIFIC EXAMPLE OF THIS PATTERN THAT ARE HEAVILY USED

PIPELINING

THE UNSUPERVISED DATA REDUCTION AND THE SUPERVISED ESTIMATOR CAN BE CHAINED IN ONE STEP SEE PIPELINE CHAINING

ESTIMATORS

614 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

PCA PRINCIPAL COMPONENT ANALYSIS

DECOMPOSITIONPCA LOOKS FOR A COMBINATION OF FEATURES THAT CAPTURE WELL THE VARIANCE OF THE ORIGINAL FEATURES  
SEEDECOMPOSING SIGNALS IN COMPONENTS MATRIX FACTORIZATION PROBLEMS

EXAMPLES

- FACES RECOGNITION EXAMPLE USING EIGENFACES AND SVMs

RANDOM PROJECTIONS

THE MODULE RANDOMPROJECTION PROVIDES SEVERAL TOOLS FOR DATA REDUCTION BY RANDOM PROJECTIONS SEE THE  
RELEVANT SECTION OF THE DOCUMENTATION RANDOM PROJECTION

EXAMPLES

- THE JOHNSONLINDENSTRAUSS BOUND FOR EMBEDDING WITH RANDOM PROJECTIONS

FEATURE AGGLOMERATION

CLUSTERFEATUREAGGLOMERATION APPLIES HIERARCHICAL CLUSTERING TO GROUP TOGETHER FEATURES THAT BEHAVE SIM  
ILARLY

EXAMPLES

- FEATURE AGGLOMERATION VS UNIVARIATE SELECTION

- FEATURE AGGLOMERATION

FEATURE SCALING

NOTE THAT IF FEATURES HAVE VERY DIFFERENT SCALING OR STATISTICAL PROPERTIES CLUSTERFEATUREAGGLOMERATION  
MAY NOT BE ABLE TO CAPTURE THE LINKS BETWEEN RELATED FEATURES USING A PREPROCESSINGSTANDARDSCALER  
CAN BE USEFUL IN THESE SETTINGS

356 RANDOM PROJECTION

THESKLEARNRANDOMPROJECTION MODULE IMPLEMENTS A SIMPLE AND COMPUTATIONALLY EFFICIENT WAY TO REDUCE  
THE DIMENSIONALITY OF THE DATA BY TRADING A CONTROLLED AMOUNT OF ACCURACY AS ADDITIONAL VARIANCE FOR FASTER PROCESSING  
TIMES AND SMALLER MODEL SIZES THIS MODULE IMPLEMENTS TWO TYPES OF UNSTRUCTURED RANDOM MATRIX GAUSSIAN RANDOM  
MATRIX ANDSPARSE RANDOM MATRIX

THE DIMENSIONS AND DISTRIBUTION OF RANDOM PROJECTIONS MATRICES ARE CONTROLLED SO AS TO PRESERVE THE PAIRWISE DISTANCES  
BETWEEN ANY TWO SAMPLES OF THE DATASET THUS RANDOM PROJECTION IS A SUITABLE APPROXIMATION TECHNIQUE FOR DISTANCE  
BASED METHOD

35 DATASET TRANSFORMATIONS 615

SCIKITLEARN USER GUIDE RELEASE 0213

REFERENCES

- SANJOY DASGUPTA 2000 EXPERIMENTS WITH RANDOM PROJECTION IN PROCEEDINGS OF THE SIXTEENTH CONFERENCE ON UNCERTAINTY IN ARTIFICIAL INTELLIGENCE UAI'00 CRAIG BOUTILIER AND MOISÉS GOLDSZMIDT EDS MORGAN KAUFMANN PUBLISHERS INC SAN FRANCISCO CA USA 143151
- ELLA BINGHAM AND HEIKKI MANNILA 2001 RANDOM PROJECTION IN DIMENSIONALITY REDUCTION APPLICATIONS TO IMAGE AND TEXT DATA IN PROCEEDINGS OF THE SEVENTH ACM SIGKDD INTERNATIONAL CONFERENCE ON KNOWLEDGE DISCOVERY AND DATA MINING KDD '01 ACM NEW YORK NY USA 245250

THE JOHNSONLINDENSTRAUSS LEMMA

THE MAIN THEORETICAL RESULT BEHIND THE EFFICIENCY OF RANDOM PROJECTION IS THE JOHNSONLINDENSTRAUSS LEMMA QUOTING WIKIPEDIA

IN MATHEMATICS THE JOHNSONLINDENSTRAUSS LEMMA IS A RESULT CONCERNING LOWDISTORTION EMBEDDINGS OF POINTS FROM HIGHDIMENSIONAL INTO LOWDIMENSIONAL EUCLIDEAN SPACE THE LEMMA STATES THAT A SMALL SET OF POINTS IN A HIGHDIMENSIONAL SPACE CAN BE EMBEDDED INTO A SPACE OF MUCH LOWER DIMENSION IN SUCH A WAY THAT DISTANCES BETWEEN THE POINTS ARE NEARLY PRESERVED THE MAP USED FOR THE EMBEDDING IS AT LEAST LIPSCHITZ AND CAN EVEN BE TAKEN TO BE AN ORTHOGONAL PROJECTION

KNOWING ONLY THE NUMBER OF SAMPLES THE SKLEARNRANDOMPROJECTION

JOHNSONLINDENSTRAUSSMINDIM ESTIMATES CONSERVATIVELY THE MINIMAL SIZE OF THE RANDOM SUBSPACE TO GUARANTEE A BOUNDED DISTORTION INTRODUCED BY THE RANDOM PROJECTION

FROM SKLEARNRANDOMPROJECTION IMPORT JOHNSONLINDENSTRAUSSMINDIM

JOHNSONLINDENSTRAUSSMINDIMNSAMPLES1E6 EPS05

663

JOHNSONLINDENSTRAUSSMINDIMNSAMPLES1E6 EPS05 01 001

ARRAY 663 11841 1112658

JOHNSONLINDENSTRAUSSMINDIMNSAMPLES1E4 1E5 1E6 EPS01

ARRAY 7894 9868 11841

EXAMPLE

- SEE THE JOHNSONLINDENSTRAUSS BOUND FOR EMBEDDING WITH RANDOM PROJECTIONS FOR A THEORETICAL EXPLICATION ON THE JOHNSONLINDENSTRAUSS LEMMA AND AN EMPIRICAL VALIDATION USING SPARSE RANDOM MATRICES

REFERENCES

- SANJOY DASGUPTA AND ANUPAM GUPTA 1999 AN ELEMENTARY PROOF OF THE JOHNSONLINDENSTRAUSS LEMMA GAUSSIAN RANDOM PROJECTION

THESKLEARNRANDOMPROJECTIONGAUSSIANRANDOMPROJECTION REDUCES THE DIMENSIONALITY BY PRO

JECTING THE ORIGINAL INPUT SPACE ON A RANDOMLY GENERATED MATRIX WHERE COMPONENTS ARE DRAWN FROM THE FOLLOWING DISTRIBUTION ¶01

¶¶¶¶¶¶¶¶¶¶

HERE A SMALL EXCERPT WHICH ILLUSTRATES HOW TO USE THE GAUSSIAN RANDOM PROJECTION TRANSFORMER

616 CHAPTER 3 USER GUIDE



SCIKITLEARN USER GUIDE RELEASE 0213

```
import numpy as np
from sklearn import randomprojection
X = np.random.rand(100, 10000)
transformer = randomprojection.GaussianRandomProjection
X_new = transformer.fit_transform(X)
X_new.shape
100 3947
```

Sparse Random Projection

The `sklearn.randomprojection.SparseRandomProjection` reduces the dimensionality by projecting the original input space using a sparse random matrix. Sparse random matrices are an alternative to dense Gaussian random projection matrix that guarantees similar embedding quality while being much more memory efficient and allowing faster computation of the projected data. If we define  $s = 1/\text{density}$ , the elements of the random matrix are drawn from

$$\begin{cases} \sqrt{s} & \text{with probability } 1 - 1/s \\ 0 & \text{with probability } 1 - 1/s \end{cases}$$

where `components` is the size of the projected subspace. By default, the density of non-zero elements is set to the minimum density as recommended by Ping Li et al.  $1/\sqrt{\text{features}}$ .

Here a small excerpt which illustrates how to use the `SparseRandomProjection` transformer:

```
import numpy as np
from sklearn import randomprojection
X = np.random.rand(100, 10000)
transformer = randomprojection.SparseRandomProjection
X_new = transformer.fit_transform(X)
X_new.shape
100 3947
```

References

- D. Achlioptas, 2003, Database-friendly random projections, *Journal of Computer and System Sciences* 66 (2003) 671–687
- Ping Li, Trevor J. Hastie and Kenneth W. Church, 2006, Very sparse random projections, in *Proceedings of the 12th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '06, ACM, New York, NY, USA, 287–296

357 Kernel Approximation

This submodule contains functions that approximate the feature mappings that correspond to certain kernels as they are used for example in support vector machines. See *Support Vector Machines*. The following feature functions perform nonlinear transformations of the input which can serve as a basis for linear classification or other algorithms. The advantage of using approximate explicit feature maps compared to the kernel trick which makes use of feature maps implicitly is that explicit mappings can be better suited for online learning and can significantly reduce the cost of learning with very large datasets. Standard kernelized SVMs do not scale well to large datasets but using an



SCIKITLEARN USER GUIDE RELEASE 0213

APPROXIMATE KERNEL MAP IT IS POSSIBLE TO USE MUCH MORE EFFICIENT LINEAR SVMs IN PARTICULAR THE COMBINATION OF KERNEL MAP APPROXIMATIONS WITH SGDCLASSIFIER CAN MAKE NONLINEAR LEARNING ON LARGE DATASETS POSSIBLE

SINCE THERE HAS NOT BEEN MUCH EMPIRICAL WORK USING APPROXIMATE EMBEDDINGS IT IS ADVISABLE TO COMPARE RESULTS AGAINST EXACT KERNEL METHODS WHEN POSSIBLE

SEE ALSO

POLYNOMIAL REGRESSION EXTENDING LINEAR MODELS WITH BASIS FUNCTIONS FOR AN EXACT POLYNOMIAL TRANSFORMATION

NYSTROEM METHOD FOR KERNEL APPROXIMATION

THE NYSTROEM METHOD AS IMPLEMENTED IN NYSTROEM IS A GENERAL METHOD FOR LOWRANK APPROXIMATIONS OF KERNELS IT ACHIEVES THIS BY ESSENTIALLY SUBSAMPLING THE DATA ON WHICH THE KERNEL IS EVALUATED BY DEFAULT NYSTROEM USES THE RBF KERNEL BUT IT CAN USE ANY KERNEL FUNCTION OR A PRECOMPUTED KERNEL MATRIX THE NUMBER OF SAMPLES USED WHICH IS ALSO THE DIMENSIONALITY OF THE FEATURES COMPUTED IS GIVEN BY THE PARAMETER NCOMPONENTS

RADIAL BASIS FUNCTION KERNEL

THE RBFSAMPLER CONSTRUCTS AN APPROXIMATE MAPPING FOR THE RADIAL BASIS FUNCTION KERNEL ALSO KNOWN AS RANDOM KITCHEN SINKS RR2007 THIS TRANSFORMATION CAN BE USED TO EXPLICITLY MODEL A KERNEL MAP PRIOR TO APPLYING A LINEAR ALGORITHM FOR EXAMPLE A LINEAR SVM

```
FROM SKLEARNKERNELAPPROXIMATION IMPORT RBFSAMPLER
FROM SKLEARNLINEARMODEL IMPORT SGDCLASSIFIER
X 0 0 1 1 1 0 0 1
Y 0 0 1 1
RBFFEATURE RBFSAMPLERGAMMA1 RANDOMSTATE1
XFEATURES RBFFEATUREFITTRANSFORMX
CLF SGDCLASSIFIERMAXITER5
CLFFITXFEATURES Y
SGDCLASSIFIERALPHA00001 AVERAGEFALSE CLASSWEIGHTNONE
EARLYSTOPPINGFALSE EPSILON01 ETA000 FITINTERCEPTTRUE
L1RATIO015 LEARNINGRATEOPTIMAL LOSSHINGE MAXITER5
NITERNOCHANGE5 NJOBSNONE PENALTYL2 POWERT05
RANDOMSTATENONE SHUFFLETRUE TOL0001 VALIDATIONFRACTION01
VERBOSE0 WARMSTARTFALSE
CLFSCOREXFEATURES Y
10
```

THE MAPPING RELIES ON A MONTE CARLO APPROXIMATION TO THE KERNEL VALUES THE FIT FUNCTION PERFORMS THE MONTE CARLO SAMPLING WHEREAS THE TRANSFORM METHOD PERFORMS THE MAPPING OF THE DATA BECAUSE OF THE INHERENT RANDOMNESS OF THE PROCESS RESULTS MAY VARY BETWEEN DIFFERENT CALLS TO THE FIT FUNCTION

THE FIT FUNCTION TAKES TWO ARGUMENTS NCOMPONENTS WHICH IS THE TARGET DIMENSIONALITY OF THE FEATURE TRANSFORM AND GAMMA THE PARAMETER OF THE RBF KERNEL A HIGHER NCOMPONENTS WILL RESULT IN A BETTER APPROXIMATION OF THE KERNEL AND WILL YIELD RESULTS MORE SIMILAR TO THOSE PRODUCED BY A KERNEL SVM NOTE THAT “FITTING” THE FEATURE FUNCTION DOES NOT ACTUALLY DEPEND ON THE DATA GIVEN TO THE FIT FUNCTION ONLY THE DIMENSIONALITY OF THE DATA IS USED DETAILS ON THE METHOD CAN BE FOUND IN RR2007

FOR A GIVEN VALUE OF NCOMPONENTS RBFSAMPLER IS OFTEN LESS ACCURATE AS NYSTROEM RBFSAMPLER IS CHEAPER TO COMPUTE THOUGH MAKING USE OF LARGER FEATURE SPACES MORE EFFICIENT

EXAMPLES

35 DATASET TRANSFORMATIONS 619

FIG 39 COMPARING AN EXACT RBF KERNEL LEFT WITH THE APPROXIMATION RIGHT

•EXPLICIT FEATURE MAP APPROXIMATION FOR RBF KERNELS

ADDITIVE CHI SQUARED KERNEL

THE ADDITIVE CHI SQUARED KERNEL IS A KERNEL ON HISTOGRAMS OFTEN USED IN COMPUTER VISION

THE ADDITIVE CHI SQUARED KERNEL AS USED HERE IS GIVEN BY

$$K(x,y) = \sum_{i=1}^n \frac{1}{2} \left( \frac{x_i - y_i}{\sigma_i} \right)^2$$

THIS IS NOT EXACTLY THE SAME AS SKLEARNMETRICSADDITIVECHI2KERNEL THE AUTHORS OF VZ2010 PREFER THE VERSION ABOVE AS IT IS ALWAYS POSITIVE DEFINITE SINCE THE KERNEL IS ADDITIVE IT IS POSSIBLE TO TREAT ALL COMPONENTS SEPARATELY FOR EMBEDDING THIS MAKES IT POSSIBLE TO SAMPLE THE FOURIER TRANSFORM IN REGULAR INTERVALS INSTEAD OF APPROXIMATING USING MONTE CARLO SAMPLING

THE CLASSADDITIVECHI2SAMPLER IMPLEMENTS THIS COMPONENT WISE DETERMINISTIC SAMPLING EACH COMPONENT IS SAMPLED  $\frac{1}{\sigma_i}$  TIMES YIELDING  $2 \times \frac{1}{\sigma_i}$  DIMENSIONS PER INPUT DIMENSION THE MULTIPLE OF TWO STEMS FROM THE REAL AND COMPLEX PART OF THE FOURIER TRANSFORM IN THE LITERATURE  $\frac{1}{\sigma_i}$  IS USUALLY CHOSEN TO BE 1 OR 2 TRANSFORMING THE DATASET TO  $2 \times \sum \frac{1}{\sigma_i}$  DIMENSIONS

THE APPROXIMATE FEATURE MAP PROVIDED BY ADDITIVECHI2SAMPLER CAN BE COMBINED WITH THE APPROXIMATE FEATURE MAP PROVIDED BY RBFSAMPLER TO YIELD AN APPROXIMATE FEATURE MAP FOR THE EXPONENTIATED CHI SQUARED KERNEL SEE THEVZ2010 FOR DETAILS AND VVZ2010 FOR COMBINATION WITH THE RBFSAMPLER

SKEWED CHI SQUARED KERNEL

THE SKEWED CHI SQUARED KERNEL IS GIVEN BY

$$K(x,y) = \sum_{i=1}^n \frac{1}{2} \left( \frac{x_i - y_i}{\sigma_i} \right)^2 \exp \left( -\frac{(x_i - y_i)^2}{\sigma_i^2} \right)$$

IT HAS PROPERTIES THAT ARE SIMILAR TO THE EXPONENTIATED CHI SQUARED KERNEL OFTEN USED IN COMPUTER VISION BUT ALLOWS FOR A SIMPLE MONTE CARLO APPROXIMATION OF THE FEATURE MAP

SCIKITLEARN USER GUIDE RELEASE 0213

THE USAGE OF THE SKEWEDCHI2SAMPLER IS THE SAME AS THE USAGE DESCRIBED ABOVE FOR THE RBFSAMPLER THE ONLY DIFFERENCE IS IN THE FREE PARAMETER THAT IS CALLED  $\gamma$  FOR A MOTIVATION FOR THIS MAPPING AND THE MATHEMATICAL DETAILS SEE LS2010

MATHEMATICAL DETAILS

KERNEL METHODS LIKE SUPPORT VECTOR MACHINES OR KERNELIZED PCA RELY ON A PROPERTY OF REPRODUCING KERNEL HILBERT SPACES FOR ANY POSITIVE DEFINITE KERNEL FUNCTION  $\kappa$  A SO CALLED MERCER KERNEL IT IS GUARANTEED THAT THERE EXISTS A MAPPING  $\phi$  INTO A HILBERT SPACE  $\mathcal{H}$  SUCH THAT

$$\kappa(x, y) = \langle \phi(x), \phi(y) \rangle$$

WHERE  $\langle \cdot, \cdot \rangle$  DENOTES THE INNER PRODUCT IN THE HILBERT SPACE

IF AN ALGORITHM SUCH AS A LINEAR SUPPORT VECTOR MACHINE OR PCA RELIES ONLY ON THE SCALAR PRODUCT OF DATA POINTS  $x, y$  ONE MAY USE THE VALUE OF  $\kappa(x, y)$  WHICH CORRESPONDS TO APPLYING THE ALGORITHM TO THE MAPPED DATA POINTS  $\phi(x), \phi(y)$  THE ADVANTAGE OF USING  $\kappa$  IS THAT THE MAPPING  $\phi$  NEVER HAS TO BE CALCULATED EXPLICITLY ALLOWING FOR ARBITRARY LARGE FEATURES EVEN INFINITE

ONE DRAWBACK OF KERNEL METHODS IS THAT IT MIGHT BE NECESSARY TO STORE MANY KERNEL VALUES  $\kappa(x_i, x_j)$  DURING OPTIMIZATION IF A KERNELIZED CLASSIFIER IS APPLIED TO NEW DATA  $x_{new}$   $\kappa(x_{new}, x_i)$  NEEDS TO BE COMPUTED TO MAKE PREDICTIONS POSSIBLY FOR MANY DIFFERENT  $x_i$  IN THE TRAINING SET

THE CLASSES IN THIS SUBMODULE ALLOW TO APPROXIMATE THE EMBEDDING  $\phi$  THEREBY WORKING EXPLICITLY WITH THE REPRESENTATIONS  $\phi(x)$  WHICH OBIATES THE NEED TO APPLY THE KERNEL OR STORE TRAINING EXAMPLES

REFERENCES

358 PAIRWISE METRICS AFFINITIES AND KERNELS

THE SKLEARNMETRICSPAIRWISE SUBMODULE IMPLEMENTS UTILITIES TO EVALUATE PAIRWISE DISTANCES OR AFFINITY OF SETS OF SAMPLES

THIS MODULE CONTAINS BOTH DISTANCE METRICS AND KERNELS A BRIEF SUMMARY IS GIVEN ON THE TWO HERE

DISTANCE METRICS ARE FUNCTIONS  $d(a, b)$  SUCH THAT  $d(a, b) \geq d(a, c)$  IF OBJECTS  $a$  AND  $b$  ARE CONSIDERED “MORE SIMILAR” THAN OBJECTS  $a$  AND  $c$  TWO OBJECTS EXACTLY ALIKE WOULD HAVE A DISTANCE OF ZERO ONE OF THE MOST POPULAR EXAMPLES IS EUCLIDEAN DISTANCE TO BE A ‘TRUE’ METRIC IT MUST OBEY THE FOLLOWING FOUR CONDITIONS

- 1  $d(a, b) \geq 0$  FOR ALL  $a, b$
- 2  $d(a, b) = 0$  IF AND ONLY IF  $a = b$  POSITIVE DEFINITENESS
- 3  $d(a, b) = d(b, a)$  SYMMETRY
- 4  $d(a, c) \leq d(a, b) + d(b, c)$  THE TRIANGLE INEQUALITY

KERNELS ARE MEASURES OF SIMILARITY IE  $\kappa(a, b) \geq \kappa(a, c)$  IF OBJECTS  $a$  AND  $b$  ARE CONSIDERED “MORE SIMILAR” THAN OBJECTS  $a$  AND  $c$  A KERNEL MUST ALSO BE POSITIVE SEMIDEFINITE

THERE ARE A NUMBER OF WAYS TO CONVERT BETWEEN A DISTANCE METRIC AND A SIMILARITY MEASURE SUCH AS A KERNEL LET  $d$  BE THE DISTANCE AND  $\kappa$  BE THE KERNEL

$$\kappa(x, y) = \frac{1}{2} (d(x, y)^2 + d(x, x)^2 + d(y, y)^2)$$

15 NPEXP  $\gamma$  WHERE ONE HEURISTIC FOR CHOOSING  $\gamma$  IS  $\frac{1}{2 \text{NUMFEATURES}}$

25  $\frac{1}{2} \text{D} \cdot \text{NPMAXD}$

35 DATASET TRANSFORMATIONS 621

SCIKITLEARN USER GUIDE RELEASE 0213

THE DISTANCES BETWEEN THE ROW VECTORS OF X AND THE ROW VECTORS OF Y CAN BE EVALUATED USING PAIRWISEDISTANCES  
IFY IS OMITTED THE PAIRWISE DISTANCES OF THE ROW VECTORS OF X ARE CALCULATED SIMILARLY PAIRWISE  
PAIRWISEKERNELS CAN BE USED TO CALCULATE THE KERNEL BETWEEN X AND Y USING DIFFERENT KERNEL FUNCTIONS SEE  
THE API REFERENCE FOR MORE DETAILS

```
import numpy as np
from sklearn.metrics import pairwise_distances
from sklearn.metrics.pairwise import pairwise_kernels
```

```
X = np.array([2, 3, 3, 5, 5, 8])
Y = np.array([1, 0, 2, 1])
pairwise_distances(X, Y, metric='manhattan')
array([[4, 2],
       [7, 5],
       [12, 10],
       [3, 0],
       [8, 5],
       [0, 3],
       [8, 5],
       [0, 3]])
pairwise_kernels(X, Y, metric='linear')
array([[2, 7],
       [3, 11],
       [5, 18],
       [3, 11],
       [5, 18],
       [3, 11],
       [5, 18],
       [3, 11]])
```

COSINE SIMILARITY

COSINESIMILARITY COMPUTES THE L2NORMALIZED DOT PRODUCT OF VECTORS THAT IS IF  $x$  AND  $y$  ARE ROW VECTORS THEIR  
COSINE SIMILARITY  $s(x, y)$  IS DEFINED AS

$$s(x, y) = \frac{x \cdot y}{\|x\|_2 \|y\|_2}$$

THIS IS CALLED COSINE SIMILARITY BECAUSE EUCLIDEAN L2 NORMALIZATION PROJECTS THE VECTORS ONTO THE UNIT SPHERE AND  
THEIR DOT PRODUCT IS THEN THE COSINE OF THE ANGLE BETWEEN THE POINTS DENOTED BY THE VECTORS

THIS KERNEL IS A POPULAR CHOICE FOR COMPUTING THE SIMILARITY OF DOCUMENTS REPRESENTED AS TFIDF VECTORS

COSINESIMILARITY ACCEPTSSCIPYSPARSE MATRICES NOTE THAT THE TFIDF FUNCTIONALITY IN SKLEARN  
FEATUREEXTRACTIONTEXT CAN PRODUCE NORMALIZED VECTORS IN WHICH CASE COSINESIMILARITY IS EQUIV  
ALENT TO LINEARKERNEL ONLY SLOWER

REFERENCES

- CD MANNING P RAGHAVAN AND H SCHÜTZE 2008 INTRODUCTION TO INFORMATION RETRIEVAL CAMBRIDGE UNI  
VERSITY PRESS <https://nlp.stanford.edu/IR-book/html/htmledition/the-vector-space-model-for-scoring-1.html>

LINEAR KERNEL  
THE FUNCTION LINEARKERNEL COMPUTES THE LINEAR KERNEL THAT IS A SPECIAL CASE OF POLYNOMIALKERNEL WITH  
DEGREE=1 AND COEF=0 HOMOGENEOUS IF X AND Y ARE COLUMN VECTORS THEIR LINEAR KERNEL IS

$$s(x, y) = x \cdot y$$

POLYNOMIAL KERNEL

THE FUNCTION POLYNOMIALKERNEL COMPUTES THE DEGREED POLYNOMIAL KERNEL BETWEEN TWO VECTORS THE POLYNOMIAL KERNEL REPRESENTS THE SIMILARITY BETWEEN TWO VECTORS CONCEPTUALLY THE POLYNOMIAL KERNELS CONSIDERS NOT ONLY THE SIMILARITY BETWEEN VECTORS UNDER THE SAME DIMENSION BUT ALSO ACROSS DIMENSIONS WHEN USED IN MACHINE LEARNING ALGORITHMS THIS ALLOWS TO ACCOUNT FOR FEATURE INTERACTION

THE POLYNOMIAL KERNEL IS DEFINED AS

$$K(x,y) = \sum_{i=0}^d (x \cdot y)^i$$

WHERE

- $x,y$  ARE THE INPUT VECTORS
  - $d$  IS THE KERNEL DEGREE
- IF  $d=0$  THE KERNEL IS SAID TO BE HOMOGENEOUS

SIGMOID KERNEL

THE FUNCTION SIGMOIDKERNEL COMPUTES THE SIGMOID KERNEL BETWEEN TWO VECTORS THE SIGMOID KERNEL IS ALSO KNOWN AS HYPERBOLIC TANGENT OR MULTILAYER PERCEPTRON BECAUSE IN THE NEURAL NETWORK FIELD IT IS OFTEN USED AS NEURON ACTIVATION FUNCTION IT IS DEFINED AS

$$K(x,y) = \tanh^2(\frac{1}{2}x \cdot y)$$

WHERE

- $x,y$  ARE THE INPUT VECTORS
- $\gamma$  IS KNOWN AS SLOPE
- $\gamma_0$  IS KNOWN AS INTERCEPT

RBFB KERNEL

THE FUNCTION RBFKERNEL COMPUTES THE RADIAL BASIS FUNCTION RBF KERNEL BETWEEN TWO VECTORS THIS KERNEL IS DEFINED AS

$$K(x,y) = \exp(-\gamma \|x - y\|^2)$$

WHERE  $x,y$  ARE THE INPUT VECTORS IF  $\gamma = \frac{1}{2\sigma^2}$  THE KERNEL IS KNOWN AS THE GAUSSIAN KERNEL OF VARIANCE  $\sigma^2$

LAPLACIAN KERNEL

THE FUNCTION LAPLACIANKERNEL IS A VARIANT ON THE RADIAL BASIS FUNCTION KERNEL DEFINED AS

$$K(x,y) = \exp(-\gamma \|x - y\|)$$

WHERE  $x,y$  ARE THE INPUT VECTORS AND  $\|x - y\|$  IS THE MANHATTAN DISTANCE BETWEEN THE INPUT VECTORS IT HAS PROVEN USEFUL IN ML APPLIED TO NOISELESS DATA SEE EG MACHINE LEARNING FOR QUANTUM MECHANICS IN A NUTSHELL

SCIKITLEARN USER GUIDE RELEASE 0213

CHISQUARED KERNEL

THE CHISQUARED KERNEL IS A VERY POPULAR CHOICE FOR TRAINING NONLINEAR SVMs IN COMPUTER VISION APPLICATIONS IT CAN BE COMPUTED USING CHI2KERNEL AND THEN PASSED TO AN SKLEARNsvmsvc WITHKERNELPRECOMPUTED

```
FROM SKLEARNsvm import SVC
FROM SKLEARNmetricspairwise import CHI2KERNEL
X 0 1 1 0 2 8 7 3
Y 0 1 0 1
K CHI2KERNELX GAMMA5
K
ARRAY1 036787944 089483932 058364548
036787944 1 051341712 083822343
089483932 051341712 1 07768366
058364548 083822343 07768366 1
SVM SVCKERNELPRECOMPUTEDFITK Y
SVMpredictK
ARRAY0 1 0 1
```

IT CAN ALSO BE DIRECTLY USED AS THE KERNEL ARGUMENT

```
SVM SVCKERNELCHI2KERNELFITX Y
SVMpredictX
ARRAY0 1 0 1
```

THE CHI SQUARED KERNEL IS GIVEN BY

$$K(x, y) = \frac{1}{2} \left( \|x\|^2 + \|y\|^2 - \|x - y\|^2 \right)$$

THE DATA IS ASSUMED TO BE NONNEGATIVE AND IS OFTEN NORMALIZED TO HAVE AN L1NORM OF ONE THE NORMALIZATION IS RATIONALIZED WITH THE CONNECTION TO THE CHI SQUARED DISTANCE WHICH IS A DISTANCE BETWEEN DISCRETE PROBABILITY DISTRIBUTIONS

THE CHI SQUARED KERNEL IS MOST COMMONLY USED ON HISTOGRAMS BAGS OF VISUAL WORDS

REFERENCES

- ZHANG J AND MARSZALEK M AND LAZEBNIK S AND SCHMID C LOCAL FEATURES AND KERNELS FOR CLASSIFICATION OF TEXTURE AND OBJECT CATEGORIES A COMPREHENSIVE STUDY INTERNATIONAL JOURNAL OF COMPUTER VISION 2007 [HTTPSRESEARCHMICROSOFTCOMENUSUMPEOPLEMANIKPROJECTSTRADEOFFPAPERSZHANGIJCV06PDF](https://research.microsoft.com/en-us/projects/trafford/papers/ZhangIJCV06.pdf)

359 TRANSFORMING THE PREDICTION TARGET Y

THESE ARE TRANSFORMERS THAT ARE NOT INTENDED TO BE USED ON FEATURES ONLY ON SUPERVISED LEARNING TARGETS SEE ALSO TRANSFORMING TARGET IN REGRESSION IF YOU WANT TO TRANSFORM THE PREDICTION TARGET FOR LEARNING BUT EVALUATE THE MODEL IN THE ORIGINAL UNTRANSFORMED SPACE

LABEL BINARIZATION

LABELBINARIZER IS A UTILITY CLASS TO HELP CREATE A LABEL INDICATOR MATRIX FROM A LIST OF MULTICLASS LABELS

624 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

FROM SKLEARN IMPORT PREPROCESSING

LB PREPROCESSINGLABELBINARIZER

LBFIT1 2 6 4 2

LABELBINARIZERNEGLABEL0 POSLABEL1 SPARSEOUTPUTFALSE

LBCLASSES

ARRAY1 2 4 6

LBTRANSFORM1 6

ARRAY1 0 0 0

0 0 0 1

FOR MULTIPLE LABELS PER INSTANCE USE MULTILABELBINARIZER

LB PREPROCESSINGMULTILABELBINARIZER

LBFITTRANSFORM1 2 3

ARRAY1 1 0

0 0 1

LBCLASSES

ARRAY1 2 3

LABEL ENCODING

LABELENCODER IS A UTILITY CLASS TO HELP NORMALIZE LABELS SUCH THAT THEY CONTAIN ONLY VALUES BETWEEN 0 AND NCLASSES

1 THIS IS SOMETIMES USEFUL FOR WRITING EFFICIENT CYTHON ROUTINES LABELENCODER CAN BE USED AS FOLLOWS

FROM SKLEARN IMPORT PREPROCESSING

LE PREPROCESSINGLABELENCODER

LEFIT1 2 2 6

LABELENCODER

LECLASSES

ARRAY1 2 6

LETRANSFORM1 1 2 6

ARRAY0 0 1 2

LEINVERSETRANSFORM0 0 1 2

ARRAY1 1 2 6

IT CAN ALSO BE USED TO TRANSFORM NONNUMERICAL LABELS AS LONG AS THEY ARE HASHABLE AND COMPARABLE TO NUMERICAL LABELS

LE PREPROCESSINGLABELENCODER

LEFITPARIS PARIS TOKYO AMSTERDAM

LABELENCODER

LISTLECLASSES

AMSTERDAM PARIS TOKYO

LETRANSFORMTOKYO TOKYO PARIS

ARRAY2 2 1

LISTLEINVERSETRANSFORM2 2 1

TOKYO TOKYO PARIS

36 DATASET LOADING UTILITIES

THESKLEARNDATASETS PACKAGE EMBEDS SOME SMALL TOY DATASETS AS INTRODUCED IN THE GETTING STARTED SECTION

THIS PACKAGE ALSO FEATURES HELPERS TO FETCH LARGER DATASETS COMMONLY USED BY THE MACHINE LEARNING COMMUNITY TO BENCHMARK ALGORITHMS ON DATA THAT COMES FROM THE ‘REAL WORLD’

36 DATASET LOADING UTILITIES 625

SCIKITLEARN USER GUIDE RELEASE 0213

TO EVALUATE THE IMPACT OF THE SCALE OF THE DATASET NSAMPLES ANDNFEATURES WHILE CONTROLLING THE STATISTICAL PROPERTIES OF THE DATA TYPICALLY THE CORRELATION AND INFORMATIVENESS OF THE FEATURES IT IS ALSO POSSIBLE TO GENERATE SYNTHETIC DATA

361 GENERAL DATASET API

THERE ARE THREE MAIN KINDS OF DATASET INTERFACES THAT CAN BE USED TO GET DATASETS DEPENDING ON THE DESIRED TYPE OF DATASET

THE DATASET LOADERS THEY CAN BE USED TO LOAD SMALL STANDARD DATASETS DESCRIBED IN THE TOY DATASETS SECTION

THE DATASET FETCHERS THEY CAN BE USED TO DOWNLOAD AND LOAD LARGER DATASETS DESCRIBED IN THE REAL WORLD DATASETS SECTION

BOTH LOADERS AND FETCHERS FUNCTIONS RETURN A DICTIONARYLIKE OBJECT HOLDING AT LEAST TWO ITEMS AN ARRAY OF SHAPE NSAMPLES NFEATURES WITH KEYDATA EXCEPT FOR 20NEWSGROUPS AND A NUMPY ARRAY OF LENGTH NSAMPLES CONTAINING THE TARGET VALUES WITH KEY TARGET

IT’S ALSO POSSIBLE FOR ALMOST ALL OF THESE FUNCTION TO CONSTRAIN THE OUTPUT TO BE A TUPLE CONTAINING ONLY THE DATA AND THE TARGET BY SETTING THE RETURNXY PARAMETER TO TRUE

THE DATASETS ALSO CONTAIN A FULL DESCRIPTION IN THEIR DESCR ATTRIBUTE AND SOME CONTAIN FEATURENAMES AND TARGETNAMES SEE THE DATASET DESCRIPTIONS BELOW FOR DETAILS

THE DATASET GENERATION FUNCTIONS THEY CAN BE USED TO GENERATE CONTROLLED SYNTHETIC DATASETS DESCRIBED IN THE GENERATED DATASETS SECTION

THESE FUNCTIONS RETURN A TUPLE X Y CONSISTING OF A NSAMPLES NFEATURES NUMPY ARRAY XAND AN ARRAY OF LENGTHNSAMPLES CONTAINING THE TARGETS Y

IN ADDITION THERE ARE ALSO MISCELLANEOUS TOOLS TO LOAD DATASETS OF OTHER FORMATS OR FROM OTHER LOCATIONS DESCRIBED IN THELOADING OTHER DATASETS SECTION

362 TOY DATASETS

SCIKITLEARN COMES WITH A FEW SMALL STANDARD DATASETS THAT DO NOT REQUIRE TO DOWNLOAD ANY FILE FROM SOME EXTERNAL WEBSITE

THEY CAN BE LOADED USING THE FOLLOWING FUNCTIONS

LOADBOSTON RETURNXY LOAD AND RETURN THE BOSTON HOUSEPRICES DATASET REGRES SION

LOADIRIS RETURNXY LOAD AND RETURN THE IRIS DATASET CLASSIFICATION

LOADDIABETES RETURNXY LOAD AND RETURN THE DIABETES DATASET REGRESSION

LOADDIGITS NCLASS RETURNXY LOAD AND RETURN THE DIGITS DATASET CLASSIFICATION

LOADLINNERUD RETURNXY LOAD AND RETURN THE LINNERUD DATASET MULTIVARIATE REGRES SION

LOADWINE RETURNXY LOAD AND RETURN THE WINE DATASET CLASSIFICATION

LOADBREASTCANCER RETURNXY LOAD AND RETURN THE BREAST CANCER WISCONSIN DATASET CLAS SIFICATION

THESE DATASETS ARE USEFUL TO QUICKLY ILLUSTRATE THE BEHAVIOR OF THE VARIOUS ALGORITHMS IMPLEMENTED IN SCIKITLEARN THEY ARE HOWEVER OFTEN TOO SMALL TO BE REPRESENTATIVE OF REAL WORLD MACHINE LEARNING TASKS

626 CHAPTER 3 USER GUIDE



SCIKITLEARN USER GUIDE RELEASE 0213  
BOSTON HOUSE PRICES DATASET  
DATA SET CHARACTERISTICS  
NUMBER OF INSTANCES 506  
NUMBER OF ATTRIBUTES 13 NUMERICCATEGORICAL PREDICTIVE MEDIAN VALUE ATTRIBUTE 14 IS USUALLY THE TARGET

- ATTRIBUTE INFORMATION IN ORDER
- CRIM PER CAPITA CRIME RATE BY TOWN
  - ZN PROPORTION OF RESIDENTIAL LAND ZONED FOR LOTS OVER 25000 SQFT
  - INDUS PROPORTION OF NONRETAIL BUSINESS ACRES PER TOWN
  - CHAS CHARLES RIVER DUMMY VARIABLE 1 IF TRACT BOUNDS RIVER 0 OTHERWISE
  - NOX NITRIC OXIDES CONCENTRATION PARTS PER 10 MILLION
  - RM AVERAGE NUMBER OF ROOMS PER DWELLING
  - AGE PROPORTION OF OWNEROCCUPIED UNITS BUILT PRIOR TO 1940
  - DIS WEIGHTED DISTANCES TO FIVE BOSTON EMPLOYMENT CENTRES
  - RAD INDEX OF ACCESSIBILITY TO RADIAL HIGHWAYS
  - TAX FULLVALUE PROPERTYTAX RATE PER 10000
  - PTRATIO PUPILTEACHER RATIO BY TOWN
  - B 1000BK 0632 WHERE BK IS THE PROPORTION OF BLACKS BY TOWN
  - LSTAT LOWER STATUS OF THE POPULATION
  - MEDV MEDIAN VALUE OF OWNEROCCUPIED HOMES IN 1000'S

MISSING ATTRIBUTE VALUES NONE  
CREATOR HARRISON D AND RUBINFELD DL  
THIS IS A COPY OF UCI ML HOUSING DATASET [HTTPSARCHIVEICSUCIEDUMLMACHINELEARNINGDATABASESHOUSING](https://archive.ics.uc.edu/ml/machine-learning-databases/housing)  
THIS DATASET WAS TAKEN FROM THE STATLIB LIBRARY WHICH IS MAINTAINED AT CARNEGIE MELLON UNIVERSITY  
THE BOSTON HOUSEPRICE DATA OF HARRISON D AND RUBINFELD DL 'HEDONIC PRICES AND THE DEMAND FOR CLEAN AIR' J ENVIRON ECONOMICS MANAGEMENT VOL5 81102 1978 USED IN BELSLEY KUH WELSCH 'REGRESSION DIAGNOSTICS' WILEY 1980 NB VARIOUS TRANSFORMATIONS ARE USED IN THE TABLE ON PAGES 244261 OF THE LATTER  
THE BOSTON HOUSEPRICE DATA HAS BEEN USED IN MANY MACHINE LEARNING PAPERS THAT ADDRESS REGRESSION PROBLEMS  
REFERENCES  
• BELSLEY KUH WELSCH 'REGRESSION DIAGNOSTICS IDENTIFYING INFLUENTIAL DATA AND SOURCES OF COLLINEARITY' WILEY 1980 244261  
• QUINLANR 1993 COMBINING INSTANCEBASED AND MODELBASED LEARNING IN PROCEEDINGS ON THE TENTH INTERNATIONAL CONFERENCE OF MACHINE LEARNING 236243 UNIVERSITY OF MASSACHUSETTS AMHERST MORGAN KAUFMANN  
36 DATASET LOADING UTILITIES 627

SCIKITLEARN USER GUIDE RELEASE 0213

IRIS PLANTS DATASET

DATA SET CHARACTERISTICS

NUMBER OF INSTANCES 150 50 IN EACH OF THREE CLASSES

NUMBER OF ATTRIBUTES 4 NUMERIC PREDICTIVE ATTRIBUTES AND THE CLASS

ATTRIBUTE INFORMATION

- SEPAL LENGTH IN CM
- SEPAL WIDTH IN CM
- PETAL LENGTH IN CM
- PETAL WIDTH IN CM
- CLASS

-IRISSETOSA

-IRISVERSICOLOUR

-IRISVIRGINICA

SUMMARY STATISTICS

SEPAL LENGTH 43 79 584 083 07826

SEPAL WIDTH 20 44 305 043 04194

PETAL LENGTH 10 69 376 176 09490 HIGH

PETAL WIDTH 01 25 120 076 09565 HIGH

MISSING ATTRIBUTE VALUES NONE

CLASS DISTRIBUTION 333 FOR EACH OF 3 CLASSES

CREATOR RA FISHER

DONOR MICHAEL MARSHALL MARSHALLPLUIOARCNASAGOV

DATE JULY 1988

THE FAMOUS IRIS DATABASE FIRST USED BY SIR RA FISHER THE DATASET IS TAKEN FROM FISHER’S PAPER NOTE THAT IT’S THE SAME AS IN R BUT NOT AS IN THE UCI MACHINE LEARNING REPOSITORY WHICH HAS TWO WRONG DATA POINTS

THIS IS PERHAPS THE BEST KNOWN DATABASE TO BE FOUND IN THE PATTERN RECOGNITION LITERATURE FISHER’S PAPER IS A CLASSIC IN THE FIELD AND IS REFERENCED FREQUENTLY TO THIS DAY SEE DUDA HART FOR EXAMPLE THE DATA SET CONTAINS 3 CLASSES OF 50 INSTANCES EACH WHERE EACH CLASS REFERS TO A TYPE OF IRIS PLANT ONE CLASS IS LINEARLY SEPARABLE FROM THE OTHER 2 THE LATTER ARE NOT LINEARLY SEPARABLE FROM EACH OTHER

REFERENCES

- FISHER RA “THE USE OF MULTIPLE MEASUREMENTS IN TAXONOMIC PROBLEMS” ANNUAL EUGENICS 7 PART II 179188 1936 ALSO IN “CONTRIBUTIONS TO MATHEMATICAL STATISTICS” JOHN WILEY NY 1950
- DUDA RO HART PE 1973 PATTERN CLASSIFICATION AND SCENE ANALYSIS Q327D83 JOHN WILEY SONS ISBN 0471223611 SEE PAGE 218
- DASARATHY BV 1980 “NOSING AROUND THE NEIGHBORHOOD A NEW SYSTEM STRUCTURE AND CLASSIFICATION RULE FOR RECOGNITION IN PARTIALLY EXPOSED ENVIRONMENTS” IEEE TRANSACTIONS ON PATTERN ANALYSIS AND MACHINE INTELLIGENCE V OL PAMI2 NO 1 6771 628 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

- GATES GW 1972 “THE REDUCED NEAREST NEIGHBOR RULE” IEEE TRANSACTIONS ON INFORMATION THEORY MAY 1972 431433

- SEE ALSO 1988 MLC PROCEEDINGS 5464 CHEESEMAN ET AL”S AUTOCLASS II CONCEPTUAL CLUSTERING SYSTEM FINDS 3 CLASSES IN THE DATA

- MANY MANY MORE

DIABETES DATASET

TEN BASELINE VARIABLES AGE SEX BODY MASS INDEX AVERAGE BLOOD PRESSURE AND SIX BLOOD SERUM MEASUREMENTS WERE OBTAINED FOR EACH OF N 442 DIABETES PATIENTS AS WELL AS THE RESPONSE OF INTEREST A QUANTITATIVE MEASURE OF DISEASE PROGRESSION ONE YEAR AFTER BASELINE

DATA SET CHARACTERISTICS

NUMBER OF INSTANCES 442

NUMBER OF ATTRIBUTES FIRST 10 COLUMNS ARE NUMERIC PREDICTIVE VALUES

TARGET COLUMN 11 IS A QUANTITATIVE MEASURE OF DISEASE PROGRESSION ONE YEAR AFTER BASELINE

ATTRIBUTE INFORMATION

- AGE
- SEX
- BODY MASS INDEX
- AVERAGE BLOOD PRESSURE

- S1

- S2

- S3

- S4

- S5

- S6

NOTE EACH OF THESE 10 FEATURE VARIABLES HAVE BEEN MEAN CENTERED AND SCALED BY THE STANDARD DEVIATION TIMES NSAMPLES IE THE SUM OF SQUARES OF EACH COLUMN TOTALS 1

SOURCE URL [HTTPSWWW4STATNCSUEDUBOOSVARSELECTDIABETESHTML](https://www4.stat.ncsu.edu/boos/varselect/diabetes.html)

FOR MORE INFORMATION SEE BRADLEY EFRON TREVOR HASTIE IAIN JOHNSTONE AND ROBERT TIBSHIRANI 2004 “LEAST ANGLE REGRESSION” ANNALS OF STATISTICS WITH DISCUSSION 407499 [HTTPSWEBSTANFORDEDUHASTIEPAPERSLARS LEASTANGLE2002PDF](https://web.stanford.edu/hastie/papers/lars/leastangle2002.pdf)

OPTICAL RECOGNITION OF HANDWRITTEN DIGITS DATASET

DATA SET CHARACTERISTICS

NUMBER OF INSTANCES 5620

NUMBER OF ATTRIBUTES 64

ATTRIBUTE INFORMATION 8X8 IMAGE OF INTEGER PIXELS IN THE RANGE 016

36 DATASET LOADING UTILITIES 629

SCIKITLEARN USER GUIDE RELEASE 0213  
MISSING ATTRIBUTE VALUES NONE  
CREATOR  
5 ALPAYDIN ALPAYDIN “ BOUNEDUTR  
DATE JULY 1998  
THIS IS A COPY OF THE TEST SET OF THE UCI ML HANDWRITTEN DIGITS DATASETS [HTTPSARCHIVEICSUCIEDUMLDATASETSOPTICAL  
RECOGNITIONOFHANDWRITTENDIGITS](https://archive.ics.uciedu/ml/dataset/optical_recognition_of_handwritten_digits)  
THE DATA SET CONTAINS IMAGES OF HANDWRITTEN DIGITS 10 CLASSES WHERE EACH CLASS REFERS TO A DIGIT  
PREPROCESSING PROGRAMS MADE AVAILABLE BY NIST WERE USED TO EXTRACT NORMALIZED BITMAPS OF HANDWRITTEN DIGITS FROM  
A PREPRINTED FORM FROM A TOTAL OF 43 PEOPLE 30 CONTRIBUTED TO THE TRAINING SET AND DIFFERENT 13 TO THE TEST SET 32X32  
BITMAPS ARE DIVIDED INTO NONOVERLAPPING BLOCKS OF 4X4 AND THE NUMBER OF ON PIXELS ARE COUNTED IN EACH BLOCK THIS  
GENERATES AN INPUT MATRIX OF 8X8 WHERE EACH ELEMENT IS AN INTEGER IN THE RANGE 016 THIS REDUCES DIMENSIONALITY AND  
GIVES INVARIANCE TO SMALL DISTORTIONS  
FOR INFO ON NIST PREPROCESSING ROUTINES SEE M D GARRIS J L BLUE G T CANDELA D L DIMMICK J GEIST P J  
GROTHER S A JANET AND C L WILSON NIST FORMBASED HANDPRINT RECOGNITION SYSTEM NISTIR 5469 1994  
REFERENCES  
• C KAYNAK 1995 METHODS OF COMBINING MULTIPLE CLASSIFIERS AND THEIR APPLICATIONS TO HANDWRITTEN DIGIT  
RECOGNITION MSC THESIS INSTITUTE OF GRADUATE STUDIES IN SCIENCE AND ENGINEERING BOGAZICI UNIVERSITY  
• 5 ALPAYDIN C KAYNAK 1998 CASCADING CLASSIFIERS KYBERNETIKA  
• KEN TANG AND PONNUTHURAI N SUGANTHAN AND XI YAO AND A KAI QIN LINEAR DIMENSIONALITYREDUCTION USING  
RELEVANCE WEIGHTED LDA SCHOOL OF ELECTRICAL AND ELECTRONIC ENGINEERING NANYANG TECHNOLOGICAL UNIVERSITY  
2005  
• CLAUDIO GENTILE A NEW APPROXIMATE MAXIMAL MARGIN CLASSIFICATION ALGORITHM NIPS 2000  
LINNERRUD DATASET  
DATA SET CHARACTERISTICS  
NUMBER OF INSTANCES 20  
NUMBER OF ATTRIBUTES 3  
MISSING ATTRIBUTE VALUES NONE  
THE LINNERRUD DATASET CONSTAINS TWO SMALL DATASET  
•PHYSIOLOGICAL CSV CONTAINING 20 OBSERVATIONS ON 3 EXERCISE VARIABLES WEIGHT WAIST AND PULSE  
•EXERCISE CSV CONTAINING 20 OBSERVATIONS ON 3 PHYSIOLOGICAL VARIABLES CHINS SITUPS AND JUMPS  
REFERENCES  
• TENENHAUS M 1998 LA REGRESSION PLS THEORIE ET PRATIQUE PARIS EDITIONS TECHNIC  
WINE RECOGNITION DATASET  
DATA SET CHARACTERISTICS  
630 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213  
NUMBER OF INSTANCES 178 50 IN EACH OF THREE CLASSES  
NUMBER OF ATTRIBUTES 13 NUMERIC PREDICTIVE ATTRIBUTES AND THE CLASS  
ATTRIBUTE INFORMATION

- ALCOHOL
- MALIC ACID
- ASH
- ALCALINITY OF ASH
- MAGNESIUM
- TOTAL PHENOLS
- FLAVANOIDS
- NONFLAVANOID PHENOLS
- PROANTHOCYANINS
- COLOR INTENSITY
- HUE
- OD280OD315 OF DILUTED WINES
- PROLINE

•CLASS  
-CLASS0  
-CLASS1  
-CLASS2

SUMMARY STATISTICS  
ALCOHOL 110 148 130 08  
MALIC ACID 074 580 234 112  
ASH 136 323 236 027  
ALCALINITY OF ASH 106 300 195 33  
MAGNESIUM 700 1620 997 143  
TOTAL PHENOLS 098 388 229 063  
FLAVANOIDS 034 508 203 100  
NONFLAVANOID PHENOLS 013 066 036 012  
PROANTHOCYANINS 041 358 159 057  
COLOUR INTENSITY 13 130 51 23  
HUE 048 171 096 023  
OD280OD315 OF DILUTED WINES 127 400 261 071  
PROLINE 278 1680 746 315  
MISSING ATTRIBUTE VALUES NONE  
CLASS DISTRIBUTION CLASS0 59 CLASS1 71 CLASS2 48  
CREATOR RA FISHER  
DONOR MICHAEL MARSHALL MARSHALLPLUIOARCNASAGOV  
36 DATASET LOADING UTILITIES 631

SCIKITLEARN USER GUIDE RELEASE 0213

DATE JULY 1988

THIS IS A COPY OF UCI ML WINE RECOGNITION DATASETS [HTTPSARCHIVEICSUCIEDUMLMACHINELEARNINGDATABASESWINE](https://archive.ics.uciedu/ml/machine-learning-databases/wine/wine.data)

THE DATA IS THE RESULTS OF A CHEMICAL ANALYSIS OF WINES GROWN IN THE SAME REGION IN ITALY BY THREE DIFFERENT CULTIVATORS THERE ARE THIRTEEN DIFFERENT MEASUREMENTS TAKEN FOR DIFFERENT CONSTITUENTS FOUND IN THE THREE TYPES OF WINE ORIGINAL OWNERS

FORINA M ET AL PARVUS AN EXTENDIBLE PACKAGE FOR DATA EXPLORATION CLASSIFICATION AND CORRELATION INSTITUTE OF PHARMACEUTICAL AND FOOD ANALYSIS AND TECHNOLOGIES VIA BRIGATA SALERNO 16147 GENOA ITALY

CITATION  
LICHMAN M 2013 UCI MACHINE LEARNING REPOSITORY [HTTPSARCHIVEICSUCIEDUML](https://archive.ics.uciedu/ml) IRVINE CA UNIVERSITY OF CALIFORNIA SCHOOL OF INFORMATION AND COMPUTER SCIENCE

REFERENCES  
1 S AEBERHARD D COOMANS AND O DE VEL COMPARISON OF CLASSIFIERS IN HIGH DIMENSIONAL SETTINGS TECH REP NO 9202 1992 DEPT OF COMPUTER SCIENCE AND DEPT OF MATHEMATICS AND STATISTICS JAMES COOK UNIVERSITY OF NORTH QUEENSLAND ALSO SUBMITTED TO TECHNOMETRICS  
THE DATA WAS USED WITH MANY OTHERS FOR COMPARING VARIOUS CLASSIFIERS THE CLASSES ARE SEPARABLE THOUGH ONLY RDA HAS ACHIEVED 100 CORRECT CLASSIFICATION RDA 100 QDA 994 LDA 989 1NN 961 ZTRANSFORMED DATA ALL RESULTS USING THE LEAVEONEOUT TECHNIQUE

2 S AEBERHARD D COOMANS AND O DE VEL “THE CLASSIFICATION PERFORMANCE OF RDA” TECH REP NO 9201 1992 DEPT OF COMPUTER SCIENCE AND DEPT OF MATHEMATICS AND STATISTICS JAMES COOK UNIVERSITY OF NORTH QUEENSLAND ALSO SUBMITTED TO JOURNAL OF CHEMOMETRICS  
BREAST CANCER WISCONSIN DIAGNOSTIC DATASET

DATA SET CHARACTERISTICS  
NUMBER OF INSTANCES 569  
NUMBER OF ATTRIBUTES 30 NUMERIC PREDICTIVE ATTRIBUTES AND THE CLASS  
ATTRIBUTE INFORMATION

- RADIUS MEAN OF DISTANCES FROM CENTER TO POINTS ON THE PERIMETER
  - TEXTURE STANDARD DEVIATION OF GRAYSCALE VALUES
  - PERIMETER
  - AREA
  - SMOOTHNESS LOCAL VARIATION IN RADIUS LENGTHS
  - COMPACTNESS PERIMETER<sup>2</sup> AREA 10
  - CONCAVITY SEVERITY OF CONCAVE PORTIONS OF THE CONTOUR
  - CONCAVE POINTS NUMBER OF CONCAVE PORTIONS OF THE CONTOUR
  - SYMMETRY
  - FRACTAL DIMENSION “COASTLINE APPROXIMATION” 1
- 632 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213  
THE MEAN STANDARD ERROR AND “WORST” OR LARGEST MEAN OF THE THREE LARGEST VALUES OF THESE  
FEATURES WERE COMPUTED FOR EACH IMAGE RESULTING IN 30 FEATURES FOR INSTANCE FIELD 3 IS  
MEAN RADIUS FIELD 13 IS RADIUS SE FIELD 23 IS WORST RADIUS  
•CLASS  
-WDBCMALIGNANT  
-WDBCBENIGN  
SUMMARY STATISTICS  
RADIUS MEAN 6981 2811  
TEXTURE MEAN 971 3928  
PERIMETER MEAN 4379 1885  
AREA MEAN 1435 25010  
SMOOTHNESS MEAN 0053 0163  
COMPACTNESS MEAN 0019 0345  
CONCAVITY MEAN 00 0427  
CONCAVE POINTS MEAN 00 0201  
SYMMETRY MEAN 0106 0304  
FRACTAL DIMENSION MEAN 005 0097  
RADIUS STANDARD ERROR 0112 2873  
TEXTURE STANDARD ERROR 036 4885  
PERIMETER STANDARD ERROR 0757 2198  
AREA STANDARD ERROR 6802 5422  
SMOOTHNESS STANDARD ERROR 0002 0031  
COMPACTNESS STANDARD ERROR 0002 0135  
CONCAVITY STANDARD ERROR 00 0396  
CONCAVE POINTS STANDARD ERROR 00 0053  
SYMMETRY STANDARD ERROR 0008 0079  
FRACTAL DIMENSION STANDARD ERROR 0001 003  
RADIUS WORST 793 3604  
TEXTURE WORST 1202 4954  
PERIMETER WORST 5041 2512  
AREA WORST 1852 42540  
SMOOTHNESS WORST 0071 0223  
COMPACTNESS WORST 0027 1058  
CONCAVITY WORST 00 1252  
CONCAVE POINTS WORST 00 0291  
SYMMETRY WORST 0156 0664  
FRACTAL DIMENSION WORST 0055 0208  
MISSING ATTRIBUTE VALUES NONE  
CLASS DISTRIBUTION 212 MALIGNANT 357 BENIGN  
CREATOR DR WILLIAM H WOLBERG W NICK STREET OLVI L MANGASARIAN  
DONOR NICK STREET  
DATE NOVEMBER 1995  
THIS IS A COPY OF UCI ML BREAST CANCER WISCONSIN DIAGNOSTIC DATASETS HTTPSGOOGLEU2UWZ2  
FEATURES ARE COMPUTED FROM A DIGITIZED IMAGE OF A FINE NEEDLE ASPIRATE FNA OF A BREAST MASS THEY DESCRIBE CHARAC  
TERISTICS OF THE CELL NUCLEI PRESENT IN THE IMAGE  
36 DATASET LOADING UTILITIES 633

SCIKITLEARN USER GUIDE RELEASE 0213

SEPARATING PLANE DESCRIBED ABOVE WAS OBTAINED USING MULTISURFACE METHODTREE MSMT K P BENNETT “DECISION TREE CONSTRUCTION VIA LINEAR PROGRAMMING” PROCEEDINGS OF THE 4TH MIDWEST ARTIFICIAL INTELLIGENCE AND COGNITIVE SCIENCE SOCIETY PP 97101 1992 A CLASSIFICATION METHOD WHICH USES LINEAR PROGRAMMING TO CONSTRUCT A DECISION TREE RELEVANT FEATURES WERE SELECTED USING AN EXHAUSTIVE SEARCH IN THE SPACE OF 14 FEATURES AND 13 SEPARATING PLANES THE ACTUAL LINEAR PROGRAM USED TO OBTAIN THE SEPARATING PLANE IN THE 3DIMENSIONAL SPACE IS THAT DESCRIBED IN K P BENNETT AND O L MANGASARIAN “ROBUST LINEAR PROGRAMMING DISCRIMINATION OF TWO LINEARLY INSEPARABLE SETS” OPTIMIZATION METHODS AND SOFTWARE 1 1992 2334

THIS DATABASE IS ALSO AVAILABLE THROUGH THE UW CS FTP SERVER  
FTP FTPCSWISCEDU CD MATHPROGCPDATASETMACHINELEARNWDBC

REFERENCES

- WN STREET WH WOLBERG AND OL MANGASARIAN NUCLEAR FEATURE EXTRACTION FOR BREAST TUMOR DIAGNOSIS ISTSPIE 1993 INTERNATIONAL SYMPOSIUM ON ELECTRONIC IMAGING SCIENCE AND TECHNOLOGY VOLUME 1905 PAGES 861870 SAN JOSE CA 1993
- OL MANGASARIAN WN STREET AND WH WOLBERG BREAST CANCER DIAGNOSIS AND PROGNOSIS VIA LINEAR PROGRAMMING OPERATIONS RESEARCH 434 PAGES 570577 JULYAUGUST 1995
- WH WOLBERG WN STREET AND OL MANGASARIAN MACHINE LEARNING TECHNIQUES TO DIAGNOSE BREAST CANCER FROM FINENEEDLE ASPIRATES CANCER LETTERS 77 1994 163171

363 REAL WORLD DATASETS

SCIKITLEARN PROVIDES TOOLS TO LOAD LARGER DATASETS DOWNLOADING THEM IF NECESSARY  
THEY CAN BE LOADED USING THE FOLLOWING FUNCTIONS

FETCHOLIVETTIFACES DATAHOME SHUFFLE    LOAD THE OLIVETTI FACES DATASET FROM ATT CLASSIFICATION

FETCH20NEWSGROUPS DATAHOME SUBSET    LOAD THE FILENAMES AND DATA FROM THE 20 NEWSGROUPS DATASET CLASSIFICATION

FETCH20NEWSGROUPSVECTORIZED SUBSET    LOAD THE 20 NEWSGROUPS DATASET AND VECTORIZE IT INTO TOKEN COUNTS CLASSIFICATION

FETCHLFWPEOPLE DATAHOME FUNNELED    LOAD THE LABELED FACES IN THE WILD LFW PEOPLE DATASET CLASSIFICATION

FETCHLFWPAIRS SUBSET DATAHOME    LOAD THE LABELED FACES IN THE WILD LFW PAIRS DATASET CLASSIFICATION

FETCHCOVTYPE DATAHOME    LOAD THE COVERTYPE DATASET CLASSIFICATION

FETCHRCV1 DATAHOME SUBSET    LOAD THE RCV1 MULTILABEL DATASET CLASSIFICATION

FETCHKDDCUP99 SUBSET DATAHOME SHUFFLE    LOAD THE KDDCUP99 DATASET CLASSIFICATION

FETCHCALIFORNIAHOUSING DATAHOME    LOAD THE CALIFORNIA HOUSING DATASET REGRESSION

THE OLIVETTI FACES DATASET

THIS DATASET CONTAINS A SET OF FACE IMAGES TAKEN BETWEEN APRIL 1992 AND APRIL 1994 AT ATT LABORATORIES CAM BRIDGE THE SKLEARNDATASETSFETCHOLIVETTIFACES FUNCTION IS THE DATA FETCHING CACHING FUNCTION THAT DOWNLOADS THE DATA ARCHIVE FROM ATT AS DESCRIBED ON THE ORIGINAL WEBSITE

THERE ARE TEN DIFFERENT IMAGES OF EACH OF 40 DISTINCT SUBJECTS FOR SOME SUBJECTS THE IMAGES WERE TAKEN

634 CHAPTER 3 USER GUIDE



SCIKITLEARN USER GUIDE RELEASE 0213

AT DIFFERENT TIMES VARYING THE LIGHTING FACIAL EXPRESSIONS OPEN CLOSED EYES SMILING NOT SMILING AND FACIAL DETAILS GLASSES NO GLASSES ALL THE IMAGES WERE TAKEN AGAINST A DARK HOMOGENEOUS BACKGROUND WITH THE SUBJECTS IN AN UPRIGHT FRONTAL POSITION WITH TOLERANCE FOR SOME SIDE MOVEMENT

DATA SET CHARACTERISTICS

CLASSES 40

SAMPLES TOTAL 400

DIMENSIONALITY 4096

FEATURES REAL BETWEEN 0 AND 1

THE IMAGE IS QUANTIZED TO 256 GREY LEVELS AND STORED AS UNSIGNED 8BIT INTEGERS THE LOADER WILL CONVERT THESE TO FLOATING POINT VALUES ON THE INTERVAL 0 1 WHICH ARE EASIER TO WORK WITH FOR MANY ALGORITHMS

THE “TARGET” FOR THIS DATABASE IS AN INTEGER FROM 0 TO 39 INDICATING THE IDENTITY OF THE PERSON PICTURED HOWEVER WITH ONLY 10 EXAMPLES PER CLASS THIS RELATIVELY SMALL DATASET IS MORE INTERESTING FROM AN UNSUPERVISED OR SEMISUPERVISED PERSPECTIVE

THE ORIGINAL DATASET CONSISTED OF 92 X 112 WHILE THE VERSION AVAILABLE HERE CONSISTS OF 64X64 IMAGES WHEN USING THESE IMAGES PLEASE GIVE CREDIT TO ATT LABORATORIES CAMBRIDGE

THE 20 NEWSGROUPS TEXT DATASET

THE 20 NEWSGROUPS DATASET COMPRISES AROUND 18000 NEWSGROUPS POSTS ON 20 TOPICS SPLIT IN TWO SUBSETS ONE FOR TRAINING OR DEVELOPMENT AND THE OTHER ONE FOR TESTING OR FOR PERFORMANCE EVALUATION THE SPLIT BETWEEN THE TRAIN AND TEST SET IS BASED UPON A MESSAGES POSTED BEFORE AND AFTER A SPECIFIC DATE

THIS MODULE CONTAINS TWO LOADERS THE FIRST ONE SKLEARNDATASETSFETCH20NEWSGROUPS RETURNS A LIST OF THE RAW TEXTS THAT CAN BE FED TO TEXT FEATURE EXTRACTORS SUCH AS SKLEARNFEATUREEXTRACTIONTEXT COUNTVECTORIZER WITH CUSTOM PARAMETERS SO AS TO EXTRACT FEATURE VECTORS THE SECOND ONE SKLEARN DATASETSFETCH20NEWSGROUPSVECTORIZED RETURNS READYTOUSE FEATURES IE IT IS NOT NECESSARY TO USE A FEATURE EXTRACTOR

DATA SET CHARACTERISTICS

CLASSES 20

SAMPLES TOTAL 18846

DIMENSIONALITY 1

FEATURES TEXT

USAGE

THESKLEARNDATASETSFETCH20NEWSGROUPS FUNCTION IS A DATA FETCHING CACHING FUNCTIONS THAT DOWN LOADS THE DATA ARCHIVE FROM THE ORIGINAL 20 NEWSGROUPS WEBSITE EXTRACTS THE ARCHIVE CONTENTS IN THE SCIKITLEARNDATA20NEWSHOME FOLDER AND CALLS THE SKLEARNDATASETSLOADFILES ON EITHER THE TRAINING OR TESTING SET FOLDER OR BOTH OF THEM

FROM SKLEARNDATASETS IMPORT FETCH20NEWSGROUPS

NEWSGROUPSTRAIN FETCH20NEWSGROUPSSUBSETTRAIN

FROM PPRINT IMPORT PPRINT

PPRINTLISTNEWSGROUPSTRAINTARGETNAMES

ALTATHEISM

COMPGRAPHICS

36 DATASET LOADING UTILITIES 635

SCIKITLEARN USER GUIDE RELEASE 0213  
COMPOSMSWINDOWSMISC  
COMPSYSIBMPCHARDWARE  
COMPSYSMACHARDWARE  
COMPWINDOWSX  
MISCFORSALE  
RECAUTOS  
RECMOTORCYCLES  
RECSPORTBASEBALL  
RECSPORTHOCKEY  
SCICRYPT  
SCIELECTRONICS  
SCIMED  
SCISPACE  
SOCRELIGIONCHRISTIAN  
TALKPOLITICSGUNS  
TALKPOLITICSMIDEAST  
TALKPOLITICSMISC  
TALKRELIGIONMISC  
THE REAL DATA LIES IN THE FILENAMES ANDTARGET ATTRIBUTES THE TARGET ATTRIBUTE IS THE INTEGER INDEX OF THE CATEGORY  
NEWSGROUPSTRAINFILENAMESSHAPE  
11314  
NEWSGROUPSTRAINTARGETSHAPE  
11314  
NEWSGROUPSTRAINTARGET10  
ARRAY 7 4 4 1 14 16 13 3 2 4  
IT IS POSSIBLE TO LOAD ONLY A SUBSELECTION OF THE CATEGORIES BY PASSING THE LIST OF THE CATEGORIES TO LOAD TO THE SKLEARN  
DATASETSFETCH20NEWSGROUPS FUNCTION  
CATS ALTATHEISM SCISPACE  
NEWSGROUPSTRAIN FETCH20NEWSGROUPSSUBSETTRAIN CATEGORIESCATS  
LISTNEWSGROUPSTRAINTARGETNAMES  
ALTATHEISM SCISPACE  
NEWSGROUPSTRAINFILENAMESSHAPE  
1073  
NEWSGROUPSTRAINTARGETSHAPE  
1073  
NEWSGROUPSTRAINTARGET10  
ARRAY0 1 1 1 0 1 1 0 0 0  
CONVERTING TEXT TO VECTORS  
IN ORDER TO FEED PREDICTIVE OR CLUSTERING MODELS WITH THE TEXT DATA ONE FIRST NEED TO TURN THE TEXT INTO VECTORS  
OF NUMERICAL VALUES SUITABLE FOR STATISTICAL ANALYSIS THIS CAN BE ACHIEVED WITH THE UTILITIES OF THE SKLEARN  
FEATUREEXTRACTIONTEXT AS DEMONSTRATED IN THE FOLLOWING EXAMPLE THAT EXTRACT TFIDF VECTORS OF UNIGRAM  
TOKENS FROM A SUBSET OF 20NEWS  
FROM SKLEARNFEATUREEXTRACTIONTEXT IMPORT TFIDFVECTORIZER  
CATEGORIES ALTATHEISM TALKRELIGIONMISC  
COMPGRAPHICS SCISPACE  
NEWSGROUPSTRAIN FETCH20NEWSGROUPSSUBSETTRAIN  
CATEGORIESCATEGORIES  
VECTORIZER TFIDFVECTORIZER  
636 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

VECTORS VECTORIZERFITTRANSFORMNEWSGROUPSTRAINDATA  
VECTORSSHAPE  
2034 34118

THE EXTRACTED TFIDF VECTORS ARE VERY SPARSE WITH AN AVERAGE OF 159 NONZERO COMPONENTS BY SAMPLE IN A MORE THAN 30000DIMENSIONAL SPACE LESS THAN 5 NONZERO FEATURES

VECTORSNNZ FLOATVECTORSSHAPE0  
15901327

SKLEARNDATASETSFETCH20NEWSGROUPSVECTORIZED IS A FUNCTION WHICH RETURNS READYTOUSE TOKEN  
COUNTS FEATURES INSTEAD OF FILE NAMES  
FILTERING TEXT FOR MORE REALISTIC TRAINING

IT IS EASY FOR A CLASSIFIER TO OVERFIT ON PARTICULAR THINGS THAT APPEAR IN THE 20 NEWSGROUPS DATA SUCH AS NEWSGROUP HEADERS MANY CLASSIFIERS ACHIEVE VERY HIGH FSCORES BUT THEIR RESULTS WOULD NOT GENERALIZE TO OTHER DOCUMENTS THAT AREN'T FROM THIS WINDOW OF TIME

FOR EXAMPLE LET’S LOOK AT THE RESULTS OF A MULTINOMIAL NAIVE BAYES CLASSIFIER WHICH IS FAST TO TRAIN AND ACHIEVES A DECENT FSCORE

```
FROM SKLEARNNAIVEBAYES IMPORT MULTINOMIALNB
FROM SKLEARN IMPORT METRICS
NEWSGROUPSTEST FETCH20NEWSGROUPSSUBSETTEST
CATEGORIESCATEGORIES
VECTORSTEST VECTORIZERTRANSFORMNEWSGROUPSTESTDATA
CLF MULTINOMIALNBALPHA01
CLFFITVECTORS NEWSGROUPSTRAINTARGET
MULTINOMIALNBALPHA001 CLASSPRIORNONE FITPRIORTRUE
PRED CLFPREDICTVECTORSTEST
METRICSF1SCORENEWSGROUPSTESTTTARGET PRED AVERAGEMACRO
088213
```

THE EXAMPLE CLASSIFICATION OF TEXT DOCUMENTS USING SPARSE FEATURES SHUFFLES THE TRAINING AND TEST DATA INSTEAD OF SEGMENTING BY TIME AND IN THAT CASE MULTINOMIAL NAIVE BAYES GETS A MUCH HIGHER FSCORE OF 088 ARE YOU SUSPICIOUS YET OF WHAT’S GOING ON INSIDE THIS CLASSIFIER

LET’S TAKE A LOOK AT WHAT THE MOST INFORMATIVE FEATURES ARE

```
IMPORT NUMPY AS NP
DEF SHOWTOP10CLASSIFIER VECTORIZER CATEGORIES
FEATURENAMES NPASARRAYVECTORIZERGETFEATURENAMES
FOR I CATEGORY INENUMERATECATEGORIES
TOP10 NPARGSORTCLASSIFIERCOEFI10
PRINTSS CATEGORY JOINFEATURENAMESTOP10
```

SHOWTOP10CLF VECTORIZER NEWSGROUPSTRAINTARGETNAMES  
ALTATHEISM EDU IT AND IN YOU THAT IS OF TO THE  
COMPGRAPHICS EDU IN GRAPHICS IT IS FOR AND OF TO THE  
SCISPACE EDU IT THAT IS IN AND SPACE TO OF THE  
TALKRELIGIONMISC NOT IT YOU IN IS THAT AND TO OF THE  
YOU CAN NOW SEE MANY THINGS THAT THESE FEATURES HAVE OVERFIT TO

36 DATASET LOADING UTILITIES 637

SCIKITLEARN USER GUIDE RELEASE 0213

- ALMOST EVERY GROUP IS DISTINGUISHED BY WHETHER HEADERS SUCH AS NNTPPOSTINGHOST AND DISTRIBUTION APPEAR MORE OR LESS OFTEN
- ANOTHER SIGNIFICANT FEATURE INVOLVES WHETHER THE SENDER IS AFFILIATED WITH A UNIVERSITY AS INDICATED EITHER BY THEIR HEADERS OR THEIR SIGNATURE
- THE WORD “ARTICLE” IS A SIGNIFICANT FEATURE BASED ON HOW OFTEN PEOPLE QUOTE PREVIOUS POSTS LIKE THIS “IN ARTICLE ARTICLE ID NAME EMAIL ADDRESS WROTE”
- OTHER FEATURES MATCH THE NAMES AND EMAIL ADDRESSES OF PARTICULAR PEOPLE WHO WERE POSTING AT THE TIME

WITH SUCH AN ABUNDANCE OF CLUES THAT DISTINGUISH NEWSGROUPS THE CLASSIFIERS BARELY HAVE TO IDENTIFY TOPICS FROM TEXT AT ALL AND THEY ALL PERFORM AT THE SAME HIGH LEVEL

FOR THIS REASON THE FUNCTIONS THAT LOAD 20 NEWSGROUPS DATA PROVIDE A PARAMETER CALLED REMOVE TELLING IT WHAT KINDS OF INFORMATION TO STRIP OUT OF EACH FILE REMOVE SHOULD BE A TUPLE CONTAINING ANY SUBSET OF HEADERS FOOTERS QUOTES TELLING IT TO REMOVE HEADERS SIGNATURE BLOCKS AND QUOTATION BLOCKS RESPECTIVELY

```
NEWSGROUPSTEST  FETCH20NEWSGROUPSSUBSETTEST
REMOVEHEADERS FOOTERS QUOTES
CATEGORIESCATEGORIES
VECTORSTEST  VECTORIZERTRANSFORMNEWSGROUPSTESTDATA
PRED  CLFPREDICTVECTORSTEST
METRICSF1SCOREPRED NEWSGROUPSTESTTARGET AVERAGEMACRO
077310
```

THIS CLASSIFIER LOST OVER A LOT OF ITS FSCORE JUST BECAUSE WE REMOVED METADATA THAT HAS LITTLE TO DO WITH TOPIC CLASSIFICATION IT LOSES EVEN MORE IF WE ALSO STRIP THIS METADATA FROM THE TRAINING DATA

```
NEWSGROUPSTRAIN  FETCH20NEWSGROUPSSUBSETTRAIN
REMOVEHEADERS FOOTERS QUOTES
CATEGORIESCATEGORIES
VECTORS  VECTORIZERFITTRANSFORMNEWSGROUPSTRAINDATA
CLF  MULTINOMIALNBALPHA01
CLFFITVECTORS NEWSGROUPSTRAINTARGET
MULTINOMIALNBALPHA001 CLASSPRIORNONE FITPRIORTRUE
VECTORSTEST  VECTORIZERTRANSFORMNEWSGROUPSTESTDATA
PRED  CLFPREDICTVECTORSTEST
METRICSF1SCORENEWSGROUPSTESTTARGET PRED AVERAGEMACRO
076995
```

SOME OTHER CLASSIFIERS COPE BETTER WITH THIS HARDER VERSION OF THE TASK TRY RUNNING SAMPLE PIPELINE FOR TEXT FEATURE EXTRACTION AND EVALUATION WITH AND WITHOUT THE FILTER OPTION TO COMPARE THE RESULTS

RECOMMENDATION

WHEN EVALUATING TEXT CLASSIFIERS ON THE 20 NEWSGROUPS DATA YOU SHOULD STRIP NEWSGROUPRELATED METADATA IN SCIKITLEARN YOU CAN DO THIS BY SETTING REMOVEHEADERS FOOTERS QUOTES THE FSCORE WILL BE LOWER BECAUSE IT IS MORE REALISTIC

EXAMPLES

- SAMPLE PIPELINE FOR TEXT FEATURE EXTRACTION AND EVALUATION
- CLASSIFICATION OF TEXT DOCUMENTS USING SPARSE FEATURES

638 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

THE LABELED FACES IN THE WILD FACE RECOGNITION DATASET

THIS DATASET IS A COLLECTION OF JPEG PICTURES OF FAMOUS PEOPLE COLLECTED OVER THE INTERNET ALL DETAILS ARE AVAILABLE ON THE OFFICIAL WEBSITE

HTTPVISWWWCSUMASSEDULFW

EACH PICTURE IS CENTERED ON A SINGLE FACE THE TYPICAL TASK IS CALLED FACE VERIFICATION GIVEN A PAIR OF TWO PICTURES A BINARY CLASSIFIER MUST PREDICT WHETHER THE TWO IMAGES ARE FROM THE SAME PERSON

AN ALTERNATIVE TASK FACE RECOGNITION OR FACE IDENTIFICATION IS GIVEN THE PICTURE OF THE FACE OF AN UNKNOWN PERSON IDENTIFY THE NAME OF THE PERSON BY REFERRING TO A GALLERY OF PREVIOUSLY SEEN PICTURES OF IDENTIFIED PERSONS

BOTH FACE VERIFICATION AND FACE RECOGNITION ARE TASKS THAT ARE TYPICALLY PERFORMED ON THE OUTPUT OF A MODEL TRAINED TO PERFORM FACE DETECTION THE MOST POPULAR MODEL FOR FACE DETECTION IS CALLED VIOLA JONES AND IS IMPLEMENTED IN THE OPENCV LIBRARY THE LFW FACES WERE EXTRACTED BY THIS FACE DETECTOR FROM VARIOUS ONLINE WEBSITES

DATA SET CHARACTERISTICS

CLASSES 5749

SAMPLES TOTAL 13233

DIMENSIONALITY 5828

FEATURES REAL BETWEEN 0 AND 255

USAGE

SCIKITLEARN PROVIDES TWO LOADERS THAT WILL AUTOMATICALLY DOWNLOAD CACHE PARSE THE METADATA FILES DECODE THE JPEG AND CONVERT THE INTERESTING SLICES INTO MEMMAPPED NUMPY ARRAYS THIS DATASET SIZE IS MORE THAN 200 MB THE FIRST LOAD TYPICALLY TAKES MORE THAN A COUPLE OF MINUTES TO FULLY DECODE THE RELEVANT PART OF THE JPEG FILES INTO NUMPY ARRAYS IF THE DATASET HAS BEEN LOADED ONCE THE FOLLOWING TIMES THE LOADING TIMES LESS THAN 200MS BY USING A MEMMAPPED VERSION MEMOIZED ON THE DISK IN THE SCIKITLEARN DATALFWHOME FOLDER USING JOBLIB

THE FIRST LOADER IS USED FOR THE FACE IDENTIFICATION TASK A MULTICLASS CLASSIFICATION TASK HENCE SUPERVISED LEARNING

```
from sklearn.datasets import fetch_lfw_people
lfw_people = fetch_lfw_people(min_faces_per_person=70, resize=0.4)

# Names of all 1000 people
names = lfw_people.target_names
print(names)
```

ARIEL SHARON

COLIN POWELL

DONALD RUMSFELD

GEORGE W BUSH

GERHARD SCHROEDER

HUGO CHAVEZ

TONY BLAIR

THE DEFAULT SLICE IS A RECTANGULAR SHAPE AROUND THE FACE REMOVING MOST OF THE BACKGROUND

lfw\_people.data.dtype

lfw\_people.data.shape

1288 1850

36 DATASET LOADING UTILITIES 639

SCIKITLEARN USER GUIDE RELEASE 0213

FWPEOPLEIMAGESHAPE  
1288 50 37

EACH OF THE 1140 FACES IS ASSIGNED TO A SINGLE PERSON ID IN THE TARGET ARRAY

FWPEOPLETARGETSHAPE  
1288

LISTFWPEOPLETARGET10  
5 6 3 1 0 1 3 4 3 0

THE SECOND LOADER IS TYPICALLY USED FOR THE FACE VERIFICATION TASK EACH SAMPLE IS A PAIR OF TWO PICTURE BELONGING OR NOT TO THE SAME PERSON

FROM SKLEARNDATASETS IMPORT FETCHLFWPAIRS

LFWPAIRSTRAIN FETCHLFWPAIRSSUBSETTRAIN

LISTLFWPAIRSTRAINTARGETNAMES

DIFFERENT PERSONS SAME PERSON

LFWPAIRSTRAINPAIRSSHAPE  
2200 2 62 47

LFWPAIRSTRAINDATASHAPE  
2200 5828

LFWPAIRSTRAINTARGETSHAPE  
2200

BOTH FOR THE SKLEARNDATASETSFETCHLFWPEOPLE ANDSKLEARNDATASETSFETCHLFWPAIRS

FUNCTION IT IS POSSIBLE TO GET AN ADDITIONAL DIMENSION WITH THE RGB COLOR CHANNELS BY PASSING COLORTRUE IN THAT CASE THE SHAPE WILL BE 2200 2 62 47 3

THESKLEARNDATASETSFETCHLFWPAIRS DATASETS IS SUBDIVIDED INTO 3 SUBSETS THE DEVELOPMENT TRAIN SET THE DEVELOPMENT TEST SET AND AN EVALUATION 10FOLDS SET MEANT TO COMPUTE PERFORMANCE METRICS USING A 10FOLDS CROSS VALIDATION SCHEME

REFERENCES

- LABELED FACES IN THE WILD A DATABASE FOR STUDYING FACE RECOGNITION IN UNCONSTRAINED ENVIRONMENTS GARY B HUANG MANU RAMESH TAMARA BERG AND ERIK LEARNEDMILLER UNIVERSITY OF MASSACHUSETTS AMHERST

TECHNICAL REPORT 0749 OCTOBER 2007

EXAMPLES

FACES RECOGNITION EXAMPLE USING EIGENFACES AND SVMs

FOREST COVERTYPES

THE SAMPLES IN THIS DATASET CORRESPOND TO 30×30M PATCHES OF FOREST IN THE US COLLECTED FOR THE TASK OF PREDICTING EACH PATCH'S COVER TYPE IE THE DOMINANT SPECIES OF TREE THERE ARE SEVEN COVERTYPES MAKING THIS A MULTICLASS CLASSIFICATION PROBLEM EACH SAMPLE HAS 54 FEATURES DESCRIBED ON THE DATASET'S HOMEPAGE SOME OF THE FEATURES ARE BOOLEAN INDICATORS WHILE OTHERS ARE DISCRETE OR CONTINUOUS MEASUREMENTS

640 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

DATA SET CHARACTERISTICS

CLASSES 7

SAMPLES TOTAL 581012

DIMENSIONALITY 54

FEATURES INT

SKLEARNDATASETSFETCHCOVTYPE WILL LOAD THE COVERTYPE DATASET IT RETURNS A DICTIONARYLIKE OBJECT WITH THE FEATURE MATRIX IN THE DATA MEMBER AND THE TARGET VALUES IN TARGET THE DATASET WILL BE DOWNLOADED FROM THE WEB IF NECESSARY

RCV1 DATASET

REUTERS CORPUS V OLUME I RCV1 IS AN ARCHIVE OF OVER 800000 MANUALLY CATEGORIZED NEWSWIRE STORIES MADE AVAILABLE BY REUTERS LTD FOR RESEARCH PURPOSES THE DATASET IS EXTENSIVELY DESCRIBED IN1

DATA SET CHARACTERISTICS

CLASSES 103

SAMPLES TOTAL 804414

DIMENSIONALITY 47236

FEATURES REAL BETWEEN 0 AND 1

SKLEARNDATASETSFETCHRCV1 WILL LOAD THE FOLLOWING VERSION RCV1V2 VECTORS FULL SETS TOPICS MULTILABELS

FROM SKLEARNDATASETS IMPORT FETCHRCV1

RCV1 FETCHRCV1

IT RETURNS A DICTIONARYLIKE OBJECT WITH THE FOLLOWING ATTRIBUTES

DATA THE FEATURE MATRIX IS A SCIPY CSR SPARSE MATRIX WITH 804414 SAMPLES AND 47236 FEATURES NONZERO VALUES CONTAINS COSINENORMALIZED LOG TFIDF VECTORS A NEARLY CHRONOLOGICAL SPLIT IS PROPOSED IN1 THE FIRST 23149 SAMPLES ARE THE TRAINING SET THE LAST 781265 SAMPLES ARE THE TESTING SET THIS FOLLOWS THE OFFICIAL LYRL2004 CHRONOLOGICAL SPLIT THE ARRAY HAS 016 OF NON ZERO VALUES

RCV1DATASHAPE

804414 47236

TARGET THE TARGET VALUES ARE STORED IN A SCIPY CSR SPARSE MATRIX WITH 804414 SAMPLES AND 103 CATEGORIES EACH SAMPLE HAS A VALUE OF 1 IN ITS CATEGORIES AND 0 IN OTHERS THE ARRAY HAS 315 OF NON ZERO VALUES

RCV1TARGETSHAPE

804414 103

SAMPLEID EACH SAMPLE CAN BE IDENTIFIED BY ITS ID RANGING WITH GAPS FROM 2286 TO 810596

RCV1SAMPLEID3

ARRAY2286 2287 2288 DTYPEUINT32

1LEWIS D D YANG Y ROSE T G LI F 2004 RCV1 A NEW BENCHMARK COLLECTION FOR TEXT CATEGORIZATION RESEARCH THE JOURNAL OF MACHINE LEARNING RESEARCH 5 361397

36 DATASET LOADING UTILITIES 641

SCIKITLEARN USER GUIDE RELEASE 0213

TARGETNAMES THE TARGET VALUES ARE THE TOPICS OF EACH SAMPLE EACH SAMPLE BELONGS TO AT LEAST ONE TOPIC AND TO UP TO 17 TOPICS THERE ARE 103 TOPICS EACH REPRESENTED BY A STRING THEIR CORPUS FREQUENCIES SPAN FIVE ORDERS OF MAGNITUDE FROM 5 OCCURRENCES FOR 'GMIL' TO 381327 FOR 'CCAT'

RCV1TARGETNAMES3TOLIST

E11 ECAT M11

THE DATASET WILL BE DOWNLOADED FROM THE RCV1 HOMEPAGE IF NECESSARY THE COMPRESSED SIZE IS ABOUT 656 MB

REFERENCES

KDDCUP 99 DATASET

THE KDD CUP '99 DATASET WAS CREATED BY PROCESSING THE TCPDUMP PORTIONS OF THE 1998 DARPA INTRUSION DETECTION SYSTEM IDS EVALUATION DATASET CREATED BY MIT LINCOLN LAB 1 THE ARTIFICIAL DATA DESCRIBED ON THE DATASET'S HOMEPAGE WAS GENERATED USING A CLOSED NETWORK AND HANDINJECTED ATTACKS TO PRODUCE A LARGE NUMBER OF DIFFERENT TYPES OF ATTACK WITH NORMAL ACTIVITY IN THE BACKGROUND AS THE INITIAL GOAL WAS TO PRODUCE A LARGE TRAINING SET FOR SUPERVISED LEARNING ALGORITHMS THERE IS A LARGE PROPORTION 80% OF ABNORMAL DATA WHICH IS UNREALISTIC IN REAL WORLD AND INAPPROPRIATE FOR UNSUPERVISED ANOMALY DETECTION WHICH AIMS AT DETECTING 'ABNORMAL' DATA IE

1 QUALITATIVELY DIFFERENT FROM NORMAL DATA

2 IN LARGE MINORITY AMONG THE OBSERVATIONS

WE THUS TRANSFORM THE KDD DATA SET INTO TWO DIFFERENT DATA SETS SA AND SF

SA IS OBTAINED BY SIMPLY SELECTING ALL THE NORMAL DATA AND A SMALL PROPORTION OF ABNORMAL DATA TO GIVES AN ANOMALY PROPORTION OF 1

SF IS OBTAINED AS IN 2 BY SIMPLY PICKING UP THE DATA WHOSE ATTRIBUTE LOGGEDIN IS POSITIVE THUS FOCUSING ON THE INTRUSION ATTACK WHICH GIVES A PROPORTION OF 0.3 OF ATTACK

HTTP AND SMTP ARE TWO SUBSETS OF SF CORRESPONDING WITH THIRD FEATURE EQUAL TO 'HTTP' RESP TO 'SMTP'

GENERAL KDD STRUCTURE

SAMPLES TOTAL 4898431

DIMENSIONALITY 41

FEATURES DISCRETE INT OR CONTINUOUS FLOAT

TARGETS STR 'NORMAL' OR NAME OF THE ANOMALY TYPE

SA STRUCTURE

SAMPLES TOTAL 976158

DIMENSIONALITY 41

FEATURES DISCRETE INT OR CONTINUOUS FLOAT

TARGETS STR 'NORMAL' OR NAME OF THE ANOMALY TYPE

SF STRUCTURE

SAMPLES TOTAL 699691

DIMENSIONALITY 4

FEATURES DISCRETE INT OR CONTINUOUS FLOAT

TARGETS STR 'NORMAL' OR NAME OF THE ANOMALY TYPE

642 CHAPTER 3 USER GUIDE



SCIKITLEARN USER GUIDE RELEASE 0213

HTTP STRUCTURE

SAMPLES TOTAL 619052

DIMENSIONALITY 3

FEATURES DISCRETE INT OR CONTINUOUS FLOAT

TARGETS STR 'NORMAL' OR NAME OF THE ANOMALY TYPE

SMTP STRUCTURE

SAMPLES TOTAL 95373

DIMENSIONALITY 3

FEATURES DISCRETE INT OR CONTINUOUS FLOAT

TARGETS STR 'NORMAL' OR NAME OF THE ANOMALY TYPE

SKLEARNDATASETSFETCHKDDCUP99 WILL LOAD THE KDDCUP99 DATASET IT RETURNS A DICTIONARYLIKE OBJECT WITH THE FEATURE MATRIX IN THE DATA MEMBER AND THE TARGET VALUES IN TARGET THE DATASET WILL BE DOWNLOADED FROM THE WEB IF NECESSARY

CALIFORNIA HOUSING DATASET

DATA SET CHARACTERISTICS

NUMBER OF INSTANCES 20640

NUMBER OF ATTRIBUTES 8 NUMERIC PREDICTIVE ATTRIBUTES AND THE TARGET

ATTRIBUTE INFORMATION

- MEDINC MEDIAN INCOME IN BLOCK
- HOUSEAGE MEDIAN HOUSE AGE IN BLOCK
- AVEROOMS AVERAGE NUMBER OF ROOMS
- AVEBEDRMS AVERAGE NUMBER OF BEDROOMS
- POPULATION BLOCK POPULATION
- AVEOCCUP AVERAGE HOUSE OCCUPANCY
- LATITUDE HOUSE BLOCK LATITUDE
- LONGITUDE HOUSE BLOCK LONGITUDE

MISSING ATTRIBUTE VALUES NONE

THIS DATASET WAS OBTAINED FROM THE STATLIB REPOSITORY HTTPLIBSTATCMUEDUDATASETS

THE TARGET VARIABLE IS THE MEDIAN HOUSE VALUE FOR CALIFORNIA DISTRICTS

THIS DATASET WAS DERIVED FROM THE 1990 US CENSUS USING ONE ROW PER CENSUS BLOCK GROUP A BLOCK GROUP IS THE SMALLEST GEOGRAPHICAL UNIT FOR WHICH THE US CENSUS BUREAU PUBLISHES SAMPLE DATA A BLOCK GROUP TYPICALLY HAS A POPULATION OF 600 TO 3000 PEOPLE

IT CAN BE DOWNLOADEDLOADED USING THE SKLEARNDATASETSFETCHCALIFORNIAHOUSING FUNCTION

REFERENCES

36 DATASET LOADING UTILITIES 643

SCIKITLEARN USER GUIDE RELEASE 0213

• PACE R KELLEY AND RONALD BARRY SPARSE SPATIAL AUTOREGRESSIONS STATISTICS AND PROBABILITY LETTERS 33  
1997 291297

364 GENERATED DATASETS

IN ADDITION SCIKITLEARN INCLUDES VARIOUS RANDOM SAMPLE GENERATORS THAT CAN BE USED TO BUILD ARTIFICIAL DATASETS OF CONTROLLED SIZE AND COMPLEXITY

GENERATORS FOR CLASSIFICATION AND CLUSTERING

THESE GENERATORS PRODUCE A MATRIX OF FEATURES AND CORRESPONDING DISCRETE TARGETS

SINGLE LABEL

BOTHMAKEBLOBS ANDMAKECLASSIFICATION CREATE MULTICLASS DATASETS BY ALLOCATING EACH CLASS ONE OR MORE NORMALLYDISTRIBUTED CLUSTERS OF POINTS MAKEBLOBS PROVIDES GREATER CONTROL REGARDING THE CENTERS AND STANDARD DEVIATIONS OF EACH CLUSTER AND IS USED TO DEMONSTRATE CLUSTERING MAKECLASSIFICATION SPECIALISES IN INTRODUCING NOISE BY WAY OF CORRELATED REDUNDANT AND UNINFORMATIVE FEATURES MULTIPLE GAUSSIAN CLUSTERS PER CLASS AND LINEAR TRANSFORMATIONS OF THE FEATURE SPACE

MAKEGAUSSIANQUANTILES DIVIDES A SINGLE GAUSSIAN CLUSTER INTO NEAREQUALSIZE CLASSES SEPARATED BY CONCENTRIC HYPERSPHERES MAKEHASTIE102 GENERATES A SIMILAR BINARY 10DIMENSIONAL PROBLEM

MAKECIRCLES ANDMAKEMOONS GENER

ATE 2D BINARY CLASSIFICATION DATASETS THAT ARE CHALLENGING TO CERTAIN ALGORITHMS EG CENTROIDBASED CLUSTERING OR LINEAR CLASSIFICATION INCLUDING OPTIONAL GAUSSIAN NOISE THEY ARE USEFUL FOR VISUALISATION MAKECIRCLES PRODUCES GAUSSIAN DATA WITH A SPHERICAL DECISION BOUNDARY FOR BINARY CLASSIFICATION WHILE MAKEMOONS PRODUCES TWO INTERLEAVING HALF CIRCLES

644 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

MULTILABEL

MAKEMULTILABELCLASSIFICATION GENERATES RANDOM SAMPLES WITH MULTIPLE LABELS REFLECTING A BAG OF WORDS DRAWN FROM A MIXTURE OF TOPICS THE NUMBER OF TOPICS FOR EACH DOCUMENT IS DRAWN FROM A POISSON DISTRIBUTION AND THE TOPICS THEMSELVES ARE DRAWN FROM A FIXED RANDOM DISTRIBUTION SIMILARLY THE NUMBER OF WORDS IS DRAWN FROM POISSON WITH WORDS DRAWN FROM A MULTINOMIAL WHERE EACH TOPIC DEFINES A PROBABILITY DISTRIBUTION OVER WORDS SIMPLIFICATIONS WITH RESPECT TO TRUE BAGOFWORDS MIXTURES INCLUDE

- PERTOPIC WORD DISTRIBUTIONS ARE INDEPENDENTLY DRAWN WHERE IN REALITY ALL WOULD BE AFFECTED BY A SPARSE BASE DISTRIBUTION AND WOULD BE CORRELATED
- FOR A DOCUMENT GENERATED FROM MULTIPLE TOPICS ALL TOPICS ARE WEIGHTED EQUALLY IN GENERATING ITS BAG OF WORDS
- DOCUMENTS WITHOUT LABELS WORDS AT RANDOM RATHER THAN FROM A BASE DISTRIBUTION

BICLUSTERING

MAKEBICCLUSTERS SHAPE NCLUSTERS NOISE    GENERATE AN ARRAY WITH CONSTANT BLOCK DIAGONAL STRUCTURE FOR BICLUSTERING

MAKECHECKERBOARD SHAPE NCLUSTERS    GENERATE AN ARRAY WITH BLOCK CHECKERBOARD STRUCTURE FOR BI CLUSTERING

GENERATORS FOR REGRESSION

MAKEREGRESSION PRODUCES REGRESSION TARGETS AS AN OPTIONALLYSPARSE RANDOM LINEAR COMBINATION OF RANDOM FEATURES WITH NOISE ITS INFORMATIVE FEATURES MAY BE UNCORRELATED OR LOW RANK FEW FEATURES ACCOUNT FOR MOST OF THE VARIANCE

OTHER REGRESSION GENERATORS GENERATE FUNCTIONS DETERMINISTICALLY FROM RANDOMIZED FEATURES

MAKESPARSEUNCORRELATED PRODUCES A TARGET AS A LINEAR COMBINATION OF FOUR FEATURES WITH FIXED COEFFICIENTS OTHERS ENCODE EXPLICITLY NONLINEAR RELATIONS MAKEFRIEDMAN1 IS RELATED BY POLYNOMIAL AND SINE TRANSFORMS MAKEFRIEDMAN2 INCLUDES FEATURE MULTIPLICATION AND RECIPROCATATION AND MAKEFRIEDMAN3 IS SIMILAR WITH AN ARCTAN TRANSFORMATION ON THE TARGET

GENERATORS FOR MANIFOLD LEARNING

MAKESCURVE NSAMPLES NOISE RANDOMSTATE GENERATE AN S CURVE DATASET

MAKESWISSROLL NSAMPLES NOISE RANDOMSTATE GENERATE A SWISS ROLL DATASET

36 DATASET LOADING UTILITIES 645

SCIKITLEARN USER GUIDE RELEASE 0213

GENERATORS FOR DECOMPOSITION

MAKELOWRANKMATRIX NSAMPLES      GENERATE A MOSTLY LOW RANK MATRIX WITH BELLSHAPED SINGULAR VALUES

MAKESPARSECODEDSIGNAL NSAMPLES      GENERATE A SIGNAL AS A SPARSE COMBINATION OF DICTIONARY ELEMENTS

MAKESPDMATRIX NNDIM RANDOMSTATE GENERATE A RANDOM SYMMETRIC POSITIVEDEFINITE MATRIX

MAKESPARSESPDMATRIX DIM ALPHA      GENERATE A SPARSE SYMMETRIC DEFINITE POSITIVE MATRIX

365 LOADING OTHER DATASETS

SAMPLE IMAGES

SCIKITLEARN ALSO EMBED A COUPLE OF SAMPLE JPEG IMAGES PUBLISHED UNDER CREATIVE COMMONS LICENSE BY THEIR AUTHORS

THOSE IMAGES CAN BE USEFUL TO TEST ALGORITHMS AND PIPELINE ON 2D DATA

LOADSAMPLEIMAGES    LOAD SAMPLE IMAGES FOR IMAGE MANIPULATION

LOADSAMPLEIMAGE IMAGENAME LOAD THE NUMPY ARRAY OF A SINGLE SAMPLE IMAGE

WARNING THE DEFAULT CODING OF IMAGES IS BASED ON THE UINT8 DTYPE TO SPARE MEMORY OFTEN MACHINE LEARNING ALGORITHMS WORK BEST IF THE INPUT IS CONVERTED TO A FLOATING POINT REPRESENTATION FIRST ALSO IF YOU PLAN TO USE MATPLOTLIBPYPLTSHOW DON'T FORGET TO SCALE TO THE RANGE 0 1 AS DONE IN THE FOLLOWING EXAMPLE

EXAMPLES

- COLOR QUANTIZATION USING KMEANS

DATASETS IN SVMLIGHT   LIBSVM FORMAT

SCIKITLEARN INCLUDES UTILITY FUNCTIONS FOR LOADING DATASETS IN THE SVMLIGHT   LIBSVM FORMAT IN THIS FORMAT EACH LINE TAKES THE FORM LABEL FEATUREIDFEATUREVALUE FEATUREIDFEATUREVALUE

THIS FORMAT IS ESPECIALLY SUITABLE FOR SPARSE DATASETS IN THIS MODULE SCIPY SPARSE CSR MATRICES ARE USED FOR X AND NUMPY ARRAYS ARE USED FOR Y

YOU MAY LOAD A DATASET LIKE AS FOLLOWS

FROM SKLEARN DATASETS IMPORT LOADSVMLIGHTFILE

XTRAIN YTRAIN    LOADSVMLIGHTFILEPATHTOTRAIN DATASETTXT

SCIKITLEARN USER GUIDE RELEASE 0213  
YOU MAY ALSO LOAD TWO OR MORE DATASETS AT ONCE  
XTRAIN YTRAIN XTEST YTEST LOADSVMLIGHTFILES  
PATHTOTRAINDATASETTXT PATHTOTESTDATASETTXT

IN THIS CASE XTRAIN ANDXTEST ARE GUARANTEED TO HAVE THE SAME NUMBER OF FEATURES ANOTHER WAY TO ACHIEVE THE  
SAME RESULT IS TO FIX THE NUMBER OF FEATURES  
XTEST YTEST LOADSVMLIGHTFILE  
PATHTOTESTDATASETTXT NFEATURESXTRAINSHAPE1

RELATED LINKS  
PUBLIC DATASETS IN SVMLIGHT LIBSVM FORMAT HTTPSWWWCSIENTUEDUTWCJLINLIBSVMTOOLS DATASETS  
FASTER APICOMPATIBLE IMPLEMENTATION HTTPSGITHUBCOMMBLONDELSVMLIGHTLOADER  
DOWNLOADING DATASETS FROM THE OPENMLORG REPOSITORY  
OPENMLORG IS A PUBLIC REPOSITORY FOR MACHINE LEARNING DATA AND EXPERIMENTS THAT ALLOWS EVERYBODY TO UPLOAD OPEN  
DATASETS  
THESKLEARNDATASETS PACKAGE IS ABLE TO DOWNLOAD DATASETS FROM THE REPOSITORY USING THE FUNCTION SKLEARN  
DATASETSFETCHOPENML  
FOR EXAMPLE TO DOWNLOAD A DATASET OF GENE EXPRESSIONS IN MICE BRAINS  
FROM SKLEARNDATASETS IMPORT FETCHOPENML  
MICE FETCHOPENMLNAMEMICEPROTEIN VERSION4  
TO FULLY SPECIFY A DATASET YOU NEED TO PROVIDE A NAME AND A VERSION THOUGH THE VERSION IS OPTIONAL SEE DATASET  
VERSIONS BELOW THE DATASET CONTAINS A TOTAL OF 1080 EXAMPLES BELONGING TO 8 DIFFERENT CLASSES  
MICE DATASHAPE  
1080 77  
MICE TARGETSHAPE  
1080  
NPUNIQUEMICE TARGET  
ARRAYCCSM CCSS CSCM CSCS TCSM TCSS TSCM TSCS  
↪ DTYPEOBJECT  
YOU CAN GET MORE INFORMATION ON THE DATASET BY LOOKING AT THE DESCR ANDDETAILS ATTRIBUTES  
PRINTMICEDESCR  
AUTHOR CLARA HIGUERA KATHELEEN J GARDINER KRZYSZTOF J CIOS  
SOURCE UCIHTTPSSARCHIVEICSUCIEDUMLDATASETSMICEPROTEINEXPRESSION  
↪2015  
PLEASE CITE HIGUERA C GARDINER KJ CIOS KJ 2015 SELFORGANIZING  
FEATURE MAPS IDENTIFY PROTEINS CRITICAL TO LEARNING IN A MOUSE MODEL OF DOWN  
SYNDROME PLOS ONE 106 E0129126  
MICEDETAILS  
ID 40966 NAME MICEPROTEIN VERSION 4 FORMAT ARFF  
UPLOADDATE 20171108T160015 LICENCE PUBLIC  
36 DATASET LOADING UTILITIES 647

SCIKITLEARN USER GUIDE RELEASE 0213

URL [HTTPSWWWOPENMLORGDATAV1DOWNLOAD17928620MICEPROTEINARFF](https://www.openml.org/data/17928620/miceprotein.arff)

FILEID 17928620 DEFAULTTARGETATTRIBUTE CLASS

ROWIDATTRIBUTE MOUSEID

IGNOREATTRIBUTE GENOTYPE TREATMENT BEHAVIOR

TAG OPENMLCC18 STUDY135 STUDY98 STUDY99

VISIBILITY PUBLIC STATUS ACTIVE

MD5CHECKSUM 3C479A6885BFA0438971388283A1CE32

THEDESCR CONTAINS A FREETEXT DESCRIPTION OF THE DATA WHILE DETAILS CONTAINS A DICTIONARY OF METADATA STORED BY OPENML LIKE THE DATASET ID FOR MORE DETAILS SEE THE OPENML DOCUMENTATION THE DATAID OF THE MICE PROTEIN DATASET IS 40966 AND YOU CAN USE THIS OR THE NAME TO GET MORE INFORMATION ON THE DATASET ON THE OPENML WEBSITE

MICEURL

[HTTPSWWWOPENMLORGDATAV1DOWNLOAD17928620MICEPROTEINARFF](https://www.openml.org/data/40966)

THEDATAID ALSO UNIQUELY IDENTIFIES A DATASET FROM OPENML

MICE FETCHOPENMLDATAID40966

MICEDetails

ID 4550 NAME MICEPROTEIN VERSION 1 FORMAT ARFF

CREATOR

UPLOADDATE 20160217T143249 LICENCE PUBLIC URL

[HTTPSWWWOPENMLORGDATAV1DOWNLOAD1804243MICEPROTEINARFF](https://www.openml.org/data/1804243/miceprotein.arff) FILEID

1804243 DEFAULTTARGETATTRIBUTE CLASS CITATION HIGUERA C

GARDINER KJ CIO S KJ 2015 SELFORGANIZING FEATURE MAPS IDENTIFY PROTEINS

CRITICAL TO LEARNING IN A MOUSE MODEL OF DOWN SYNDROME PLOS ONE 106

E0129126 WEB LINK JOURNALPONE0129126 TAG OPENML100 STUDY14

STUDY34 VISIBILITY PUBLIC STATUS ACTIVE MD5CHECKSUM

3C479A6885BFA0438971388283A1CE32

DATASET VERSIONS

A DATASET IS UNIQUELY SPECIFIED BY ITS DATAID BUT NOT NECESSARILY BY ITS NAME SEVERAL DIFFERENT “VERSIONS” OF A DATASET WITH THE SAME NAME CAN EXIST WHICH CAN CONTAIN ENTIRELY DIFFERENT DATASETS IF A PARTICULAR VERSION OF A DATASET HAS BEEN FOUND TO CONTAIN SIGNIFICANT ISSUES IT MIGHT BE DEACTIVATED USING A NAME TO SPECIFY A DATASET WILL YIELD THE EARLIEST VERSION OF A DATASET THAT IS STILL ACTIVE THAT MEANS THAT FETCHOPENMLNAMEMICEPROTEIN CAN YIELD DIFFERENT RESULTS AT DIFFERENT TIMES IF EARLIER VERSIONS BECOME INACTIVE YOU CAN SEE THAT THE DATASET WITH DATAID 40966 THAT WE FETCHED ABOVE IS THE VERSION 1 OF THE “MICEPROTEIN” DATASET

MICEDetailsVERSION

1

IN FACT THIS DATASET ONLY HAS ONE VERSION THE IRIS DATASET ON THE OTHER HAND HAS MULTIPLE VERSIONS

IRIS FETCHOPENMLNAMEIRIS

IRISDetailsVERSION

1

IRISDetailsID

61

IRIS61 FETCHOPENMLDATAID61

IRIS61DetailsVERSION

1

IRIS61DetailsID

61

648 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

IRIS969 FETCHOPENMLDATAID969

IRIS969DETAILSVERSION

3

IRIS969DETAILSID

969

SPECIFYING THE DATASET BY THE NAME “IRIS” YIELDS THE LOWEST VERSION VERSION 1 WITH THE DATAID 61 TO MAKE SURE YOU ALWAYS GET THIS EXACT DATASET IT IS SAFEST TO SPECIFY IT BY THE DATASET DATAID THE OTHER DATASET WITH DATAID 969 IS VERSION 3 VERSION 2 HAS BECOME INACTIVE AND CONTAINS A BINARIZED VERSION OF THE DATA

NPUNIQUEIRIS969TARGET

ARRAYN P DTYPESOBJECT

YOU CAN ALSO SPECIFY BOTH THE NAME AND THE VERSION WHICH ALSO UNIQUELY IDENTIFIES THE DATASET

IRISVERSION3 FETCHOPENMLNAMEIRIS VERSION3

IRISVERSION3DETAILSVERSION

3

IRISVERSION3DETAILSID

969

REFERENCES

• VANSCHOREN VAN RIJN BISCHL AND TORGIO “OPENML NETWORKED SCIENCE IN MACHINE LEARNING” ACM SIGKDD

EXPLORATIONS NEWSLETTER 152 4960 2014

LOADING FROM EXTERNAL DATASETS

SCIKITLEARN WORKS ON ANY NUMERIC DATA STORED AS NUMPY ARRAYS OR SCIPY SPARSE MATRICES OTHER TYPES THAT ARE CONVERTIBLE

TO NUMERIC ARRAYS SUCH AS PANDAS DATAFRAME ARE ALSO ACCEPTABLE

HERE ARE SOME RECOMMENDED WAYS TO LOAD STANDARD COLUMNAR DATA INTO A FORMAT USABLE BY SCIKITLEARN

• PANDASIO PROVIDES TOOLS TO READ DATA FROM COMMON FORMATS INCLUDING CSV EXCEL JSON AND SQL DATAFRAMES

MAY ALSO BE CONSTRUCTED FROM LISTS OF TUPLES OR DICTS PANDAS HANDLES HETEROGENEOUS DATA SMOOTHLY AND PROVIDES

TOOLS FOR MANIPULATION AND CONVERSION INTO A NUMERIC ARRAY SUITABLE FOR SCIKITLEARN

• SCIPYIO SPECIALIZES IN BINARY FORMATS OFTEN USED IN SCIENTIFIC COMPUTING CONTEXT SUCH AS MAT AND ARFF

• NUMPYROUTINESIO FOR STANDARD LOADING OF COLUMNAR DATA INTO NUMPY ARRAYS

• SCIKITLEARN’S DATASETSLOADSVMLIGHTFILE FOR THE SVMLIGHT OR LIBSVM SPARSE FORMAT

• SCIKITLEARN’S DATASETSLOADFILES FOR DIRECTORIES OF TEXT FILES WHERE THE NAME OF EACH DIRECTORY IS THE

NAME OF EACH CATEGORY AND EACH FILE INSIDE OF EACH DIRECTORY CORRESPONDS TO ONE SAMPLE FROM THAT CATEGORY

FOR SOME MISCELLANEOUS DATA SUCH AS IMAGES VIDEOS AND AUDIO YOU MAY WISH TO REFER TO

• SKIMAGEIO OR IMAGEIO FOR LOADING IMAGES AND VIDEOS INTO NUMPY ARRAYS

• SCIPYIOWAVFILEREAD FOR READING WA V FILES INTO A NUMPY ARRAY

CATEGORICAL OR NOMINAL FEATURES STORED AS STRINGS COMMON IN PANDAS DATAFRAMES WILL NEED CONVERTING TO

NUMERICAL FEATURES USING SKLEARNPREPROCESSINGONEHOTENCODER ORSKLEARNPREPROCESSING

ORDINALENCODER OR SIMILAR SEE PREPROCESSING DATA

36 DATASET LOADING UTILITIES 649

SCIKITLEARN USER GUIDE RELEASE 0213

NOTE IF YOU MANAGE YOUR OWN NUMERICAL DATA IT IS RECOMMENDED TO USE AN OPTIMIZED FILE FORMAT SUCH AS HDF5 TO REDUCE DATA LOAD TIMES VARIOUS LIBRARIES SUCH AS H5PY PYTABLES AND PANDAS PROVIDES A PYTHON INTERFACE FOR READING AND WRITING DATA IN THAT FORMAT

37 COMPUTING WITH SCIKITLEARN

371 STRATEGIES TO SCALE COMPUTATIONALLY BIGGER DATA

FOR SOME APPLICATIONS THE AMOUNT OF EXAMPLES FEATURES OR BOTH ANDOR THE SPEED AT WHICH THEY NEED TO BE PROCESSED ARE CHALLENGING FOR TRADITIONAL APPROACHES IN THESE CASES SCIKITLEARN HAS A NUMBER OF OPTIONS YOU CAN CONSIDER TO MAKE YOUR SYSTEM SCALE

SCALING WITH INSTANCES USING OUTFCORE LEARNING

OUTFCORE OR “EXTERNAL MEMORY” LEARNING IS A TECHNIQUE USED TO LEARN FROM DATA THAT CANNOT FIT IN A COMPUTER’S MAIN MEMORY RAM

HERE IS A SKETCH OF A SYSTEM DESIGNED TO ACHIEVE THIS GOAL

1 A WAY TO STREAM INSTANCES

2 A WAY TO EXTRACT FEATURES FROM INSTANCES

3 AN INCREMENTAL ALGORITHM

STREAMING INSTANCES

BASICALLY 1 MAY BE A READER THAT YIELDS INSTANCES FROM FILES ON A HARD DRIVE A DATABASE FROM A NETWORK STREAM ETC HOWEVER DETAILS ON HOW TO ACHIEVE THIS ARE BEYOND THE SCOPE OF THIS DOCUMENTATION

EXTRACTING FEATURES

2 COULD BE ANY RELEVANT WAY TO EXTRACT FEATURES AMONG THE DIFFERENT FEATURE EXTRACTION METHODS SUPPORTED BY SCIKIT LEARN HOWEVER WHEN WORKING WITH DATA THAT NEEDS VECTORIZATION AND WHERE THE SET OF FEATURES OR VALUES IS NOT KNOWN IN ADVANCE ONE SHOULD TAKE EXPLICIT CARE A GOOD EXAMPLE IS TEXT CLASSIFICATION WHERE UNKNOWN TERMS ARE LIKELY TO BE FOUND DURING TRAINING IT IS POSSIBLE TO USE A STATEFUL VECTORIZER IF MAKING MULTIPLE PASSES OVER THE DATA IS REASONABLE FROM AN APPLICATION POINT OF VIEW OTHERWISE ONE CAN TURN UP THE DIFFICULTY BY USING A STATELESS FEATURE EXTRACTOR CURRENTLY THE PREFERRED WAY TO DO THIS IS TO USE THE SOCALLED HASHING TRICK AS IMPLEMENTED BY SKLEARN FEATUREEXTRACTIONFEATUREHASHER FOR DATASETS WITH CATEGORICAL VARIABLES REPRESENTED AS LIST OF PYTHON DICTS ORSKLEARNFEATUREEXTRACTIONTEXTHASHINGVECTORIZER FOR TEXT DOCUMENTS

INCREMENTAL LEARNING

FINALLY FOR 3 WE HAVE A NUMBER OF OPTIONS INSIDE SCIKITLEARN ALTHOUGH NOT ALL ALGORITHMS CAN LEARN INCREMENTALLY IE WITHOUT SEEING ALL THE INSTANCES AT ONCE ALL ESTIMATORS IMPLEMENTING THE PARTIALFIT API ARE CANDIDATES

ACTUALLY THE ABILITY TO LEARN INCREMENTALLY FROM A MINIBATCH OF INSTANCES SOMETIMES CALLED “ONLINE LEARNING” IS KEY TO OUTFCORE LEARNING AS IT GUARANTEES THAT AT ANY GIVEN TIME THERE WILL BE ONLY A SMALL AMOUNT OF INSTANCES IN THE



SCIKITLEARN USER GUIDE RELEASE 0213

MAIN MEMORY CHOOSING A GOOD SIZE FOR THE MINIBATCH THAT BALANCES RELEVANCY AND MEMORY FOOTPRINT COULD INVOLVE SOME TUNING1

HERE IS A LIST OF INCREMENTAL ESTIMATORS FOR DIFFERENT TASKS

- CLASSIFICATION
- SKLEARNNAIVEBAYESMULTINOMIALNB
- SKLEARNNAIVEBAYESBERNOULLINB
- SKLEARNLINEARMODELPERCEPTRON
- SKLEARNLINEARMODELSGDCLASSIFIER
- SKLEARNLINEARMODELPASSIVEAGGRESSIVECLASSIFIER
- SKLEARNNEURALNETWORKMLPCLASSIFIER
- REGRESSION
- SKLEARNLINEARMODELSGDSREGRESSOR
- SKLEARNLINEARMODELPASSIVEAGGRESSIVEREGRESSOR
- SKLEARNNEURALNETWORKMLPREGRESSOR
- CLUSTERING
- SKLEARNCLUSTERMINIBATCHKMEANS
- SKLEARNCLUSTERBIRCH
- DECOMPOSITION FEATURE EXTRACTION
- SKLEARNDECOMPOSITIONMINIBATCHDICTIONARYLEARNING
- SKLEARNDECOMPOSITIONINCREMENTALPCA
- SKLEARNDECOMPOSITIONLATENTDIRICHLETALLOCATION
- PREPROCESSING
- SKLEARNPREPROCESSINGSTANDARDSCALER
- SKLEARNPREPROCESSINGMINMAXSCALER
- SKLEARNPREPROCESSINGMAXABSSCALER

FOR CLASSIFICATION A SOMEWHAT IMPORTANT THING TO NOTE IS THAT ALTHOUGH A STATELESS FEATURE EXTRACTION ROUTINE MAY BE ABLE TO COPE WITH NEWUNSEEN ATTRIBUTES THE INCREMENTAL LEARNER ITSELF MAY BE UNABLE TO COPE WITH NEWUNSEEN TARGETS CLASSES IN THIS CASE YOU HAVE TO PASS ALL THE POSSIBLE CLASSES TO THE FIRST PARTIALFIT CALL USING THE CLASSES PARAMETER

ANOTHER ASPECT TO CONSIDER WHEN CHOOSING A PROPER ALGORITHM IS THAT NOT ALL OF THEM PUT THE SAME IMPORTANCE ON EACH EXAMPLE OVER TIME NAMELY THE PERCEPTRON IS STILL SENSITIVE TO BADLY LABELED EXAMPLES EVEN AFTER MANY EXAMPLES WHEREAS THE SGDANDPASSIVEAGGRESSIVE FAMILIES ARE MORE ROBUST TO THIS KIND OF ARTIFACTS CONVERSELY THE LATTER ALSO TEND TO GIVE LESS IMPORTANCE TO REMARKABLY DIFFERENT YET PROPERLY LABELED EXAMPLES WHEN THEY COME LATE IN THE STREAM AS THEIR LEARNING RATE DECREASES OVER TIME

1DEPENDING ON THE ALGORITHM THE MINIBATCH SIZE CAN INFLUENCE RESULTS OR NOT SGD PASSIVEAGGRESSIVE AND DISCRETE NAIVEBAYES / ONLINE AND ARE NOT AFFECTED BY BATCH SIZE CONVERSELY MINIBATCHKMEANS CONVERGENCE RATE IS AFFECTED BY THE BATCH SIZE ALSO IT CAN VARY DRAMATICALLY WITH BATCH SIZE

37 COMPUTING WITH SCIKITLEARN 651

SCIKITLEARN USER GUIDE RELEASE 0213  
EXAMPLES  
FINALLY WE HAVE A FULLFLEDGED EXAMPLE OF OUTOFCORE CLASSIFICATION OF TEXT DOCUMENTS IT IS AIMED AT PROVIDING A  
STARTING POINT FOR PEOPLE WANTING TO BUILD OUTOFCORE LEARNING SYSTEMS AND DEMONSTRATES MOST OF THE NOTIONS DISCUSSED  
ABOVE  
FURTHERMORE IT ALSO SHOWS THE EVOLUTION OF THE PERFORMANCE OF DIFFERENT ALGORITHMS WITH THE NUMBER OF PROCESSED  
EXAMPLES  
NOW LOOKING AT THE COMPUTATION TIME OF THE DIFFERENT PARTS WE SEE THAT THE VECTORIZATION IS MUCH MORE EXPENSIVE  
THAN LEARNING ITSELF FROM THE DIFFERENT ALGORITHMS MULTINOMIALNB IS THE MOST EXPENSIVE BUT ITS OVERHEAD CAN BE  
MITIGATED BY INCREASING THE SIZE OF THE MINIBATCHES EXERCISE CHANGE MINIBATCHSIZE TO 100 AND 10000 IN THE  
PROGRAM AND COMPARE  
652 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

NOTES

372 COMPUTATIONAL PERFORMANCE

FOR SOME APPLICATIONS THE PERFORMANCE MAINLY LATENCY AND THROUGHPUT AT PREDICTION TIME OF ESTIMATORS IS CRUCIAL IT MAY ALSO BE OF INTEREST TO CONSIDER THE TRAINING THROUGHPUT BUT THIS IS OFTEN LESS IMPORTANT IN A PRODUCTION SETUP WHERE IT OFTEN TAKES PLACE OFFLINE

WE WILL REVIEW HERE THE ORDERS OF MAGNITUDE YOU CAN EXPECT FROM A NUMBER OF SCIKITLEARN ESTIMATORS IN DIFFERENT CONTEXTS AND PROVIDE SOME TIPS AND TRICKS FOR OVERCOMING PERFORMANCE BOTTLENECKS

PREDICTION LATENCY IS MEASURED AS THE ELAPSED TIME NECESSARY TO MAKE A PREDICTION EG IN MICROSECONDS LATENCY IS OFTEN VIEWED AS A DISTRIBUTION AND OPERATIONS ENGINEERS OFTEN FOCUS ON THE LATENCY AT A GIVEN PERCENTILE OF THIS DISTRIBUTION EG THE 90 PERCENTILE

PREDICTION THROUGHPUT IS DEFINED AS THE NUMBER OF PREDICTIONS THE SOFTWARE CAN DELIVER IN A GIVEN AMOUNT OF TIME EG IN PREDICTIONS PER SECOND

AN IMPORTANT ASPECT OF PERFORMANCE OPTIMIZATION IS ALSO THAT IT CAN HURT PREDICTION ACCURACY INDEED SIMPLER MODELS EG LINEAR INSTEAD OF NONLINEAR OR WITH FEWER PARAMETERS OFTEN RUN FASTER BUT ARE NOT ALWAYS ABLE TO TAKE INTO ACCOUNT THE SAME EXACT PROPERTIES OF THE DATA AS MORE COMPLEX ONES

PREDICTION LATENCY

ONE OF THE MOST STRAIGHTFORWARD CONCERNS ONE MAY HAVE WHEN USINGCHOOSING A MACHINE LEARNING TOOLKIT IS THE LATENCY AT WHICH PREDICTIONS CAN BE MADE IN A PRODUCTION ENVIRONMENT

THE MAIN FACTORS THAT INFLUENCE THE PREDICTION LATENCY ARE

1 NUMBER OF FEATURES

37 COMPUTING WITH SCIKITLEARN 653

SCIKITLEARN USER GUIDE RELEASE 0213

2 INPUT DATA REPRESENTATION AND SPARSITY

3 MODEL COMPLEXITY

4 FEATURE EXTRACTION

A LAST MAJOR PARAMETER IS ALSO THE POSSIBILITY TO DO PREDICTIONS IN BULK OR ONEATATIME MODE

BULK VERSUS ATOMIC MODE

IN GENERAL DOING PREDICTIONS IN BULK MANY INSTANCES AT THE SAME TIME IS MORE EFFICIENT FOR A NUMBER OF REASONS

BRANCHING PREDICTABILITY CPU CACHE LINEAR ALGEBRA LIBRARIES OPTIMIZATIONS ETC HERE WE SEE ON A SETTING WITH FEW

FEATURES THAT INDEPENDENTLY OF ESTIMATOR CHOICE THE BULK MODE IS ALWAYS FASTER AND FOR SOME OF THEM BY 1 TO 2 ORDERS

OF MAGNITUDE

654 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

TO BENCHMARK DIFFERENT ESTIMATORS FOR YOUR CASE YOU CAN SIMPLY CHANGE THE NFEATURES PARAMETER IN THIS EXAMPLE  
PREDICTION LATENCY THIS SHOULD GIVE YOU AN ESTIMATE OF THE ORDER OF MAGNITUDE OF THE PREDICTION LATENCY  
CONFIGURING SCIKITLEARN FOR REDUCED VALIDATION OVERHEAD

SCIKITLEARN DOES SOME VALIDATION ON DATA THAT INCREASES THE OVERHEAD PER CALL TO PREDICT AND SIMILAR FUNCTIONS  
IN PARTICULAR CHECKING THAT FEATURES ARE FINITE NOT NAN OR INFINITE INVOLVES A FULL PASS OVER THE DATA IF YOU EN  
SURE THAT YOUR DATA IS ACCEPTABLE YOU MAY SUPPRESS CHECKING FOR FINITENESS BY SETTING THE ENVIRONMENT VARIABLE  
SKLEARNASSUMEFINITE TO A NONEMPTY STRING BEFORE IMPORTING SCIKITLEARN OR CONFIGURE IT IN PYTHON WITH  
SKLEARNSETCONFIG FOR MORE CONTROL THAN THESE GLOBAL SETTINGS A CONFIGCONTEXT ALLOWS YOU TO SET THIS  
CONFIGURATION WITHIN A SPECIFIED CONTEXT

```
import sklearn
with sklearn.config_context(assume_finite=True):
    pass # do learning/prediction here with reduced validation
```

NOTE THAT THIS WILL AFFECT ALL USES OF SKLEARNUTILSASSERTALLFINITE WITHIN THE CONTEXT

INFLUENCE OF THE NUMBER OF FEATURES  
OBVIOUSLY WHEN THE NUMBER OF FEATURES INCREASES SO DOES THE MEMORY CONSUMPTION OF EACH EXAMPLE INDEED FOR A  
MATRIX OF  $n$  INSTANCES WITH  $m$  FEATURES THE SPACE COMPLEXITY IS IN  $O(nm)$  FROM A COMPUTING PERSPECTIVE IT ALSO  
MEANS THAT THE NUMBER OF BASIC OPERATIONS EG MULTIPLICATIONS FOR VECTOR MATRIX PRODUCTS IN LINEAR MODELS INCREASES  
TOO HERE IS A GRAPH OF THE EVOLUTION OF THE PREDICTION LATENCY WITH THE NUMBER OF FEATURES

SCIKITLEARN USER GUIDE RELEASE 0213

OVERALL YOU CAN EXPECT THE PREDICTION TIME TO INCREASE AT LEAST LINEARLY WITH THE NUMBER OF FEATURES NONLINEAR CASES CAN HAPPEN DEPENDING ON THE GLOBAL MEMORY FOOTPRINT AND ESTIMATOR

INFLUENCE OF THE INPUT DATA REPRESENTATION

SCIPY PROVIDES SPARSE MATRIX DATA STRUCTURES WHICH ARE OPTIMIZED FOR STORING SPARSE DATA THE MAIN FEATURE OF SPARSE FORMATS IS THAT YOU DON'T STORE ZEROS SO IF YOUR DATA IS SPARSE THEN YOU USE MUCH LESS MEMORY A NONZERO VALUE IN A SPARSE CSR OR CSC REPRESENTATION WILL ONLY TAKE ON AVERAGE ONE 32BIT INTEGER POSITION THE 64 BIT FLOATING POINT VALUE AN ADDITIONAL 32BIT PER ROW OR COLUMN IN THE MATRIX USING SPARSE INPUT ON A DENSE OR SPARSE LINEAR MODEL CAN SPEEDUP PREDICTION BY QUITE A BIT AS ONLY THE NON ZERO VALUED FEATURES IMPACT THE DOT PRODUCT AND THUS THE MODEL PREDICTIONS HENCE IF YOU HAVE 100 NON ZEROS IN 1E6 DIMENSIONAL SPACE YOU ONLY NEED 100 MULTIPLY AND ADD OPERATION INSTEAD OF 1E6

CALCULATION OVER A DENSE REPRESENTATION HOWEVER MAY LEVERAGE HIGHLY OPTIMISED VECTOR OPERATIONS AND MULTITHREADING IN BLAS AND TENDS TO RESULT IN FEWER CPU CACHE MISSES SO THE SPARSITY SHOULD TYPICALLY BE QUITE HIGH 10 NONZEROS MAX TO BE CHECKED DEPENDING ON THE HARDWARE FOR THE SPARSE INPUT REPRESENTATION TO BE FASTER THAN THE DENSE INPUT REPRESENTATION ON A MACHINE WITH MANY CPUS AND AN OPTIMIZED BLAS IMPLEMENTATION

HERE IS SAMPLE CODE TO TEST THE SPARSITY OF YOUR INPUT

```
DEFSPARSITYRATIOX
RETURN10 NPCOUNTNONZEROX FLOATXSHAPE0 XSHAPE1
PRINTINPUT SPARSITY RATIO SPARSITYRATIOX
```

AS A RULE OF THUMB YOU CAN CONSIDER THAT IF THE SPARSITY RATIO IS GREATER THAN 90 YOU CAN PROBABLY BENEFIT FROM SPARSE FORMATS CHECK SCIPY'S SPARSE MATRIX FORMATS DOCUMENTATION FOR MORE INFORMATION ON HOW TO BUILD OR CONVERT YOUR DATA TO SPARSE MATRIX FORMATS MOST OF THE TIME THE CSR ANDCSC FORMATS WORK BEST

656 CHAPTER 3 USER GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

INFLUENCE OF THE MODEL COMPLEXITY

GENERALLY SPEAKING WHEN MODEL COMPLEXITY INCREASES PREDICTIVE POWER AND LATENCY ARE SUPPOSED TO INCREASE INCREASING PREDICTIVE POWER IS USUALLY INTERESTING BUT FOR MANY APPLICATIONS WE WOULD BETTER NOT INCREASE PREDICTION LATENCY TOO MUCH WE WILL NOW REVIEW THIS IDEA FOR DIFFERENT FAMILIES OF SUPERVISED MODELS

FORSKLEARNLINEARMODEL EG LASSO ELASTICNET SGDCLASSIFIERREGRESSOR RIDGE RIDGECLASSIFIER PAS

SIVEAGGRESSIVECLASSIFIERREGRESSOR LINEARSVC LOGISTICREGRESSION THE DECISION FUNCTION THAT IS APPLIED AT PREDICTION TIME IS THE SAME A DOT PRODUCT SO LATENCY SHOULD BE EQUIVALENT

HERE IS AN EXAMPLE USING SKLEARNLINEARMODELSTOCHASTICGRADIENTSGDCLASSIFIER WITH THE

ELASTICNET PENALTY THE REGULARIZATION STRENGTH IS GLOBALLY CONTROLLED BY THE ALPHA PARAMETER WITH A SUFFICIENTLY HIGHALPHA ONE CAN THEN INCREASE THE L1RATIO PARAMETER OF ELASTICNET TO ENFORCE VARIOUS LEVELS OF SPARSITY

IN THE MODEL COEFFICIENTS HIGHER SPARSITY HERE IS INTERPRETED AS LESS MODEL COMPLEXITY AS WE NEED FEWER COEFFICIENTS TO DESCRIBE IT FULLY OF COURSE SPARSITY INFLUENCES IN TURN THE PREDICTION TIME AS THE SPARSE DOTPRODUCT TAKES TIME ROUGHLY PROPORTIONAL TO THE NUMBER OF NONZERO COEFFICIENTS

FOR THESKLEARNNSVM FAMILY OF ALGORITHMS WITH A NONLINEAR KERNEL THE LATENCY IS TIED TO THE NUMBER OF SUPPORT VECTORS THE FEWER THE FASTER LATENCY AND THROUGHPUT SHOULD ASYMPTOTICALLY GROW LINEARLY WITH THE NUMBER OF SUPPORT VECTORS IN A SVC OR SVR MODEL THE KERNEL WILL ALSO INFLUENCE THE LATENCY AS IT IS USED TO COMPUTE THE PROJECTION OF THE INPUT VECTOR ONCE PER SUPPORT VECTOR IN THE FOLLOWING GRAPH THE NUPARAMETER OF SKLEARNNSVM CLASSESNUVR WAS USED TO INFLUENCE THE NUMBER OF SUPPORT VECTORS

SCIKITLEARN USER GUIDE RELEASE 0213  
FOR SKLEARN ENSEMBLE OF TREES EG RANDOM FOREST GBT EXTRA TREES ETC THE NUMBER OF TREES AND THEIR  
DEPTH PLAY THE MOST IMPORTANT ROLE LATENCY AND THROUGHPUT SHOULD SCALE LINEARLY WITH THE NUMBER OF TREES IN  
THIS CASE WE USED DIRECTLY THE NESTIMATORS PARAMETER OF SKLEARN ENSEMBLE GRADIENT BOOSTING  
GRADIENT BOOSTING REGRESSOR  
IN ANY CASE BE WARNED THAT DECREASING MODEL COMPLEXITY CAN HURT ACCURACY AS MENTIONED ABOVE FOR INSTANCE A NON  
LINEARLY SEPARABLE PROBLEM CAN BE HANDLED WITH A SPEEDY LINEAR MODEL BUT PREDICTION POWER WILL VERY LIKELY SUFFER IN  
THE PROCESS  
658 CHAPTER 3 USER GUIDE



SCIKITLEARN USER GUIDE RELEASE 0213

FEATURE EXTRACTION LATENCY

MOST SCIKITLEARN MODELS ARE USUALLY PRETTY FAST AS THEY ARE IMPLEMENTED EITHER WITH COMPILED CYTHON EXTENSIONS OR OPTIMIZED COMPUTING LIBRARIES ON THE OTHER HAND IN MANY REAL WORLD APPLICATIONS THE FEATURE EXTRACTION PROCESS IE TURNING RAW DATA LIKE DATABASE ROWS OR NETWORK PACKETS INTO NUMPY ARRAYS GOVERNS THE OVERALL PREDICTION TIME FOR EXAMPLE ON THE REUTERS TEXT CLASSIFICATION TASK THE WHOLE PREPARATION READING AND PARSING SGML FILES TOKENIZING THE TEXT AND HASHING IT INTO A COMMON VECTOR SPACE IS TAKING 100 TO 500 TIMES MORE TIME THAN THE ACTUAL PREDICTION CODE DEPENDING ON THE CHOSEN MODEL

IN MANY CASES IT IS THUS RECOMMENDED TO CAREFULLY TIME AND PROFILE YOUR FEATURE EXTRACTION CODE AS IT MAY BE A GOOD PLACE TO START OPTIMIZING WHEN YOUR OVERALL LATENCY IS TOO SLOW FOR YOUR APPLICATION

PREDICTION THROUGHPUT

ANOTHER IMPORTANT METRIC TO CARE ABOUT WHEN SIZING PRODUCTION SYSTEMS IS THE THROUGHPUT IE THE NUMBER OF PREDICTIONS YOU CAN MAKE IN A GIVEN AMOUNT OF TIME HERE IS A BENCHMARK FROM THE PREDICTION LATENCY EXAMPLE THAT MEASURES THIS QUANTITY FOR A NUMBER OF ESTIMATORS ON SYNTHETIC DATA

37 COMPUTING WITH SCIKITLEARN 659

SCIKITLEARN USER GUIDE RELEASE 0213

THESE THROUGHPUTS ARE ACHIEVED ON A SINGLE PROCESS AN OBVIOUS WAY TO INCREASE THE THROUGHPUT OF YOUR APPLICATION IS TO SPAWN ADDITIONAL INSTANCES USUALLY PROCESSES IN PYTHON BECAUSE OF THE GIL THAT SHARE THE SAME MODEL ONE MIGHT ALSO ADD MACHINES TO SPREAD THE LOAD A DETAILED EXPLANATION ON HOW TO ACHIEVE THIS IS BEYOND THE SCOPE OF THIS DOCUMENTATION THOUGH

TIPS AND TRICKS

LINEAR ALGEBRA LIBRARIES

AS SCIKITLEARN RELIES HEAVILY ON NUMPYSCIPY AND LINEAR ALGEBRA IN GENERAL IT MAKES SENSE TO TAKE EXPLICIT CARE OF THE VERSIONS OF THESE LIBRARIES BASICALLY YOU OUGHT TO MAKE SURE THAT NUMPY IS BUILT USING AN OPTIMIZED BLAS LAPACK LIBRARY

NOT ALL MODELS BENEFIT FROM OPTIMIZED BLAS AND LAPACK IMPLEMENTATIONS FOR INSTANCE MODELS BASED ON RANDOMIZED DECISION TREES TYPICALLY DO NOT RELY ON BLAS CALLS IN THEIR INNER LOOPS NOR DO KERNEL SVMs SVC SVR NUSVC NUSVR ON THE OTHER HAND A LINEAR MODEL IMPLEMENTED WITH A BLAS DGEMM CALL VIA NUMPYDOT WILL TYPICALLY BENEFIT HUGELY FROM A TUNED BLAS IMPLEMENTATION AND LEAD TO ORDERS OF MAGNITUDE SPEEDUP OVER A NONOPTIMIZED BLAS

YOU CAN DISPLAY THE BLAS LAPACK IMPLEMENTATION USED BY YOUR NUMPY SCIPY SCIKITLEARN INSTALL WITH THE FOLLOWING COMMANDS

```
from numpy.distutils.system_info import get_info
```

```
print(get_info('blas_opt'))
```

```
print(get_info('lapack_opt'))
```

OPTIMIZED BLAS LAPACK IMPLEMENTATIONS INCLUDE

- ATLAS NEED HARDWARE SPECIFIC TUNING BY REBUILDING ON THE TARGET MACHINE

SCIKITLEARN USER GUIDE RELEASE 0213

- OPENBLAS
  - MKL
  - APPLE ACCELERATE AND VECLIB FRAMEWORKS OSX ONLY
- MORE INFORMATION CAN BE FOUND ON THE SCIPY INSTALL PAGE AND IN THIS BLOG POST FROM DANIEL NOURI WHICH HAS SOME NICE STEP BY STEP INSTALL INSTRUCTIONS FOR DEBIAN UBUNTU

LIMITING WORKING MEMORY

SOME CALCULATIONS WHEN IMPLEMENTED USING STANDARD NUMPY VECTORIZED OPERATIONS INVOLVE USING A LARGE AMOUNT OF TEMPORARY MEMORY THIS MAY POTENTIALLY EXHAUST SYSTEM MEMORY WHERE COMPUTATIONS CAN BE PERFORMED IN FIXED MEMORY CHUNKS WE ATTEMPT TO DO SO AND ALLOW THE USER TO HINT AT THE MAXIMUM SIZE OF THIS WORKING MEMORY DE FAULTING TO 1GB USING SKLEARNSETCONFIG ORCONFIGCONTEXT THE FOLLOWING SUGGESTS TO LIMIT TEMPORARY

WORKING MEMORY TO 128 MIB

IMPORT SKLEARN

WITH SKLEARNCONFIGCONTEXTWORKINGMEMORY128

PASS DO CHUNKED WORK HERE

AN EXAMPLE OF A CHUNKED OPERATION ADHERING TO THIS SETTING IS METRICPAIRWISEDISTANCESCHUNKED WHICH FACILITATES COMPUTING ROWWISE REDUCTIONS OF A PAIRWISE DISTANCE MATRIX

MODEL COMPRESSION

MODEL COMPRESSION IN SCIKITLEARN ONLY CONCERNS LINEAR MODELS FOR THE MOMENT IN THIS CONTEXT IT MEANS THAT WE WANT TO CONTROL THE MODEL SPARSITY IE THE NUMBER OF NONZERO COORDINATES IN THE MODEL VECTORS IT IS GENERALLY A GOOD IDEA TO COMBINE MODEL SPARSITY WITH SPARSE INPUT DATA REPRESENTATION

HERE IS SAMPLE CODE THAT ILLUSTRATES THE USE OF THE SPARSIFY METHOD

CLF SGDREGRESSORPENALTYELASTICNET L1RATIO0.25

CLFFITXTRAIN YTRAINSPARSIFY

CLFPREDICTXTEST

IN THIS EXAMPLE WE PREFER THE ELASTICNET PENALTY AS IT IS OFTEN A GOOD COMPROMISE BETWEEN MODEL COMPACTNESS AND PREDICTION POWER ONE CAN ALSO FURTHER TUNE THE L1RATIO PARAMETER IN COMBINATION WITH THE REGULARIZATION STRENGTHALPHA TO CONTROL THIS TRADEOFF

A TYPICAL BENCHMARK ON SYNTHETIC DATA YIELDS A 30 DECREASE IN LATENCY WHEN BOTH THE MODEL AND INPUT ARE SPARSE WITH 0.000024 AND 0.027400 NONZERO COEFFICIENTS RATIO RESPECTIVELY YOUR MILEAGE MAY VARY DEPENDING ON THE SPARSITY AND SIZE OF YOUR DATA AND MODEL FURTHERMORE SPARSIFYING CAN BE VERY USEFUL TO REDUCE THE MEMORY USAGE OF PREDICTIVE MODELS DEPLOYED ON PRODUCTION SERVERS

MODEL RESHAPING

MODEL RESHAPING CONSISTS IN SELECTING ONLY A PORTION OF THE AVAILABLE FEATURES TO FIT A MODEL IN OTHER WORDS IF A MODEL DISCARDS FEATURES DURING THE LEARNING PHASE WE CAN THEN STRIP THOSE FROM THE INPUT THIS HAS SEVERAL BENEFITS FIRSTLY IT REDUCES MEMORY AND THEREFORE TIME OVERHEAD OF THE MODEL ITSELF IT ALSO ALLOWS TO DISCARD EXPLICIT FEATURE SELECTION COMPONENTS IN A PIPELINE ONCE WE KNOW WHICH FEATURES TO KEEP FROM A PREVIOUS RUN FINALLY IT CAN HELP REDUCE PROCESSING TIME AND IO USAGE UPSTREAM IN THE DATA ACCESS AND FEATURE EXTRACTION LAYERS BY NOT COLLECTING AND BUILDING FEATURES THAT ARE DISCARDED BY THE MODEL FOR INSTANCE IF THE RAW DATA COME FROM A DATABASE IT CAN MAKE IT POSSIBLE TO WRITE SIMPLER AND FASTER QUERIES OR REDUCE IO USAGE BY MAKING THE QUERIES RETURN LIGHTER RECORDS AT THE 37 COMPUTING WITH SCIKITLEARN 661

SCIKITLEARN USER GUIDE RELEASE 0213  
MOMENT RESHAPING NEEDS TO BE PERFORMED MANUALLY IN SCIKITLEARN IN THE CASE OF SPARSE INPUT PARTICULARLY IN CSR  
FORMAT IT IS GENERALLY SUFFICIENT TO NOT GENERATE THE RELEVANT FEATURES LEAVING THEIR COLUMNS EMPTY  
LINKS

- SCIKITLEARN DEVELOPER PERFORMANCE DOCUMENTATION
  - SCIPY SPARSE MATRIX FORMATS DOCUMENTATION
- 373 PARALLELISM RESOURCE MANAGEMENT AND CONFIGURATION  
PARALLEL AND DISTRIBUTED COMPUTING

SCIKITLEARN USES THE JOBLIB LIBRARY TO ENABLE PARALLEL COMPUTING INSIDE ITS ESTIMATORS SEE THE JOBLIB DOCUMENTATION FOR  
THE SWITCHES TO CONTROL PARALLEL COMPUTING

NOTE THAT BY DEFAULT SCIKITLEARN USES ITS EMBEDDED VENDORED VERSION OF JOBLIB A CONFIGURATION SWITCH DOCUMENTED  
BELOW CONTROLS THIS BEHAVIOR  
CONFIGURATION SWITCHES

PYTHON RUNTIME  
SKLEARNSETCONFIG CONTROLS THE FOLLOWING BEHAVIORS  
ASSUMEFINITE USED TO SKIP VALIDATION WHICH ENABLES FASTER COMPUTATIONS BUT MAY LEAD TO SEGMENTATION  
FAULTS IF THE DATA CONTAINS NANS  
WORKINGMEMORY THE OPTIMAL SIZE OF TEMPORARY ARRAYS USED BY SOME ALGORITHMS  
ENVIRONMENT VARIABLES

THESE ENVIRONMENT VARIABLES SHOULD BE SET BEFORE IMPORTING SCIKITLEARN  
SKLEARNSITEJOBLIB WHEN THIS ENVIRONMENT VARIABLE IS SET TO A NON ZERO VALUE SCIKITLEARN USES  
THE SITE JOBLIB RATHER THAN ITS VENDORED VERSION CONSEQUENTLY JOBLIB MUST BE INSTALLED FOR SCIKITLEARN  
TO RUN NOTE THAT USING THE SITE JOBLIB IS AT YOUR OWN RISKS THE VERSIONS OF SCIKITLEARN AND JOBLIB  
NEED TO BE COMPATIBLE CURRENTLY JOBLIB 011 IS SUPPORTED IN ADDITION DUMPS FROM JOBLIBMEMORY  
MIGHT BE INCOMPATIBLE AND YOU MIGHT LOOSE SOME CACHES AND HAVE TO REDOWNLOAD SOME DATASETS  
DEPRECATED SINCE VERSION 021 AS OF VERSION 021 THIS PARAMETER HAS NO EFFECT VENDORED JOBLIB WAS  
REMOVED AND SITE JOBLIB IS ALWAYS USED

SKLEARNASSUMEFINITE SETS THE DEFAULT VALUE FOR THE ASSUMEFINITE ARGUMENT OF  
SKLEARNSETCONFIG  
SKLEARNWORKINGMEMORY SETS THE DEFAULT VALUE FOR THE LIMITING WORKING MEMORY ARGUMENT  
OFSKLEARNSETCONFIG  
SKLEARNSEED SETS THE SEED OF THE GLOBAL RANDOM GENERATOR WHEN RUNNING THE TESTS FOR REPRODUCIBIL  
ITY  
SKLEARNSKIPNETWORKTESTS WHEN THIS ENVIRONMENT VARIABLE IS SET TO A NON ZERO VALUE THE  
TESTS THAT NEED NETWORK ACCESS ARE SKIPPED

CHAPTER  
FOUR  
GLOSSARY OF COMMON TERMS AND API ELEMENTS  
THIS GLOSSARY HOPES TO DEFINITELY REPRESENT THE TACIT AND EXPLICIT CONVENTIONS APPLIED IN SCIKITLEARN AND ITS API WHILE PROVIDING A REFERENCE FOR USERS AND CONTRIBUTORS IT AIMS TO DESCRIBE THE CONCEPTS AND EITHER DETAIL THEIR CORRESPONDING API OR LINK TO OTHER RELEVANT PARTS OF THE DOCUMENTATION WHICH DO SO BY LINKING TO GLOSSARY ENTRIES FROM THE API REFERENCE AND USER GUIDE WE MAY MINIMIZE REDUNDANCY AND INCONSISTENCY  
WE BEGIN BY LISTING GENERAL CONCEPTS AND ANY THAT DIDN'T FIT ELSEWHERE BUT MORE SPECIFIC SETS OF RELATED TERMS ARE LISTED BELOW CLASS APIS AND ESTIMATOR TYPES TARGET TYPES METHODS PARAMETERS ATTRIBUTES DATA AND SAMPLE PROPERTIES  
41 GENERAL CONCEPTS  
1D  
1D ARRAY ONEDIMENSIONAL ARRAY A NUMPY ARRAY WHOSE SHAPE HAS LENGTH 1 A VECTOR  
2D  
2D ARRAY TWODIMENSIONAL ARRAY A NUMPY ARRAY WHOSE SHAPE HAS LENGTH 2 OFTEN REPRESENTS A MATRIX  
API REFERS TO BOTH THE SPECIFIC INTERFACES FOR ESTIMATORS IMPLEMENTED IN SCIKITLEARN AND THE GENERALIZED CONVENTIONS ACROSS TYPES OF ESTIMATORS AS DESCRIBED IN THIS GLOSSARY AND OVERVIEWED IN THE CONTRIBUTOR DOCUMENTATION  
THE SPECIFIC INTERFACES THAT CONSTITUTE SCIKITLEARN'S PUBLIC API ARE LARGELY DOCUMENTED IN API REFERENCE HOW EVER WE LESS FORMALLY CONSIDER ANYTHING AS PUBLIC API IF NONE OF THE IDENTIFIERS REQUIRED TO ACCESS IT BEGINS WITH  
WE GENERALLY TRY TO MAINTAIN BACKWARDS COMPATIBILITY FOR ALL OBJECTS IN THE PUBLIC API  
PRIVATE API INCLUDING FUNCTIONS MODULES AND METHODS BEGINNING ARE NOT ASSURED TO BE STABLE  
ARRAYLIKE THE MOST COMMON DATA FORMAT FOR INPUT TO SCIKITLEARN ESTIMATORS AND FUNCTIONS ARRAYLIKE IS ANY TYPE OBJECT FOR WHICH NUMPYASARRAY WILL PRODUCE AN ARRAY OF APPROPRIATE SHAPE USUALLY 1 OR 2DIMENSIONAL OF APPROPRIATE DTYPE USUALLY NUMERIC  
THIS INCLUDES  
• A NUMPY ARRAY  
• A LIST OF NUMBERS  
• A LIST OF LENGTHK LISTS OF NUMBERS FOR SOME FIXED LENGTH K  
• APANDASDATAFRAME WITH ALL COLUMNS NUMERIC  
• A NUMERIC PANDASSERIES  
IT EXCLUDES  
• ASPARSE MATRIX  
• AN ITERATOR  
663

SCIKITLEARN USER GUIDE RELEASE 0213

• A GENERATOR

NOTE THAT OUTPUT FROM SCIKITLEARN ESTIMATORS AND FUNCTIONS EG PREDICTIONS SHOULD GENERALLY BE ARRAYS OR SPARSE MATRICES OR LISTS THEREOF AS IN MULTIOUTPUT TREEDECISIONTREECLASSIFIER 'SPREDICTPROBA' AN ESTIMATOR WHERE PREDICT RETURNS A LIST OR A PANDASSERIES IS NOT VALID

ATTRIBUTE

ATTRIBUTES WE MOSTLY USE ATTRIBUTE TO REFER TO HOW MODEL INFORMATION IS STORED ON AN ESTIMATOR DURING FITTING ANY PUBLIC ATTRIBUTE STORED ON AN ESTIMATOR INSTANCE IS REQUIRED TO BEGIN WITH AN ALPHABETIC CHARACTER AND END IN A SINGLE UNDERSCORE IF IT IS SET IN FITORPARTIALFIT THESE ARE WHAT IS DOCUMENTED UNDER AN ESTIMATOR'S ATTRIBUTES DOCUMENTATION THE INFORMATION STORED IN ATTRIBUTES IS USUALLY EITHER SUFFICIENT STATISTICS USED FOR PREDICTION OR TRANSFORMATION TRANSDUCTIVE OUTPUTS SUCH AS LABELS OREMBEDDING OR DIAGNOSTIC DATA SUCH AS FEATUREIMPORTANCES

COMMON ATTRIBUTES ARE LISTED BELOW

A PUBLIC ATTRIBUTE MAY HAVE THE SAME NAME AS A CONSTRUCTOR PARAMETER WITH APPENDED THIS IS USED TO STORE A VALIDATED OR ESTIMATED VERSION OF THE USER'S INPUT FOR EXAMPLE DECOMPOSITIONPCA IS CONSTRUCTED WITH ANNCOMPONENTS PARAMETER FROM THIS TOGETHER WITH OTHER PARAMETERS AND THE DATA PCA ESTIMATES THE ATTRIBUTENCOMPONENTS

FURTHER PRIVATE ATTRIBUTES USED IN PREDICTIONTRANSFORMATIONETC MAY ALSO BE SET WHEN FITTING THESE BEGIN WITH A SINGLE UNDERSCORE AND ARE NOT ASSURED TO BE STABLE FOR PUBLIC ACCESS

A PUBLIC ATTRIBUTE ON AN ESTIMATOR INSTANCE THAT DOES NOT END IN AN UNDERSCORE SHOULD BE THE STORED UNMODIFIED VALUE OF AN INIT PARAMETER OF THE SAME NAME BECAUSE OF THIS EQUIVALENCE THESE ARE DOCUMENTED UNDER AN ESTIMATOR'S PARAMETERS DOCUMENTATION

BACKWARDS COMPATIBILITY WE GENERALLY TRY TO MAINTAIN BACKWARDS COMPATIBILITY IE INTERFACES AND BEHAVIORS MAY BE EXTENDED BUT NOT CHANGED OR REMOVED FROM RELEASE TO RELEASE BUT THIS COMES WITH SOME EXCEPTIONS PUBLIC API ONLY THE BEHAVIOUR OF OBJECTS ACCESSED THROUGH PRIVATE IDENTIFIERS THOSE BEGINNING MAY BE CHANGED ARBITRARILY BETWEEN VERSIONS

AS DOCUMENTED WE WILL GENERALLY ASSUME THAT THE USERS HAVE ADHERED TO THE DOCUMENTED PARAMETER TYPES AND RANGES IF THE DOCUMENTATION ASKS FOR A LIST AND THE USER GIVES A TUPLE WE DO NOT ASSURE CONSISTENT BEHAVIOR FROM VERSION TO VERSION

DEPRECATION BEHAVIORS MAY CHANGE FOLLOWING A DEPRECATION PERIOD USUALLY TWO RELEASES LONG WARNINGS ARE ISSUED USING PYTHON'S WARNINGS MODULE

KEYWORD ARGUMENTS WE MAY SOMETIMES ASSUME THAT ALL OPTIONAL PARAMETERS OTHER THAN X AND Y TO FIT AND SIMILAR METHODS ARE PASSED AS KEYWORD ARGUMENTS ONLY AND MAY BE POSITIONALLY REORDERED

BUG FIXES AND ENHANCEMENTS BUG FIXES AND - LESS OFTEN - ENHANCEMENTS MAY CHANGE THE BEHAVIOR OF ESTIMATORS INCLUDING THE PREDICTIONS OF AN ESTIMATOR TRAINED ON THE SAME DATA AND RANDOMSTATE WHEN THIS HAPPENS WE ATTEMPT TO NOTE IT CLEARLY IN THE CHANGELOG

SERIALIZATION WE MAKE NO ASSURANCES THAT PICKLING AN ESTIMATOR IN ONE VERSION WILL ALLOW IT TO BE UNPICKLED TO AN EQUIVALENT MODEL IN THE SUBSEQUENT VERSION FOR ESTIMATORS IN THE SKLEARN PACKAGE WE ISSUE A WARNING WHEN THIS UNPICKLING IS ATTEMPTED EVEN IF IT MAY HAPPEN TO WORK SEE SECURITY MAINTAINABILITY LIMITATIONS

UTILSESTIMATORCHECKSCHECKESTIMATOR WE PROVIDE LIMITED BACKWARDS COMPATIBILITY ASSUR

ANCES FOR THE ESTIMATOR CHECKS WE MAY ADD EXTRA REQUIREMENTS ON ESTIMATORS TESTED WITH THIS FUNCTION USUALLY WHEN THESE WERE INFORMALLY ASSUMED BUT NOT FORMALLY TESTED

DESPITE THIS INFORMAL CONTRACT WITH OUR USERS THE SOFTWARE IS PROVIDED AS IS AS STATED IN THE LICENCE WHEN A RELEASE INADVERTENTLY INTRODUCES CHANGES THAT ARE NOT BACKWARDS COMPATIBLE THESE ARE KNOWN AS SOFTWARE REGRESSIONS

CALLABLE A FUNCTION CLASS OR AN OBJECT WHICH IMPLEMENTS THE CALL METHOD ANYTHING THAT RETURNS TRUE WHEN THE ARGUMENT OF CALLABLE

664 CHAPTER 4 GLOSSARY OF COMMON TERMS AND API ELEMENTS

SCIKITLEARN USER GUIDE RELEASE 0213

CATEGORICAL FEATURE A CATEGORICAL OR NOMINAL FEATURE IS ONE THAT HAS A FINITE SET OF DISCRETE VALUES ACROSS THE POPULATION OF DATA THESE ARE COMMONLY REPRESENTED AS COLUMNS OF INTEGERS OR STRINGS STRINGS WILL BE REJECTED BY MOST SCIKITLEARN ESTIMATORS AND INTEGERS WILL BE TREATED AS ORDINAL OR COUNTVALUED FOR THE USE WITH MOST ESTIMATORS CATEGORICAL VARIABLES SHOULD BE ONEHOT ENCODED NOTABLE EXCEPTIONS INCLUDE TREEBASED MODELS SUCH AS RANDOM FORESTS AND GRADIENT BOOSTING MODELS THAT OFTEN WORK BETTER AND FASTER WITH INTEGERCODED CATEGORICAL VARIABLES ORDINALENCODER HELPS ENCODING STRINGVALUED CATEGORICAL FEATURES AS ORDINAL INTEGERS AND ONEHOTENCODER CAN BE USED TO ONEHOT ENCODE CATEGORICAL FEATURES SEE ALSO ENCODING CATEGORICAL FEATURES AND THE CATEGORICALENCODING PACKAGE FOR TOOLS RELATED TO ENCODING CATEGORICAL FEATURES

CLONE  
CLONED TO COPY AN ESTIMATOR INSTANCE AND CREATE A NEW ONE WITH IDENTICAL PARAMETERS BUT WITHOUT ANY FITTED ATTRIBUTES USINGCLONE

WHENFIT IS CALLED A METAESTIMATOR USUALLY CLONES A WRAPPED ESTIMATOR INSTANCE BEFORE FITTING THE CLONED INSTANCE EXCEPTIONS FOR LEGACY REASONS INCLUDE PIPELINE ANDFEATUREUNION  
COMMON TESTS THIS REFERS TO THE TESTS RUN ON ALMOST EVERY ESTIMATOR CLASS IN SCIKITLEARN TO CHECK THEY COMPLY WITH BASIC API CONVENTIONS THEY ARE AVAILABLE FOR EXTERNAL USE THROUGH UTILSESTIMATORCHECKS  
CHECKESTIMATOR WITH MOST OF THE IMPLEMENTATION IN SKLEARNUTILSESTIMATORCHECKSPY

NOTE SOME EXCEPTIONS TO THE COMMON TESTING REGIME ARE CURRENTLY HARDCODED INTO THE LIBRARY BUT WE HOPE TO REPLACE THIS BY MARKING EXCEPTIONAL BEHAVIOURS ON THE ESTIMATOR USING SEMANTIC ESTIMATOR TAGS  
DEPRECATION WE USE DEPRECATION TO SLOWLY VIOLATE OUR BACKWARDS COMPATIBILITY ASSURANCES USUALLY TO TO

- CHANGE THE DEFAULT VALUE OF A PARAMETER OR
  - REMOVE A PARAMETER ATTRIBUTE METHOD CLASS ETC
- WE WILL ORDINARILY ISSUE A WARNING WHEN A DEPRECATED ELEMENT IS USED ALTHOUGH THERE MAY BE LIMITATIONS TO THIS FOR INSTANCE WE WILL RAISE A WARNING WHEN SOMEONE SETS A PARAMETER THAT HAS BEEN DEPRECATED BUT MAY NOT WHEN THEY ACCESS THAT PARAMETER’S ATTRIBUTE ON THE ESTIMATOR INSTANCE

SEE THE CONTRIBUTORS’ GUIDE  
DIMENSIONALITY MAY BE USED TO REFER TO THE NUMBER OF FEATURES IE NFEATURES OR COLUMNS IN A 2D FEATURE MATRIX DIMENSIONS ARE HOWEVER ALSO USED TO REFER TO THE LENGTH OF A NUMPY ARRAY’S SHAPE DISTINGUISHING A 1D ARRAY FROM A 2D MATRIX

DOCSTRING THE EMBEDDED DOCUMENTATION FOR A MODULE CLASS FUNCTION ETC USUALLY IN CODE AS A STRING AT THE BEGINNING OF THE OBJECT’S DEFINITION AND ACCESSIBLE AS THE OBJECT’S DOC ATTRIBUTE  
WE TRY TO ADHERE TO PEP257 AND FOLLOW NUMPYDOC CONVENTIONS

DOUBLE UNDERSCORE  
DOUBLE UNDERSCORE NOTATION WHEN SPECIFYING PARAMETER NAMES FOR NESTED ESTIMATORS MAY BE USED TO SEPARATE BETWEEN PARENT AND CHILD IN SOME CONTEXTS THE MOST COMMON USE IS WHEN SETTING PARAMETERS THROUGH A META ESTIMATOR WITH SETPARAMS AND HENCE IN SPECIFYING A SEARCH GRID IN PARAMETER SEARCH SEE PARAMETER IT IS ALSO USED INPIPELINEPIPELINEFIT FOR PASSING SAMPLE PROPERTIES TO THEFIT METHODS OF ESTIMATORS IN THE

PIPELINE  
DTYPE  
DATA TYPE NUMPY ARRAYS ASSUME A HOMOGENEOUS DATA TYPE THROUGHOUT AVAILABLE IN THE DTYPE ATTRIBUTE OF AN ARRAY OR SPARSE MATRIX WE GENERALLY ASSUME SIMPLE DATA TYPES FOR SCIKITLEARN DATA FLOAT OR INTEGER WE MAY SUPPORT OBJECT OR STRING DATA TYPES FOR ARRAYS BEFORE ENCODING OR VECTORIZING OUR ESTIMATORS DO NOT WORK WITH STRUCT ARRAYS FOR INSTANCE

TODO MENTION EFFICIENCY AND PRECISION ISSUES CASTING POLICY  
41 GENERAL CONCEPTS 665

SCIKITLEARN USER GUIDE RELEASE 0213

DUCK TYPING WE TRY TO APPLY DUCK TYPING TO DETERMINE HOW TO HANDLE SOME INPUT VALUES EG CHECKING WHETHER A GIVEN ESTIMATOR IS A CLASSIFIER THAT IS WE AVOID USING ISINSTANCE WHERE POSSIBLE AND RELY ON THE PRESENCE OR ABSENCE OF ATTRIBUTES TO DETERMINE AN OBJECT’S BEHAVIOUR SOME NUANCE IS REQUIRED WHEN FOLLOWING THIS APPROACH

- FOR SOME ESTIMATORS AN ATTRIBUTE MAY ONLY BE AVAILABLE ONCE IT IS FITTED FOR INSTANCE WE CANNOT A PRIORI DETERMINE IF PREDICTPROBA IS AVAILABLE IN A GRID SEARCH WHERE THE GRID INCLUDES ALTERNATING BETWEEN A PROBABILISTIC AND A NONPROBABILISTIC PREDICTOR IN THE FINAL STEP OF THE PIPELINE IN THE FOLLOWING WE CAN ONLY DETERMINE IF CLF IS PROBABILISTIC AFTER FITTING IT ON SOME DATA

```
FROM SKLEARNMODELSELECTION IMPORT GRIDSEARCHCV
FROM SKLEARNLINEARMODEL IMPORT SGDCLASSIFIER
CLF = GRIDSEARCHCVSGDCLASSIFIER
PARAMGRIDLOSS LOG HINGE
```

THIS MEANS THAT WE CAN ONLY CHECK FOR DUCKTYPED ATTRIBUTES AFTER FITTING AND THAT WE MUST BE CAREFUL TO MAKE METAESTIMATORS ONLY PRESENT ATTRIBUTES ACCORDING TO THE STATE OF THE UNDERLYING ESTIMATOR AFTER FITTING

- CHECKING IF AN ATTRIBUTE IS PRESENT USING HASATTR IS IN GENERAL JUST AS EXPENSIVE AS GETTING THE ATTRIBUTE GETATTR OR DOT NOTATION IN SOME CASES GETTING THE ATTRIBUTE MAY INDEED BE EXPENSIVE EG FOR SOME IMPLEMENTATIONS OF FEATUREIMPORTANCES WHICH MAY SUGGEST THIS IS AN API DESIGN FLAW SO CODE WHICH DOESHASATTR FOLLOWED BY GETATTR SHOULD BE AVOIDED GETATTR WITHIN A TRYEXCEPT BLOCK IS PREFERRED

- FOR DETERMINING SOME ASPECTS OF AN ESTIMATOR’S EXPECTATIONS OR SUPPORT FOR SOME FEATURE WE USE ESTIMATOR TAGS INSTEAD OF DUCK TYPING

EARLY STOPPING THIS CONSISTS IN STOPPING AN ITERATIVE OPTIMIZATION METHOD BEFORE THE CONVERGENCE OF THE TRAINING LOSS TO AVOID OVERFITTING THIS IS GENERALLY DONE BY MONITORING THE GENERALIZATION SCORE ON A VALIDATION SET WHEN AVAILABLE IT IS ACTIVATED THROUGH THE PARAMETER EARLYSTOPPING OR BY SETTING A POSITIVE NITERNOCHANGE ESTIMATOR INSTANCE WE SOMETIMES USE THIS TERMINOLOGY TO DISTINGUISH AN ESTIMATOR CLASS FROM A CONSTRUCTED INSTANCE FOR EXAMPLE IN THE FOLLOWING CLS IS AN ESTIMATOR CLASS WHILE EST1 AND EST2 ARE INSTANCES

```
CLS = RANDOMFORESTCLASSIFIER
EST1 = CLS
EST2 = RANDOMFORESTCLASSIFIER
```

EXAMPLES WE TRY TO GIVE EXAMPLES OF BASIC USAGE FOR MOST FUNCTIONS AND CLASSES IN THE API

- AS DOCTESTS IN THEIR DOCSTRINGS IE WITHIN THE SKLEARN LIBRARY CODE ITSELF
- AS EXAMPLES IN THE EXAMPLE GALLERY RENDERED USING SPHINXGALLERY FROM SCRIPTS IN THE EXAMPLES DIRECTORY EXEMPLIFYING KEY FEATURES OR PARAMETERS OF THE ESTIMATORFUNCTION THESE SHOULD ALSO BE REFERENCED FROM THE USER GUIDE
- SOMETIMES IN THE USER GUIDE BUILT FROM DOC ALONGSIDE A TECHNICAL DESCRIPTION OF THE ESTIMATOR

EVALUATION METRIC

EVALUATION METRICS EVALUATION METRICS GIVE A MEASURE OF HOW WELL A MODEL PERFORMS WE MAY USE THIS TERM SPECIFICALLY TO REFER TO THE FUNCTIONS IN METRICS DISREGARDING METRICSPAIRWISE AS DISTINCT FROM THE SCORE METHOD AND THE SCORING API USED IN CROSS VALIDATION SEE MODEL EVALUATION QUANTIFYING THE QUALITY OF PREDICTIONS

THESE FUNCTIONS USUALLY ACCEPT A GROUND TRUTH OR THE RAW DATA WHERE THE METRIC EVALUATES CLUSTERING WITHOUT A GROUND TRUTH AND A PREDICTION BE IT THE OUTPUT OF PREDICT YPRED OF PREDICTPROBA YPROBA OR OF AN ARBITRARY SCORE FUNCTION INCLUDING DECISIONFUNCTION YSCORE FUNCTIONS ARE USUALLY NAMED TO END WITH SCORE IF A GREATER SCORE INDICATES A BETTER MODEL AND LOSS IF A LESSER SCORE INDICATES A BETTER MODEL THIS DIVERSITY OF INTERFACE MOTIVATES THE SCORING API



SCIKITLEARN USER GUIDE RELEASE 0213

NOTE THAT SOME ESTIMATORS CAN CALCULATE METRICS THAT ARE NOT INCLUDED IN METRICS AND ARE ESTIMATORSPECIFIC NOTABLY MODEL LIKELIHOODS

ESTIMATOR TAGS A PROPOSED FEATURE EG 8022 BY WHICH THE CAPABILITIES OF AN ESTIMATOR ARE DESCRIBED THROUGH A SET OF SEMANTIC TAGS THIS WOULD ENABLE SOME RUNTIME BEHAVIORS BASED ON ESTIMATOR INSPECTION BUT IT ALSO ALLOWS EACH ESTIMATOR TO BE TESTED FOR APPROPRIATE INVARIANCES WHILE BEING EXCEPTED FROM OTHER COMMON TESTS SOME ASPECTS OF ESTIMATOR TAGS ARE CURRENTLY DETERMINED THROUGH THE DUCK TYPING OF METHODS LIKE PREDICTPROBA AND THROUGH SOME SPECIAL ATTRIBUTES ON ESTIMATOR OBJECTS

ESTIMATORTYPE THIS STRINGVALUED ATTRIBUTE IDENTIFIES AN ESTIMATOR AS BEING A CLASSIFIER REGRESSOR ETC IT IS SET BY MIXINS SUCH AS BASECLASSIFIERMIXIN BUT NEEDS TO BE MORE EXPLICITLY ADOPTED ON A META ESTIMATOR ITS VALUE SHOULD USUALLY BE CHECKED BY WAY OF A HELPER SUCH AS BASEISCLASSIFIER PAIRWISE THIS BOOLEAN ATTRIBUTE INDICATES WHETHER THE DATA X PASSED TOFIT AND SIMILAR METHODS CONSISTS OF PAIRWISE MEASURES OVER SAMPLES RATHER THAN A FEATURE REPRESENTATION FOR EACH SAMPLE IT IS USUALLY TRUE WHERE AN ESTIMATOR HAS A METRIC ORAFFINITY ORKERNEL PARAMETER WITH VALUE 'PRECOMPUTED' ITS PRIMARY PURPOSE IS THAT WHEN A METAESTIMATOR EXTRACTS A SUBSAMPLE OF DATA INTENDED FOR A PAIRWISE ESTIMATOR THE DATA NEEDS TO BE INDEXED ON BOTH AXES WHILE OTHER DATA IS INDEXED ONLY ON THE FIRST AXIS

FEATURE

FEATURES

FEATURE VECTOR IN THE ABSTRACT A FEATURE IS A FUNCTION IN ITS MATHEMATICAL SENSE MAPPING A SAMPLED OBJECT TO A NUMERIC OR CATEGORICAL QUANTITY "FEATURE" IS ALSO COMMONLY USED TO REFER TO THESE QUANTITIES BEING THE INDIVIDUAL ELEMENTS OF A VECTOR REPRESENTING A SAMPLE IN A DATA MATRIX FEATURES ARE REPRESENTED AS COLUMNS EACH COLUMN CONTAINS THE RESULT OF APPLYING A FEATURE FUNCTION TO A SET OF SAMPLES

ELSEWHERE FEATURES ARE KNOWN AS ATTRIBUTES PREDICTORS REGRESSORS OR INDEPENDENT VARIABLES

NEARLY ALL ESTIMATORS IN SCIKITLEARN ASSUME THAT FEATURES ARE NUMERIC FINITE AND NOT MISSING EVEN WHEN THEY HAVE SEMANTICALLY DISTINCT DOMAINS AND DISTRIBUTIONS CATEGORICAL ORDINAL COUNTVALUED REALVALUED INTERVAL SEE ALSO CATEGORICAL FEATURE ANDMISSING VALUES

NFEATURES INDICATES THE NUMBER OF FEATURES IN A DATASET

FITTING CALLING FITORFITTRANSFORM FITPREDICT ETC ON AN ESTIMATOR

FITTED THE STATE OF AN ESTIMATOR AFTER FITTING

THERE IS NO CONVENTIONAL PROCEDURE FOR CHECKING IF AN ESTIMATOR IS FITTED HOWEVER AN ESTIMATOR THAT IS NOT FITTED

- SHOULD RAISE EXCEPTIONSNOTFITTEDERROR WHEN A PREDICTION METHOD PREDICT TRANSFORM ETC IS

CALLED UTILSVALIDATIONCHECKISFITTED IS USED INTERNALLY FOR THIS PURPOSE

- SHOULD NOT HAVE ANY ATTRIBUTES BEGINNING WITH AN ALPHABETIC CHARACTER AND ENDING WITH AN UNDERSCORE

NOTE THAT A DESCRIPTOR FOR THE ATTRIBUTE MAY STILL BE PRESENT ON THE CLASS BUT HASATTR SHOULD RETURN FALSE

FUNCTION WE PROVIDE AD HOC FUNCTION INTERFACES FOR MANY ALGORITHMS WHILE ESTIMATOR CLASSES PROVIDE A MORE CONSISTENT INTERFACE

IN PARTICULAR SCIKITLEARN MAY PROVIDE A FUNCTION INTERFACE THAT FITS A MODEL TO SOME DATA AND RETURNS THE LEARNT

MODEL PARAMETERS AS IN LINEARMODELENETPATH FOR TRANSDUCTIVE MODELS THIS ALSO RETURNS THE EM

BEDDING OR CLUSTER LABELS AS IN MANIFOLDSPECTRALEMMBEDDING ORCLUSTERDBSCAN MANY PREPRO

CESSING TRANSFORMERS ALSO PROVIDE A FUNCTION INTERFACE AKIN TO CALLING FITTRANSFORM AS INPREPROCESSING

MAXABSSCALE USERS SHOULD BE CAREFUL TO AVOID DATA LEAKAGE WHEN MAKING USE OF THESE FITTRANSFORM

EQUIVALENT FUNCTIONS

WE DO NOT HAVE A STRICT POLICY ABOUT WHEN TO OR WHEN NOT TO PROVIDE FUNCTION FORMS OF ESTIMATORS BUT MAINTAINERS

SHOULD CONSIDER CONSISTENCY WITH EXISTING INTERFACES AND WHETHER PROVIDING A FUNCTION WOULD LEAD USERS ASTRAY

FROM BEST PRACTICES AS REGARDS DATA LEAKAGE ETC

41 GENERAL CONCEPTS 667

SCIKITLEARN USER GUIDE RELEASE 0213

GALLERY SEEEXAMPLES

HYPERPARAMETER

HYPERPARAMETER SEEPARAMETER

IMPUTE

IMPUTATION MOST MACHINE LEARNING ALGORITHMS REQUIRE THAT THEIR INPUTS HAVE NO MISSING VALUES AND WILL NOT WORK IF THIS REQUIREMENT IS VIOLATED ALGORITHMS THAT ATTEMPT TO FILL IN OR IMPUTE MISSING VALUES ARE REFERRED TO AS IMPUTATION ALGORITHMS

INDEXABLE ANARRAYLIKE SPARSE MATRIX PANDAS DATAFRAME OR SEQUENCE USUALLY A LIST

INDUCTION

INDUCTIVE INDUCTIVE CONTRASTED WITH TRANSDUCTIVE MACHINE LEARNING BUILDS A MODEL OF SOME DATA THAT CAN THEN BE APPLIED TO NEW INSTANCES MOST ESTIMATORS IN SCIKITLEARN ARE INDUCTIVE HAVING PREDICT ANDOR TRANSFORM METHODS

JOBLIB A PYTHON LIBRARY [HTTPS://JOBLIBREADTHEDOCSIO](https://joblib.readthedocs.io) USED IN SCIKITLEARN TO FACILITE SIMPLE PARALLELISM AND CACHING

JOBLIB IS ORIENTED TOWARDS EFFICIENTLY WORKING WITH NUMPY ARRAYS SUCH AS THROUGH USE OF MEMORY MAPPING SEE PARALLEL AND DISTRIBUTED COMPUTING FOR MORE INFORMATION

LABEL INDICATOR MATRIX

MULTILABEL INDICATOR MATRIX

MULTILABEL INDICATOR MATRICES THE FORMAT USED TO REPRESENT MULTILABEL DATA WHERE EACH ROW OF A 2D ARRAY OR SPARSE MATRIX CORRESPONDS TO A SAMPLE EACH COLUMN CORRESPONDS TO A CLASS AND EACH ELEMENT IS 1 IF THE SAMPLE IS LABELED WITH THE CLASS AND 0 IF NOT

LEAKAGE

DATA LEAKAGE A PROBLEM IN CROSS VALIDATION WHERE GENERALIZATION PERFORMANCE CAN BE OVERESTIMATED SINCE KNOWLEDGE OF THE TEST DATA WAS INADVERTENTLY INCLUDED IN TRAINING A MODEL THIS IS A RISK FOR INSTANCE WHEN APPLYING ATRANSFORMER TO THE ENTIRETY OF A DATASET RATHER THAN EACH TRAINING PORTION IN A CROSS VALIDATION SPLIT

WE AIM TO PROVIDE INTERFACES SUCH AS PIPELINE ANDMODELSELECTION THAT SHIELD THE USER FROM DATA LEAKAGE

MEMMMAPPING

MEMORY MAP

MEMORY MAPPING A MEMORY EFFICIENCY STRATEGY THAT KEEPS DATA ON DISK RATHER THAN COPYING IT INTO MAIN MEMORY

MEMORY MAPS CAN BE CREATED FOR ARRAYS THAT CAN BE READ WRITTEN OR BOTH USING NUMPYMEMMAP WHEN USING JOBLIB TO PARALLELIZE OPERATIONS IN SCIKITLEARN IT MAY AUTOMATICALLY MEMMAP LARGE ARRAYS TO REDUCE MEMORY DUPLICATION OVERHEAD IN MULTIPROCESSING

MISSING VALUES MOST SCIKITLEARN ESTIMATORS DO NOT WORK WITH MISSING VALUES WHEN THEY DO EG IN IMPUTE

SIMPLEIMPUTER NAN IS THE PREFERRED REPRESENTATION OF MISSING VALUES IN FLOAT ARRAYS IF THE ARRAY HAS INTEGER DTYPE NAN CANNOT BE REPRESENTED FOR THIS REASON WE SUPPORT SPECIFYING ANOTHER MISSINGVALUES VALUE WHEN IMPUTATION OR LEARNING CAN BE PERFORMED IN INTEGER SPACE UNLABELED DATA IS A SPECIAL CASE OF MISSING VALUES IN THE TARGET

NFEATURES THE NUMBER OF FEATURES

NOUTPUTS THE NUMBER OF OUTPUTS IN THE TARGET

NSAMPLES THE NUMBER OF SAMPLES

NTARGETS SYNONYM FOR NOUTPUTS

NARRATIVE DOCS

668 CHAPTER 4 GLOSSARY OF COMMON TERMS AND API ELEMENTS

SCIKITLEARN USER GUIDE RELEASE 0213

NARRATIVE DOCUMENTATION AN ALIAS FOR USER GUIDE IE DOCUMENTATION WRITTEN IN DOCMODULES UNLIKE THE API REFERENCE PROVIDED THROUGH DOCSTRINGS THE USER GUIDE AIMS TO

- GROUP TOOLS PROVIDED BY SCIKITLEARN TOGETHER THEMATICALLY OR IN TERMS OF USAGE
- MOTIVATE WHY SOMEONE WOULD USE EACH PARTICULAR TOOL OFTEN THROUGH COMPARISON
- PROVIDE BOTH INTUITIVE AND TECHNICAL DESCRIPTIONS OF TOOLS
- PROVIDE OR LINK TO EXAMPLES OF USING KEY FEATURES OF A TOOL

NPA SHORTHAND FOR NUMPY DUE TO THE CONVENTIONAL IMPORT STATEMENT

IMPORT NUMPY AS NP

ONLINE LEARNING WHERE A MODEL IS ITERATIVELY UPDATED BY RECEIVING EACH BATCH OF GROUND TRUTH TARGETS SOON AFTER MAKING PREDICTIONS ON CORRESPONDING BATCH OF DATA INTRINSICALLY THE MODEL MUST BE USABLE FOR PREDICTION AFTER EACH BATCH SEE PARTIALFIT

OUTOFCORE AN EFFICIENCY STRATEGY WHERE NOT ALL THE DATA IS STORED IN MAIN MEMORY AT ONCE USUALLY BY PERFORMING LEARNING ON BATCHES OF DATA SEE PARTIALFIT

OUTPUTS INDIVIDUAL SCALARCATEGORICAL VARIABLES PER SAMPLE IN THE TARGET FOR EXAMPLE IN MULTILABEL CLASSIFICATION EACH POSSIBLE LABEL CORRESPONDS TO A BINARY OUTPUT ALSO CALLED RESPONSES TASKS OR TARGETS SEE MULTICLASS MULTIOUTPUT ANDCONTINUOUS MULTIOUTPUT

PAIR A TUPLE OF LENGTH TWO

PARAMETER

PARAMETERS

PARAM

PARAMS WE MOSTLY USE PARAMETER TO REFER TO THE ASPECTS OF AN ESTIMATOR THAT CAN BE SPECIFIED IN ITS CONSTRUCTION FOR EXAMPLEMAXDEPTH ANDRANDOMSTATE ARE PARAMETERS OF RANDOMFORESTCLASSIFIER PARAMETERS

TO AN ESTIMATOR'S CONSTRUCTOR ARE STORED UNMODIFIED AS ATTRIBUTES ON THE ESTIMATOR INSTANCE AND CONVENTIONALLY START WITH AN ALPHABETIC CHARACTER AND END WITH AN ALPHANUMERIC CHARACTER EACH ESTIMATOR'S CONSTRUCTOR PARAM ETERS ARE DESCRIBED IN THE ESTIMATOR'S DOCSTRING

WE DO NOT USE PARAMETERS IN THE STATISTICAL SENSE WHERE PARAMETERS ARE VALUES THAT SPECIFY A MODEL AND CAN BE ESTIMATED FROM DATA WHAT WE CALL PARAMETERS MIGHT BE WHAT STATISTICIANS CALL HYPERPARAMETERS TO THE MODEL ASPECTS FOR CONFIGURING MODEL STRUCTURE THAT ARE OFTEN NOT DIRECTLY LEARNT FROM DATA HOWEVER OUR PARAMETERS ARE ALSO USED TO PRESCRIBE MODELING OPERATIONS THAT DO NOT AFFECT THE LEARNT MODEL SUCH AS NJOBS FOR CONTROLLING PARALLELISM

WHEN TALKING ABOUT THE PARAMETERS OF A METAESTIMATOR WE MAY ALSO BE INCLUDING THE PARAMETERS OF THE ESTIMATORS WRAPPED BY THE METAESTIMATOR ORDINARILY THESE NESTED PARAMETERS ARE DENOTED BY USING ADOUBLE UNDERSCORE TO SEPARATE BETWEEN THE ESTIMATORASPARAMETER AND ITS PARAMETER THUS CLF

BAGGINGCLASSIFIERBASEESTIMATORDECISIONTREECLASSIFIERMAXDEPTH3 HAS

A DEEP PARAMETER BASEESTIMATORMAXDEPTH WITH VALUE 3 WHICH IS ACCESSIBLE WITH CLF

BASEESTIMATORMAXDEPTH ORCLFGETPARAMSBASEESTIMATORMAXDEPTH

THE LIST OF PARAMETERS AND THEIR CURRENT VALUES CAN BE RETRIEVED FROM AN ESTIMATOR INSTANCE USING ITS GETPARAMS METHOD

BETWEEN CONSTRUCTION AND FITTING PARAMETERS MAY BE MODIFIED USING SETPARAMS TO ENABLE THIS PARAMETERS ARE NOT ORDINARILY VALIDATED OR ALTERED WHEN THE ESTIMATOR IS CONSTRUCTED OR WHEN EACH PARAMETER IS SET PARAMETER VALIDATION IS PERFORMED WHEN FITIS CALLED

COMMON PARAMETERS ARE LISTED BELOW

PAIRWISE METRIC

41 GENERAL CONCEPTS 669

SCIKITLEARN USER GUIDE RELEASE 0213

PAIRWISE METRICS IN ITS BROAD SENSE A PAIRWISE METRIC DEFINES A FUNCTION FOR MEASURING SIMILARITY OR DISSIMILARITY BETWEEN TWO SAMPLES WITH EACH ORDINARILY REPRESENTED AS A FEATURE VECTOR WE PARTICULARLY PROVIDE IMPLEMENTATIONS OF DISTANCE METRICS AS WELL AS IMPROPER METRICS LIKE COSINE DISTANCE THROUGH METRICS PAIRWISEDISTANCES AND OF KERNEL FUNCTIONS A CONSTRAINED CLASS OF SIMILARITY FUNCTIONS IN METRICS PAIRWISEKERNELS THESE CAN COMPUTE PAIRWISE DISTANCE MATRICES THAT ARE SYMMETRIC AND HENCE STORE DATA REDUNDANTLY

SEE ALSO PRECOMPUTED ANDMETRIC

NOTE THAT FOR MOST DISTANCE METRICS WE RELY ON IMPLEMENTATIONS FROM SCIPYSPATIALDISTANCE BUT MAY REIMPLEMENT FOR EFFICIENCY IN OUR CONTEXT THE NEIGHBORS MODULE ALSO DUPLICATES SOME METRIC IMPLEMENTATIONS FOR INTEGRATION WITH EFFICIENT BINARY TREE SEARCH DATA STRUCTURES

PDA SHORTHAND FOR PANDAS DUE TO THE CONVENTIONAL IMPORT STATEMENT

IMPORT PANDAS AS PD

PRECOMPUTED WHERE ALGORITHMS RELY ON PAIRWISE METRICS AND CAN BE COMPUTED FROM PAIRWISE METRICS ALONE WE OFTEN ALLOW THE USER TO SPECIFY THAT THE XPROVIDED IS ALREADY IN THE PAIRWISE DISSIMILARITY SPACE RATHER THAN IN A FEATURE SPACE THAT IS WHEN PASSED TO FIT IT IS A SQUARE SYMMETRIC MATRIX WITH EACH VECTOR INDICATING DISSIMILARITY TO EVERY SAMPLE AND WHEN PASSED TO PREDICTIONTRANSFORMATION METHODS EACH ROW CORRESPONDS TO A TESTING SAMPLE AND EACH COLUMN TO A TRAINING SAMPLE

USE OF PRECOMPUTED X IS USUALLY INDICATED BY SETTING A METRIC AFFINITY ORKERNEL PARAMETER TO THE STRING 'PRECOMPUTED' AN ESTIMATOR SHOULD MARK ITSELF AS BEING PAIRWISE IF THIS IS THE CASE RECTANGULAR DATA THAT CAN BE REPRESENTED AS A MATRIX WITH SAMPLES ON THE FIRST AXIS AND A FIXED FINITE SET OF FEATURES ON THE SECOND IS CALLED RECTANGULAR

THIS TERM EXCLUDES SAMPLES WITH NONVECTORIAL STRUCTURE SUCH AS TEXT AN IMAGE OF ARBITRARY SIZE A TIME SERIES OF ARBITRARY LENGTH A SET OF VECTORS ETC THE PURPOSE OF A VECTORIZER IS TO PRODUCE RECTANGULAR FORMS OF SUCH DATA SAMPLE

SAMPLES WE USUALLY USE THIS TERM AS A NOUN TO INDICATE A SINGLE FEATURE VECTOR ELSEWHERE A SAMPLE IS CALLED AN INSTANCE DATA POINT OR OBSERVATION NSAMPLES INDICATES THE NUMBER OF SAMPLES IN A DATASET BEING THE NUMBER OF ROWS IN A DATA ARRAY X

SAMPLE PROPERTY

SAMPLE PROPERTIES A SAMPLE PROPERTY IS DATA FOR EACH SAMPLE EG AN ARRAY OF LENGTH NSAMPLES PASSED TO AN ESTIMATOR METHOD OR A SIMILAR FUNCTION ALONGSIDE BUT DISTINCT FROM THE FEATURES X AND TARGET Y THE MOST PROMINENT EXAMPLE IS SAMPLEWEIGHT SEE OTHERS AT DATA AND SAMPLE PROPERTIES

AS OF VERSION 019 WE DO NOT HAVE A CONSISTENT APPROACH TO HANDLING SAMPLE PROPERTIES AND THEIR ROUTING IN METAESTIMATORS THOUGH AFITPARAMS PARAMETER IS OFTEN USED

SCIKITLEARNCONTRIB A VENUE FOR PUBLISHING SCIKITLEARNCOMPATIBLE LIBRARIES THAT ARE BROADLY AUTHORIZED BY THE CORE DEVELOPERS AND THE CONTRIB COMMUNITY BUT NOT MAINTAINED BY THE CORE DEVELOPER TEAM SEE HTTPS

SCIKITLEARNCONTRIBGITHUBIO

SCIKITLEARN ENHANCEMENT PROPOSALS

SLEP

SLEPS CHANGES TO THE API PRINCIPLES AND CHANGES TO DEPENDENCIES OR SUPPORTED VERSIONS HAPPEN VIA A SLEP AND FOLLOWS THE DECISIONMAKING PROCESS OUTLINED IN SCIKITLEARN GOVERNANCE AND DECISIONMAKING FOR ALL VOTES A PROPOSAL MUST HAVE BEEN MADE PUBLIC AND DISCUSSED BEFORE THE VOTE SUCH PROPOSAL MUST BE A CONSOLIDATED DOCUMENT IN THE FORM OF A 'SCIKITLEARN ENHANCEMENT PROPOSAL' SLEP RATHER THAN A LONG DISCUSSION ON AN ISSUE A SLEP MUST BE SUBMITTED AS A PULLREQUEST TO ENHANCEMENT PROPOSALS USING THE SLEP TEMPLATE

SEMISUPERVISED

670 CHAPTER 4 GLOSSARY OF COMMON TERMS AND API ELEMENTS

SCIKITLEARN USER GUIDE RELEASE 0213

SEMISUPERVISED LEARNING

SEMISUPERVISED LEARNING WHERE THE EXPECTED PREDICTION LABEL OR GROUND TRUTH IS ONLY AVAILABLE FOR SOME SAMPLES PROVIDED AS TRAINING DATA WHEN FITTING THE MODEL WE CONVENTIONALLY APPLY THE LABEL 1 TO UNLABELED SAMPLES IN SEMISUPERVISED CLASSIFICATION

SPARSE MATRIX A REPRESENTATION OF TWO DIMENSIONAL NUMERIC DATA THAT IS MORE MEMORY EFFICIENT THE CORRESPONDING DENSE NUMPY ARRAY WHERE ALMOST ALL ELEMENTS ARE ZERO WE USE THE SCIPYSPARSE FRAMEWORK WHICH PROVIDES SEVERAL UNDERLYING SPARSE DATA REPRESENTATIONS OR FORMATS SOME FORMATS ARE MORE EFFICIENT THAN OTHERS FOR PARTICULAR TASKS AND WHEN A PARTICULAR FORMAT PROVIDES ESPECIAL BENEFIT WE TRY TO DOCUMENT THIS FACT IN SCIKIT LEARN PARAMETER DESCRIPTIONS

SOME SPARSE MATRIX FORMATS NOTABLY CSR CSC COO AND LIL DISTINGUISH BETWEEN IMPLICIT AND EXPLICIT ZEROS EXPLICIT ZEROS ARE STORED IE THEY CONSUME MEMORY IN A DATA ARRAY IN THE DATA STRUCTURE WHILE IMPLICIT ZEROS CORRESPOND TO EVERY ELEMENT NOT OTHERWISE DEFINED IN EXPLICIT STORAGE

TWO SEMANTICS FOR SPARSE MATRICES ARE USED IN SCIKITLEARN

MATRIX SEMANTICS THE SPARSE MATRIX IS INTERPRETED AS AN ARRAY WITH IMPLICIT AND EXPLICIT ZEROS BEING INTERPRETED AS THE NUMBER 0 THIS IS THE INTERPRETATION MOST OFTEN ADOPTED EG WHEN SPARSE MATRICES ARE USED FOR FEATURE MATRICES OR MULTILABEL INDICATOR MATRICES

GRAPH SEMANTICS AS WITH SCIPYSPARSECSGRAPH EXPLICIT ZEROS ARE INTERPRETED AS THE NUMBER 0 BUT IMPLICIT ZEROS INDICATE A MASKED OR ABSENT VALUE SUCH AS THE ABSENCE OF AN EDGE BETWEEN TWO VERTICES OF A GRAPH WHERE AN EXPLICIT VALUE INDICATES AN EDGE'S WEIGHT THIS INTERPRETATION IS ADOPTED TO REPRESENT CONNECTIVITY IN CLUSTERING IN REPRESENTATIONS OF NEAREST NEIGHBORHOODS EG NEIGHBORS KNEIGHBORSGRAPH AND FOR PRECOMPUTED DISTANCE REPRESENTATION WHERE ONLY DISTANCES IN THE NEIGHBORHOOD OF EACH POINT ARE REQUIRED

WHEN WORKING WITH SPARSE MATRICES WE ASSUME THAT IT IS SPARSE FOR A GOOD REASON AND AVOID WRITING CODE THAT DENSIFIES A USER PROVIDED SPARSE MATRIX INSTEAD MAINTAINING SPARSITY OR RAISING AN ERROR IF NOT POSSIBLE IE IF AN ESTIMATOR DOES NOT CANNOT SUPPORT SPARSE MATRICES

SUPERVISED

SUPERVISED LEARNING LEARNING WHERE THE EXPECTED PREDICTION LABEL OR GROUND TRUTH IS AVAILABLE FOR EACH SAMPLE WHEN FITTING THE MODEL PROVIDED AS Y THIS IS THE APPROACH TAKEN IN A CLASSIFIER OR REGRESSOR AMONG OTHER ESTIMATORS

TARGET

TARGETS THE DEPENDENT VARIABLE IN SUPERVISED AND SEMISUPERVISED LEARNING PASSED AS Y TO AN ESTIMATOR'S FIT METHOD ALSO KNOWN AS DEPENDENT VARIABLE OUTCOME VARIABLE RESPONSE VARIABLE GROUND TRUTH OR LABEL SCIKITLEARN WORKS WITH TARGETS THAT HAVE MINIMAL STRUCTURE A CLASS FROM A FINITE SET A FINITE REAL VALUED NUMBER MULTIPLE CLASSES OR MULTIPLE NUMBERS SEE TARGET TYPES

TRANSDUCTION

TRANSDUCTIVE A TRANSDUCTIVE CONTRASTED WITH INDUCTIVE MACHINE LEARNING METHOD IS DESIGNED TO MODEL A SPECIFIC DATASET BUT NOT TO APPLY THAT MODEL TO UNSEEN DATA EXAMPLES INCLUDE MANIFOLD TSNE CLUSTER AGGLOMERATIVE CLUSTERING AND NEIGHBORS LOCAL OUTLIER FACTOR

UNLABELED

UNLABELED DATA SAMPLES WITH AN UNKNOWN GROUND TRUTH WHEN FITTING EQUIVALENTLY MISSING VALUES IN THE TARGET SEE ALSO SEMISUPERVISED AND UNSUPERVISED LEARNING

UNSUPERVISED

UNSUPERVISED LEARNING LEARNING WHERE THE EXPECTED PREDICTION LABEL OR GROUND TRUTH IS NOT AVAILABLE FOR EACH SAMPLE WHEN FITTING THE MODEL AS IN CLUSTERERS AND OUTLIER DETECTORS UNSUPERVISED ESTIMATORS IGNORE ANY Y PASSED TO FIT

41 GENERAL CONCEPTS 671

42 CLASS APIS AND ESTIMATOR TYPES

CLASSIFIER

CLASSIFIERS ASUPERVISED ORSEMISUPERVISED PREDICTOR WITH A FINITE SET OF DISCRETE POSSIBLE OUTPUT VALUES

A CLASSIFIER SUPPORTS MODELING SOME OF BINARY MULTICLASS MULTILABEL ORMULTICLASS MULTIOUTPUT TARGETS WITHIN SCIKITLEARN ALL CLASSIFIERS SUPPORT MULTICLASS CLASSIFICATION DEFAULTING TO USING A ONEVSREST STRATEGY OVER THE BINARY CLASSIFICATION PROBLEM

CLASSIFIERS MUST STORE A CLASSES ATTRIBUTE AFTER FITTING AND USUALLY INHERIT FROM BASECLASSIFIERMIXIN

WHICH SETS THEIR ESTIMATOR TYPE ATTRIBUTE

A CLASSIFIER CAN BE DISTINGUISHED FROM OTHER ESTIMATORS WITH ISCLASSIFIER

A CLASSIFIER MUST IMPLEMENT

- FIT
- PREDICT
- SCORE

IT MAY ALSO BE APPROPRIATE TO IMPLEMENT DECISIONFUNCTION PREDICTPROBA ANDPREDICTLOGPROBA

CLUSTERER

CLUSTERERS AUNSUPERVISED PREDICTOR WITH A FINITE SET OF DISCRETE OUTPUT VALUES

A CLUSTERER USUALLY STORES LABELS AFTER FITTING AND MUST DO SO IF IT IS TRANSDUCTIVE

A CLUSTERER MUST IMPLEMENT

- FIT
- FITPREDICT IFTRANSDUCTIVE
- PREDICT IFINDUCTIVE

DENSITY ESTIMATOR

TODO

ESTIMATOR AN OBJECT WHICH MANAGES THE ESTIMATION AND DECODING OF A MODEL THE MODEL IS ESTIMATED AS A DETERMINISTIC FUNCTION OF

- PARAMETERS PROVIDED IN OBJECT CONSTRUCTION OR WITH SETPARAMS
- THE GLOBAL NUMPYRANDOM RANDOM STATE IF THE ESTIMATOR’S RANDOMSTATE PARAMETER IS SET TO NONE AND
- ANY DATA OR SAMPLE PROPERTIES PASSED TO THE MOST RECENT CALL TO FITFITTRANSFORM ORFITPREDICT OR DATA

SIMILARLY PASSED IN A SEQUENCE OF CALLS TO PARTIALFIT

THE ESTIMATED MODEL IS STORED IN PUBLIC AND PRIVATE ATTRIBUTES ON THE ESTIMATOR INSTANCE FACILITATING DECODING

THROUGH PREDICTION AND TRANSFORMATION METHODS

ESTIMATORS MUST PROVIDE A FITMETHOD AND SHOULD PROVIDE SETPARAMS ANDGETPARAMS ALTHOUGH THESE ARE USUALLY

PROVIDED BY INHERITANCE FROM BASEBASEESTIMATOR

THE CORE FUNCTIONALITY OF SOME ESTIMATORS MAY ALSO BE AVAILABLE AS A FUNCTION

FEATURE EXTRACTOR

FEATURE EXTRACTORS ATRANSFORMER WHICH TAKES INPUT WHERE EACH SAMPLE IS NOT REPRESENTED AS AN ARRAYLIKE OBJECT OF FIXED LENGTH AND PRODUCES AN ARRAYLIKE OBJECT OF FEATURES FOR EACH SAMPLE AND THUS A 2DIMENSIONAL ARRAYLIKE FOR A SET OF SAMPLES IN OTHER WORDS IT LOSSILY MAPS A NONRECTANGULAR DATA REPRESENTATION INTO RECTANGULAR DATA

SCIKITLEARN USER GUIDE RELEASE 0213

FEATURE EXTRACTORS MUST IMPLEMENT AT LEAST

- FIT
- TRANSFORM
- GETFEATURENAMES

METAESTIMATOR

METAESTIMATORS

METAESTIMATOR

METAESTIMATORS ANESTIMATOR WHICH TAKES ANOTHER ESTIMATOR AS A PARAMETER EXAMPLES INCLUDE PIPELINE

PIPELINE MODELSELECTIONGRIDSEARCHCV FEATURESELECTIONSELECTFROMMODEL AND

ENSEMBLEBAGGINGCLASSIFIER

IN A METAESTIMATOR’S FITMETHOD ANY CONTAINED ESTIMATORS SHOULD BE CLONED BEFORE THEY ARE FIT ALTHOUGH FIXME

PIPELINE AND FEATUREUNION DO NOT DO THIS CURRENTLY AN EXCEPTION TO THIS IS THAT AN ESTIMATOR MAY EXPLIC

ITLY DOCUMENT THAT IT ACCEPTS A PREFITTED ESTIMATOR EG USING PREFITTRUE INFEATURESELECTION

SELECTFROMMODEL ONE KNOWN ISSUE WITH THIS IS THAT THE PREFITTED ESTIMATOR WILL LOSE ITS MODEL IF THE

METAESTIMATOR IS CLONED A METAESTIMATOR SHOULD HAVE FIT CALLED BEFORE PREDICTION EVEN IF ALL CONTAINED

ESTIMATORS ARE PREFITTED

IN CASES WHERE A METAESTIMATOR’S PRIMARY BEHAVIORS EG PREDICT ORTRANSFORM IMPLEMENTATION ARE FUNCTIONS

OF PREDICTIONTRANSFORMATION METHODS OF THE PROVIDED BASE ESTIMATOR OR MULTIPLE BASE ESTIMATORS A META

ESTIMATOR SHOULD PROVIDE AT LEAST THE STANDARD METHODS PROVIDED BY THE BASE ESTIMATOR IT MAY NOT BE POSSIBLE

TO IDENTIFY WHICH METHODS ARE PROVIDED BY THE UNDERLYING ESTIMATOR UNTIL THE METAESTIMATOR HAS BEEN FITTED

SEE ALSO DUCK TYPING FOR WHICH UTILSMETAESTIMATORSIFDELEGATEHASMETHOD MAY HELP IT

SHOULD ALSO PROVIDE OR MODIFY THE ESTIMATOR TAGS ANDCLASSES ATTRIBUTE PROVIDED BY THE BASE ESTIMATOR

METAESTIMATORS SHOULD BE CAREFUL TO VALIDATE DATA AS MINIMALLY AS POSSIBLE BEFORE PASSING IT TO AN UNDERLYING

ESTIMATOR THIS SAVES COMPUTATION TIME AND MAY FOR INSTANCE ALLOW THE UNDERLYING ESTIMATOR TO EASILY WORK

WITH DATA THAT IS NOT RECTANGULAR

OUTLIER DETECTOR

OUTLIER DETECTORS ANUNSUPERVISED BINARY PREDICTOR WHICH MODELS THE DISTINCTION BETWEEN CORE AND OUTLYING SAMPLES

OUTLIER DETECTORS MUST IMPLEMENT

- FIT
- FITPREDICT IFTRANSDUCTIVE
- PREDICT IFINDUCTIVE

INDUCTIVE OUTLIER DETECTORS MAY ALSO IMPLEMENT DECISIONFUNCTION TO GIVE A NORMALIZED INLIER SCORE WHERE OUTLIERS

HAVE SCORE BELOW 0 SCORESAMPLES MAY PROVIDE AN UNNORMALIZED SCORE PER SAMPLE

PREDICTOR

PREDICTORS ANESTIMATOR SUPPORTING PREDICT ANDOR FITPREDICT THIS ENCOMPASSES CLASSIFIER REGRESSOR OUTLIER DETEC

TORANDCLUSTERER

IN STATISTICS “PREDICTORS” REFERS TO FEATURES

REGRESSOR

REGRESSORS ASUPERVISED ORSEMISUPERVISED PREDICTOR WITH CONTINUOUS OUTPUT VALUES

REGRESSORS USUALLY INHERIT FROM BASEREGRESSORMIXIN WHICH SETS THEIR ESTIMATORATYPE ATTRIBUTE

A REGRESSOR CAN BE DISTINGUISHED FROM OTHER ESTIMATORS WITH ISREGRESSOR

42 CLASS APIS AND ESTIMATOR TYPES 673

SCIKITLEARN USER GUIDE RELEASE 0213

A REGRESSOR MUST IMPLEMENT

- FIT
- PREDICT
- SCORE

TRANSFORMER

TRANSFORMERS AN ESTIMATOR SUPPORTING TRANSFORM ANDOR FITTRANSFORM A PURELY TRANSDUCTIVE TRANSFORMER SUCH AS MANIFOLDTSNE MAY NOT IMPLEMENT TRANSFORM

VECTORIZER

VECTORIZERS SEEFEATURE EXTRACTOR

THERE ARE FURTHER APIS SPECIFICALLY RELATED TO A SMALL FAMILY OF ESTIMATORS SUCH AS

CROSSVALIDATION SPLITTER

CV SPLITTER

CROSSVALIDATION GENERATOR A NONESTIMATOR FAMILY OF CLASSES USED TO SPLIT A DATASET INTO A SEQUENCE OF TRAIN AND TEST

PORTIONS SEE CROSSVALIDATION EVALUATING ESTIMATOR PERFORMANCE BY PROVIDING SPLIT ANDGETNSPLITS METH

ODS NOTE THAT UNLIKE ESTIMATORS THESE DO NOT HAVE FITMETHODS AND DO NOT PROVIDE SETPARAMS ORGETPARAMS

PARAMETER VALIDATION MAY BE PERFORMED IN INIT

CROSSVALIDATION ESTIMATOR AN ESTIMATOR THAT HAS BUILTIN CROSSVALIDATION CAPABILITIES TO AUTOMATICALLY SELECT THE BEST

HYPERPARAMETERS SEE THE USER GUIDE SOME EXAMPLE OF CROSSVALIDATION ESTIMATORS ARE ELASTICNETCV AND

LOGISTICREGRESSIONCV CROSSVALIDATION ESTIMATORS ARE NAMED ESTIMATORCV AND TEND TO BE ROUGHLY

EQUIVALENT TO GRIDSEARCHCVESTIMATOR THE ADVANTAGE OF USING A CROSSVALIDATION ESTIMATOR

OVER THE CANONICAL ESTIMATOR CLASS ALONG WITH GRID SEARCH IS THAT THEY CAN TAKE ADVANTAGE OF WARMSTARTING

BY REUSING PRECOMPUTED RESULTS IN THE PREVIOUS STEPS OF THE CROSSVALIDATION PROCESS THIS GENERALLY LEADS TO

SPEED IMPROVEMENTS AN EXCEPTION IS THE RIDGECV CLASS WHICH CAN INSTEAD PERFORM EFFICIENT LEAVEONEOUT

CV

SCORER A NONESTIMATOR CALLABLE OBJECT WHICH EVALUATES AN ESTIMATOR ON GIVEN TEST DATA RETURNING A NUMBER UNLIKE

EVALUATION METRICS A GREATER RETURNED NUMBER MUST CORRESPOND WITH A BETTER SCORE SEE THE SCORING PARAMETER

DEFINING MODEL EVALUATION RULES

FURTHER EXAMPLES

- NEIGHBORSDISTANCEMETRIC
- GAUSSIANPROCESSKERNELSKERNEL
- TREECRITERION

43 TARGET TYPES

BINARY A CLASSIFICATION PROBLEM CONSISTING OF TWO CLASSES A BINARY TARGET MAY REPRESENTED AS FOR A MULTICLASS

PROBLEM BUT WITH ONLY TWO LABELS A BINARY DECISION FUNCTION IS REPRESENTED AS A 1D ARRAY

SEMANTICALLY ONE CLASS IS OFTEN CONSIDERED THE “POSITIVE” CLASS UNLESS OTHERWISE SPECIFIED EG USING POSLABEL

INEVALUATION METRICS WE CONSIDER THE CLASS LABEL WITH THE GREATER VALUE NUMERICALLY OR LEXICOGRAPHICALLY AS

THE POSITIVE CLASS OF LABELS 0 1 1 IS THE POSITIVE CLASS OF 1 2 2 IS THE POSITIVE CLASS OF ‘NO’ ‘YES’ ‘YES’

IS THE POSITIVE CLASS OF ‘NO’ ‘YES’ ‘NO’ IS THE POSITIVE CLASS THIS AFFECTS THE OUTPUT OF DECISIONFUNCTION FOR

INSTANCE

NOTE THAT A DATASET SAMPLED FROM A MULTICLASS YOR A CONTINUOUS YMAY APPEAR TO BE BINARY

674 CHAPTER 4 GLOSSARY OF COMMON TERMS AND API ELEMENTS



SCIKITLEARN USER GUIDE RELEASE 0213

TYPEOFTARGET WILL RETURN 'BINARY' FOR BINARY INPUT OR A SIMILAR ARRAY WITH ONLY A SINGLE CLASS PRESENT  
CONTINUOUS A REGRESSION PROBLEM WHERE EACH SAMPLE'S TARGET IS A FINITE FLOATING POINT NUMBER REPRESENTED AS A  
1DIMENSIONAL ARRAY OF FLOATS OR SOMETIMES INTS  
TYPEOFTARGET WILL RETURN 'CONTINUOUS' FOR CONTINUOUS INPUT BUT IF THE DATA IS ALL INTEGERS IT WILL BE  
IDENTIFIED AS 'MULTICLASS'

CONTINUOUS MULTIOUTPUT  
MULTIOUTPUT CONTINUOUS A REGRESSION PROBLEM WHERE EACH SAMPLE'S TARGET CONSISTS OF NOUTPUTS OUTPUTS EACH  
ONE A FINITE FLOATING POINT NUMBER FOR A FIXED INT NOUTPUTS 1 IN A PARTICULAR DATASET  
CONTINUOUS MULTIOUTPUT TARGETS ARE REPRESENTED AS MULTIPLE CONTINUOUS TARGETS HORIZONTALLY STACKED INTO AN ARRAY  
OF SHAPENSAMPLES NOUTPUTS  
TYPEOFTARGET WILL RETURN 'CONTINUOUSMULTIOUTPUT' FOR CONTINUOUS MULTIOUTPUT INPUT BUT IF THE DATA IS ALL  
INTEGERS IT WILL BE IDENTIFIED AS 'MULTICLASSMULTIOUTPUT'

MULTICLASS A CLASSIFICATION PROBLEM CONSISTING OF MORE THAN TWO CLASSES A MULTICLASS TARGET MAY BE REPRESENTED AS  
A 1DIMENSIONAL ARRAY OF STRINGS OR INTEGERS A 2D COLUMN VECTOR OF INTEGERS IE A SINGLE OUTPUT IN MULTIOUTPUT  
TERMS IS ALSO ACCEPTED  
WE DO NOT OFFICIALLY SUPPORT OTHER ORDERABLE HASHABLE OBJECTS AS CLASS LABELS EVEN IF ESTIMATORS MAY HAPPEN TO  
WORK WHEN GIVEN CLASSIFICATION TARGETS OF SUCH TYPE  
FOR SEMISUPERVISED CLASSIFICATION UNLABELED SAMPLES SHOULD HAVE THE SPECIAL LABEL 1 IN Y  
WITHIN SCKITLEARN ALL ESTIMATORS SUPPORTING BINARY CLASSIFICATION ALSO SUPPORT MULTICLASS CLASSIFICATION USING  
ONEVSREST BY DEFAULT  
APREPROCESSINGLABELENCODER HELPS TO CANONICALIZE MULTICLASS TARGETS AS INTEGERS  
TYPEOFTARGET WILL RETURN 'MULTICLASS' FOR MULTICLASS INPUT THE USER MAY ALSO WANT TO HANDLE 'BINARY'  
INPUT IDENTICALLY TO 'MULTICLASS'

MULTICLASS MULTIOUTPUT  
MULTIOUTPUT MULTICLASS A CLASSIFICATION PROBLEM WHERE EACH SAMPLE'S TARGET CONSISTS OF NOUTPUTS OUTPUTS EACH  
A CLASS LABEL FOR A FIXED INT NOUTPUTS 1 IN A PARTICULAR DATASET EACH OUTPUT HAS A FIXED SET OF AVAILABLE  
CLASSES AND EACH SAMPLE IS LABELLED WITH A CLASS FOR EACH OUTPUT AN OUTPUT MAY BE BINARY OR MULTICLASS AND IN  
THE CASE WHERE ALL OUTPUTS ARE BINARY THE TARGET IS MULTILABEL  
MULTICLASS MULTIOUTPUT TARGETS ARE REPRESENTED AS MULTIPLE MULTICLASS TARGETS HORIZONTALLY STACKED INTO AN ARRAY  
OF SHAPENSAMPLES NOUTPUTS  
XXX FOR SIMPLICITY WE MAY NOT ALWAYS SUPPORT STRING CLASS LABELS FOR MULTICLASS MULTIOUTPUT AND INTEGER CLASS  
LABELS SHOULD BE USED  
MULTIOUTPUT PROVIDES ESTIMATORS WHICH ESTIMATE MULTIOUTPUT PROBLEMS USING MULTIPLE SINGLEOUTPUT ESTIMA  
TORS THIS MAY NOT FULLY ACCOUNT FOR DEPENDENCIES AMONG THE DIFFERENT OUTPUTS WHICH METHODS NATIVELY HANDLING  
THE MULTIOUTPUT CASE EG DECISION TREES NEAREST NEIGHBORS NEURAL NETWORKS MAY DO BETTER  
TYPEOFTARGET WILL RETURN 'MULTICLASSMULTIOUTPUT' FOR MULTICLASS MULTIOUTPUT INPUT  
MULTILABEL AMULTICLASS MULTIOUTPUT TARGET WHERE EACH OUTPUT IS BINARY THIS MAY BE REPRESENTED AS A 2D DENSE  
ARRAY OR SPARSE MATRIX OF INTEGERS SUCH THAT EACH COLUMN IS A SEPARATE BINARY TARGET WHERE POSITIVE LABELS ARE  
INDICATED WITH 1 AND NEGATIVE LABELS ARE USUALLY 1 OR 0 SPARSE MULTILABEL TARGETS ARE NOT SUPPORTED EVERYWHERE  
THAT DENSE MULTILABEL TARGETS ARE SUPPORTED  
SEMANTICALLY A MULTILABEL TARGET CAN BE THOUGHT OF AS A SET OF LABELS FOR EACH SAMPLE WHILE NOT USED INTER  
NALLYPREPROCESSINGMULTILABELBINARIZER IS PROVIDED AS A UTILITY TO CONVERT FROM A LIST OF SETS  
REPRESENTATION TO A 2D ARRAY OR SPARSE MATRIX ONEHOT ENCODING A MULTICLASS TARGET WITH PREPROCESSING  
LABELBINARIZER TURNS IT INTO A MULTILABEL PROBLEM

43 TARGET TYPES 675

SCIKITLEARN USER GUIDE RELEASE 0213

TYPEOFTARGET WILL RETURN ‘MULTILABELINDICATOR’ FOR MULTILABEL INPUT WHETHER SPARSE OR DENSE

MULTIOUTPUT

MULTIOUTPUT A TARGET WHERE EACH SAMPLE HAS MULTIPLE CLASSIFICATIONREGRESSION LABELS SEE MULTICLASS MULTIOUTPUT ANDCONTINUOUS MULTIOUTPUT WE DO NOT CURRENTLY SUPPORT MODELLING MIXED CLASSIFICATION AND REGRESSION TARGETS

44 METHODS

DECISIONFUNCTION IN A FITTED CLASSIFIER OROUTLIER DETECTOR PREDICTS A “SOFT” SCORE FOR EACH SAMPLE IN RELATION TO EACH CLASS RATHER THAN THE “HARD” CATEGORICAL PREDICTION PRODUCED BY PREDICT ITS INPUT IS USUALLY ONLY SOME OBSERVED DATA X

IF THE ESTIMATOR WAS NOT ALREADY FITTED CALLING THIS METHOD SHOULD RAISE A EXCEPTIONSNOTFITTEDERROR

OUTPUT CONVENTIONS

BINARY CLASSIFICATION A 1DIMENSIONAL ARRAY WHERE VALUES STRICTLY GREATER THAN ZERO INDICATE THE POSITIVE CLASS IE THE LAST CLASS IN CLASSES

MULTICLASS CLASSIFICATION A 2DIMENSIONAL ARRAY WHERE THE ROWWISE ARGMAXIMUM IS THE PREDICTED CLASS COLUMNS ARE ORDERED ACCORDING TO CLASSES

MULTILABEL CLASSIFICATION SCIKITLEARN IS INCONSISTENT IN ITS REPRESENTATION OF MULTILABEL DECISION FUNCTIONS SOME ESTIMATORS REPRESENT IT LIKE MULTICLASS MULTIOUTPUT IE A LIST OF 2D ARRAYS EACH WITH TWO COLUMNS OTHERS REPRESENT IT WITH A SINGLE 2D ARRAY WHOSE COLUMNS CORRESPOND TO THE INDIVIDUAL BINARY CLASSIFICATION DECISIONS THE LATTER REPRESENTATION IS AMBIGUOUSLY IDENTICAL TO THE MULTICLASS CLASSIFICATION FORMAT THOUGH ITS SEMANTICS DIFFER IT SHOULD BE INTERPRETED LIKE IN THE BINARY CASE BY THRESHOLDING AT 0

TODO THIS GIST HIGHLIGHTS THE USE OF THE DIFFERENT FORMATS FOR MULTILABEL

MULTIOUTPUT CLASSIFICATION A LIST OF 2D ARRAYS CORRESPONDING TO EACH MULTICLASS DECISION FUNCTION

OUTLIER DETECTION A 1DIMENSIONAL ARRAY WHERE A VALUE GREATER THAN OR EQUAL TO ZERO INDICATES AN INLIER

FIT THEFIT METHOD IS PROVIDED ON EVERY ESTIMATOR IT USUALLY TAKES SOME SAMPLESXTARGETSYIF THE MODEL IS SUPERVISED AND POTENTIALLY OTHER SAMPLE PROPERTIES SUCH AS SAMPLEWEIGHT IT SHOULD

- CLEAR ANY PRIOR ATTRIBUTES STORED ON THE ESTIMATOR UNLESS WARMSTART IS USED
- VALIDATE AND INTERPRET ANY PARAMETERS IDEALLY RAISING AN ERROR IF INVALID
- VALIDATE THE INPUT DATA
- ESTIMATE AND STORE MODEL ATTRIBUTES FROM THE ESTIMATED PARAMETERS AND PROVIDED DATA AND
- RETURN THE NOW FITTED ESTIMATOR TO FACILITATE METHOD CHAINING

TARGET TYPES DESCRIBES POSSIBLE FORMATS FOR Y

FITPREDICT USED ESPECIALLY FOR UNSUPERVISED TRANSDUCTIVE ESTIMATORS THIS FITS THE MODEL AND RETURNS THE PRE DITIONS SIMILAR TO PREDICT ON THE TRAINING DATA IN CLUSTERERS THESE PREDICTIONS ARE ALSO STORED IN THE LABELS ATTRIBUTE AND THE OUTPUT OF FITPREDICTX IS USUALLY EQUIVALENT TO FITXPREDICTX THE PA

RAMETERS TO FITPREDICT ARE THE SAME AS THOSE TO FIT

FITTRANSFORM A METHOD ON TRANSFORMERS WHICH FITS THE ESTIMATOR AND RETURNS THE TRANSFORMED TRAINING DATA IT TAKES PARAMETERS AS IN FITAND ITS OUTPUT SHOULD HAVE THE SAME SHAPE AS CALLING FITX

TRANSFORMX THERE ARE NONETHELESS RARE CASES WHERE FITTRANSFORMX ANDFITX

TRANSFORMX DO NOT RETURN THE SAME VALUE WHEREIN TRAINING DATA NEEDS TO BE HANDLED DIFFERENTLY DUE TO MODEL BLENDING IN STACKED ENSEMBLES FOR INSTANCE SUCH CASES SHOULD BE CLEARLY DOCUMENTED

TRANSDUCTIVE TRANSFORMERS MAY ALSO PROVIDE FITTRANSFORM BUT NOT TRANSFORM

676 CHAPTER 4 GLOSSARY OF COMMON TERMS AND API ELEMENTS

SCIKITLEARN USER GUIDE RELEASE 0213

ONE REASON TO IMPLEMENT FITTRANSFORM IS THAT PERFORMING FIT ANDTRANSFORM SEPARATELY WOULD BE LESS EFFICIENT THAN TOGETHER BASETRANSFORMERMIXIN PROVIDES A DEFAULT IMPLEMENTATION PROVIDING A CONSISTENT INTERFACE ACROSS TRANSFORMERS WHERE FITTRANSFORM IS OR IS NOT SPECIALISED ININDUCTIVE LEARNING - WHERE THE GOAL IS TO LEARN A GENERALISED MODEL THAT CAN BE APPLIED TO NEW DATA - USERS SHOULD BE CAREFUL NOT TO APPLY FITTRANSFORM TO THE ENTIRETY OF A DATASET IE TRAINING AND TEST DATA TOGETHER BEFORE FURTHER MODELLING AS THIS RESULTS IN DATA LEAKAGE

GETFEATURENAMES PRIMARILY FOR FEATURE EXTRACTORS BUT ALSO USED FOR OTHER TRANSFORMERS TO PROVIDE STRING NAMES FOR EACH COLUMN IN THE OUTPUT OF THE ESTIMATOR'S TRANSFORM METHOD IT OUTPUTS A LIST OF STRINGS AND MAY TAKE A LIST OF STRINGS AS INPUT CORRESPONDING TO THE NAMES OF INPUT COLUMNS FROM WHICH OUTPUT COLUMN NAMES CAN BE GENERATED BY DEFAULT INPUT FEATURES ARE NAMED X0 X1

GETNSPLITS ON A CV SPLITTER NOT AN ESTIMATOR RETURNS THE NUMBER OF ELEMENTS ONE WOULD GET IF ITERATING THROUGH THE RETURN VALUE OF SPLIT GIVEN THE SAME PARAMETERS TAKES THE SAME PARAMETERS AS SPLIT GETPARAMS GETS ALL PARAMETERS AND THEIR VALUES THAT CAN BE SET USING SETPARAMS A PARAMETER DEEP CAN BE USED WHEN SET TO FALSE TO ONLY RETURN THOSE PARAMETERS NOT INCLUDING IE NOT DUE TO INDIRECTION VIA CONTAINED ESTIMATORS

MOST ESTIMATORS ADOPT THE DEFINITION FROM BASEBASEESTIMATOR WHICH SIMPLY ADOPTS THE PARAMETERS DEFINED FORINIT PIPELINEPIPELINE AMONG OTHERS REIMPLEMENTS GETPARAMS TO DECLARE THE ESTIMATORS NAMED IN ITS STEPS PARAMETERS AS THEMSELVES BEING PARAMETERS

PARTIALFIT FACILITATES FITTING AN ESTIMATOR IN AN ONLINE FASHION UNLIKE FIT REPEATEDLY CALLING PARTIALFIT DOES NOT CLEAR THE MODEL BUT UPDATES IT WITH RESPECT TO THE DATA PROVIDED THE PORTION OF DATA PROVIDED TO PARTIALFIT MAY BE CALLED A MINIBATCH EACH MINIBATCH MUST BE OF CONSISTENT SHAPE ETC IN ITERATIVE ESTIMATORS PARTIALFIT OFTEN ONLY PERFORMS A SINGLE ITERATION

PARTIALFIT MAY ALSO BE USED FOR OUTFOCORE LEARNING ALTHOUGH USUALLY LIMITED TO THE CASE WHERE LEARNING CAN BE PERFORMED ONLINE IE THE MODEL IS USABLE AFTER EACH PARTIALFIT AND THERE IS NO SEPARATE PROCESSING NEEDED TO FINALIZE THE MODEL CLUSTERBIRCH INTRODUCES THE CONVENTION THAT CALLING PARTIALFITX WILL PRODUCE A MODEL THAT IS NOT FINALIZED BUT THE MODEL CAN BE FINALIZED BY CALLING PARTIALFIT IE WITHOUT PASSING A FURTHER MINIBATCH

GENERALLY ESTIMATOR PARAMETERS SHOULD NOT BE MODIFIED BETWEEN CALLS TO PARTIALFIT ALTHOUGH PARTIALFIT SHOULD VALIDATE THEM AS WELL AS THE NEW MINIBATCH OF DATA IN CONTRAST WARMSTART IS USED TO REPEATEDLY FIT THE SAME ESTIMATOR WITH THE SAME DATA BUT VARYING PARAMETERS LIKEFITPARTIALFIT SHOULD RETURN THE ESTIMATOR OBJECT

TO CLEAR THE MODEL A NEW ESTIMATOR SHOULD BE CONSTRUCTED FOR INSTANCE WITH BASECLONE PREDICT MAKES A PREDICTION FOR EACH SAMPLE USUALLY ONLY TAKING XAS INPUT BUT SEE UNDER REGRESSOR OUTPUT CONVENTIONS BELOW IN A CLASSIFIER ORREGRESSOR THIS PREDICTION IS IN THE SAME TARGET SPACE USED IN FITTING EG ONE OF 'RED' 'AMBER' 'GREEN' IF THE YIN FITTING CONSISTED OF THESE STRINGS DESPITE THIS EVEN WHEN YPASSED TO FIT IS A LIST OR OTHER ARRAYLIKE THE OUTPUT OF PREDICT SHOULD ALWAYS BE AN ARRAY OR SPARSE MATRIX IN A CLUSTERER OR OUTLIER DETECTOR THE PREDICTION IS AN INTEGER

IF THE ESTIMATOR WAS NOT ALREADY FITTED CALLING THIS METHOD SHOULD RAISE A EXCEPTIONSNOTFITTEDERROR OUTPUT CONVENTIONS

CLASSIFIER AN ARRAY OF SHAPE NSAMPLESNSAMPLES NOUTPUTS MULTILABEL DATA MAY BE REPRESENTED AS A SPARSE MATRIX IF A SPARSE MATRIX WAS USED IN FITTING EACH ELEMENT SHOULD BE ONE OF THE VALUES IN THE CLASSIFIER'S CLASSES ATTRIBUTE

CLUSTERER AN ARRAY OF SHAPE NSAMPLES WHERE EACH VALUE IS FROM 0 TO NCLUSTERS 1 IF THE CORRESPONDING SAMPLE IS CLUSTERED AND 1 IF THE SAMPLE IS NOT CLUSTERED AS IN CLUSTERDBSCAN

OUTLIER DETECTOR AN ARRAY OF SHAPE NSAMPLES WHERE EACH VALUE IS 1 FOR AN OUTLIER AND 1 OTHERWISE 44 METHODS 677

SCIKITLEARN USER GUIDE RELEASE 0213

REGRESSOR A NUMERIC ARRAY OF SHAPE NSAMPLES USUALLY FLOAT64 SOME REGRESSORS HAVE EXTRA OPTIONS IN THEIRPREDICT METHOD ALLOWING THEM TO RETURN STANDARD DEVIATION RETURNSTDTRUE OR COVARIANCE RETURNCOVTRUE RELATIVE TO THE PREDICTED VALUE IN THIS CASE THE RETURN VALUE IS A TUPLE OF ARRAYS CORRESPONDING TO PREDICTION MEAN STD COV AS REQUIRED  
PREDICTLOGPROBA THE NATURAL LOGARITHM OF THE OUTPUT OF PREDICTPROBA PROVIDED TO FACILITATE NUMERICAL STABILITY

PREDICTPROBA A METHOD IN CLASSIFIERS ANDCLUSTERERS THAT ARE ABLE TO RETURN PROBABILITY ESTIMATES FOR EACH CLASSCLUSTER ITS INPUT IS USUALLY ONLY SOME OBSERVED DATA X

IF THE ESTIMATOR WAS NOT ALREADY FITTED CALLING THIS METHOD SHOULD RAISE A EXCEPTIONSNOTFITTEDERROR  
OUTPUT CONVENTIONS ARE LIKE THOSE FOR DECISIONFUNCTION EXCEPT IN THE BINARY CLASSIFICATION CASE WHERE ONE COLUMN IS OUTPUT FOR EACH CLASS WHILE DECISIONFUNCTION OUTPUTS A 1D ARRAY FOR BINARY AND MULTICLASS PREDICTIONS EACH ROW SHOULD ADD TO 1

LIKE OTHER METHODS PREDICTPROBA SHOULD ONLY BE PRESENT WHEN THE ESTIMATOR CAN MAKE PROBABILISTIC PREDICTIONS SEE DUCK TYPING THIS MEANS THAT THE PRESENCE OF THE METHOD MAY DEPEND ON ESTIMATOR PARAMETERS EG INLINEARMODELSGDCLASSIFIER OR TRAINING DATA EG IN MODELSELECTIONGRIDSEARCHCV

AND MAY ONLY APPEAR AFTER FITTING

SCORE A METHOD ON AN ESTIMATOR USUALLY A PREDICTOR WHICH EVALUATES ITS PREDICTIONS ON A GIVEN DATASET AND RETURNS A SINGLE NUMERICAL SCORE A GREATER RETURN VALUE SHOULD INDICATE BETTER PREDICTIONS ACCURACY IS USED FOR CLASSIFIERS AND R2 FOR REGRESSORS BY DEFAULT

IF THE ESTIMATOR WAS NOT ALREADY FITTED CALLING THIS METHOD SHOULD RAISE A EXCEPTIONSNOTFITTEDERROR  
SOME ESTIMATORS IMPLEMENT A CUSTOM ESTIMATORSPECIFIC SCORE FUNCTION OFTEN THE LIKELIHOOD OF THE DATA UNDER THE MODEL

SCORESAMPLES TODO

IF THE ESTIMATOR WAS NOT ALREADY FITTED CALLING THIS METHOD SHOULD RAISE A EXCEPTIONSNOTFITTEDERROR  
SETPARAMS AVAILABLE IN ANY ESTIMATOR TAKES KEYWORD ARGUMENTS CORRESPONDING TO KEYS IN GETPARAMS EACH IS PROVIDED A NEW VALUE TO ASSIGN SUCH THAT CALLING GETPARAMS AFTERSETPARAMS WILL REFLECT THE CHANGED PARAMETERS MOST ESTIMATORS USE THE IMPLEMENTATION IN BASEBASEESTIMATOR WHICH HANDLES NESTED PARAMETERS AND OTHERWISE SETS THE PARAMETER AS AN ATTRIBUTE ON THE ESTIMATOR THE METHOD IS OVERRIDDEN IN PIPELINE PIPELINE AND RELATED ESTIMATORS

SPLIT ON A CV SPLITTER NOT AN ESTIMATOR THIS METHOD ACCEPTS PARAMETERS XYGROUPS WHERE ALL MAY BE OPTIONAL AND RETURNS AN ITERATOR OVER TRAINIDX TESTIDX PAIRS EACH OF TRAINTESTIDX IS A 1D INTEGER ARRAY WITH VALUES FROM 0 FROM XSHAPE0 1 OF ANY LENGTH SUCH THAT NO VALUES APPEAR IN BOTH SOME TRAINIDX AND ITS CORRESPONDING TESTIDX

TRANSFORM IN A TRANSFORMER TRANSFORMS THE INPUT USUALLY ONLY X INTO SOME TRANSFORMED SPACE CONVENTIONALLY NOTATED AS XT OUTPUT IS AN ARRAY OR SPARSE MATRIX OF LENGTH NSAMPLES AND WITH NUMBER OF COLUMNS FIXED AFTER FITTING

IF THE ESTIMATOR WAS NOT ALREADY FITTED CALLING THIS METHOD SHOULD RAISE A EXCEPTIONSNOTFITTEDERROR  
45 PARAMETERS

THESE COMMON PARAMETER NAMES SPECIFICALLY USED IN ESTIMATOR CONSTRUCTION SEE CONCEPT PARAMETER SOMETIMES ALSO APPEAR AS PARAMETERS OF FUNCTIONS OR NONESTIMATOR CONSTRUCTORS

CLASSWEIGHT USED TO SPECIFY SAMPLE WEIGHTS WHEN FITTING CLASSIFIERS AS A FUNCTION OF THE TARGET CLASS WHERE SAMPLEWEIGHT IS ALSO SUPPORTED AND GIVEN IT IS MULTIPLIED BY THE CLASSWEIGHT CONTRIBUTION SIMILARLY

678 CHAPTER 4 GLOSSARY OF COMMON TERMS AND API ELEMENTS

SCIKITLEARN USER GUIDE RELEASE 0213

WHERECLASSWEIGHT IS USED IN A MULTIOUTPUT INCLUDING MULTILABEL TASKS THE WEIGHTS ARE MULTIPLIED ACROSS OUTPUTS IE COLUMNS OF Y

BY DEFAULT ALL SAMPLES HAVE EQUAL WEIGHT SUCH THAT CLASSES ARE EFFECTIVELY WEIGHTED BY THEIR THEIR PREVALENCE IN THE TRAINING DATA THIS COULD BE ACHIEVED EXPLICITLY WITH CLASSWEIGHTLABEL1 1 LABEL2 1

FOR ALL CLASS LABELS

MORE GENERALLY CLASSWEIGHT IS SPECIFIED AS A DICT MAPPING CLASS LABELS TO WEIGHTS CLASSLABEL WEIGHT SUCH THAT EACH SAMPLE OF THE NAMED CLASS IS GIVEN THAT WEIGHT

CLASSWEIGHTBALANCED CAN BE USED TO GIVE ALL CLASSES EQUAL WEIGHT BY GIVING EACH SAMPLE A

WEIGHT INVERSELY RELATED TO ITS CLASS’S PREVALENCE IN THE TRAINING DATA NSAMPLES NCLASSES NP

BINCOUNTY CLASS WEIGHTS WILL BE USED DIFFERENTLY DEPENDING ON THE ALGORITHM FOR LINEAR MODELS SUCH

AS LINEAR SVM OR LOGISTIC REGRESSION THE CLASS WEIGHTS WILL ALTER THE LOSS FUNCTION BY WEIGHTING THE LOSS OF EACH

SAMPLE BY ITS CLASS WEIGHT FOR TREEBASED ALGORITHMS THE CLASS WEIGHTS WILL BE USED FOR REWEIGHTING THE SPLITTING

CRITERION NOTE HOWEVER THAT THIS REBALANCING DOES NOT TAKE THE WEIGHT OF SAMPLES IN EACH CLASS INTO ACCOUNT

FOR MULTIOUTPUT CLASSIFICATION A LIST OF DICTS IS USED TO SPECIFY WEIGHTS FOR EACH OUTPUT FOR EXAMPLE FOR FOUR

CLASS MULTILABEL CLASSIFICATION WEIGHTS SHOULD BE 0 1 1 1 0 1 1 5 0 1

1 1 0 1 1 1 INSTEAD OF11 25 31 41

THECLASSWEIGHT PARAMETER IS VALIDATED AND INTERPRETED WITH UTILSCOMPUTECLASSWEIGHT

CV DETERMINES A CROSS VALIDATION SPLITTING STRATEGY AS USED IN CROSSVALIDATION BASED ROUTINES CV

IS ALSO AVAILABLE IN ESTIMATORS SUCH AS MULTIOUTPUTCLASSIFIERCHAIN ORCALIBRATION

CALIBRATEDCLASSIFIERCV WHICH USE THE PREDICTIONS OF ONE ESTIMATOR AS TRAINING DATA FOR ANOTHER TO

NOT OVERFIT THE TRAINING SUPERVISION

POSSIBLE INPUTS FOR CVARE USUALLY

- AN INTEGER SPECIFYING THE NUMBER OF FOLDS IN KFOLD CROSS VALIDATION KFOLD WILL BE STRATIFIED OVER CLASSES

IF THE ESTIMATOR IS A CLASSIFIER DETERMINED BY BASEISCLASSIFIER AND THE TARGETS MAY REPRESENT A

BINARY OR MULTICLASS BUT NOT MULTIOUTPUT CLASSIFICATION PROBLEM DETERMINED BY UTILSMULTICLASS

TYPEOFTARGET

- ACROSSVALIDATION SPLITTER INSTANCE REFER TO THE USER GUIDE FOR SPLITTERS AVAILABLE WITHIN SCIKITLEARN

- AN ITERABLE YIELDING TRAINTEST SPLITS

WITH SOME EXCEPTIONS ESPECIALLY WHERE NOT USING CROSS VALIDATION AT ALL IS AN OPTION THE DEFAULT IS 3FOLD AND

WILL CHANGE TO 5FOLD IN VERSION 022

CVVALUES ARE VALIDATED AND INTERPRETED WITH UTILSCHECKCV

KERNEL TODO

MAXITER FOR ESTIMATORS INVOLVING ITERATIVE OPTIMIZATION THIS DETERMINES THE MAXIMUM NUMBER OF ITERA

TIONS TO BE PERFORMED IN FIT IFMAXITER ITERATIONS ARE RUN WITHOUT CONVERGENCE A EXCEPTIONS

CONVERGENCEWARNING SHOULD BE RAISED NOTE THAT THE INTERPRETATION OF “A SINGLE ITERATION” IS INCONSISTENT

ACROSS ESTIMATORS SOME BUT NOT ALL USE IT TO MEAN A SINGLE EPOCH IE A PASS OVER EVERY SAMPLE IN THE DATA

FIXME PERHAPS WE SHOULD HAVE SOME COMMON TESTS ABOUT THE RELATIONSHIP BETWEEN CONVERGENCEWARNING AND

MAXITER

MEMORY SOME ESTIMATORS MAKE USE OF JOBLIBMEMORY TO STORE PARTIAL SOLUTIONS DURING FITTING THUS WHEN FIT

IS CALLED AGAIN THOSE PARTIAL SOLUTIONS HAVE BEEN MEMOIZED AND CAN BE REUSED

AMEMORY PARAMETER CAN BE SPECIFIED AS A STRING WITH A PATH TO A DIRECTORY OR A JOBLIBMEMORY INSTANCE

OR AN OBJECT WITH A SIMILAR INTERFACE IE A CACHE METHOD CAN BE USED

MEMORY VALUES ARE VALIDATED AND INTERPRETED WITH UTILSVALIDATIONCHECKMEMORY

45 PARAMETERS 679

SCIKITLEARN USER GUIDE RELEASE 0213

METRIC AS A PARAMETER THIS IS THE SCHEME FOR DETERMINING THE DISTANCE BETWEEN TWO DATA POINTS SEE METRICS PAIRWISEDISTANCES IN PRACTICE FOR SOME ALGORITHMS AN IMPROPER DISTANCE METRIC ONE THAT DOES NOT OBEY THE TRIANGLE INEQUALITY SUCH AS COSINE DISTANCE MAY BE USED

XXX HIERARCHICAL CLUSTERING USES AFFINITY WITH THIS MEANING

WE ALSO USE METRIC TO REFER TO EVALUATION METRICS BUT AVOID USING THIS SENSE AS A PARAMETER NAME

NCOMPONENTS THE NUMBER OF FEATURES WHICH A TRANSFORMER SHOULD TRANSFORM THE INPUT INTO SEE COMPONENTS FOR THE SPECIAL CASE OF AFFINE PROJECTION

NITERNOCHANGE NUMBER OF ITERATIONS WITH NO IMPROVEMENT TO WAIT BEFORE STOPPING THE ITERATIVE PROCEDURE THIS IS ALSO KNOWN AS A PATIENCE PARAMETER IT IS TYPICALLY USED WITH EARLY STOPPING TO AVOID STOPPING TOO EARLY

NJOBS THIS IS USED TO SPECIFY HOW MANY CONCURRENT PROCESSESTHREADS SHOULD BE USED FOR PARALLELIZED ROUTINES SCIKITLEARN USES ONE PROCESSOR FOR ITS PROCESSING BY DEFAULT ALTHOUGH IT ALSO MAKES USE OF NUMPY WHICH MAY BE CONFIGURED TO USE A THREADED NUMERICAL PROCESSOR LIBRARY LIKE MKL SEE FAQ

NJOBS IS AN INT SPECIFYING THE MAXIMUM NUMBER OF CONCURRENTLY RUNNING JOBS IF SET TO 1 ALL CPUS ARE USED IF 1 IS GIVEN NO JOBLIB LEVEL PARALLELISM IS USED AT ALL WHICH IS USEFUL FOR DEBUGGING EVEN WITH NJOBS 1 PARALLELISM MAY OCCUR DUE TO NUMERICAL PROCESSING LIBRARIES SEE FAQ FOR NJOBS BELOW 1 NCPUS 1

NJOBS ARE USED THUS FOR NJOBS 2 ALL CPUS BUT ONE ARE USED

NJOBSNONE MEANS UNSET IT WILL GENERALLY BE INTERPRETED AS NJOBS1 UNLESS THE CURRENT JOBLIB PARALLEL BACKEND CONTEXT SPECIFIES OTHERWISE

THE USE OFNJOBS BASED PARALLELISM IN ESTIMATORS VARIES

- MOST OFTEN PARALLELISM HAPPENS IN FITTING BUT SOMETIMES PARALLELISM HAPPENS IN PREDICTION EG IN RANDOM FORESTS
- SOME PARALLELISM USES A MULTITHREADING BACKEND BY DEFAULT SOME A MULTIPROCESSING BACKEND IT IS POSSIBLE TO OVERRIDE THE DEFAULT BACKEND BY USING SKLEARNUTILSPARALLELBACKEND
- WHETHER PARALLEL PROCESSING IS HELPFUL AT IMPROVING RUNTIME DEPENDS ON MANY FACTORS AND IT'S USUALLY A GOOD IDEA TO EXPERIMENT RATHER THAN ASSUMING THAT INCREASING THE NUMBER OF JOBS IS ALWAYS A GOOD THING IT CAN BE HIGHLY DETRIMENTAL TO PERFORMANCE TO RUN MULTIPLE COPIES OF SOME ESTIMATORS OR FUNCTIONS IN PARALLEL

NESTED USES OF NJOBS BASED PARALLELISM WITH THE SAME BACKEND WILL RESULT IN AN EXCEPTION SO GRIDSEARCHCVONEVSRESTCLASSIFIERSVC NJOBS2 NJOBS2 WON'T WORK

WHENNJOBS IS NOT 1 THE ESTIMATOR BEING PARALLELIZED MUST BE PICKLABLE THIS MEANS FOR INSTANCE THAT LAMBDAS CANNOT BE USED AS ESTIMATOR PARAMETERS

RANDOMSTATE WHENEVER RANDOMIZATION IS PART OF A SCIKITLEARN ALGORITHM A RANDOMSTATE PARAMETER MAY BE PROVIDED TO CONTROL THE RANDOM NUMBER GENERATOR USED NOTE THAT THE MERE PRESENCE OF RANDOMSTATE DOESN'T MEAN THAT RANDOMIZATION IS ALWAYS USED AS IT MAY BE DEPENDENT ON ANOTHER PARAMETER EG SHUFFLE BEING SET

RANDOMSTATE 'S VALUE MAY BE

- NONE DEFAULT USE THE GLOBAL RANDOM STATE FROM NUMPYRANDOM
- AN INTEGER USE A NEW RANDOM NUMBER GENERATOR SEEDED BY THE GIVEN INTEGER TO MAKE A RANDOMIZED ALGORITHM DETERMINISTIC IE RUNNING IT MULTIPLE TIMES WILL PRODUCE THE SAME RESULT AN ARBITRARY INTEGER

RANDOMSTATE CAN BE USED HOWEVER IT MAY BE WORTHWHILE CHECKING THAT YOUR RESULTS ARE STABLE ACROSS A NUMBER OF DIFFERENT DISTINCT RANDOM SEEDS POPULAR INTEGER RANDOM SEEDS ARE 0 AND 42

ANUMPYRANDOMRANDOMSTATE INSTANCE USE THE PROVIDED RANDOM STATE ONLY AFFECTING OTHER USERS OF THE SAME RANDOM STATE INSTANCE CALLING FIT MULTIPLE TIMES WILL REUSE THE SAME INSTANCE AND WILL PRODUCE DIFFERENT RESULTS

680 CHAPTER 4 GLOSSARY OF COMMON TERMS AND API ELEMENTS

SCIKITLEARN USER GUIDE RELEASE 0213

UTILSCHECKRANDOMSTATE IS USED INTERNALLY TO VALIDATE THE INPUT RANDOMSTATE AND RETURN A RANDOMSTATE INSTANCE

SCORING SPECIFIES THE SCORE FUNCTION TO BE MAXIMIZED USUALLY BY CROSS VALIDATION OR - IN SOME CASES - MULTIPLE SCORE FUNCTIONS TO BE REPORTED THE SCORE FUNCTION CAN BE A STRING ACCEPTED BY METRICSGETSCORER OR A CALLABLE SCORER NOT TO BE CONFUSED WITH AN EVALUATION METRIC AS THE LATTER HAVE A MORE DIVERSE API SCORING MAY ALSO BE SET TO NONE IN WHICH CASE THE ESTIMATOR'S SCORE METHOD IS USED SEE THE SCORING PARAMETER DEFINING MODEL EVALUATION RULES IN THE USER GUIDE

WHERE MULTIPLE METRICS CAN BE EVALUATED SCORING MAY BE GIVEN EITHER AS A LIST OF UNIQUE STRINGS OR A DICT WITH NAMES AS KEYS AND CALLABLES AS VALUES NOTE THAT THIS DOES NOTSPECIFY WHICH SCORE FUNCTION IS TO BE MAXIMISED AND ANOTHER PARAMETER SUCH AS REFIT MAY BE USED FOR THIS PURPOSE

THESCORING PARAMETER IS VALIDATED AND INTERPRETED USING METRICSCHECKSCORING

VERBOSE LOGGING IS NOT HANDLED VERY CONSISTENTLY IN SCIKITLEARN AT PRESENT BUT WHEN IT IS PROVIDED AS AN OPTION THEVERBOSE PARAMETER IS USUALLY AVAILABLE TO CHOOSE NO LOGGING SET TO FALSE ANY TRUE VALUE SHOULD ENABLE SOME LOGGING BUT LARGER INTEGERS EG ABOVE 10 MAY BE NEEDED FOR FULL VERBOSITY VERBOSE LOGS ARE USUALLY PRINTED TO STANDARD OUTPUT ESTIMATORS SHOULD NOT PRODUCE ANY OUTPUT ON STANDARD OUTPUT WITH THE DEFAULT VERBOSE SETTING

WARMSTART WHEN FITTING AN ESTIMATOR REPEATEDLY ON THE SAME DATASET BUT FOR MULTIPLE PARAMETER VALUES SUCH AS TO FIND THE VALUE MAXIMIZING PERFORMANCE AS IN GRID SEARCH IT MAY BE POSSIBLE TO REUSE ASPECTS OF THE MODEL LEARNT FROM THE PREVIOUS PARAMETER VALUE SAVING TIME WHEN WARMSTART IS TRUE THE EXISTING FITTED MODEL ATTRIBUTES ARE USED TO INITIALISE THE NEW MODEL IN A SUBSEQUENT CALL TO FIT

NOTE THAT THIS IS ONLY APPLICABLE FOR SOME MODELS AND SOME PARAMETERS AND EVEN SOME ORDERS OF PARAMETER VALUES FOR EXAMPLE WARMSTART MAY BE USED WHEN BUILDING RANDOM FORESTS TO ADD MORE TREES TO THE FOREST INCREASING NESTIMATORS BUT NOT TO REDUCE THEIR NUMBER

PARTIALFIT ALSO RETAINS THE MODEL BETWEEN CALLS BUT DIFFERS WITH WARMSTART THE PARAMETERS CHANGE AND THE DATA IS MOREORLESS CONSTANT ACROSS CALLS TO FIT WITHPARTIALFIT THE MINIBATCH OF DATA CHANGES AND MODEL PARAMETERS STAY FIXED

THERE ARE CASES WHERE YOU WANT TO USE WARMSTART TO FIT ON DIFFERENT BUT CLOSELY RELATED DATA FOR EXAM PLE ONE MAY INITIALLY FIT TO A SUBSET OF THE DATA THEN FINETUNE THE PARAMETER SEARCH ON THE FULL DATASET FOR CLASSIFICATION ALL DATA IN A SEQUENCE OF WARMSTART CALLS TOFIT MUST INCLUDE SAMPLES FROM EACH CLASS

46 ATTRIBUTES

SEE CONCEPT ATTRIBUTE

CLASSES A LIST OF CLASS LABELS KNOWN TO THE CLASSIFIER MAPPING EACH LABEL TO A NUMERICAL INDEX USED IN THE MODEL REPRESENTATION OUR OUTPUT FOR INSTANCE THE ARRAY OUTPUT FROM PREDICTPROBA HAS COLUMNS ALIGNED WITH CLASSES FOR MULTIOUTPUT CLASSIFIERS CLASSES SHOULD BE A LIST OF LISTS WITH ONE CLASS LISTING FOR EACH OUTPUT FOR EACH OUTPUT THE CLASSES SHOULD BE SORTED NUMERICALLY OR LEXICOGRAPHICALLY FOR STRINGS CLASSES AND THE MAPPING TO INDICES IS OFTEN MANAGED WITH PREPROCESSINGLABELENCODER

COMPONENTS AN AFFINE TRANSFORMATION MATRIX OF SHAPE NCOMPONENTS NFEATURES USED IN MANY LIN EARTRANSFORMERS WHERE NCOMPONENTS IS THE NUMBER OF OUTPUT FEATURES AND NFEATURES IS THE NUMBER OF INPUT FEATURES

SEE ALSO COMPONENTS WHICH IS A SIMILAR ATTRIBUTE FOR LINEAR PREDICTORS

COEF THE WEIGHTCOEFFICIENT MATRIX OF A GENERALISED LINEAR MODEL PREDICTOR OF SHAPENFEATURES FOR BINARY CLASSIFICATION AND SINGLEOUTPUT REGRESSION NCLASSES NFEATURES FOR MULTICLASS CLASSIFICATION AND NTARGETS NFEATURES FOR MULTIOUTPUT REGRESSION NOTE THIS DOES NOT INCLUDE THE INTERCEPT OR BIAS TERM WHICH IS STORED IN INTERCEPT

46 ATTRIBUTES 681

SCIKITLEARN USER GUIDE RELEASE 0213

WHEN AVAILABLE FEATUREIMPORTANCES IS NOT USUALLY PROVIDED AS WELL BUT CAN BE CALCULATED AS THE NORM OF EACH FEATURE’S ENTRY IN COEF

SEE ALSO COMPONENTS WHICH IS A SIMILAR ATTRIBUTE FOR LINEAR TRANSFORMERS

EMBEDDING AN EMBEDDING OF THE TRAINING DATA IN MANIFOLD LEARNING ESTIMATORS WITH SHAPE NSAMPLES

NCOMPONENTS IDENTICAL TO THE OUTPUT OF FITTRANSFORM SEE ALSO LABELS

NITER THE NUMBER OF ITERATIONS ACTUALLY PERFORMED WHEN FITTING AN ITERATIVE ESTIMATOR THAT MAY STOP UPON CONVERGENCE SEE ALSO MAXITER

FEATUREIMPORTANCES A VECTOR OF SHAPE NFEATURES AVAILABLE IN SOME PREDICTORS TO PROVIDE A RELATIVE MEASURE OF THE IMPORTANCE OF EACH FEATURE IN THE PREDICTIONS OF THE MODEL

LABELS A VECTOR CONTAINING A CLUSTER LABEL FOR EACH SAMPLE OF THE TRAINING DATA IN CLUSTERERS IDENTICAL TO THE OUTPUT OFFITPREDICT SEE ALSO EMBEDDING

47 DATA AND SAMPLE PROPERTIES

SEE CONCEPT SAMPLE PROPERTY

GROUPS USED IN CROSS VALIDATION ROUTINES TO IDENTIFY SAMPLES WHICH ARE CORRELATED EACH VALUE IS AN IDENTIFIER SUCH THAT IN A SUPPORTING CV SPLITTER SAMPLES FROM SOME GROUPS VALUE MAY NOT APPEAR IN BOTH A TRAINING SET AND ITS CORRESPONDING TEST SET SEE CROSSVALIDATION ITERATORS FOR GROUPED DATA

SAMPLEWEIGHT A RELATIVE WEIGHT FOR EACH SAMPLE INTUITIVELY IF ALL WEIGHTS ARE INTEGERS A WEIGHTED MODEL OR SCORE SHOULD BE EQUIVALENT TO THAT CALCULATED WHEN REPEATING THE SAMPLE THE NUMBER OF TIMES SPECIFIED IN THE WEIGHT WEIGHTS MAY BE SPECIFIED AS FLOATS SO THAT SAMPLE WEIGHTS ARE USUALLY EQUIVALENT UP TO A CONSTANT POSITIVE SCALING FACTOR

FIXME IS THIS INTERPRETATION ALWAYS THE CASE IN PRACTICE WE HAVE NO COMMON TESTS

SOME ESTIMATORS SUCH AS DECISION TREES SUPPORT NEGATIVE WEIGHTS FIXME THIS FEATURE OR ITS ABSENCE MAY NOT BE TESTED OR DOCUMENTED IN MANY ESTIMATORS

THIS IS NOT ENTIRELY THE CASE WHERE OTHER PARAMETERS OF THE MODEL CONSIDER THE NUMBER OF SAMPLES IN A REGION AS WITHMINSAMPLES INCLUSTERDBSCAN IN THIS CASE A COUNT OF SAMPLES BECOMES TO A SUM OF THEIR WEIGHTS

IN CLASSIFICATION SAMPLE WEIGHTS CAN ALSO BE SPECIFIED AS A FUNCTION OF CLASS WITH THE CLASSWEIGHT ESTIMATOR PARAMETER

XDENOTES DATA THAT IS OBSERVED AT TRAINING AND PREDICTION TIME USED AS INDEPENDENT VARIABLES IN LEARNING THE NOTATION IS UPPERCASE TO DENOTE THAT IT IS ORDINARILY A MATRIX SEE RECTANGULAR WHEN A MATRIX EACH SAMPLE MAY BE REPRESENTED BY A FEATURE VECTOR OR A VECTOR OF PRECOMPUTED DISSIMILARITY WITH EACH TRAINING SAMPLE X MAY ALSO NOT BE A MATRIX AND MAY REQUIRE A FEATURE EXTRACTOR OR APAIRWISE METRIC TO TURN IT INTO ONE BEFORE LEARNING A MODEL

XT SHORTHAND FOR “TRANSFORMED X”

Y

YDENOTES DATA THAT MAY BE OBSERVED AT TRAINING TIME AS THE DEPENDENT VARIABLE IN LEARNING BUT WHICH IS UNAVAILABLE AT PREDICTION TIME AND IS USUALLY THE TARGET OF PREDICTION THE NOTATION MAY BE UPPERCASE TO DENOTE THAT IT IS A MATRIX REPRESENTING MULTIOUTPUT TARGETS FOR INSTANCE BUT USUALLY WE USE YAND SOMETIMES DO SO EVEN WHEN MULTIPLE OUTPUTS ARE ASSUMED

682 CHAPTER 4 GLOSSARY OF COMMON TERMS AND API ELEMENTS



CHAPTER  
FIVE  
EXAMPLES  
51 MISCELLANEOUS EXAMPLES  
MISCELLANEOUS AND INTRODUCTORY EXAMPLES FOR SCIKITLEARN  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
511 COMPACT ESTIMATOR REPRESENTATIONS  
THIS EXAMPLE ILLUSTRATES THE USE OF THE PRINTCHANGEDONLY GLOBAL PARAMETER  
SETTING PRINTCHANGEDONLY TO TRUE WILL ALTERATE THE REPRESENTATION OF ESTIMATORS TO ONLY SHOW THE PARAMETERS THAT HAVE BEEN SET TO NONDEFAULT VALUES THIS CAN BE USED TO HAVE MORE COMPACT REPRESENTATIONS  
OUT  
DEFAULT REPRESENTATION  
LOGISTICREGRESSIONC10 CLASSWEIGHTNONE DUALFALSE FITINTERCEPTTRUE  
INTERCEPTSCALING1 L1RATIONONE MAXITER100  
MULTICLASSWARN NJOBSNONE PENALTYL1  
RANDOMSTATENONE SOLVERWARN TOL00001 VERBOSE0  
WARMSTARTFALSE  
WITH CHANGEDONLY OPTION  
LOGISTICREGRESSIONPENALTYL1  
PRINTDOC  
FROM SKLEARNLINEARMODEL IMPORT LOGISTICREGRESSION  
FROM SKLEARN IMPORT SETCONFIG  
LR LOGISTICREGRESSIONPENALTYL1  
PRINTDEFAULT REPRESENTATION  
PRINTLR  
LOGISTICREGRESSIONC10 CLASSWEIGHTNONE DUALFALSE FITINTERCEPTTRUE  
INTERCEPTSCALING1 L1RATIONONE MAXITER100  
683

SCIKITLEARN USER GUIDE RELEASE 0213  
MULTICLASSWARN NJOBSNONE PENALTYL1  
RANDOMSTATENONE SOLVERWARN TOL00001 VERBOSE0  
WARMSTARTFALSE  
SETCONFIGPRINTCHANGEDONLYTRUE  
PRINTNWITH CHANGEDONLY OPTION  
PRINTLR  
LOGISTICREGRESSIONPENALTYL1  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0003 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
512 ISOTONIC REGRESSION  
AN ILLUSTRATION OF THE ISOTONIC REGRESSION ON GENERATED DATA THE ISOTONIC REGRESSION FINDS A NONDECREASING APPROX  
IMATION OF A FUNCTION WHILE MINIMIZING THE MEAN SQUARED ERROR ON THE TRAINING DATA THE BENEFIT OF SUCH A MODEL IS  
THAT IT DOES NOT ASSUME ANY FORM FOR THE TARGET FUNCTION SUCH AS LINEARITY FOR COMPARISON A LINEAR REGRESSION IS ALSO  
PRESENTED  
PRINTDOC  
684 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
AUTHOR NELLE VAROQUAUX NELLEVAROQUAUXGMAILCOM
ALEXANDRE GRAMFORT ALEXANDREGGRAMFORTINRIAFR
LICENSE BSD
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM MATPLOTLIBCOLLECTIONS IMPORT LINECOLLECTION
FROM SKLEARNLINEARMODEL IMPORT LINEARREGRESSION
FROM SKLEARNISOTONIC IMPORT ISOTONICREGRESSION
FROM SKLEARNUTILS IMPORT CHECKRANDOMSTATE
N 100
X NPARANGEN
RS CHECKRANDOMSTATE0
Y RSRANDINT50 50 SIZEN 50 NPLOG1PNPARANGEN

FIT ISOTONICREGRESSION AND LINEARREGRESSION MODELS
IR ISOTONICREGRESSION
Y IRFITTRANSFORMX Y
LR LINEARREGRESSION
LRFITX NPNEWAXIS Y X NEEDS TO BE 2D FOR LINEARREGRESSION

PLOT RESULT
SEGMENTS I YI I YI FORIINRANGEN
LC LINECOLLECTIONSEGMENTS ZORDER0
LCSETARRAYNPONESLENY
LCSETLINEWIDTHHSNPFULLN 05
FIG PLTFigure
PLTPLOTX Y R MARKERSIZE12
PLTPLOTX Y B MARKERSIZE12
PLTPLOTX LRPREDICTX NPNEWAXIS B
PLTGCAADDCOLLECTIONLC
PLTLEGENDDATA ISOTONIC FIT LINEAR FIT LOCLOWER RIGHT
PLTTITLEISOTONIC REGRESSION
PLTSHOW
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0056 SECONDS
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE
513 FACE COMPLETION WITH A MULTIOUTPUT ESTIMATORS
THIS EXAMPLE SHOWS THE USE OF MULTIOUTPUT ESTIMATOR TO COMPLETE IMAGES THE GOAL IS TO PREDICT THE LOWER HALF OF A
FACE GIVEN ITS UPPER HALF
THE FIRST COLUMN OF IMAGES SHOWS TRUE FACES THE NEXT COLUMNS ILLUSTRATE HOW EXTREMELY RANDOMIZED TREES K NEAREST
NEIGHBORS LINEAR REGRESSION AND RIDGE REGRESSION COMPLETE THE LOWER HALF OF THOSE FACES
51 MISCELLANEOUS EXAMPLES 685
```

SCIKITLEARN USER GUIDE RELEASE 0213  
OUT  
DOWNLOADING OLIVETTI FACES FROM [HTTPSNDOWNLOADERFIGSHARECOMFILES5976027](https://download.figshare.com/files/5976027) TO  
'→HOMECIRCLECISCIKITLEARNDATA  
686 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
PRINTDOC
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARNDATASETS IMPORT FETCHOLIVETTIFACES
FROM SKLEARNUTILSVALIDATION IMPORT CHECKRANDOMSTATE
FROM SKLEARNENSEMBLE IMPORT EXTRATREESREGRESSOR
FROM SKLEARNNEIGHBORS IMPORT KNEIGHBORSREGRESSOR
FROM SKLEARNLINEARMODEL IMPORT LINEARREGRESSION
FROM SKLEARNLINEARMODEL IMPORT RIDGECV
LOAD THE FACES DATASETS
DATA FETCHOLIVETTIFACES
TARGETS DATATARGET
DATA DATAIMAGESRESHAPELENDATAIMAGES 1
TRAIN DATATARGETS 30
TEST DATATARGETS 30 TEST ON INDEPENDENT PEOPLE
TEST ON A SUBSET OF PEOPLE
NFACES 5
RNG CHECKRANDOMSTATE4
FACEIDS RNGRANDINTTESTSHAPE0 SIZENFACES
TEST TESTFACEIDS
NPIXELS DATASHAPE1
UPPER HALF OF THE FACES
XTRAIN TRAIN NPIXELS 1 2
LOWER HALF OF THE FACES
YTRAIN TRAIN NPIXELS 2
XTEST TEST NPIXELS 1 2
YTEST TEST NPIXELS 2
FIT ESTIMATORS
ESTIMATORS
EXTRA TREES EXTRATREESREGRESSORNESTIMATORS10 MAXFEATURES32
RANDOMSTATE0
KNN KNEIGHBORSREGRESSOR
LINEAR REGRESSION LINEARREGRESSION
RIDGE RIDGECV

YTESTPREDICT DICT
FORNAME ESTIMATOR INESTIMATORSITEMS
ESTIMATORFITXTRAIN YTRAIN
YTESTPREDICTNAME ESTIMATORPREDICTXTEST
PLOT THE COMPLETED FACES
IMAGESHAPE 64 64
NCOLS 1 LENESTIMATORS
PLTFIGUREFIGSIZE2 NCOLS 226 NFACES
PLTSUPTITLEFACE COMPLETION WITH MULTIOUTPUT ESTIMATORS SIZE16
FORIINRANGENFACES
TRUEFACE NPHSTACKXTESTI YTESTI
51 MISCELLANEOUS EXAMPLES 687
```

SCIKITLEARN USER GUIDE RELEASE 0213

IFI

SUB PLTSUBPLOTNFACES NCOLS I NCOLS 1

ELSE

SUB PLTSUBPLOTNFACES NCOLS I NCOLS 1

TITLETRUE FACES

SUBAXISOFF

SUBIMSHOWTRUEFACERESHAPEIMAGESHAPE

CMAPPLTCMGRAY

INTERPOLATIONNEAREST

FORJ ESTINENUMERATESORTEDESTIMATORS

COMPLETEDFACE NPHSTACKXTESTI YTESTPREDICTESTI

IFI

SUB PLTSUBPLOTNFACES NCOLS I NCOLS 2 J

ELSE

SUB PLTSUBPLOTNFACES NCOLS I NCOLS 2 J

TITLEEST

SUBAXISOFF

SUBIMSHOWCOMPLETEDFACERESHAPEIMAGESHAPE

CMAPPLTCMGRAY

INTERPOLATIONNEAREST

PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 3749 SECONDS

NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE

514 MULTILABEL CLASSIFICATION

THIS EXAMPLE SIMULATES A MULTILABEL DOCUMENT CLASSIFICATION PROBLEM THE DATASET IS GENERATED RANDOMLY BASED ON THE FOLLOWING PROCESS

- PICK THE NUMBER OF LABELS N POISSONNLABELS
- N TIMES CHOOSE A CLASS C C MULTINOMIALTHETA
- PICK THE DOCUMENT LENGTH K POISSONLENGTH
- K TIMES CHOOSE A WORD W MULTINOMIALTHETAC

IN THE ABOVE PROCESS REJECTION SAMPLING IS USED TO MAKE SURE THAT N IS MORE THAN 2 AND THAT THE DOCUMENT LENGTH IS NEVER ZERO LIKEWISE WE REJECT CLASSES WHICH HAVE ALREADY BEEN CHOSEN THE DOCUMENTS THAT ARE ASSIGNED TO BOTH CLASSES ARE PLOTTED SURROUNDED BY TWO COLORED CIRCLES

THE CLASSIFICATION IS PERFORMED BY PROJECTING TO THE FIRST TWO PRINCIPAL COMPONENTS FOUND BY PCA AND CCA FOR VISUAL ISATION PURPOSES FOLLOWED BY USING THE SKLEARNMULTICLASSONEVSRESTCLASSIFIER METACCLASSIFIER USING TWO SVCS WITH LINEAR KERNELS TO LEARN A DISCRIMINATIVE MODEL FOR EACH CLASS NOTE THAT PCA IS USED TO PERFORM AN UNSUPERVISED DIMENSIONALITY REDUCTION WHILE CCA IS USED TO PERFORM A SUPERVISED ONE

NOTE IN THE PLOT “UNLABELED SAMPLES” DOES NOT MEAN THAT WE DON’T KNOW THE LABELS AS IN SEMISUPERVISED LEARNING BUT THAT THE SAMPLES SIMPLY DO NOTHAVE A LABEL

688 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
PRINTDOC
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARNDATASETS IMPORT MAKEMULTILABELCLASSIFICATION
FROM SKLEARNMULTICLASS IMPORT ONEVSRESTCLASSIFIER
FROM SKLEARN SVM IMPORT SVC
FROM SKLEARN DECOMPOSITION IMPORT PCA
FROM SKLEARN CROSS DECOMPOSITION IMPORT CCA
DEF PLOT HYPERPLANE CLF MINX MAXX LINESTYLE LABEL
GET THE SEPARATING HYPERPLANE
W CLFCOEF0
A W0 W1
XX NPLINSPACE MINX 5 MAXX 5 MAKE SURE THE LINE IS LONG ENOUGH
YY AXX CLFINTERCEPT0 W1
PLT PLOT XX YY LINESTYLE LABEL LABEL
DEF PLOT SUBFIGURE X Y SUBPLOT TITLE TRANSFORM
IF TRANSFORM PCA
X PCANCOMPONENTS2FIT TRANSFORM X
ELIF TRANSFORM CCA
X CCANCOMPONENTS2FIT X Y TRANSFORM X
51 MISCELLANEOUS EXAMPLES 689
```

SCIKITLEARN USER GUIDE RELEASE 0213  
ELSE  
RAISEVALUEERROR  
MINX NPMINX 0  
MAXX NPMAXX 0  
MINY NPMINX 1  
MAXY NPMAXX 1  
CLASSIF ONEVSRESTCLASSIFIERSVCKERNELLINEAR  
CLASSIFFITX Y  
PLTSUBPLOT2 2 SUBPLOT  
PLTTITLETITLE  
ZEROCLASS NPWHEREY 0  
ONECLASS NPWHEREY 1  
PLTSCATTERX 0 X 1 S40 CGRAY EDGECOLORS0 0 0  
PLTSCATTERXZEROCLASS 0 XZEROCLASS 1 S160 EDGECOLORSB  
FACECOLORSNONE LINEWIDTHS2 LABELCLASS 1  
PLTSCATTERXONECLASS 0 XONECLASS 1 S80 EDGECOLORSORANGE  
FACECOLORSNONE LINEWIDTHS2 LABELCLASS 2  
PLOTHYPERPLANECLASSIFESTIMATORS0 MINX MAXX K  
BOUNDARY NFOR CLASS 1  
PLOTHYPERPLANECLASSIFESTIMATORS1 MINX MAXX K  
BOUNDARY NFOR CLASS 2  
PLXTICKS  
PLTYTICKS  
PLTXLIMMINX 5 MAXX MAXX 5 MAXX  
PLTYLIMMINY 5 MAXY MAXY 5 MAXY  
IFSUBPLOT 2  
PLTXLABELFIRST PRINCIPAL COMPONENT  
PLTYLABELSECOND PRINCIPAL COMPONENT  
PLTLEGENDLOCUPPER LEFT  
PLTFIGUREFIGSIZE8 6  
X Y MAKEMULTILABELCLASSIFICATIONNCLASSES2 NLABELS1  
ALLOWUNLABELEDTRUE  
RANDOMSTATE1  
PLOTSUBFIGUREX Y 1 WITH UNLABELED SAMPLES CCA CCA  
PLOTSUBFIGUREX Y 2 WITH UNLABELED SAMPLES PCA PCA  
X Y MAKEMULTILABELCLASSIFICATIONNCLASSES2 NLABELS1  
ALLOWUNLABELEDFALSE  
RANDOMSTATE1  
PLOTSUBFIGUREX Y 3 WITHOUT UNLABELED SAMPLES CCA CCA  
PLOTSUBFIGUREX Y 4 WITHOUT UNLABELED SAMPLES PCA PCA  
PLTSUBPLOTSADJUST04 02 97 94 09 2  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0093 SECONDS  
690 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213

NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)

515 COMPARING ANOMALY DETECTION ALGORITHMS FOR OUTLIER DETECTION ON TOY DATASETS

THIS EXAMPLE SHOWS CHARACTERISTICS OF DIFFERENT ANOMALY DETECTION ALGORITHMS ON 2D DATASETS DATASETS CONTAIN ONE OR TWO MODES REGIONS OF HIGH DENSITY TO ILLUSTRATE THE ABILITY OF ALGORITHMS TO COPE WITH MULTIMODAL DATA

FOR EACH DATASET 15 OF SAMPLES ARE GENERATED AS RANDOM UNIFORM NOISE THIS PROPORTION IS THE VALUE GIVEN TO THE NU PARAMETER OF THE ONECLASSSSVM AND THE CONTAMINATION PARAMETER OF THE OTHER OUTLIER DETECTION ALGORITHMS DECISION BOUNDARIES BETWEEN INLIERS AND OUTLIERS ARE DISPLAYED IN BLACK EXCEPT FOR LOCAL OUTLIER FACTOR LOF AS IT HAS NO PREDICT METHOD TO BE APPLIED ON NEW DATA WHEN IT IS USED FOR OUTLIER DETECTION

THESKLEARNSSVMONECLASSSSVM IS KNOWN TO BE SENSITIVE TO OUTLIERS AND THUS DOES NOT PERFORM VERY WELL FOR OUTLIER DETECTION THIS ESTIMATOR IS BEST SUITED FOR NOVELTY DETECTION WHEN THE TRAINING SET IS NOT CONTAMINATED BY OUTLIERS THAT SAID OUTLIER DETECTION IN HIGHDIMENSION OR WITHOUT ANY ASSUMPTIONS ON THE DISTRIBUTION OF THE INLYING DATA IS VERY CHALLENGING AND A ONECLASS SVM MIGHT GIVE USEFUL RESULTS IN THESE SITUATIONS DEPENDING ON THE VALUE OF ITS HYPERPARAMETERS

SKLEARNCOVARIANCEELLIPTICENVELOPE ASSUMES THE DATA IS GAUSSIAN AND LEARNS AN ELLIPSE IT THUS DEGRADES WHEN THE DATA IS NOT UNIMODAL NOTICE HOWEVER THAT THIS ESTIMATOR IS ROBUST TO OUTLIERS

SKLEARNENSEMBLEISOLATIONFOREST ANDSKLEARNNEIGHBORSLOCALOUTLIERFACTOR SEEM TO PERFORM REASONABLY WELL FOR MULTIMODAL DATA SETS THE ADVANTAGE OF SKLEARNNEIGHBORS LOCALOUTLIERFACTOR OVER THE OTHER ESTIMATORS IS SHOWN FOR THE THIRD DATA SET WHERE THE TWO MODES HAVE DIFFERENT DENSITIES THIS ADVANTAGE IS EXPLAINED BY THE LOCAL ASPECT OF LOF MEANING THAT IT ONLY COMPARES THE SCORE OF ABNORMALITY OF ONE SAMPLE WITH THE SCORES OF ITS NEIGHBORS

FINALLY FOR THE LAST DATA SET IT IS HARD TO SAY THAT ONE SAMPLE IS MORE ABNORMAL THAN ANOTHER SAMPLE AS THEY ARE UNIFORMLY DISTRIBUTED IN A HYPERCUBE EXCEPT FOR THE SKLEARNSSVMONECLASSSSVM WHICH OVERFITS A LITTLE ALL ESTIMATORS PRESENT DECENT SOLUTIONS FOR THIS SITUATION IN SUCH A CASE IT WOULD BE WISE TO LOOK MORE CLOSELY AT THE SCORES OF ABNORMALITY OF THE SAMPLES AS A GOOD ESTIMATOR SHOULD ASSIGN SIMILAR SCORES TO ALL THE SAMPLES WHILE THESE EXAMPLES GIVE SOME INTUITION ABOUT THE ALGORITHMS THIS INTUITION MIGHT NOT APPLY TO VERY HIGH DIMENSIONAL DATA

FINALLY NOTE THAT PARAMETERS OF THE MODELS HAVE BEEN HERE HANDPICKED BUT THAT IN PRACTICE THEY NEED TO BE ADJUSTED IN THE ABSENCE OF LABELLED DATA THE PROBLEM IS COMPLETELY UNSUPERVISED SO MODEL SELECTION CAN BE A CHALLENGE

51 MISCELLANEOUS EXAMPLES 691

SCIKITLEARN USER GUIDE RELEASE 0213  
AUTHOR ALEXANDRE GRAMFORT ALEXANDREGRAMFORTINRIA  
ALBERT THOMAS ALBERTTHOMASTELECOMPARISTECHFR  
LICENSE BSD 3 CLAUSE  
IMPORT TIME  
IMPORT NUMPY AS NP  
IMPORT MATPLOTLIB  
IMPORT MATPLOTLIBPYPLOT AS PLT  
692 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

FROM SKLEARN IMPORT SVM

FROM SKLEARNDATASETS IMPORT MAKEMOONS MAKEBLOBS

FROM SKLEARNCOVARIANCE IMPORT ELLIPTICENVELOPE

FROM SKLEARNENSEMBLE IMPORT ISOLATIONFOREST

FROM SKLEARNNEIGHBORS IMPORT LOCALOUTLIERFACTOR

PRINTDOC

MATPLOTLIBRCPARAMSCONTOURNEGATIVELINESTYLE SOLID

EXAMPLE SETTINGS

NSAMPLES 300

OUTLIERSFRACTION 015

NOUTLIERS INTOUTLIERSFRACTION NSAMPLES

NINLIERS NSAMPLES NOUTLIERS

DEFINE OUTLIERANOMALY DETECTION METHODS TO BE COMPARED

ANOMALYALGORITHMS

ROBUST COVARIANCE ELLIPTICENVELOPECONTAMINATIONOUTLIERSFRACTION

ONECLASS SVM SVMONECLASSSSVMNUOUTLIERSFRACTION KERNELRBF

GAMMA01

ISOLATION FOREST ISOLATIONFORESTBEHAVIOURNEW

CONTAMINATIONOUTLIERSFRACTION

RANDOMSTATE42

LOCAL OUTLIER FACTOR LOCALOUTLIERFACTOR

NNEIGHBORS35 CONTAMINATIONOUTLIERSFRACTION

DEFINE DATASETS

BLOBSPARAMS DICRANDOMSTATE0 NSAMPLESNINLIERS NFEATURES2

DATASETS

MAKEBLOBSCENTERS0 0 0 0 CLUSTERSTD05

BLOBSPARAMS0

MAKEBLOBSCENTERS2 2 2 2 CLUSTERSTD05 05

BLOBSPARAMS0

MAKEBLOBSCENTERS2 2 2 2 CLUSTERSTD15 3

BLOBSPARAMS0

4MAKEMOONSNSAMPLESNSAMPLES NOISE05 RANDOMSTATE00

NPARRAY05 025

14NPRANDOMRANDOMSTATE42RANDNSAMPLES 2 05

COMPARE GIVEN CLASSIFIERS UNDER GIVEN SETTINGS

XX YY NPMESHGRIDNPLINSPACE7 7 150

NPLINSPACE7 7 150

PLTFIGUREFIGSIZELENANOMALYALGORITHMS 2 3 125

PLTSUBPLOTSADJUSTLEFT02 RIGHT98 BOTTOM001 TOP96 WSPACE05

HSPACE01

PLOTNUM 1

RNG NPRANDOMRANDOMSTATE42

FORIDATASET X INENUMERATEDDATASETS

ADD OUTLIERS

X NPCONCATENATEX RNGUNIFORMLOW6 HIGH6

SIZENOUTLIERS 2 AXIS0

FORNAME ALGORITHM INANOMALYALGORITHMS

51 MISCELLANEOUS EXAMPLES 693

SCIKITLEARN USER GUIDE RELEASE 0213

T0 TIMETIME

ALGORITHMFITX

T1 TIMETIME

PLTSUBPLOTLENDATASETS LENANOMALYALGORITHMS PLOTNUM

IFIDATASET 0

PLTTITLENAME SIZE18

  FIT THE DATA AND TAG OUTLIERS

IFNAME LOCAL OUTLIER FACTOR

YPRED ALGORITHMFITPREDICTX

ELSE

YPRED ALGORITHMFITXPREDICTX

  PLOT THE LEVELS LINES AND THE POINTS

IFNAME LOCAL OUTLIER FACTOR LOF DOES NOT IMPLEMENT PREDICT

Z ALGORITHMPPREDICTNPCXXRAVEL YYRAVEL

Z ZRESHAPEXXSHAPE

PLTCONTOURXX YY Z LEVELS0 LINEWIDTHS2 COLORSBLACK

COLORS NPARRAY377EB8 FF7F00

PLTSCATTERX 0 X 1 S10 COLORCOLORSYYPRED 1 2

PLTXLIM7 7

PLTYLIM7 7

PLTXTICKS

PLTYTICKS

PLTTEXT99 01 2FS T1 T0LSTRIPO

TRANSFORMPLTGCATTRANSAXES SIZE15

HORIZONTALALIGNMENTRIGHT

PLOTNUM 1

PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 3491 SECONDS

NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)

516 THE JOHNSONLINDENSTRAUSS BOUND FOR EMBEDDING WITH RANDOM PROJECTIONS

THE JOHNSONLINDENSTRAUSS LEMMA STATES THAT ANY HIGH DIMENSIONAL DATASET CAN BE RANDOMLY PROJECTED INTO A LOWER DIMENSIONAL EUCLIDEAN SPACE WHILE CONTROLLING THE DISTORTION IN THE PAIRWISE DISTANCES

THEORETICAL BOUNDS

THE DISTORTION INTRODUCED BY A RANDOM PROJECTION PIS ASSERTED BY THE FACT THAT PIS DEFINING AN EPSEMBEDDING WITH GOOD PROBABILITY AS DEFINED BY

$$1-\frac{1}{n}\sum_{i,j=1}^n\|u_i-u_j\|_2^2\leq\frac{1}{n}\sum_{i,j=1}^n\|v_i-v_j\|_2^2$$

WHERE U AND V ARE ANY ROWS TAKEN FROM A DATASET OF SHAPE NSAMPLES NFEATURES AND P IS A PROJECTION BY A RANDOM GAUSSIAN NO 1 MATRIX WITH SHAPE NCOMPONENTS NFEATURES OR A SPARSE ACHLIOPTAS MATRIX

THE MINIMUM NUMBER OF COMPONENTS TO GUARANTEES THE EPSEMBEDDING IS GIVEN BY

$$\frac{1}{\epsilon^2}\log\frac{1}{\delta}\leq\frac{1}{\epsilon^2}\log\frac{1}{\delta}$$

694 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

THE FIRST PLOT SHOWS THAT WITH AN INCREASING NUMBER OF SAMPLES NSAMPLES THE MINIMAL NUMBER OF DIMENSIONS NCOMPONENTS INCREASED LOGARITHMICALLY IN ORDER TO GUARANTEE AN EPSEMBEDDING  
THE SECOND PLOT SHOWS THAT AN INCREASE OF THE ADMISSIBLE DISTORTION EPS ALLOWS TO REDUCE DRASTICALLY THE MINIMAL NUMBER OF DIMENSIONS NCOMPONENTS FOR A GIVEN NUMBER OF SAMPLES NSAMPLES  
EMPIRICAL VALIDATION

WE VALIDATE THE ABOVE BOUNDS ON THE DIGITS DATASET OR ON THE 20 NEWSGROUPS TEXT DOCUMENT TFIDF WORD FREQUENCIES DATASET

- FOR THE DIGITS DATASET SOME 8X8 GRAY LEVEL PIXELS DATA FOR 500 HANDWRITTEN DIGITS PICTURES ARE RANDOMLY PROJECTED TO SPACES FOR VARIOUS LARGER NUMBER OF DIMENSIONS NCOMPONENTS
- FOR THE 20 NEWSGROUPS DATASET SOME 500 DOCUMENTS WITH 100K FEATURES IN TOTAL ARE PROJECTED USING A SPARSE RANDOM MATRIX TO SMALLER EUCLIDEAN SPACES WITH VARIOUS VALUES FOR THE TARGET NUMBER OF DIMENSIONS NCOMPONENTS

THE DEFAULT DATASET IS THE DIGITS DATASET TO RUN THE EXAMPLE ON THE TWENTY NEWSGROUPS DATASET PASS THE -TWENTY NEWSGROUPS COMMAND LINE ARGUMENT TO THIS SCRIPT  
FOR EACH VALUE OF NCOMPONENTS WE PLOT

- 2D DISTRIBUTION OF SAMPLE PAIRS WITH PAIRWISE DISTANCES IN ORIGINAL AND PROJECTED SPACES AS X AND Y AXIS RESPECTIVELY
- 1D HISTOGRAM OF THE RATIO OF THOSE DISTANCES PROJECTED ORIGINAL

WE CAN SEE THAT FOR LOW VALUES OF NCOMPONENTS THE DISTRIBUTION IS WIDE WITH MANY DISTORTED PAIRS AND A SKEWED DISTRIBUTION DUE TO THE HARD LIMIT OF ZERO RATIO ON THE LEFT AS DISTANCES ARE ALWAYS POSITIVES WHILE FOR LARGER VALUES OF NCOMPONENTS THE DISTORTION IS CONTROLLED AND THE DISTANCES ARE WELL PRESERVED BY THE RANDOM PROJECTION  
REMARKS

ACCORDING TO THE JL LEMMA PROJECTING 500 SAMPLES WITHOUT TOO MUCH DISTORTION WILL REQUIRE AT LEAST SEVERAL THOUSANDS DIMENSIONS IRRESPECTIVE OF THE NUMBER OF FEATURES OF THE ORIGINAL DATASET  
HENCE USING RANDOM PROJECTIONS ON THE DIGITS DATASET WHICH ONLY HAS 64 FEATURES IN THE INPUT SPACE DOES NOT MAKE SENSE IT DOES NOT ALLOW FOR DIMENSIONALITY REDUCTION IN THIS CASE  
ON THE TWENTY NEWSGROUPS ON THE OTHER HAND THE DIMENSIONALITY CAN BE DECREASED FROM 56436 DOWN TO 10000 WHILE REASONABLY PRESERVING PAIRWISE DISTANCES  
51 MISCELLANEOUS EXAMPLES 695



SCIKITLEARN USER GUIDE RELEASE 0213

- 

51 MISCELLANEOUS EXAMPLES 697





SCIKITLEARN USER GUIDE RELEASE 0213

- 

51 MISCELLANEOUS EXAMPLES 699

SCIKITLEARN USER GUIDE RELEASE 0213

- 700 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

- 

51 MISCELLANEOUS EXAMPLES 701



SCIKITLEARN USER GUIDE RELEASE 0213

- OUT  
EMBEDDING 500 SAMPLES WITH DIM 64 USING VARIOUS RANDOM PROJECTIONS  
PROJECTED 500 SAMPLES FROM 64 TO 300 IN 0016S  
RANDOM MATRIX WITH SIZE 0028MB  
MEAN DISTANCES RATE 097 008  
PROJECTED 500 SAMPLES FROM 64 TO 1000 IN 0048S  
RANDOM MATRIX WITH SIZE 0096MB  
MEAN DISTANCES RATE 099 005  
PROJECTED 500 SAMPLES FROM 64 TO 10000 IN 0594S  
RANDOM MATRIX WITH SIZE 0964MB  
MEAN DISTANCES RATE 101 001  
PRINTDOC  
IMPORT SYS  
FROM TIME IMPORT TIME  
IMPORT NUMPY AS NP  
IMPORT MATPLOTLIB  
IMPORT MATPLOTLIBPYPLOT AS PLT  
FROM DISTUTILSVERSION IMPORT LOOSEVERSION  
FROM SKLEARNRANDOMPROJECTION IMPORT JOHNSONLINDENSTRAUSSMINDIM  
51 MISCELLANEOUS EXAMPLES 703

SCIKITLEARN USER GUIDE RELEASE 0213

FROM SKLEARNRANDOMPROJECTION IMPORT SPARSERANDOMPROJECTION

FROM SKLEARNDATASETS IMPORT FETCH20NEWSGROUPSVECTORIZED

FROM SKLEARNDATASETS IMPORT LOADDIGITS

FROM SKLEARNMETRICSPAIRWISE IMPORT EUCLIDEANDISTANCES

  NORMED IS BEING DEPRECATED IN FAVOR OF DENSITY IN HISTOGRAMS

IFLOOSEVERSIONMATPLOTLIBVERSION 21

DENSITYPARAM DENSITY TRUE

ELSE

DENSITYPARAM NORMED TRUE

PART 1 PLOT THE THEORETICAL DEPENDENCY BETWEEN NCOMPONENTSMIN AND

NSAMPLES

RANGE OF ADMISSIBLE DISTORTIONS

EPSRANGE NPLinspace01 099 5

COLORS PLTCMBLUESNPLinspace03 10 LENEPSRANGE

RANGE OF NUMBER OF SAMPLES OBSERVATION TO EMBED

NSAMPLESRANGE NPLOGSPACE1 9 9

PLTFigure

FOREPS COLOR INZIPEPSRANGE COLORS

MINNCOMPONENTS JOHNSONLINDENSTRAUSSMINDIMNSAMPLESRANGE EPSEPS

PLTLOGLOGNSAMPLESRANGE MINNCOMPONENTS COLORCOLOR

PLTLEGENDEPS 01F EPS FOREPSINEPSRANGE LOCLOWER RIGHT

PLTXLABELNUMBER OF OBSERVATIONS TO EMBED

PLTYLABELMINIMUM NUMBER OF DIMENSIONS

PLTTITLEJOHNSONLINDENSTRAUSS BOUNDS NNSAMPLES VS NCOMPONENTS

RANGE OF ADMISSIBLE DISTORTIONS

EPSRANGE NPLinspace001 099 100

RANGE OF NUMBER OF SAMPLES OBSERVATION TO EMBED

NSAMPLESRANGE NPLOGSPACE2 6 5

COLORS PLTCMBLUESNPLinspace03 10 LENNSAMPLESRANGE

PLTFigure

FORNSAMPLES COLOR INZIPNSAMPLESRANGE COLORS

MINNCOMPONENTS JOHNSONLINDENSTRAUSSMINDIMNSAMPLES EPSEPSRANGE

PLTSEMILOGYEPSRANGE MINNCOMPONENTS COLORCOLOR

PLTLEGENDNSAMPLES D NFORNINNSAMPLESRANGE LOCUPPER RIGHT

PLTXLABELDISTORTION EPS

PLTYLABELMINIMUM NUMBER OF DIMENSIONS

PLTTITLEJOHNSONLINDENSTRAUSS BOUNDS NNCOMPONENTS VS EPS

PART 2 PERFORM SPARSE RANDOM PROJECTION OF SOME DIGITS IMAGES WHICH ARE

QUITE LOW DIMENSIONAL AND DENSE OR DOCUMENTS OF THE 20 NEWSGROUPS DATASET

WHICH IS BOTH HIGH DIMENSIONAL AND SPARSE

IFTWENTYNEWSGROUPS INSYSARGV

NEED AN INTERNET CONNECTION HENCE NOT ENABLED BY DEFAULT

DATA FETCH20NEWSGROUPSVECTORIZEDDATA500

ELSE

DATA LOADDIGITSDATA500

704 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
NSAMPLES NFEATURES DATASHAPE  
PRINTEMBEDDING DSAMPLES WITH DIM DUSING VARIOUS RANDOM PROJECTIONS  
  NSAMPLES NFEATURES  
NCOMPONENTSRANGE NPARRAY300 1000 10000  
DISTS EUCLIDEANDISTANCESDATA SQUAREDTRUERAVEL  
  SELECT ONLY NONIDENTICAL SAMPLES PAIRS  
NONZERO DISTS 0  
DISTS DISTSNONZERO  
FORNCOMPONENTS INNCOMPONENTSRANGE  
T0 TIME  
RP SPARSERANDOMPROJECTIONNCOMPONENTSNCOMPONENTS  
PROJECTEDDATA RPFITTRANSFORMDATA  
PRINTPROJECTED DSAMPLES FROM DTODIN03FS  
  NSAMPLES NFEATURES NCOMPONENTS TIME T0  
IFHASATTRRP COMPONENTS  
NBYTES RPCOMPONENTSDATANBYTES  
NBYTES RPCOMPONENTSINDICESNBYTES  
PRINTRANDOM MATRIX WITH SIZE 03FMB NBYTES 1E6  
PROJECTEDDISTS EUCLIDEANDISTANCES  
PROJECTEDDATA SQUAREDTRUERAVELNONZERO  
PLTFigure  
PLTHEXBINDISTS PROJECTEDDISTS GRIDSIZE100 CMAPPLTCMPUBU  
PLTXLABELPAIRWISE SQUARED DISTANCES IN ORIGINAL SPACE  
PLTYLABELPAIRWISE SQUARED DISTANCES IN PROJECTED SPACE  
PLTTITLEPAIRWISE DISTANCES DISTRIBUTION FOR NCOMPONENTS D  
NCOMPONENTS  
CB PLTCOLORBAR  
CBSETLABELSAMPLE PAIRS COUNTS  
RATES PROJECTEDDISTS DISTS  
PRINTMEAN DISTANCES RATE 02F02F  
  NPMEANRATES NPSTDRATES  
PLTFigure  
PLTHISTRATES BINS50 RANGE0 2 EDGECOLORK DENSITYPARAM  
PLTXLABELSQUARED DISTANCES RATE PROJECTED ORIGINAL  
PLTYLABELDISTRIBUTION OF SAMPLES PAIRS  
PLTTITLEHISTOGRAM OF PAIRWISE DISTANCE RATES FOR NCOMPONENTS D  
NCOMPONENTS  
  TODO COMPUTE THE EXPECTED VALUE OF EPS AND ADD THEM TO THE PREVIOUS PLOT  
  AS VERTICAL LINES REGION  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 1837 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
51 MISCELLANEOUS EXAMPLES 705

BOTH KERNEL RIDGE REGRESSION KRR AND SVR LEARN A NONLINEAR FUNCTION BY EMPLOYING THE KERNEL TRICK IE THEY LEARN A LINEAR FUNCTION IN THE SPACE INDUCED BY THE RESPECTIVE KERNEL WHICH CORRESPONDS TO A NONLINEAR FUNCTION IN THE ORIGINAL SPACE THEY DIFFER IN THE LOSS FUNCTIONS RIDGE VERSUS EPSILONINSENSITIVE LOSS IN CONTRAST TO SVR FITTING A KRR CAN BE DONE IN CLOSEDFORM AND IS TYPICALLY FASTER FOR MEDIUMSIZED DATASETS ON THE OTHER HAND THE LEARNED MODEL IS NONSPARSE AND THUS SLOWER THAN SVR AT PREDICTIONTIME

THIS EXAMPLE ILLUSTRATES BOTH METHODS ON AN ARTIFICIAL DATASET WHICH CONSISTS OF A SINUSOIDAL TARGET FUNCTION AND STRONG NOISE ADDED TO EVERY FIFTH DATAPOINT THE FIRST FIGURE COMPARES THE LEARNED MODEL OF KRR AND SVR WHEN BOTH COMPLEXITYREGULARIZATION AND BANDWIDTH OF THE RBF KERNEL ARE OPTIMIZED USING GRIDSEARCH THE LEARNED FUNCTIONS ARE VERY SIMILAR HOWEVER FITTING KRR IS APPROX SEVEN TIMES FASTER THAN FITTING SVR BOTH WITH GRIDSEARCH HOWEVER PREDICTION OF 100000 TARGET VALUES IS MORE THAN TREE TIMES FASTER WITH SVR SINCE IT HAS LEARNED A SPARSE MODEL USING ONLY APPROX 13 OF THE 100 TRAINING DATAPOINTS AS SUPPORT VECTORS

THE NEXT FIGURE COMPARES THE TIME FOR FITTING AND PREDICTION OF KRR AND SVR FOR DIFFERENT SIZES OF THE TRAINING SET FITTING KRR IS FASTER THAN SVR FOR MEDIUM SIZED TRAINING SETS LESS THAN 1000 SAMPLES HOWEVER FOR LARGER TRAINING SETS SVR SCALES BETTER WITH REGARD TO PREDICTION TIME SVR IS FASTER THAN KRR FOR ALL SIZES OF THE TRAINING SET BECAUSE OF THE LEARNED SPARSE SOLUTION NOTE THAT THE DEGREE OF SPARSITY AND THUS THE PREDICTION TIME DEPENDS ON THE PARAMETERS EPSILON AND C OF THE SVR

-



SCIKITLEARN USER GUIDE RELEASE 0213

- 

51 MISCELLANEOUS EXAMPLES 707

SCIKITLEARN USER GUIDE RELEASE 0213

- OUT  
SVR COMPLEXITY AND BANDWIDTH SELECTED AND MODEL FITTED IN 0389 S  
KRR COMPLEXITY AND BANDWIDTH SELECTED AND MODEL FITTED IN 0175 S  
SUPPORT VECTOR RATIO 0320  
SVR PREDICTION FOR 100000 INPUTS IN 0117 S  
KRR PREDICTION FOR 100000 INPUTS IN 0141 S  
AUTHORS JAN HENDRIK METZEN JHMINFORMATIKUNIBREMENDE  
LICENSE BSD 3 CLAUSE  
IMPORT TIME  
IMPORT NUMPY AS NP  
FROM SKLEARN SVM IMPORT SVR  
FROM SKLEARNMODELSELECTION IMPORT GRIDSEARCHCV  
FROM SKLEARNMODELSELECTION IMPORT LEARNINGCURVE  
FROM SKLEARNKERNELRIDGE IMPORT KERNELRIDGE  
IMPORT MATPLOTLIBPYPLOT AS PLT  
708 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
RNG NPRANDOMRANDOMSTATE0

GENERATE SAMPLE DATA  
X 5RNGRAND10000 1  
Y NPSINXRAVEL  
ADD NOISE TO TARGETS  
Y5 3 05 RNGRANDXSHAPE0 5  
XPLOT NPLINSPACE0 5 100000 NONE

FIT REGRESSION MODEL  
TRAINSIZ 100  
SVR GRIDSEARCHCVSVRKERNELRBF GAMMA01 CV5  
PARAMGRIDC 1E0 1E1 1E2 1E3  
GAMMA NPLOGSPACE2 2 5  
KR GRIDSEARCHCVKERNELRIDGEKERNELRBF GAMMA01 CV5  
PARAMGRIDALPHA 1E0 01 1E2 1E3  
GAMMA NPLOGSPACE2 2 5  
T0 TIMETIME  
SVRFITXTRAINSIZ YTRAINSIZ  
SVRFIT TIMETIME T0  
PRINTSVR COMPLEXITY AND BANDWIDTH SELECTED AND MODEL FITTED IN 3FS  
SVRFIT  
T0 TIMETIME  
KRFITXTRAINSIZ YTRAINSIZ  
KRFIT TIMETIME T0  
PRINTKRR COMPLEXITY AND BANDWIDTH SELECTED AND MODEL FITTED IN 3FS  
KRFIT  
SVRATIO SVRBESTESTIMATORSUPPORTSHAPE0 TRAINSIZ  
PRINTSUPPORT VECTOR RATIO 3F SVRATIO  
T0 TIMETIME  
YSVR SVRPREDICTXPLOT  
SVRPREDICT TIMETIME T0  
PRINTSVR PREDICTION FOR DINPUTS IN 3FS  
XPLOTSHAPE0 SVRPREDICT  
T0 TIMETIME  
YKR KRPREDICTXPLOT  
KRPREDICT TIMETIME T0  
PRINTKRR PREDICTION FOR DINPUTS IN 3FS  
XPLOTSHAPE0 KRPREDICT

LOOK AT THE RESULTS  
SVIND SVRBESTESTIMATORSUPPORT  
PLTSCATTERXSVIND YSVIND CR S50 LABELSVR SUPPORT VECTORS  
ZORDER2 EDGECOLORS0 0 0  
PLTSCATTERX100 Y100 CK LABELDATA ZORDER1  
EDGECOLORS0 0 0  
51 MISCELLANEOUS EXAMPLES 709

SCIKITLEARN USER GUIDE RELEASE 0213  
PLTPLOTXPLOT YSVR CR  
LABELSVR FIT 3FS PREDICT 3FS SVRFIT SVRPREDICT  
PLTPLOTXPLOT YKR CG  
LABELKRR FIT 3FS PREDICT 3FS KRFIT KRPREDICT  
PLTXLABELDATA  
PLTYLABELTARGET  
PLTTITLESVR VERSUS KERNEL RIDGE  
PLTLEGEND  
VISUALIZE TRAINING AND PREDICTION TIME  
PLTFigure  
GENERATE SAMPLE DATA  
X 5RNGRAND10000 1  
Y NPSINXRavel  
Y5 3 05 RNGRANDXSHAPE0 5  
SIZES NPLOGSPACE1 4 7ASTYPENPINT  
FORNAME ESTIMATOR INKRR KERNELRIDGEKERNELRBF ALPHA01  
GAMMA10  
SVR SVRKERNELRBF C1E1 GAMMA10ITEMS  
TRAINTIME  
TESTTIME  
FORTRAINTESTSIZE INSIZES  
T0 TIMETIME  
ESTIMATORFITXTRAINTESTSIZE YTRAINTESTSIZE  
TRAINTIMEAPPENDTIMETIME T0  
T0 TIMETIME  
ESTIMATORPREDICTXPLOT1000  
TESTTIMEAPPENDTIMETIME T0  
PLTPLOTSIZES TRaintime O COLORR IFNAME SVR ELSEG  
LABELSTRain NAME  
PLTPLOTSIZES TESTTIME O COLORR IFNAME SVR ELSEG  
LABELSTEST NAME  
PLTXSCALELOG  
PLTYSCALELOG  
PLTXLABELTRAIN SIZE  
PLTYLABELTIME SECONDS  
PLTTITLEEXECUTION TIME  
PLTLEGENDLOCBEST  
VISUALIZE LEARNING CURVES  
PLTFigure  
SVR SVRKERNELRBF C1E1 GAMMA01  
KR KERNELRIDGEKERNELRBF ALPHA01 GAMMA01  
TRAINSIZES TRAINSCORESSVR TESTSCORESSVR  
LEARNINGCURVESVR X100 Y100 TRAINSIZESNPLINSPACE01 1 10  
SCORINGNEGMEANSQUAREDERROR CV10  
TRAINSIZESABS TRAINSCORESKR TESTSCORESKR  
LEARNINGCURVEKR X100 Y100 TRAINSIZESNPLINSPACE01 1 10  
SCORINGNEGMEANSQUAREDERROR CV10  
PLTPLOTTRAINSIZES TESTSCORESSVRMEAN1 O COLORR  
LABELSVR  
PLTPLOTTRAINSIZES TESTSCORESKRMEAN1 O COLORG  
710 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

LABELKRR

PLTXLABELTRAIN SIZE

PLTYLABELMEAN SQUARED ERROR

PLTTITLELEARNING CURVES

PLTLEGENDLOCBEST

PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 13067 SECONDS

NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE

518 EXPLICIT FEATURE MAP APPROXIMATION FOR RBF KERNELS

AN EXAMPLE ILLUSTRATING THE APPROXIMATION OF THE FEATURE MAP OF AN RBF KERNEL

IT SHOWS HOW TO USE RBFSAMPLER ANDNYSTROEM TO APPROXIMATE THE FEATURE MAP OF AN RBF KERNEL FOR CLASSIFICATION WITH AN SVM ON THE DIGITS DATASET RESULTS USING A LINEAR SVM IN THE ORIGINAL SPACE A LINEAR SVM USING THE APPROXIMATE MAPPINGS AND USING A KERNELIZED SVM ARE COMPARED TIMINGS AND ACCURACY FOR VARYING AMOUNTS OF MONTE CARLO SAMPLINGS IN THE CASE OF RBFSAMPLER WHICH USES RANDOM FOURIER FEATURES AND DIFFERENT SIZED SUBSETS OF THE TRAINING SET FOR NYSTROEM FOR THE APPROXIMATE MAPPING ARE SHOWN

PLEASE NOTE THAT THE DATASET HERE IS NOT LARGE ENOUGH TO SHOW THE BENEFITS OF KERNEL APPROXIMATION AS THE EXACT SVM IS STILL REASONABLY FAST

SAMPLING MORE DIMENSIONS CLEARLY LEADS TO BETTER CLASSIFICATION RESULTS BUT COMES AT A GREATER COST THIS MEANS THERE IS A TRADEOFF BETWEEN RUNTIME AND ACCURACY GIVEN BY THE PARAMETER NCOMPONENTS NOTE THAT SOLVING THE LINEAR SVM AND ALSO THE APPROXIMATE KERNEL SVM COULD BE GREATLY ACCELERATED BY USING STOCHASTIC GRADIENT DESCENT VIA SKLEARNLINEARMODELSGDCLASSIFIER THIS IS NOT EASILY POSSIBLE FOR THE CASE OF THE KERNELIZED SVM

THE SECOND PLOT VISUALIZED THE DECISION SURFACES OF THE RBF KERNEL SVM AND THE LINEAR SVM WITH APPROXIMATE KERNEL MAPS THE PLOT SHOWS DECISION SURFACES OF THE CLASSIFIERS PROJECTED ONTO THE FIRST TWO PRINCIPAL COMPONENTS OF THE DATA THIS VISUALIZATION SHOULD BE TAKEN WITH A GRAIN OF SALT SINCE IT IS JUST AN INTERESTING SLICE THROUGH THE DECISION SURFACE IN 64 DIMENSIONS IN PARTICULAR NOTE THAT A DATAPOINT REPRESENTED AS A DOT DOES NOT NECESSARILY BE CLASSIFIED INTO THE REGION IT IS LYING IN SINCE IT WILL NOT LIE ON THE PLANE THAT THE FIRST TWO PRINCIPAL COMPONENTS SPAN

THE USAGE OF RBFSAMPLER ANDNYSTROEM IS DESCRIBED IN DETAIL IN KERNEL APPROXIMATION

51 MISCELLANEOUS EXAMPLES 711

SCIKITLEARN USER GUIDE RELEASE 0213

- 
- 

712 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
PRINTDOC  
AUTHOR GAELEVAROQUAUXGAELE.DOTVAROQUAUX.AT.NORMALESUP.DOT.ORG  
ANDREAS.MUELLER@AMUELLERAISUNIBONN.DE  
LICENSE BSD 3 CLAUSE  
STANDARD SCIENTIFIC PYTHON IMPORTS  
import matplotlib.pyplot as plt  
import numpy as np  
from time import time  
import datasets, classifiers and performance metrics  
from sklearn import datasets, svm, pipeline  
from sklearn.kernel.approximation import RBFSampler  
nystroem  
from sklearn.decomposition import PCA  
the digits dataset  
digits = datasets.load\_digits(n\_class=9)  
to apply a classifier on this data we need to flatten the image to  
turn the data in a samples feature matrix  
n\_samples, length = digits.data.shape  
data = digits.data.reshape(-1, 16)  
data = data.mean(axis=0)  
we learn the digits on the first half of the digits  
data\_train, target\_train = data[:n\_samples//2], digits.target[:n\_samples//2]  
digit\_target = digits.target[n\_samples//2:]  
now predict the value of the digit on the second half  
data\_test, target\_test = data[n\_samples//2:], digits.target[n\_samples//2:]  
digit\_target = digits.target[n\_samples//2:]  
data\_test = scaler.transform(data\_test)  
create a classifier a support vector classifier  
kernel\_svm = svm.SVC(gamma=2)  
linear\_svm = svm.LinearSVC  
create pipeline from kernel approximation  
and linear svm  
feature\_map\_fourier = RBFSampler(gamma=2, random\_state=1)  
feature\_map\_nystroem = nystroem(gamma=2, random\_state=1)  
fourier\_approx\_svm = pipeline.Pipeline([('feature\_map', feature\_map\_fourier), ('svm', svm.LinearSVC)])  
nystroem\_approx\_svm = pipeline.Pipeline([('feature\_map', feature\_map\_nystroem), ('svm', svm.LinearSVC)])  
fit and predict using linear and kernel svm  
kernel\_svm\_fit\_time = time  
kernel\_svm\_fit(data\_train, target\_train)  
kernel\_svm\_score = kernel\_svm.score(data\_test, target\_test)  
kernel\_svm\_time = time  
linear\_svm\_time = time  
51 MISCELLANEOUS EXAMPLES 713

SCIKITLEARN USER GUIDE RELEASE 0213  
LINEARSVMFITDATATRAIN TARGETSTRAIN  
LINEARSVMSCORE LINEARSVMSCOREDATATEST TARGETSTEST  
LINEARSVMTIME TIME LINEARSVMTIME  
SAMPLESIZES 30 NPARANGE1 10  
FOURIERSCORES  
NYSTROEMSCORES  
FOURIERTIMES  
NYSTROEMTIMES  
FORDINSAMPLESIZES  
FOURIERAPPROXSVMSETPARAMSFEATUREMAPNCOMPONENTSD  
NYSTROEMAPPROXSVMSETPARAMSFEATUREMAPNCOMPONENTSD  
START TIME  
NYSTROEMAPPROXSVMFITDATATRAIN TARGETSTRAIN  
NYSTROEMTIMESAPPENDTIME START  
START TIME  
FOURIERAPPROXSVMFITDATATRAIN TARGETSTRAIN  
FOURIERTIMESAPPENDTIME START  
FOURIERSCORE FOURIERAPPROXSVMSCOREDATATEST TARGETSTEST  
NYSTROEMSCORE NYSTROEMAPPROXSVMSCOREDATATEST TARGETSTEST  
NYSTROEMSCORESAPPENDNYSTROEMSCORE  
FOURIERSCORESAPPENDFOURIERSCORE  
PLOT THE RESULTS  
PLTFIGUREFIGSIZE8 8  
ACCURACY PLTSUBPLOT211  
SECOND Y AXIS FOR TIMEINGS  
TIMESCALE PLTSUBPLOT212  
ACCURACYPLOTSAMPLESIZES NYSTROEMSCORES LABELNYSTROEM APPROX KERNEL  
TIMESCALEPLOTSAMPLESIZES NYSTROEMTIMES  
LABELNYSTROEM APPROX KERNEL  
ACCURACYPLOTSAMPLESIZES FOURIERSCORES LABELFOURIER APPROX KERNEL  
TIMESCALEPLOTSAMPLESIZES FOURIERTIMES  
LABELFOURIER APPROX KERNEL  
HORIZONTAL LINES FOR EXACT RBF AND LINEAR KERNELS  
ACCURACYPLOTSAMPLESIZES0 SAMPLESIZES1  
LINEARSVMSCORE LINEARSVMSCORE LABELLINEAR SVM  
TIMESCALEPLOTSAMPLESIZES0 SAMPLESIZES1  
LINEARSVMTIME LINEARSVMTIME LABELLINEAR SVM  
ACCURACYPLOTSAMPLESIZES0 SAMPLESIZES1  
KERNELSVMSCORE KERNELSVMSCORE LABELRBF SVM  
TIMESCALEPLOTSAMPLESIZES0 SAMPLESIZES1  
KERNELSVMTIME KERNELSVMTIME LABELRBF SVM  
VERTICAL LINE FOR DATASET DIMENSIONALITY 64  
ACCURACYPLOT64 64 07 1 LABELNFEATURES  
LEGENDS AND LABELS  
ACCURACYSETTITLECLASSIFICATION ACCURACY  
TIMESCALESETTITLETRAINING TIMES  
ACCURACYSETXLIMSAMPLESIZES0 SAMPLESIZES1  
714 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213  
ACCURACYSETXTICKS  
ACCURACYSETYLIMNPMINFOURIERSCORES 1  
TIMESCALESETXLABELSAMPLING STEPS TRANSFORMED FEATURE DIMENSION  
ACCURACYSETYLABELCLASSIFICATION ACCURACY  
TIMESCALESETYLABELTRAINING TIME IN SECONDS  
ACCURACYLEGENDLOCBEST  
TIMESCALELEGENDLOCBEST  
VISUALIZE THE DECISION SURFACE PROJECTED DOWN TO THE FIRST  
TWO PRINCIPAL COMPONENTS OF THE DATASET  
PCA PCANCOMPONENTS8FITDATATRAIN  
X PCATransformDATATRAIN  
GENERATE GRID ALONG FIRST TWO PRINCIPAL COMPONENTS  
MULTIPLES NPARANGE2 2 01  
STEPS ALONG FIRST COMPONENT  
FIRST MULTIPLES NPNEWAXIS PCACOMPONENTS0  
STEPS ALONG SECOND COMPONENT  
SECOND MULTIPLES NPNEWAXIS PCACOMPONENTS1  
COMBINE  
GRID FIRSTNPNEWAXIS SECOND NPNEWAXIS  
FLATGRID GRIDRESHAPE1 DATASHAPE1  
TITLE FOR THE PLOTS  
TITLES SVC WITH RBF KERNEL  
SVC LINEAR KERNEL NWITH FOURIER RBF FEATURE MAP N  
NCOMPONENTS100  
SVC LINEAR KERNEL NWITH NYSTROEM RBF FEATURE MAP N  
NCOMPONENTS100  
PLTTIGHTLAYOUT  
PLTFIGUREFIGSIZE12 5  
PREDICT AND PLOT  
FOR CLFINENUMERATEKERNELSVM NYSTROEMAPPROXSVM  
FOURIERAPPROXSVM  
PLOT THE DECISION BOUNDARY FOR THAT WE WILL ASSIGN A COLOR TO EACH  
POINT IN THE MESH XMIN XMAXXYMIN YMAX  
PLTSUBPLOT1 3 I 1  
Z CLFPREDICTFLATGRID  
PUT THE RESULT INTO A COLOR PLOT  
Z ZRESHAPEGRIDSHAPE1  
PLTCONTOURFMULTIPLES MULTIPLES Z CMAPPLTCMPAIED  
PLTAXISOFF  
PLOT ALSO THE TRAINING POINTS  
PLTSCATTERX 0 X 1 CTARGETSTRAIN CMAPPLTCMPAIED  
EDGECOLORS0 0 0  
PLTTITLETITLES  
PLTTIGHTLAYOUT  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 2188 SECONDS  
51 MISCELLANEOUS EXAMPLES 715

SCIKITLEARN USER GUIDE RELEASE 0213

52 EXAMPLES BASED ON REAL WORLD DATASETS

APPLICATIONS TO REAL WORLD PROBLEMS WITH SOME MEDIUM SIZED DATASETS OR INTERACTIVE USER INTERFACE

NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)

521 OUTLIER DETECTION ON A REAL DATA SET

THIS EXAMPLE ILLUSTRATES THE NEED FOR ROBUST COVARIANCE ESTIMATION ON A REAL DATA SET IT IS USEFUL BOTH FOR OUTLIER DETECTION AND FOR A BETTER UNDERSTANDING OF THE DATA STRUCTURE

WE SELECTED TWO SETS OF TWO VARIABLES FROM THE BOSTON HOUSING DATA SET AS AN ILLUSTRATION OF WHAT KIND OF ANALYSIS CAN BE DONE WITH SEVERAL OUTLIER DETECTION TOOLS FOR THE PURPOSE OF VISUALIZATION WE ARE WORKING WITH TWODIMENSIONAL EXAMPLES BUT ONE SHOULD BE AWARE THAT THINGS ARE NOT SO TRIVIAL IN HIGHDIMENSION AS IT WILL BE POINTED OUT IN BOTH EXAMPLES BELOW THE MAIN RESULT IS THAT THE EMPIRICAL COVARIANCE ESTIMATE AS A NONROBUST ONE IS HIGHLY INFLUENCED BY THE HETEROGENEOUS STRUCTURE OF THE OBSERVATIONS ALTHOUGH THE ROBUST COVARIANCE ESTIMATE IS ABLE TO FOCUS ON THE MAIN MODE OF THE DATA DISTRIBUTION IT STICKS TO THE ASSUMPTION THAT THE DATA SHOULD BE GAUSSIAN DISTRIBUTED YIELDING SOME BIASED ESTIMATION OF THE DATA STRUCTURE BUT YET ACCURATE TO SOME EXTENT THE ONECLASS SVM DOES NOT ASSUME ANY PARAMETRIC FORM OF THE DATA DISTRIBUTION AND CAN THEREFORE MODEL THE COMPLEX SHAPE OF THE DATA MUCH BETTER

FIRST EXAMPLE

THE FIRST EXAMPLE ILLUSTRATES HOW ROBUST COVARIANCE ESTIMATION CAN HELP CONCENTRATING ON A RELEVANT CLUSTER WHEN AN OTHER ONE EXISTS HERE MANY OBSERVATIONS ARE CONFOUNDED INTO ONE AND BREAK DOWN THE EMPIRICAL COVARIANCE ESTIMATION OF COURSE SOME SCREENING TOOLS WOULD HAVE POINTED OUT THE PRESENCE OF TWO CLUSTERS SUPPORT VECTOR MACHINES GAUSSIAN MIXTURE MODELS UNIVARIATE OUTLIER DETECTION BUT HAD IT BEEN A HIGHDIMENSIONAL EXAMPLE NONE OF THESE COULD BE APPLIED THAT EASILY

SECOND EXAMPLE

THE SECOND EXAMPLE SHOWS THE ABILITY OF THE MINIMUM COVARIANCE DETERMINANT ROBUST ESTIMATOR OF COVARIANCE TO CONCENTRATE ON THE MAIN MODE OF THE DATA DISTRIBUTION THE LOCATION SEEMS TO BE WELL ESTIMATED ALTHOUGH THE COVARIANCE IS HARD TO ESTIMATE DUE TO THE BANANASHAPED DISTRIBUTION ANYWAY WE CAN GET RID OF SOME OUTLYING OBSERVATIONS THE ONECLASS SVM IS ABLE TO CAPTURE THE REAL DATA STRUCTURE BUT THE DIFFICULTY IS TO ADJUST ITS KERNEL BANDWIDTH PARAMETER SO AS TO OBTAIN A GOOD COMPROMISE BETWEEN THE SHAPE OF THE DATA SCATTER MATRIX AND THE RISK OF OVERFITTING THE DATA

716 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

- 

52 EXAMPLES BASED ON REAL WORLD DATASETS 717

SCIKITLEARN USER GUIDE RELEASE 0213

•

PRINTDOC  
AUTHOR VIRGILE FRITSCH VIRGILEFRITSCHINRIA.FR  
LICENSE BSD 3 CLAUSE  
IMPORT NUMPY AS NP  
FROM SKLEARNCOVARIANCE IMPORT ELLIPTICENVELOPE  
FROM SKLEARN SVM IMPORT ONECLASS SVM  
IMPORT MATPLOTLIBPYPLOT AS PLT  
IMPORT MATPLOTLIBFONTMANAGER  
FROM SKLEARN DATASETS IMPORT LOADBOSTON  
GET DATA  
X1 LOADBOSTON DATA 8 10 TWO CLUSTERS  
X2 LOADBOSTON DATA 5 12 BANANASHAPED  
DEFINE CLASSIFIERS TO BE USED  
CLASSIFIERS  
EMPIRICAL COVARIANCE ELLIPTICENVELOPE SUPPORT FRACTION 1  
CONTAMINATION 0.261  
ROBUST COVARIANCE MINIMUM COVARIANCE DETERMINANT  
ELLIPTICENVELOPE CONTAMINATION 0.261  
OC SVM ONECLASS SVM N U 0.261 GAMMA 0.05  
COLORS M G B  
LEGEND 1  
LEGEND 2  
718 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

LEARN A FRONTIER FOR OUTLIER DETECTION WITH SEVERAL CLASSIFIERS

XX1 YY1 NPMESHGRIDNPLINSPACE8 28 500 NPLINSPACE3 40 500

XX2 YY2 NPMESHGRIDNPLINSPACE3 10 500 NPLINSPACE5 45 500

FORI CLFNAME CLF INENUMERATECLASSIFIERSITEMS

PLTFigure1

CLFFITX1

Z1 CLFDECISIONFUNCTIONNPCXX1RAVEL YY1RAVEL

Z1 Z1RESHAPEXX1SHAPE

LEGEND1CLFNAME PLTCONTOUR

XX1 YY1 Z1 LEVELS0 LINEWIDTHS2 COLORSCOLORSI

PLTFigure2

CLFFITX2

Z2 CLFDECISIONFUNCTIONNPCXX2RAVEL YY2RAVEL

Z2 Z2RESHAPEXX2SHAPE

LEGEND2CLFNAME PLTCONTOUR

XX2 YY2 Z2 LEVELS0 LINEWIDTHS2 COLORSCOLORSI

LEGEND1VALUESLIST LISTLEGEND1VALUES

LEGEND1KEYSLIST LISTLEGEND1KEYS

PLOT THE RESULTS SHAPE OF THE DATA POINTS CLOUD

PLTFigure1 TWO CLUSTERS

PLTTITLEOUTLIER DETECTION ON A REAL DATA SET BOSTON HOUSING

PLTSCATTERX1 0 X1 1 COLORBLACK

BBOXARGS DICTBOXSTYLEROUND FC08

ARROWARGS DICTARROWSTYLE

PLTANNOTATESEVERAL CONFOUNDED POINTS XY24 19

XYCOORDSDATA TEXTCOORDSDATA

XYTEXT13 10 BBOXBBOXARGS ARROWPROPSARROWARGS

PLTXLIMXX1MIN XX1MAX

PLTYLIMYY1MIN YY1MAX

PLTLEGENDLEGEND1VALUESLIST0COLLECTIONS0

LEGEND1VALUESLIST1COLLECTIONS0

LEGEND1VALUESLIST2COLLECTIONS0

LEGEND1KEYSLIST0 LEGEND1KEYSLIST1 LEGEND1KEYSLIST2

LOCUPPER CENTER

PROPMATPLOTLIBFONTMANAGERFONTPROPERTIESIZE12

PLTYLABELACCESSIBILITY TO RADIAL HIGHWAYS

PLTXLABELPUPILTEACHER RATIO BY TOWN

LEGEND2VALUESLIST LISTLEGEND2VALUES

LEGEND2KEYSLIST LISTLEGEND2KEYS

PLTFigure2 BANANA SHAPE

PLTTITLEOUTLIER DETECTION ON A REAL DATA SET BOSTON HOUSING

PLTSCATTERX2 0 X2 1 COLORBLACK

PLTXLIMXX2MIN XX2MAX

PLTYLIMYY2MIN YY2MAX

PLTLEGENDLEGEND2VALUESLIST0COLLECTIONS0

LEGEND2VALUESLIST1COLLECTIONS0

LEGEND2VALUESLIST2COLLECTIONS0

LEGEND2KEYSLIST0 LEGEND2KEYSLIST1 LEGEND2KEYSLIST2

LOCUPPER CENTER

PROPMATPLOTLIBFONTMANAGERFONTPROPERTIESIZE12

PLTYLABEL LOWER STATUS OF THE POPULATION

PLTXLABELAVERAGE NUMBER OF ROOMS PER DWELLING

PLTSHOW

52 EXAMPLES BASED ON REAL WORLD DATASETS 719

SCIKITLEARN USER GUIDE RELEASE 0213

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 3436 SECONDS

NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)

522 COMPRESSIVE SENSING TOMOGRAPHY RECONSTRUCTION WITH L1 PRIOR LASSO

THIS EXAMPLE SHOWS THE RECONSTRUCTION OF AN IMAGE FROM A SET OF PARALLEL PROJECTIONS ACQUIRED ALONG DIFFERENT ANGLES SUCH A DATASET IS ACQUIRED IN COMPUTED TOMOGRAPHY CT

WITHOUT ANY PRIOR INFORMATION ON THE SAMPLE THE NUMBER OF PROJECTIONS REQUIRED TO RECONSTRUCT THE IMAGE IS OF THE ORDER OF THE LINEAR SIZE LOF THE IMAGE IN PIXELS FOR SIMPLICITY WE CONSIDER HERE A SPARSE IMAGE WHERE ONLY PIXELS ON THE BOUNDARY OF OBJECTS HAVE A NONZERO VALUE SUCH DATA COULD CORRESPOND FOR EXAMPLE TO A CELLULAR MATERIAL

NOTE HOWEVER THAT MOST IMAGES ARE SPARSE IN A DIFFERENT BASIS SUCH AS THE HAAR WAVELETS ONLY L7 PROJECTIONS ARE ACQUIRED THEREFORE IT IS NECESSARY TO USE PRIOR INFORMATION AVAILABLE ON THE SAMPLE ITS SPARSITY THIS IS AN EXAMPLE OF COMPRESSIVE SENSING

THE TOMOGRAPHY PROJECTION OPERATION IS A LINEAR TRANSFORMATION IN ADDITION TO THE DATAFIDELITY TERM CORRESPONDING TO A LINEAR REGRESSION WE PENALIZE THE L1 NORM OF THE IMAGE TO ACCOUNT FOR ITS SPARSITY THE RESULTING OPTIMIZATION PROBLEM IS CALLED THE LASSO WE USE THE CLASS SKLEARNLINEARMODELLASSO THAT USES THE COORDINATE DESCENT ALGORITHM IMPORTANTLY THIS IMPLEMENTATION IS MORE COMPUTATIONALLY EFFICIENT ON A SPARSE MATRIX THAN THE PROJECTION OPERATOR USED HERE

THE RECONSTRUCTION WITH L1 PENALIZATION GIVES A RESULT WITH ZERO ERROR ALL PIXELS ARE SUCCESSFULLY LABELED WITH 0 OR 1 EVEN IF NOISE WAS ADDED TO THE PROJECTIONS IN COMPARISON AN L2 PENALIZATION SKLEARNLINEARMODELRIDGE PRODUCES A LARGE NUMBER OF LABELING ERRORS FOR THE PIXELS IMPORTANT ARTIFACTS ARE OBSERVED ON THE RECONSTRUCTED IMAGE CONTRARY TO THE L1 PENALIZATION NOTE IN PARTICULAR THE CIRCULAR ARTIFACT SEPARATING THE PIXELS IN THE CORNERS THAT HAVE CONTRIBUTED TO FEWER PROJECTIONS THAN THE CENTRAL DISK

PRINTDOC

AUTHOR EMMANUELLE GOUILLART EMMANUELLEGOUILLARTNSUPORG

LICENSE BSD 3 CLAUSE

IMPORT NUMPY AS NP

FROM SCIPY IMPORT SPARSE

FROM SCIPY IMPORT NDIMAGE

FROM SKLEARNLINEARMODEL IMPORT LASSO

720 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
FROM SKLEARNLINEARMODEL IMPORT RIDGE  
IMPORT MATPLOTLIBPYPLOT AS PLT  
DEFWEIGHTSX DX1 ORIG0  
X NPRAVELX  
FLOORX NPFLOORX ORIG DXASTYPENPINT64  
ALPHA X ORIG FLOORX DX DX  
RETURNNPHSTACKFLOORX FLOORX 1 NPHSTACK1 ALPHA ALPHA  
DEFGENERATECENTERCOORDINATESLX  
X Y NPMGRIDLX LXASTYPENPFLOAT64  
CENTER LX 2  
X 05 CENTER  
Y 05 CENTER  
RETURNX Y  
DEFBUILDPROJECTIONOPERATORLX NDIR  
COMPUTE THE TOMOGRAPHY DESIGN MATRIX  
PARAMETERS  
  
LX INT  
LINEAR SIZE OF IMAGE ARRAY  
NDIR INT  
NUMBER OF ANGLES AT WHICH PROJECTIONS ARE ACQUIRED  
RETURNS  
  
P SPARSE MATRIX OF SHAPE NDIR LX LX 2  
  
X Y GENERATECENTERCOORDINATESLX  
ANGLES NPLinspace0 NPPI NDIR ENDPOINTFALSE  
DATAINDS WEIGHTS CAMERAINDS  
DATAUNRAVELINDICES NPARANGELX 2  
DATAUNRAVELINDICES NPHSTACKDATAUNRAVELINDICES  
DATAUNRAVELINDICES  
FORI ANGLE INENUMERATEANGLES  
XROT NPCOSANGLE X NPSINANGLE Y  
INDS W WEIGHTSXROT DX1 ORIGXMIN  
MASK NPLOGICALANDINDS 0 INDs LX  
WEIGHTS LISTWMASK  
CAMERAINDS LISTINDSMASK I LX  
DATAINDS LISTDATAUNRAVELINDICESMASK  
PROJOPERATOR SPARSECOOMATRIXWEIGHTS CAMERAINDS DATAINDS  
RETURNPROJOPERATOR  
DEFGENERATESYNTHETICDATA  
SYNTHETIC BINARY DATA  
RS NPRANDOMRANDOMSTATE0  
NPTS 36  
X Y NPOGRID0L 0L  
MASKOUTER X L 2 2 Y L 2 2 L 2 2  
52 EXAMPLES BASED ON REAL WORLD DATASETS 721

SCIKITLEARN USER GUIDE RELEASE 0213  
MASK NPZEROSL L  
POINTS L RSRAND2 NPTS  
MASKPOINTS0ASTYPENPINT POINTS1ASTYPENPINT 1  
MASK NIMAGEGAUSSIANFILTERMASK SIGMAL NPTS  
RES NPLOGICALANDMASK MASKMEAN MASKOUTER  
RETURNNPLOGICALXORRES NIMAGEBINARYEROSIONRES  
GENERATE SYNTHETIC IMAGES AND PROJECTIONS  
L 128  
PROJOPERATOR BUILDPROJECTIONOPERATORL L 7  
DATA GENERATESYNTHETICDATA  
PROJ PROJOPERATOR DATARAVEL NPNEWAXIS  
PROJ 015 NPRANDOMRANDN PROJSHAPE  
RECONSTRUCTION WITH L2 RIDGE PENALIZATION  
RGRRIDGE RIDGEALPHA02  
RGRRIDGEFITPROJOPERATOR PROJRAVEL  
RECL2 RGRRIDGECOEFFRESHAPEL L  
RECONSTRUCTION WITH L1 LASSO PENALIZATION  
THE BEST VALUE OF ALPHA WAS DETERMINED USING CROSS VALIDATION  
WITH LASSOCV  
RGRLASSO LASSOALPHA0001  
RGRLASSOFITPROJOPERATOR PROJRAVEL  
RECL1 RGRLASSOCOEFFRESHAPEL L  
PLTFIGUREFIGSIZE8 33  
PLTSUBPLOT131  
PLTIMSHOWDATA CMAPPLTCMGRAY INTERPOLATIONNEAREST  
PLTAXISOFF  
PLTTITLEORIGINAL IMAGE  
PLTSUBPLOT132  
PLTIMSHOWRECL2 CMAPPLTCMGRAY INTERPOLATIONNEAREST  
PLTTITLEL2 PENALIZATION  
PLTAXISOFF  
PLTSUBPLOT133  
PLTIMSHOWRECL1 CMAPPLTCMGRAY INTERPOLATIONNEAREST  
PLTTITLEL1 PENALIZATION  
PLTAXISOFF  
PLTSUBPLOTSADJUSTHSPACE001 WSPACE001 TOP1 BOTTOM0 LEFT0  
RIGHT1  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 9761 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
523 TOPIC EXTRACTION WITH NONNEGATIVE MATRIX FACTORIZATION AND LATENT DIRICHLET  
ALLOCATION  
THIS IS AN EXAMPLE OF APPLYING SKLEARNDECOMPOSITIONNMF ANDSKLEARNDECOMPOSITION  
LATENTDIRICHLETALLOCATION ON A CORPUS OF DOCUMENTS AND EXTRACT ADDITIVE MODELS OF THE TOPIC STRUCTURE  
722 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213

OF THE CORPUS THE OUTPUT IS A LIST OF TOPICS EACH REPRESENTED AS A LIST OF TERMS WEIGHTS ARE NOT SHOWN

NONNEGATIVE MATRIX FACTORIZATION IS APPLIED WITH TWO DIFFERENT OBJECTIVE FUNCTIONS THE FROBENIUS NORM AND THE

GENERALIZED KULLBACKLEIBLER DIVERGENCE THE LATTER IS EQUIVALENT TO PROBABILISTIC LATENT SEMANTIC INDEXING

THE DEFAULT PARAMETERS NSAMPLES NFEATURES NCOMPONENTS SHOULD MAKE THE EXAMPLE RUNNABLE IN A COUPLE OF

TENS OF SECONDS YOU CAN TRY TO INCREASE THE DIMENSIONS OF THE PROBLEM BUT BE AWARE THAT THE TIME COMPLEXITY IS

POLYNOMIAL IN NMF IN LDA THE TIME COMPLEXITY IS PROPORTIONAL TO NSAMPLES ITERATIONS

OUT

LOADING DATASET

DONE IN 7911S

EXTRACTING TFIDF FEATURES FOR NMF

DONE IN 0268S

EXTRACTING TF FEATURES FOR LDA

DONE IN 0254S

FITTING THE NMF MODEL FROBENIUS NORM WITH TFIDF FEATURES NSAMPLES2000 AND N

↪FEATURES1000

DONE IN 0406S

TOPICS IN NMF MODEL FROBENIUS NORM

TOPIC 0 JUST PEOPLE DON THINK LIKE KNOW TIME GOOD MAKE WAY REALLY SAY RIGHT VE WANT

↪DID LL NEW USE YEARS

TOPIC 1 WINDOWS USE DOS USING WINDOW PROGRAM OS DRIVERS APPLICATION HELP SOFTWARE

↪PC RUNNING MS SCREEN FILES VERSION CARD CODE WORK

TOPIC 2 GOD JESUS BIBLE FAITH CHRISTIAN CHRIST CHRISTIANS DOES HEAVEN SIN BELIEVE

↪LORD LIFE CHURCH MARY ATHEISM BELIEF HUMAN LOVE RELIGION

TOPIC 3 THANKS KNOW DOES MAIL ADVANCE HI INFO INTERESTED EMAIL ANYBODY LOOKING CARD

↪HELP LIKE APPRECIATED INFORMATION SEND LIST VIDEO NEED

TOPIC 4 CAR CARS TIRES MILES 00 NEW ENGINE INSURANCE PRICE CONDITION OIL POWER

↪SPEED GOOD 000 BRAKE YEAR MODELS USED BOUGHT

TOPIC 5 EDU SOON COM SEND UNIVERSITY INTERNET MIT FTP MAIL CC PUB ARTICLE

↪INFORMATION HOPE PROGRAM MAC EMAIL HOME CONTACT BLOOD

TOPIC 6 FILE PROBLEM FILES FORMAT WIN SOUND FTP PUB READ SAVE SITE HELP IMAGE

↪AVAILABLE CREATE COPY RUNNING MEMORY SELF VERSION

TOPIC 7 GAME TEAM GAMES YEAR WIN PLAY SEASON PLAYERS NHL RUNS GOAL HOCKEY TORONTO

↪DIVISION FLYERS PLAYER DEFENSE LEAFS BAD TEAMS

TOPIC 8 DRIVE DRIVES HARD DISK FLOPPY SOFTWARE CARD MAC COMPUTER POWER SCSI

↪CONTROLLER APPLE MB 00 PC ROM SALE PROBLEM INTERNAL

TOPIC 9 KEY CHIP CLIPPER KEYS ENCRYPTION GOVERNMENT PUBLIC USE SECURE ENFORCEMENT

↪PHONE NSA COMMUNICATIONS LAW ENCRYPTED SECURITY CLINTON USED LEGAL STANDARD

FITTING THE NMF MODEL GENERALIZED KULLBACKLEIBLER DIVERGENCE WITH TFIDF FEATURES

↪NSAMPLES2000 AND NFEATURES1000

DONE IN 1769S

TOPICS IN NMF MODEL GENERALIZED KULLBACKLEIBLER DIVERGENCE

TOPIC 0 JUST PEOPLE DON LIKE DID KNOW MAKE REALLY RIGHT THINK SAY THINGS TIME LOOK

↪WAY DIDN VE COURSE PROBABLY GOOD

TOPIC 1 HELP THANKS WINDOWS KNOW HI NEED USING DOES LOOKING ANYBODY APPRECIATED

↪CARD MAIL SOFTWARE USE INFO EMAIL FTP AVAILABLE PC

TOPIC 2 DOES GOD BELIEVE KNOW MEAN TRUE CHRISTIANS READ POINT JESUS CHRISTIAN

↪CHURCH COME PEOPLE FACT SAYS RELIGION SAY AGREE BIBLE

TOPIC 3 KNOW THANKS MAIL INTERESTED LIKE NEW JUST BIKE EMAIL EDU ADVANCE WANT

↪CONTACT REALLY LIST HEARD COM POST HEAR INFORMATION

TOPIC 4 10 NEW 30 12 20 50 11 SALE 16 15 TIME 14 OLD POWER AGO GOOD 100 GREAT OFFER

↪COST

52 EXAMPLES BASED ON REAL WORLD DATASETS 723

SCIKITLEARN USER GUIDE RELEASE 0213

TOPIC 5 NUMBER 1993 DATA SUBJECT GOVERNMENT NEW NUMBERS PROVIDE INFORMATION SPACE

↳FOLLOWING COM RESEARCH INCLUDE LARGE NOTE GROUP MAJOR TIME TALK

TOPIC 6 EDU PROBLEM FILE COM REMEMBER TRY SOON ARTICLE MIKE FILES CODE PROGRAM SUN

↳FREE SEND THINK CASES MANAGER LITTLE CALLED

TOPIC 7 GAME YEAR TEAM GAMES WORLD FACT SECOND CASE WON SAID WIN DIVISION PLAY BEST

↳CLEARLY CLAIM ALLOW EXAMPLE USED DOESN

TOPIC 8 THINK DON DRIVE HARD NEED BIT MAC MAKE SURE READ APPLE GOING COMES DISK

↳COMPUTER CASE PRETTY DRIVES SOFTWARE VE

TOPIC 9 GOOD JUST USE LIKE DOESN GOT WAY DON LL GOING DOES CHIP BETTER DOING BAD

↳KEY WANT SURE BIT CAR

FITTING LDA MODELS WITH TF FEATURES NSAMPLES2000 AND NFEATURES1000

DONE IN 31675

TOPICS IN LDA MODEL

TOPIC 0 EDU COM MAIL SEND GRAPHICS FTP PUB AVAILABLE CONTACT UNIVERSITY LIST FAQ CA

↳INFORMATION CS 1993 PROGRAM SUN UK MIT

TOPIC 1 DON LIKE JUST KNOW THINK VE WAY USE RIGHT GOOD GOING MAKE SURE LL POINT GOT

↳NEED REALLY TIME DOESN

TOPIC 2 CHRISTIAN THINK ATHEISM FAITH PITTSBURGH NEW BIBLE RADIO GAMES ALT LOT JUST

↳RELIGION LIKE BOOK READ PLAY TIME SUBJECT BELIEVE

TOPIC 3 DRIVE DISK WINDOWS THANKS USE CARD DRIVES HARD VERSION PC SOFTWARE FILE

↳USING SCSI HELP DOES NEW DOS CONTROLLER 16

TOPIC 4 HIV HEALTH AIDS DISEASE APRIL MEDICAL CARE RESEARCH 1993 LIGHT INFORMATION

↳STUDY NATIONAL SERVICE TEST LED 10 PAGE NEW DRUG

TOPIC 5 GOD PEOPLE DOES JUST GOOD DON JESUS SAY ISRAEL WAY LIFE KNOW TRUE FACT TIME

↳LAW WANT BELIEVE MAKE THINK

TOPIC 6 55 10 11 18 15 TEAM GAME 19 PERIOD PLAY 23 12 13 FLYERS 20 25 22 17 24 16

TOPIC 7 CAR YEAR JUST CARS NEW ENGINE LIKE BIKE GOOD OIL INSURANCE BETTER TIRES 000

↳THING SPEED MODEL BRAKE DRIVING PERFORMANCE

TOPIC 8 PEOPLE SAID DID JUST DIDN KNOW TIME LIKE WENT THINK CHILDREN CAME COME DON

↳TOOK YEARS SAY DEAD TOLD STARTED

TOPIC 9 KEY SPACE LAW GOVERNMENT PUBLIC USE ENCRYPTION EARTH SECTION SECURITY MOON

↳PROBE ENFORCEMENT KEYS STATES LUNAR MILITARY CRIME SURFACE TECHNOLOGY

AUTHOR OLIVIER GRISEL OLIVIERGRISELENSTAORG

LARS BUITINCK

CHYIKWEI YAU CHYIKWEIYAU@GMAIL.COM

LICENSE BSD 3 CLAUSE

FROM TIME IMPORT TIME

FROM SKLEARNFEATUREEXTRACTIONTEXT IMPORT TFIDFVECTORIZER COUNTVECTORIZER

FROM SKLEARNDECOMPOSITION IMPORT NMF LATENTDIRICHLETALLOCATION

FROM SKLEARNDATASETS IMPORT FETCH20NEWSGROUPS

NSAMPLES 2000

NFEATURES 1000

NCOMPONENTS 10

NTOPWORDS 20

DEFPRINTTOPWORDSMODEL FEATURENAMES NTOPWORDS

724 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
FORTOPICIDX TOPIC INENUMERATEMODELCOMPONENTS
MESSAGE TOPIC D TOPICIDX
MESSAGE JOINFEATURENAMESI
FORIINTOPICARGSORTNTOPWORDS 11
PRINTMESSAGE
PRINT
LOAD THE 20 NEWSGROUPS DATASET AND VECTORIZE IT WE USE A FEW HEURISTICS
TO FILTER OUT USELESS TERMS EARLY ON THE POSTS ARE STRIPPED OF HEADERS
FOOTERS AND QUOTED REPLIES AND COMMON ENGLISH WORDS WORDS OCCURRING IN
ONLY ONE DOCUMENT OR IN AT LEAST 95 OF THE DOCUMENTS ARE REMOVED
PRINTLOADING DATASET
TO TIME
DATASET FETCH20NEWSGROUPSSHUFFLETRUE RANDOMSTATE1
REMOVEHEADERS FOOTERS QUOTES
DATASAMPLES DATASETDATANSAMPLES
PRINTDONE IN 03FS TIME TO
USE TFIDF FEATURES FOR NMF
PRINTEXTRACTING TFIDF FEATURES FOR NMF
TFIDFVECTORIZER TFIDFVECTORIZERMAXDF095 MINDF2
MAXFEATURESNFEATURES
STOPWORDSENGLISH
TO TIME
TFIDF TFIDFVECTORIZERFITTRANSFORMDATASAMPLES
PRINTDONE IN 03FS TIME TO
USE TF RAW TERM COUNT FEATURES FOR LDA
PRINTEXTRACTING TF FEATURES FOR LDA
TFVECTORIZER COUNTVECTORIZERMAXDF095 MINDF2
MAXFEATURESNFEATURES
STOPWORDSENGLISH
TO TIME
TF TFVECTORIZERFITTRANSFORMDATASAMPLES
PRINTDONE IN 03FS TIME TO
PRINT
FIT THE NMF MODEL
PRINTFITTING THE NMF MODEL FROBENIUS NORM WITH TFIDF FEATURES
NSAMPLES DAND NFEATURES D
NSAMPLES NFEATURES
TO TIME
NMF NMFNCOMPONENTSNSCOMPONENTS RANDOMSTATE1
ALPHA1 L1RATIO5FITTFIDF
PRINTDONE IN 03FS TIME TO
PRINTNTOPICS IN NMF MODEL FROBENIUS NORM
TFIDFFEATURENAMES TFIDFVECTORIZERGETFEATURENAMES
PRINTTOPWORDSNMF TFIDFFEATURENAMES NTOPWORDS
FIT THE NMF MODEL
PRINTFITTING THE NMF MODEL GENERALIZED KULLBACKLEIBLER DIVERGENCE WITH
TFIDF FEATURES NSAMPLES DAND NFEATURES D
NSAMPLES NFEATURES
TO TIME
NMF NMFNCOMPONENTSNSCOMPONENTS RANDOMSTATE1
52 EXAMPLES BASED ON REAL WORLD DATASETS 725
```

SCIKITLEARN USER GUIDE RELEASE 0213  
BETALOSSKULLBACKLEIBLER SOLVERMU MAXITER1000 ALPHA1  
L1RATIO5FITTFIDF  
PRINTDONE IN 03FS TIME TO  
PRINTNTOPICS IN NMF MODEL GENERALIZED KULLBACKLEIBLER DIVERGENCE  
TFIDFFEATURENAMES TFIDFVECTORIZERGETFEATURENAMES  
PRINTTOPWORDSNMF TFIDFFEATURENAMES NTOPWORDS  
PRINTFITTING LDA MODELS WITH TF FEATURES  
NSAMPLES DAND NFEATURES D  
NSAMPLES NFEATURES  
LDA LATENTDIRICHLETALLOCATIONNNCOMPONENTSNSCOMPONENTS MAXITER5  
LEARNINGMETHODONLINE  
LEARNINGOFFSET50  
RANDOMSTATE0  
TO TIME  
LDAFITTF  
PRINTDONE IN 03FS TIME TO  
PRINTNTOPICS IN LDA MODEL  
TFFEATURENAMES TFVECTORIZERGETFEATURENAMES  
PRINTTOPWORDSLDA TFFEATURENAMES NTOPWORDS  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 13781 SECONDS  
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE  
524 FACES RECOGNITION EXAMPLE USING EIGENFACES AND SVMs  
THE DATASET USED IN THIS EXAMPLE IS A PREPROCESSED EXCERPT OF THE “LABELED FACES IN THE WILD” AKA LFW  
HTTPVISWWWCSUMASSEDULFWLFWFUNNELEDTGZ 233MB  
EXPECTED RESULTS FOR THE TOP 5 MOST REPRESENTED PEOPLE IN THE DATASET  
ARIEL SHARON 067 092 077 13  
COLIN POWELL 075 078 076 60  
DONALD RUMSFELD 078 067 072 27  
GEORGE W BUSH 086 086 086 146  
GERHARD SCHROEDER 076 076 076 25  
HUGO CHAVEZ 067 067 067 15  
TONY BLAIR 081 069 075 36  
AVG TOTAL 080 080 080 322  
726 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

•

52 EXAMPLES BASED ON REAL WORLD DATASETS 727

SCIKITLEARN USER GUIDE RELEASE 0213

•

OUT

TOTAL DATASET SIZE

NSAMPLES 1288

NFEATURES 1850

NCLASSES 7

EXTRACTING THE TOP 150 EIGENFACES FROM 966 FACES

DONE IN 0118S

PROJECTING THE INPUT DATA ON THE EIGENFACES ORTHONORMAL BASIS

DONE IN 0005S

FITTING THE CLASSIFIER TO THE TRAINING SET

DONE IN 36078S

BEST ESTIMATOR FOUND BY GRID SEARCH

SVCC10000 CLASSWEIGHTBALANCED GAMMA0005

PREDICTING PEOPLES NAMES ON THE TEST SET

DONE IN 0061S

PRECISION RECALL F1SCORE SUPPORT

728 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
ARIEL SHARON 075 046 057 13
COLIN POWELL 079 087 083 60
DONALD RUMSFELD 094 063 076 27
GEORGE W BUSH 083 098 090 146
GERHARD SCHROEDER 091 080 085 25
HUGO CHAVEZ 100 053 070 15
TONY BLAIR 096 075 084 36
ACCURACY 085 322
MACRO AVG 088 072 078 322
WEIGHTED AVG 086 085 084 322
6 2 0 5 0 0 0
1 52 0 7 0 0 0
1 3 17 6 0 0 0
0 3 0 143 0 0 0
0 1 0 3 20 0 1
0 4 0 2 1 8 0
0 1 1 6 1 0 27
FROM TIME IMPORT TIME
IMPORT LOGGING
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARNMODELSELECTION IMPORT TRAINTESTSPLIT
FROM SKLEARNMODELSELECTION IMPORT GRIDSEARCHCV
FROM SKLEARNDATASETS IMPORT FETCHLFWPEOPLE
FROM SKLEARNMETRICS IMPORT CLASSIFICATIONREPORT
FROM SKLEARNMETRICS IMPORT CONFUSIONMATRIX
FROM SKLEARNDECOMPOSITION IMPORT PCA
FROM SKLEARNSVM IMPORT SVC
PRINTDOC
    DISPLAY PROGRESS LOGS ON STDOUT
LOGGINGBASICCONFIGLEVELLOGGINGINFO FORMAT ASCTIMES MESSAGES

DOWNLOAD THE DATA IF NOT ALREADY ON DISK AND LOAD IT AS NUMPY ARRAYS
LFWPEOPLE  FETCHLFWPEOPLEMINFACESPERPERSON70 RESIZE04
INTROSPECT THE IMAGES ARRAYS TO FIND THE SHAPES FOR PLOTTING
NSAMPLES H W  LFWPEOPLEIMAGESSHAPE
FOR MACHINE LEARNING WE USE THE 2 DATA DIRECTLY AS RELATIVE PIXEL
POSITIONS INFO IS IGNORED BY THIS MODEL
X  LFWPEOPLEDATA
NFEATURES  XSHAPE1
    THE LABEL TO PREDICT IS THE ID OF THE PERSON
52 EXAMPLES BASED ON REAL WORLD DATASETS 729
```

SCIKITLEARN USER GUIDE RELEASE 0213

Y LFWPEOPLETARGET  
TARGETNAMES LFWPEOPLETARGETNAMES  
NCLASSES TARGETNAMESSHAPE0  
PRINTTOTAL DATASET SIZE  
PRINTNSAMPLES D NSAMPLES  
PRINTNFEATURES D NFEATURES  
PRINTNCLASSES D NCLASSES

SPLIT INTO A TRAINING SET AND A TEST SET USING A STRATIFIED K FOLD  
SPLIT INTO A TRAINING AND TESTING SET  
XTRAIN XTEST YTRAIN YTEST TRAINTESTSPLIT  
X Y TESTSIZE025 RANDOMSTATE42

COMPUTE A PCA EIGENFACES ON THE FACE DATASET TREATED AS UNLABELED  
DATASET UNSUPERVISED FEATURE EXTRACTION DIMENSIONALITY REDUCTION  
NCOMPONENTS 150  
PRINTEXTRACTING THE TOP DEIGENFACES FROM DFACES  
NCOMPONENTS XTRAINSHAPE0  
T0 TIME  
PCA PCANCOMPONENTSNCOMPONENTS SVDSOLVERRANDOMIZED  
WHITENTRUEFITXTRAIN  
PRINTDONE IN 03FS TIME T0  
EIGENFACES PCACOMPONENTSRESHAPENCOMPONENTS H W  
PRINTPROJECTING THE INPUT DATA ON THE EIGENFACES ORTHONORMAL BASIS  
T0 TIME  
XTRAINPCA PCATransformXTRAIN  
XTESTPCA PCATransformXTEST  
PRINTDONE IN 03FS TIME T0

TRAIN A SVM CLASSIFICATION MODEL  
PRINTFITTING THE CLASSIFIER TO THE TRAINING SET  
T0 TIME  
PARAMGRID C 1E3 5E3 1E4 5E4 1E5  
GAMMA 00001 00005 0001 0005 001 01  
CLF GRIDSEARCHCVSVCKERNELRBF CLASSWEIGHTBALANCED  
PARAMGRID CV5 IIDFALSE  
CLF CLFFITXTRAINPCA YTRAIN  
PRINTDONE IN 03FS TIME T0  
PRINTBEST ESTIMATOR FOUND BY GRID SEARCH  
PRINTCLFBESTESTIMATOR

QUANTITATIVE EVALUATION OF THE MODEL QUALITY ON THE TEST SET  
PRINTPREDICTING PEOPLES NAMES ON THE TEST SET  
730 CHAPTER 5 EXAMPLES



```
SCIKITLEARN USER GUIDE RELEASE 0213
TO TIME
YPRED CLFPREDICTXTESTPCA
PRINTDONE IN 03FS TIME TO
PRINTCLASSIFICATIONREPORTYTEST YPRED TARGETNAMESTARGETNAMES
PRINTCONFUSIONMATRIXYTEST YPRED LABELSRANGENCLASSES

QUALITATIVE EVALUATION OF THE PREDICTIONS USING MATPLOTLIB
DEFPLOTGALLERYIMAGES TITLES H W NROW3 NCOL4
HELPER FUNCTION TO PLOT A GALLERY OF PORTRAITS
PLTFIGUREFIGSIZE18 NCOL 24 NROW
PLTSUBPLOTSADJUSTBOTTOM0 LEFT01 RIGHT99 TOP90 HSPACE35
FORIINRANGENROW NCOL
PLTSUBPLOTNROW NCOL I 1
PLTIMSHOWIMAGESIRESHAPEH W CMAPPLTCMGRAY
PLTTITLETITLES I SIZE12
PLTXTICKS
PLTYTICKS
PLOT THE RESULT OF THE PREDICTION ON A PORTION OF THE TEST SET
DEFTITLEYPRED YTEST TARGETNAMES I
PREDNAME TARGETNAMESYPREDIRSPLOT 11
TRUENAME TARGETNAMESYTESTIRSPLOT 11
RETURNPREDICTED SNTRUE S PREDNAME TRUENAME
PREDICTIONTITLES TITLEYPRED YTEST TARGETNAMES I
FORIINRANGEYPREDSHAPE0
PLOTGALLERYXTEST PREDICTIONTITLES H W
PLOT THE GALLERY OF THE MOST SIGNIFICATIVE EIGENFACES
EIGENFACETITLES EIGENFACE D IFORIINRANGE EIGENFACESSHAPE0
PLOTGALLERYEIGENFACES EIGENFACETITLES H W
PLTSHOW
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 59514 SECONDS
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE
525 MODEL COMPLEXITY INFLUENCE
DEMONSTRATE HOW MODEL COMPLEXITY INFLUENCES BOTH PREDICTION ACCURACY AND COMPUTATIONAL PERFORMANCE
THE DATASET IS THE BOSTON HOUSING DATASET RESP 20 NEWSGROUPS FOR REGRESSION RESP CLASSIFICATION
FOR EACH CLASS OF MODELS WE MAKE THE MODEL COMPLEXITY VARY THROUGH THE CHOICE OF RELEVANT MODEL PARAMETERS AND
MEASURE THE INFLUENCE ON BOTH COMPUTATIONAL PERFORMANCE LATENCY AND PREDICTIVE POWER MSE OR HAMMING LOSS
52 EXAMPLES BASED ON REAL WORLD DATASETS 731
```

SCIKITLEARN USER GUIDE RELEASE 0213

- 
- 

732 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

- OUT  
BENCHMARKING SGDCLASSIFIERALPHA0001 L1RATIO025 LOSSMODIFIEDHUBER  
PENALTYELASTICNET  
COMPLEXITY 4466 HAMMING LOSS MISCLASSIFICATION RATIO 02491 PRED TIME 0  
↪021496S  
BENCHMARKING SGDCLASSIFIERALPHA0001 L1RATIO05 LOSSMODIFIEDHUBER  
PENALTYELASTICNET  
COMPLEXITY 1663 HAMMING LOSS MISCLASSIFICATION RATIO 02915 PRED TIME 0  
↪017370S  
BENCHMARKING SGDCLASSIFIERALPHA0001 L1RATIO075 LOSSMODIFIEDHUBER  
PENALTYELASTICNET  
COMPLEXITY 880 HAMMING LOSS MISCLASSIFICATION RATIO 03180 PRED TIME 0  
↪012989S  
BENCHMARKING SGDCLASSIFIERALPHA0001 L1RATIO09 LOSSMODIFIEDHUBER  
PENALTYELASTICNET  
COMPLEXITY 639 HAMMING LOSS MISCLASSIFICATION RATIO 03337 PRED TIME 0  
↪011292S  
BENCHMARKING NUSVRC10000 GAMMA30517578125E05 NU01  
COMPLEXITY 69 MSE 318139 PRED TIME 0000283S  
BENCHMARKING NUSVRC10000 GAMMA30517578125E05 NU025  
COMPLEXITY 136 MSE 256140 PRED TIME 0000506S  
BENCHMARKING NUSVRC10000 GAMMA30517578125E05  
COMPLEXITY 244 MSE 223375 PRED TIME 0000868S  
BENCHMARKING NUSVRC10000 GAMMA30517578125E05 NU075  
COMPLEXITY 351 MSE 213688 PRED TIME 0001226S  
BENCHMARKING NUSVRC10000 GAMMA30517578125E05 NU09  
COMPLEXITY 404 MSE 211033 PRED TIME 0001402S  
BENCHMARKING GRADIENTBOOSTINGREGRESSORNESTIMATORS10  
52 EXAMPLES BASED ON REAL WORLD DATASETS 733

```
SCIKITLEARN USER GUIDE RELEASE 0213
COMPLEXITY 10 MSE 290148 PRED TIME 0000093S
BENCHMARKING GRADIENTBOOSTINGREGRESSORNESTIMATORS50
COMPLEXITY 50 MSE 89630 PRED TIME 0000165S
BENCHMARKING GRADIENTBOOSTINGREGRESSOR
COMPLEXITY 100 MSE 77187 PRED TIME 0000227S
BENCHMARKING GRADIENTBOOSTINGREGRESSORNESTIMATORS200
COMPLEXITY 200 MSE 66955 PRED TIME 0000608S
BENCHMARKING GRADIENTBOOSTINGREGRESSORNESTIMATORS500
COMPLEXITY 500 MSE 71437 PRED TIME 0000776S
PRINTDOC
AUTHOR EUSTACHE DIEMERT EUSTACHEDIEMERTFR
LICENSE BSD 3 CLAUSE
IMPORT TIME
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM MPLTOOLKITSAXESGRID1PARASITEAXES IMPORT HOSTSUBPLOT
FROM MPLTOOLKITSAXISARTISTAXISLINES IMPORT AXES
FROM SCIPYSPARSECSR IMPORT CSRMATRIX
FROM SKLEARN IMPORT DATASETS
FROM SKLEARNUTILS IMPORT SHUFFLE
FROM SKLEARNMETRICS IMPORT MEANSQUAREDERROR
FROM SKLEARN SVMCLASSES IMPORT NUSVR
FROM SKLEARNENSEMBLEGRADIENTBOOSTING IMPORT GRADIENTBOOSTINGREGRESSOR
FROM SKLEARNLINEARMODELSTOCHASTICGRADIENT IMPORT SGDCLASSIFIER
FROM SKLEARNMETRICS IMPORT HAMMINGLOSS

ROUTINES
INITIALIZE RANDOM GENERATOR
NPRANDOMSEED0
DEFGENERATEDATACASE SPARSEFALSE
GENERATE REGRESSIONCLASSIFICATION DATA
BUNCH NONE
IFCASE REGRESSION
BUNCH DATASETSLOADBOSTON
ELIFCASE CLASSIFICATION
BUNCH DATASETSFETCH20NEWSGROUPSVECTORIZEDSUBSETALL
X Y SHUFFLEBUNCHDATA BUNCHTARGET
OFFSET INTXSHAPE0 08
XTRAIN YTRAIN XOFFSET YOFFSET
XTEST YTEST XOFFSET YOFFSET
734 CHAPTER 5 EXAMPLES
```

SCIKITLEARN USER GUIDE RELEASE 0213  
IFSPARSE  
XTRAIN CSRMATRIXXTTRAIN  
XTEST CSRMATRIXXTTEST  
ELSE  
XTRAIN NPARRAYXTTRAIN  
XTEST NPARRAYXTTEST  
YTEST NPARRAYYTEST  
YTRAIN NPARRAYYTRAIN  
DATA XTRAIN XTRAIN XTEST XTEST YTRAIN YTRAIN  
YTEST YTEST  
RETURNNDATA  
DEFBENCHMARKINFLUENCECONF

BENCHMARK INFLUENCE OF CHANGINGPARAM ON BOTH MSE AND LATENCY

PREDICTIONTIMES  
PREDICTIONPOWERS  
COMPLEXITIES  
FORPARAMVALUE INCONFCHANGINGPARAMVALUES  
CONFTUNEDPARAMSCONFCHANGINGPARAM PARAMVALUE  
ESTIMATOR CONFESTIMATOR CONFTUNEDPARAMS  
PRINTBENCHMARKING S ESTIMATOR  
ESTIMATORFITCONFDATAXTTRAIN CONFDATAYTRAIN  
CONFPOSTFITHOOKESTIMATOR  
COMPLEXITY CONFCOMPLEXITYCOMPUTERESTIMATOR  
COMPLEXITIESAPPENDCOMPLEXITY  
STARTTIME TIMETIME  
FORINRANGECONFNSAMPLES  
YPRED ESTIMATORPREDICTCONFDATAXTTEST  
ELAPSEDTIME TIMETIME STARTTIME FLOATCONFNSAMPLES  
PREDICTIONTIMESAPPENDELAPSEDTIME  
PREDScore CONFpredictionperformanceCOMPUTER  
CONFDATAYTEST YPRED  
PREDICTIONPOWERSAPPENDPREDScore  
PRINTCOMPLEXITY DS4F PRED TIME FSN  
COMPLEXITY CONFpredictionperformanceLABEL PREDScore  
ELAPSEDTIME  
RETURNPREDICTIONPOWERS PREDICTIONTIMES COMPLEXITIES  
DEFPLOTINFLUENCECONF MSEVALUES PREDICTIONTIMES COMPLEXITIES

PLOT INFLUENCE OF MODEL COMPLEXITY ON BOTH ACCURACY AND LATENCY

PLTFigureFIGSIZE12 6  
HOST HOSTSUBPLOT111 AXESCLASSAXES  
PLTSUBPLOTSADJUSTRIGHT075  
PAR1 HOSTTWINX  
HOSTSETXLABELMODEL COMPLEXITY S CONFCOMPLEXITYLABEL  
Y1LABEL CONFpredictionperformanceLABEL  
Y2LABEL TIME S  
HOSTSETYLABELY1LABEL  
PAR1SETYLABELY2LABEL  
P1 HOSTPLOTCOMPLEXITIES MSEVALUES B LABELprediction ERROR  
P2 PAR1PLOTCOMPLEXITIES PREDICTIONTIMES R  
LABELLATENCY  
52 EXAMPLES BASED ON REAL WORLD DATASETS 735

SCIKITLEARN USER GUIDE RELEASE 0213  
HOSTLEGENDLOCUPPER RIGHT  
HOSTAXISLEFTLABELSETCOLORP1GETCOLOR  
PAR1AXISRIGHTLABELSETCOLORP2GETCOLOR  
PLTTITLEINFLUENCE OF MODEL COMPLEXITY S CONFESTIMATORNAME  
PLTSHOW  
DEFCOUNTNONZEROCOEFFICIENTSESTIMATOR  
A ESTIMATORCOEFTOARRAY  
RETURNNPCOUNTNONZEROA

MAIN CODE  
REGRESSIONDATA GENERATEDATAREGRESSION  
CLASSIFICATIONDATA GENERATEDATACLASSIFICATION SPARSETRUE  
CONFIGURATIONS  
ESTIMATOR SGDCLASSIFIER  
TUNEDPARAMS PENALTY ELASTICNET ALPHA 0001 LOSS  
MODIFIEDHUBER FITINTERCEPT TRUE TOL 1E3  
CHANGINGPARAM L1RATIO  
CHANGINGPARAMVALUES 025 05 075 09  
COMPLEXITYLABEL NONZERO COEFFICIENTS  
COMPLEXITYCOMPUTER COUNTNONZEROCOEFFICIENTS  
PREDICTIONPERFORMANCECOMPUTER HAMMINGLOSS  
PREDICTIONPERFORMANCELABEL HAMMING LOSS MISCLASSIFICATION RATIO  
POSTFITHOOK LAMBDA X SPARSIFY  
DATA CLASSIFICATIONDATA  
NSAMPLES 30  
ESTIMATOR NUSVR  
TUNEDPARAMS C 1E3 GAMMA 2 15  
CHANGINGPARAM NU  
CHANGINGPARAMVALUES 01 025 05 075 09  
COMPLEXITYLABEL NSUPPORTVECTORS  
COMPLEXITYCOMPUTER LAMBDA X LENXSUPPORTVECTORS  
DATA REGRESSIONDATA  
POSTFITHOOK LAMBDA X  
PREDICTIONPERFORMANCECOMPUTER MEANSQUAREDERROR  
PREDICTIONPERFORMANCELABEL MSE  
NSAMPLES 30  
ESTIMATOR GRADIENTBOOSTINGREGRESSOR  
TUNEDPARAMS LOSS LS  
CHANGINGPARAM NESTIMATORS  
CHANGINGPARAMVALUES 10 50 100 200 500  
COMPLEXITYLABEL NTREES  
COMPLEXITYCOMPUTER LAMBDA X NESTIMATORS  
DATA REGRESSIONDATA  
POSTFITHOOK LAMBDA X  
PREDICTIONPERFORMANCECOMPUTER MEANSQUAREDERROR  
PREDICTIONPERFORMANCELABEL MSE  
NSAMPLES 30

FORCONFINCONFIGURATIONS  
PREDICTIONPERFORMANCES PREDICTIONTIMES COMPLEXITIES  
BENCHMARKINFLUENCECONF  
PLOTINFLUENCECONF PREDICTIONPERFORMANCES PREDICTIONTIMES  
COMPLEXITIES  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 35625 SECONDS  
736 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

NOTE [CLICK HERE](#) TO DOWNLOAD THE FULL EXAMPLE CODE

526 SPECIES DISTRIBUTION MODELING

MODELING SPECIES' GEOGRAPHIC DISTRIBUTIONS IS AN IMPORTANT PROBLEM IN CONSERVATION BIOLOGY IN THIS EXAMPLE WE MODEL THE GEOGRAPHIC DISTRIBUTION OF TWO SOUTH AMERICAN MAMMALS GIVEN PAST OBSERVATIONS AND 14 ENVIRONMENTAL VARIABLES SINCE WE HAVE ONLY POSITIVE EXAMPLES THERE ARE NO UNSUCCESSFUL OBSERVATIONS WE CAST THIS PROBLEM AS A DENSITY ESTIMATION PROBLEM AND USE THE ONECLASSSSVM PROVIDED BY THE PACKAGE SKLEARN SVM AS OUR MODELING TOOL THE DATASET IS PROVIDED BY PHILLIPS ET AL 2006 IF AVAILABLE THE EXAMPLE USES BASEMAP TO PLOT THE COAST LINES AND NATIONAL BOUNDARIES OF SOUTH AMERICA

THE TWO SPECIES ARE

- “BRADYPUS VARIEGATUS” THE BROWNTHOATED SLOTH
- “MICRORYZOMYS MINUTUS” ALSO KNOWN AS THE FOREST SMALL RICE RAT A RODENT THAT LIVES IN PERU COLOMBIA ECUADOR PERU AND VENEZUELA

REFERENCES

- “MAXIMUM ENTROPY MODELING OF SPECIES GEOGRAPHIC DISTRIBUTIONS” S J PHILLIPS R P ANDERSON R E SCHAPIRE ECOLOGICAL MODELLING 190231259 2006

52 EXAMPLES BASED ON REAL WORLD DATASETS 737

SCIKITLEARN USER GUIDE RELEASE 0213  
OUT

MODELING DISTRIBUTION OF SPECIES BRADYPUS VARIEGATUS  
FIT ONECLASSSSVM DONE  
PLOT COASTLINES FROM COVERAGE  
PREDICT SPECIES DISTRIBUTION  
AREA UNDER THE ROC CURVE 0868443

MODELING DISTRIBUTION OF SPECIES MICRORYZOMYS MINUTUS  
FIT ONECLASSSSVM DONE  
PLOT COASTLINES FROM COVERAGE  
PREDICT SPECIES DISTRIBUTION  
AREA UNDER THE ROC CURVE 0993919  
TIME ELAPSED 2018S  
AUTHORS PETER PRETTENHOFER PETERPRETTENHOFERGMAILCOM  
JAKE VANDERPLAS VANDERPLASASTROWASHINGTONEDU

LICENSE BSD 3 CLAUSE  
FROM TIME IMPORT TIME  
IMPORT NUMPY AS NP  
IMPORT MATPLOTLIBPYPLOT AS PLT  
FROM SKLEARNDATASETSBASE IMPORT BUNCH  
FROM SKLEARNDATASETS IMPORT FETCHSPECIESDISTRIBUTIONS  
FROM SKLEARNDATASETSSPECIESDISTRIBUTIONS IMPORT CONSTRUCTGRIDS  
FROM SKLEARN IMPORT SVM METRICS  
IF BASEMAP IS AVAILABLE WELL USE IT  
OTHERWISE WELL IMPROVISE LATER  
TRY  
FROM MPLTOOLKITSBASEMAP IMPORT BASEMAP  
BASEMAP TRUE  
EXCEPTIMPORTERROR  
BASEMAP FALSE  
PRINTDOC  
DEFCREATESPECIESBUNCHSPECIESNAME TRAIN TEST COVERAGES XGRID YGRID  
CREATE A BUNCH WITH INFORMATION ABOUT A PARTICULAR ORGANISM  
THIS WILL USE THE TESTTRAIN RECORD ARRAYS TO EXTRACT THE  
DATA SPECIFIC TO THE GIVEN SPECIES NAME

BUNCH BUNCHNAME JOINSPECIESNAMESPLIT2  
SPECIESNAME SPECIESNAMEENCODEASCII  
738 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213  
POINTS DICTTESTTEST TRAINTRAIN  
FORLABEL PTS INPOINTSITEMS  
CHOOSE POINTS ASSOCIATED WITH THE DESIRED SPECIES  
PTS PTSPTSSPECIES SPECIESNAME  
BUNCHPTS S LABEL PTS  
DETERMINE COVERAGE VALUES FOR EACH OF THE TRAINING TESTING POINTS  
IX NPSEARCHSORTEDXGRID PTSDD LONG  
IY NPSEARCHSORTEDYGRID PTSDD LAT  
BUNCHCOV S LABEL COVERAGES IY IXT  
RETURNBUNCH  
DEFPLOTSPECIESDISTRIBUTIONSPECIESBRADYPUSVARIEGATUSO  
MICRORYZOMYSMINUTUSO

PLOT THE SPECIES DISTRIBUTION

IFLENSPECIES 2  
PRINTNOTE WHEN MORE THAN TWO SPECIES ARE PROVIDED  
ONLY THE FIRST TWO WILL BE USED  
TO TIME  
LOAD THE COMPRESSED DATA  
DATA FETCHSPECIESDISTRIBUTIONS  
SET UP THE DATA GRID  
XGRID YGRID CONSTRUCTGRIDSDATA  
THE GRID IN XY COORDINATES  
X Y NPMESHGRIDXGRID YGRID1  
CREATE A BUNCH FOR EACH SPECIES  
BVBUNCH CREATESPECIESBUNCHSPECIES0  
DATATRAIN DATATEST  
DATACOVERAGES XGRID YGRID  
MMBUNCH CREATESPECIESBUNCHSPECIES1  
DATATRAIN DATATEST  
DATACOVERAGES XGRID YGRID  
BACKGROUND POINTS GRID COORDINATES FOR EVALUATION  
NPRANDOMSEED13  
BACKGROUNDPOINTS NPCNPRANDOMRANDINTLOW0 HIGHDATANY  
SIZE10000  
NPRANDOMRANDINTLOW0 HIGHDATANX  
SIZE10000T  
WELL MAKE USE OF THE FACT THAT COVERAGES6 HAS MEASUREMENTS AT ALL  
LAND POINTS THIS WILL HELP US DECIDE BETWEEN LAND AND WATER  
LANDREFERENCE DATACOVERAGES6  
FIT PREDICT AND PLOT FOR EACH SPECIES  
FORI SPECIES INENUMERATEBVBUNCH MMBUNCH  
PRINT80  
PRINTMODELING DISTRIBUTION OF SPECIES S SPECIESNAME  
52 EXAMPLES BASED ON REAL WORLD DATASETS 739

SCIKITLEARN USER GUIDE RELEASE 0213  
STANDARDIZE FEATURES  
MEAN SPECIESCOVTRAINMEANAXISO  
STD SPECIESCOVTRAINSTDAXISO  
TRAINCOVERSTD SPECIESCOVTRAIN MEAN STD  
FIT ONECLASSSSVM  
PRINT FIT ONECLASSSSVM END  
CLF SVMONECLASSSSVMNU01 KERNELRBF GAMMA05  
CLFFITTRAINCOVERSTD  
PRINTDONE  
PLOT MAP OF SOUTH AMERICA  
PLTSUBPLOT1 2 I 1  
IFBASEMAP  
PRINT PLOT COASTLINES USING BASEMAP  
M BASEMAPPROJECTIONCYL LLCRNRLATYMIN  
URCRNRLATYMAX LLCRNRLONXMIN  
URCRNRLONXMAX RESOLUTIONC  
MDRAWCOASTLINES  
MDRAWCOUNTRIES  
ELSE  
PRINT PLOT COASTLINES FROM COVERAGE  
PLTCONTOURX Y LANDREFERENCE  
LEVELS9998 COLORSK  
LINESTYLESSOLID  
PLXTICKS  
PLTYTICKS  
PRINT PREDICT SPECIES DISTRIBUTION  
PREDICT SPECIES DISTRIBUTION USING THE TRAINING DATA  
Z NPONESDATANY DATANX DTYENPFLOAT64  
WELL PREDICT ONLY FOR THE LAND POINTS  
IDX NPWHERELANDREFERENCE 9999  
COVERAGESLAND DATACOVERAGES IDX0 IDX1T  
PRED CLFDECISIONFUNCTIONCOVERAGESLAND MEAN STD  
Z PREDMIN  
ZIDX0 IDX1 PRED  
LEVELS NPLINSPACEZMIN ZMAX 25  
ZLANDREFERENCE 9999 9999  
PLOT CONTOURS OF THE PREDICTION  
PLTCONTOURFX Y Z LEVELSLEVELS CMAPPLTCMREDS  
PLTCOLORBARFORMAT 2F  
SCATTER TRAININGTESTING POINTS  
PLTSCATTERSPECIESPTSTRAINDD LONG SPECIESPTSTRAINDD LAT  
S22 CBLACK  
MARKER LABELTRAIN  
PLTSCATTERSPECIESPTSTESTDD LONG SPECIESPTSTESTDD LAT  
S22 CBLACK  
MARKERX LABELTEST  
PLTLEGEND  
PLTTITLESPECIESNAME  
740 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

PLTAXISEQUAL

COMPUTE AUC WITH REGARDS TO BACKGROUND POINTS

PREDBACKGROUND ZBACKGROUNDPOINTS0 BACKGROUNDPOINTS1

PREDTEST CLFDECISIONFUNCTIONSPECIESCOVTEST MEAN STD

SCORES NPRPREDTEST PREDBACKGROUND

Y NPRNPONESPREDTESTSHAPE NPZEROSPREDBACKGROUNDSHAPE

FPR TPR THRESHOLDS METRICSMROCCURVEY SCORES

ROCAUC METRICSAUCFPR TPR

PLTTEXT35 70 AUC 3F ROCAUC HARIGHT

PRINTNAREA UNDER THE ROC CURVE F ROCAUC

PRINTNTIME ELAPSED 2FS TIME TO

PLOTSPECIESDISTRIBUTION

PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 20185 SECONDS

NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE

527 VISUALIZING THE STOCK MARKET STRUCTURE

THIS EXAMPLE EMPLOYS SEVERAL UNSUPERVISED LEARNING TECHNIQUES TO EXTRACT THE STOCK MARKET STRUCTURE FROM VARIATIONS IN HISTORICAL QUOTES  
THE QUANTITY THAT WE USE IS THE DAILY VARIATION IN QUOTE PRICE QUOTES THAT ARE LINKED TEND TO COFLUCTUATE DURING A DAY  
LEARNING A GRAPH STRUCTURE

WE USE SPARSE INVERSE COVARIANCE ESTIMATION TO FIND WHICH QUOTES ARE CORRELATED CONDITIONALLY ON THE OTHERS SPECIFICALLY SPARSE INVERSE COVARIANCE GIVES US A GRAPH THAT IS A LIST OF CONNECTION FOR EACH SYMBOL THE SYMBOLS THAT IT IS CONNECTED TOO ARE THOSE USEFUL TO EXPLAIN ITS FLUCTUATIONS  
CLUSTERING

WE USE CLUSTERING TO GROUP TOGETHER QUOTES THAT BEHAVE SIMILARLY HERE AMONGST THE VARIOUS CLUSTERING TECHNIQUES AVAILABLE IN THE SCIKITLEARN WE USE AFFINITY PROPAGATION AS IT DOES NOT ENFORCE EQUALSIZE CLUSTERS AND IT CAN CHOOSE AUTOMATICALLY THE NUMBER OF CLUSTERS FROM THE DATA  
NOTE THAT THIS GIVES US A DIFFERENT INDICATION THAN THE GRAPH AS THE GRAPH REFLECTS CONDITIONAL RELATIONS BETWEEN VARIABLES WHILE THE CLUSTERING REFLECTS MARGINAL PROPERTIES VARIABLES CLUSTERED TOGETHER CAN BE CONSIDERED AS HAVING A SIMILAR IMPACT AT THE LEVEL OF THE FULL STOCK MARKET  
EMBEDDING IN 2D SPACE

FOR VISUALIZATION PURPOSES WE NEED TO LAY OUT THE DIFFERENT SYMBOLS ON A 2D CANVAS FOR THIS WE USE MANIFOLD LEARNING TECHNIQUES TO RETRIEVE 2D EMBEDDING  
52 EXAMPLES BASED ON REAL WORLD DATASETS 741

SCIKITLEARN USER GUIDE RELEASE 0213

VISUALIZATION

THE OUTPUT OF THE 3 MODELS ARE COMBINED IN A 2D GRAPH WHERE NODES REPRESENTS THE STOCKS AND EDGES THE

- CLUSTER LABELS ARE USED TO DEFINE THE COLOR OF THE NODES
- THE SPARSE COVARIANCE MODEL IS USED TO DISPLAY THE STRENGTH OF THE EDGES
- THE 2D EMBEDDING IS USED TO POSITION THE NODES IN THE PLAN

THIS EXAMPLE HAS A FAIR AMOUNT OF VISUALIZATIONRELATED CODE AS VISUALIZATION IS CRUCIAL HERE TO DISPLAY THE GRAPH ONE OF THE CHALLENGE IS TO POSITION THE LABELS MINIMIZING OVERLAP FOR THIS WE USE AN HEURISTIC BASED ON THE DIRECTION OF THE NEAREST NEIGHBOR ALONG EACH AXIS

OUT

CLUSTER 1 APPLE AMAZON YAHOO

CLUSTER 2 COMCAST CABLEVISION TIME WARNER

CLUSTER 3 CONOCOPHILLIPS CHEVRON TOTAL VALERO ENERGY EXXON

CLUSTER 4 CISCO DELL HP IBM MICROSOFT SAP TEXAS INSTRUMENTS

CLUSTER 5 BOEING GENERAL DYNAMICS NORTHROP GRUMMAN RAYTHEON

CLUSTER 6 AIG AMERICAN EXPRESS BANK OF AMERICA CATERPILLAR CVS DUPONT DE

↳NEMOURS FORD GENERAL ELECTRICS GOLDMAN SACHS HOME DEPOT JPMORGAN CHASE

↳MARRIOTT 3M RYDER WELLS FARGO WALMART

CLUSTER 7 MCDONALDS

742 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
CLUSTER 8 GLAXOSMITHKLINE NOVARTIS PFIZER SANOFIAVENTIS UNILEVER
CLUSTER 9 KELLOGG COCA COLA PEPSI
CLUSTER 10 COLGATEPALMOLIVE KIMBERLYCLARK PROCTER GAMBLE
CLUSTER 11 CANON HONDA NAVISTAR SONY TOYOTA XEROX
AUTHOR GAEL VAROQUAUX GAELVAROQUAUXNORMALESUPORG
LICENSE BSD 3 CLAUSE
IMPORT SYS
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM MATPLOTLIBCOLLECTIONS IMPORT LINECOLLECTION
IMPORT PANDAS AS PD
FROM SKLEARN IMPORT CLUSTER COVARIANCE MANIFOLD
PRINTDOC
```

```
RETRIEVE THE DATA FROM INTERNET
THE DATA IS FROM 2003 2008 THIS IS REASONABLY CALM NOT TOO LONG AGO SO
THAT WE GET HIGHTECH FIRMS AND BEFORE THE 2008 CRASH THIS KIND OF
HISTORICAL DATA CAN BE OBTAINED FOR FROM APIS LIKE THE QUANDLCOM AND
ALPHAVANTAGECO ONES
SYMBOLDICT
TOT TOTAL
XOM EXXON
CVX CHEVRON
COP CONOCOPHILLIPS
VLO VALERO ENERGY
MSFT MICROSOFT
IBM IBM
TWX TIME WARNER
CMCSA COMCAST
CVC CABLEVISION
YHOO YAHOO
DELL DELL
HPQ HP
AMZN AMAZON
TM TOYOTA
CAJ CANON
SNE SONY
F FORD
HMC HONDA
NAV NAVISTAR
NOC NORTHROP GRUMMAN
BA BOEING
KO COCA COLA
52 EXAMPLES BASED ON REAL WORLD DATASETS 743
```

SCIKITLEARN USER GUIDE RELEASE 0213

MMM 3M  
MCD MCDONALD S  
PEP PEPSI  
K KELLOGG  
UN UNILEVER  
MAR MARRIOTT  
PG PROCTER GAMBLE  
CL COLGATEPALMOLIVE  
GE GENERAL ELECTRICS  
WFC WELLS FARGO  
JPM JPMORGAN CHASE  
AIG AIG  
AXP AMERICAN EXPRESS  
BAC BANK OF AMERICA  
GS GOLDMAN SACHS  
AAPL APPLE  
SAP SAP  
CSCO CISCO  
TXN TEXAS INSTRUMENTS  
XRX XEROX  
WMT WALMART  
HD HOME DEPOT  
GSK GLAXOSMITHKLINE  
PFE PFIZER  
SNY SANOFIAVENTIS  
NVS NOVARTIS  
KMB KIMBERLYCLARK  
R RYDER  
GD GENERAL DYNAMICS  
RTN RAYTHEON  
CVS CVS  
CAT CATERPILLAR  
DD DUPONT DE NEMOURS

SYMBOLS NAMES NPARRAYSORTEDSYMBOLDICTITEMST  
QUOTES  
FORSYMBOLINSYMBOLS  
PRINTFETCHING QUOTE HISTORY FOR R SYMBOL FILESYSSTDERR  
URL HTTPSRAGHUBUSERCONTENTCOMSCIKITLEARNEXAMPLESDATA  
MASTERFINANCIALDATACSV  
QUOTESAPPENDPDREADCSVURLFORMATSYMBOL  
CLOSEPRICES NPVSTACKQCLOSE FORQINQUOTES  
OPENPRICES NPVSTACKQOPEN FORQINQUOTES  
THE DAILY VARIATIONS OF THE QUOTES ARE WHAT CARRY MOST INFORMATION  
VARIATION CLOSEPRICES OPENPRICES

LEARN A GRAPHICAL STRUCTURE FROM THE CORRELATIONS  
EDGEMODEL COVARIANCEGRAPHICALLASSOCVCV5  
STANDARDIZE THE TIME SERIES USING CORRELATIONS RATHER THAN COVARIANCE  
IS MORE EFFICIENT FOR STRUCTURE RECOVERY  
744 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
X VARIATIONCOPYT  
X XSTDAXISO  
EDGEMODELFITX

CLUSTER USING AFFINITY PROPAGATION  
LABELS CLUSTERAFFINITYPROPAGATIONEDGEMODELCOVARIANCE  
NLABELS LABELSMAX  
FORIINRANGENLABELS 1  
PRINTCLUSTER IS I 1 JOINNAMESLABELS I

FIND A LOWDIMENSION EMBEDDING FOR VISUALIZATION FIND THE BEST POSITION OF  
THE NODES THE STOCKS ON A 2D PLANE  
WE USE A DENSE EIGENSOLVER TO ACHIEVE REPRODUCIBILITY ARPACK IS  
INITIATED WITH RANDOM VECTORS THAT WE DONT CONTROL IN ADDITION WE  
USE A LARGE NUMBER OF NEIGHBORS TO CAPTURE THE LARGESCALE STRUCTURE  
NODEPOSITIONMODEL MANIFOLDLOCALLYLINEAREMBEDDING  
NCOMPONENTS2 EIGENSOLVERDENSE NNEIGHBORS6  
EMBEDDING NODEPOSITIONMODELFITTRANSFORMXTT

VISUALIZATION  
PLTFigure1 FACECOLORW FIGSIZE10 8  
PLTCLF  
AX PLTAXES0 0 1 1  
PLTAXISOFF  
DISPLAY A GRAPH OF THE PARTIAL CORRELATIONS  
PARTIALCORRELATIONS EDGEMODELPRECISIONCOPY  
D 1 NPSQRTNPDIAGPARTIALCORRELATIONS  
PARTIALCORRELATIONS D  
PARTIALCORRELATIONS D NPNEWAXIS  
NONZERO NPABSNPTRIUPARTIALCORRELATIONS K1 002  
PLOT THE NODES USING THE COORDINATES OF OUR EMBEDDING  
PLTSCATTEREMBEDDING0 EMBEDDING1 S100 D2 CLABELS  
CMAPPLTCMNIPYSPECTRAL  
PLOT THE EDGES  
STARTIDX ENDIDX NPWHERE NONZERO  
A SEQUENCE OF LINE0LINE1LINE2 WHERE  
LINEN X0 Y0 X1 Y1 XM YM  
SEGMENTS EMBEDDING START EMBEDDING STOP  
FORSTART STOP INZIPSTARTIDX ENDIDX  
VALUES NPABSPARTIALCORRELATIONSNONZERO  
LC LINECOLLECTIONSEGMENTS  
ZORDER0 CMAPPLTCMHOTR  
NORMPLTNORMALIZE0 7 VALUESMAX  
LCSETARRAYVALUES  
LCSETLINEWIDTHS15 VALUES  
AXADDCOLLECTIONLC  
ADD A LABEL TO EACH NODE THE CHALLENGE HERE IS THAT WE WANT TO  
52 EXAMPLES BASED ON REAL WORLD DATASETS 745

SCIKITLEARN USER GUIDE RELEASE 0213  
POSITION THE LABELS TO AVOID OVERLAP WITH OTHER LABELS  
FORINDEX NAME LABEL X Y INENUMERATE  
ZIPNAMES LABELS EMBEDDINGT  
DX X EMBEDDING0  
DXINDEX 1  
DY Y EMBEDDING1  
DYINDEX 1  
THISDX DXNPARGMINNPABSDY  
THISDY DYNPARGMINNPABSDX  
IFTHISDX 0  
HORIZONTALALIGNMENT LEFT  
X X 002  
ELSE  
HORIZONTALALIGNMENT RIGHT  
X X 002  
IFTHISDY 0  
VERTICALALIGNMENT BOTTOM  
Y Y 002  
ELSE  
VERTICALALIGNMENT TOP  
Y Y 002  
PLTTEXTX Y NAME SIZE10  
HORIZONTALALIGNMENTHORIZONTALALIGNMENT  
VERTICALALIGNMENTVERTICALALIGNMENT  
BBOXDICTFACECOLORW  
EDGECOLORPLTCMNIPYSPECTRALLABEL FLOATNLABELS  
ALPHA6  
PLTXLIMEMBEDDING0MIN 15 EMBEDDING0PTP  
EMBEDDING0MAX 10 EMBEDDING0PTP  
PLTYLIMEMBEDDING1MIN 03 EMBEDDING1PTP  
EMBEDDING1MAX 03 EMBEDDING1PTP  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 4857 SECONDS  
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE  
528 WIKIPEDIA PRINCIPAL EIGENVECTOR  
A CLASSICAL WAY TO ASSERT THE RELATIVE IMPORTANCE OF VERTICES IN A GRAPH IS TO COMPUTE THE PRINCIPAL EIGENVECTOR OF THE  
ADJACENCY MATRIX SO AS TO ASSIGN TO EACH VERTEX THE VALUES OF THE COMPONENTS OF THE FIRST EIGENVECTOR AS A CENTRALITY  
SCORE  
HTTPSENWIKIPEDIAORGWIKIEIGENVECTORCENTRALITY  
ON THE GRAPH OF WEBPAGES AND LINKS THOSE VALUES ARE CALLED THE PAGERANK SCORES BY GOOGLE  
THE GOAL OF THIS EXAMPLE IS TO ANALYZE THE GRAPH OF LINKS INSIDE WIKIPEDIA ARTICLES TO RANK ARTICLES BY RELATIVE IMPORTANCE  
ACCORDING TO THIS EIGENVECTOR CENTRALITY  
THE TRADITIONAL WAY TO COMPUTE THE PRINCIPAL EIGENVECTOR IS TO USE THE POWER ITERATION METHOD  
HTTPSENWIKIPEDIAORGWIKIPOWERITERATION  
746 CHAPTER 5 EXAMPLES



```
SCIKITLEARN USER GUIDE RELEASE 0213
HERE THE COMPUTATION IS ACHIEVED THANKS TO MARTINSSON'S RANDOMIZED SVD ALGORITHM IMPLEMENTED IN SCIKITLEARN
THE GRAPH DATA IS FETCHED FROM THE DBPEDIA DUMPS DBPEDIA IS AN EXTRACTION OF THE LATENT STRUCTURED DATA OF THE
WIKIPEDIA CONTENT
AUTHOR OLIVIER GRISEL OLIVIERGRISELENSTAORG
LICENSE BSD 3 CLAUSE
FROM BZ2 IMPORT BZ2FILE
IMPORT OS
FROM DATETIME IMPORT DATETIME
FROM PPRINT IMPORT PPRINT
FROM TIME IMPORT TIME
IMPORT NUMPY AS NP
FROM SCIPY IMPORT SPARSE
FROM JOBLIB IMPORT MEMORY
FROM SKLEARNDECOMPOSITION IMPORT RANDOMIZEDSVD
FROM URLLIBREQUEST IMPORT URLOPEN
PRINTDOC

WHERE TO DOWNLOAD THE DATA IF NOT ALREADY ON DISK
REDIRECTSURL HTTPDOWNLOADSDBPEDIAORG351ENREDIRECTSENNTBZ2
REDIRECTSFILENAME REDIRECTSURLRSPLIT 11
PAGELINKSURL HTTPDOWNLOADSDBPEDIAORG351ENPAGELINKSENNTBZ2
PAGELINKSFILENAME PAGELINKSURLRSPLIT 11
RESOURCES
REDIRECTSURL REDIRECTSFILENAME
PAGELINKSURL PAGELINKSFILENAME

FORURL FILENAME INRESOURCES
IF NOTOSPATHEXISTSFILENAME
PRINTDOWNLOADING DATA FROM S PLEASE WAIT URL
OPENER URLOPENURL
OPENFILENAME WBWRITEOPENERREAD
PRINT

LOADING THE REDIRECT FILES
MEMORY MEMORYCACHEDIR
DEFINDEXREDIRECTS INDEXMAP K
FIND THE INDEX OF AN ARTICLE NAME AFTER REDIRECT RESOLUTION
K REDIRECTSGETK K
RETURNINDEXMAPSETDEFAULTK LENINDEXMAP
52 EXAMPLES BASED ON REAL WORLD DATASETS 747
```

SCIKITLEARN USER GUIDE RELEASE 0213  
DBPEDIARESOURCEPREFIXLEN LENHTTPDBPEDIAORGRESOURCE  
SHORTNAMESLICE SLICEDBPEDIARESOURCEPREFIXLEN 1 1  
DEFSHORTNAMEENTURI  
REMOVE THE AND URI MARKERS AND THE COMMON URI PREFIX  
RETURNNTURISHORTNAMESLICE  
DEFGETREDIRECTSREDIRECTSFILENAME  
PARSE THE REDIRECTIONS AND BUILD A TRANSITIVELY CLOSED MAP OUT OF IT  
REDIRECTS  
PRINTPARSING THE NT REDIRECT FILE  
FORL LINE INENUMERATEBZ2FILEREDIRECTSFILENAME  
SPLIT LINESPLIT  
IFLENSPLIT 4  
PRINTIGNORING MALFORMED LINE LINE  
CONTINUE  
REDIRECTSSHORTNAMESPLIT0 SHORTNAMESPLIT2  
IFL 1000000 0  
PRINTS LINE 08D DATETIMENOWISOFORMAT L  
COMPUTE THE TRANSITIVE CLOSURE  
PRINTCOMPUTING THE TRANSITIVE CLOSURE OF THE REDIRECT RELATION  
FORL SOURCE INENUMERATEREDIRECTSKEYS  
TRANSITIVETARGET NONE  
TARGET REDIRECTSSOURCE  
SEEN SOURCE  
WHILETRUE  
TRANSITIVETARGET TARGET  
TARGET REDIRECTSGETTARGET  
IFTARGETISNONEORTARGETINSEEN  
BREAK  
SEENADDTARGET  
REDIRECTSSOURCE TRANSITIVETARGET  
IFL 1000000 0  
PRINTS LINE 08D DATETIMENOWISOFORMAT L  
RETURNREDIRECTS  
DISABLING JOBLIB AS THE PICKLING OF LARGE DICTS SEEMS MUCH TOO SLOW  
MEMORYCACHE  
DEFGETADJACENCYMATRIXREDIRECTSFILENAME PAGELINKSFILENAME LIMITNONE  
EXTRACT THE ADJACENCY GRAPH AS A SCIPY SPARSE MATRIX  
REDIRECTS ARE RESOLVED FIRST  
RETURNS X THE SCIPY SPARSE ADJACENCY MATRIX REDIRECTS AS PYTHON  
DICT FROM ARTICLE NAMES TO ARTICLE NAMES AND INDEXMAP A PYTHON DICT  
FROM ARTICLE NAMES TO PYTHON INT ARTICLE INDEXES  
  
PRINTCOMPUTING THE REDIRECT MAP  
REDIRECTS GETREDIRECTSREDIRECTSFILENAME  
PRINTCOMPUTING THE INTEGER INDEX MAP  
INDEXMAP DICT  
748 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

LINKS LIST

FORL LINE INENUMERATEBZ2FILEPAGELINKSFILENAME

SPLIT LINESPLIT

IFLENSPLIT 4

PRINTIGNORING MALFORMED LINE LINE

CONTINUE

I INDEXREDIRECTS INDEXMAP SHORTNAMESPLIT0

J INDEXREDIRECTS INDEXMAP SHORTNAMESPLIT2

LINKSAPPENDI J

IFL 1000000 0

PRINTS LINE 08D DATETIMENOWISOFORMAT L

IFLIMITIS NOTNONEANDL LIMIT 1

BREAK

PRINTCOMPUTING THE ADJACENCY MATRIX

X SPARSELILMATRIXLENINDEXMAP LENINDEXMAP DTYPENPFLOAT32

FORI JINLINKS

XI J 10

DELLINKS

PRINTCONVERTING TO CSR REPRESENTATION

X XTOCSR

PRINTCSR CONVERSION DONE

RETURNX REDIRECTS INDEXMAP

STOP AFTER 5M LINKS TO MAKE IT POSSIBLE TO WORK IN RAM

X REDIRECTS INDEXMAP GETADJACENCYMATRIX

REDIRECTSFILENAME PAGELINKSFILENAME LIMIT5000000

NAMES I NAME FORNAME I ININDEXMAPITEMS

PRINTCOMPUTING THE PRINCIPAL SINGULAR VECTORS USING RANDOMIZEDSVD

T0 TIME

U S V RANDOMIZEDSVDX 5 NITER3

PRINTDONE IN 03FS TIME T0

PRINT THE NAMES OF THE WIKIPEDIA RELATED STRONGEST COMPONENTS OF THE

PRINCIPAL SINGULAR VECTOR WHICH SHOULD BE SIMILAR TO THE HIGHEST EIGENVECTOR

PRINTTOP WIKIPEDIA PAGES ACCORDING TO PRINCIPAL SINGULAR VECTORS

PPRINTNAMESI FORIINNPABSUT0ARGSORT10

PPRINTNAMESI FORIINNPABSV0ARGSORT10

DEFCENTRALITYSCORESX ALPHA085 MAXITER100 TOL1E10

POWER ITERATION COMPUTATION OF THE PRINCIPAL EIGENVECTOR

THIS METHOD IS ALSO KNOWN AS GOOGLE PAGERANK AND THE IMPLEMENTATION

IS BASED ON THE ONE FROM THE NETWORKX PROJECT BSD LICENSED TOO

WITH COPYRIGHTS BY

ARIC HAGBERG HAGBERGLANLGOV

DAN SCHULT DSCHULTCOLGATEEDU

PIETER SWART SWARTLANLGOV

N XSHAPE0

X XCOPY

INCOMINGCOUNTS NPASARRAYXSUMAXIS1RAVEL

52 EXAMPLES BASED ON REAL WORLD DATASETS 749

```
SCIKITLEARN USER GUIDE RELEASE 0213
PRINTNORMALIZING THE GRAPH
FORIININCOMINGCOUNTSNONZERO0
XDATAIXINDPTRIXINDPTRI 1 10 INCOMINGCOUNTSI
DANGLE NPASARRAYNPWHEREINPISCLOSEXSUMAXIS1 0
10 N 0RAVEL
SCORES NPFULLN 1 N DTYPENPFLOAT32 INITIAL GUESS
FORIINRANGEMAXITER
PRINTPOWER ITERATION D I
PREVSCORES SCORES
SCORES ALPHA SCORES X NPDOTDANGLE PREVSCORES
1 ALPHA PREVSCORESSUM N
CHECK CONVERGENCE NORMALIZED LINF NORM
SCORESMAX NPABSSCORESMAX
IFSCORESMAX 00
SCORESMAX 10
ERR NPABSSCORES PREVSCORESMAX SCORESMAX
PRINTERROR 06F ERR
IFERR N TOL
RETURNSCORES
RETURNSCORES
PRINTCOMPUTING PRINCIPAL EIGENVECTOR SCORE USING A POWER ITERATION METHOD
T0 TIME
SCORES CENTRALITYSCORESX MAXITER100 TOL1E10
PRINTDONE IN 03FS TIME T0
PPRINTNAMESI FORIINNPAABSSCORESARGSORT10
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0000 SECONDS
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE
529 LIBSVM GUI
A SIMPLE GRAPHICAL FRONTEND FOR LIBSVM MAINLY INTENDED FOR DIDACTIC PURPOSES YOU CAN CREATE DATA POINTS BY POINT
AND CLICK AND VISUALIZE THE DECISION REGION INDUCED BY DIFFERENT KERNELS AND PARAMETER SETTINGS
TO CREATE POSITIVE EXAMPLES CLICK THE LEFT MOUSE BUTTON TO CREATE NEGATIVE EXAMPLES CLICK THE RIGHT BUTTON
IF ALL EXAMPLES ARE FROM THE SAME CLASS IT USES A ONECLASS SVM
PRINTDOC
AUTHOR PETER PRETTENHOER PETERPRETTENHOFERGMAILCOM

LICENSE BSD 3 CLAUSE
IMPORT MATPLOTLIB
MATPLOTLIBUSETKAGG
FROM MATPLOTLIBBACKENDSBACKENDTKAGG IMPORT FIGURECANVASTKAGG
FROM MATPLOTLIBBACKENDSBACKENDTKAGG IMPORT NAVIGATIONTOOLBAR2TKAGG
FROM MATPLOTLIBFIGURE IMPORT FIGURE
FROM MATPLOTLIBCONTOUR IMPORT CONTOURSET
750 CHAPTER 5 EXAMPLES
```

```
SCIKITLEARN USER GUIDE RELEASE 0213
IMPORT SYS
IMPORT NUMPY AS NP
IMPORT TKINTER AS TK
FROM SKLEARN IMPORT SVM
FROM SKLEARNDATASETS IMPORT DUMPSVMLIGHTFILE
YMIN YMAX 50 50
XMIN XMAX 50 50
CLASS MODEL
THE MODEL WHICH HOLD THE DATA IT IMPLEMENTS THE
OBSERVABLE IN THE OBSERVER PATTERN AND NOTIFIES THE
REGISTERED OBSERVERS ON CHANGE EVENT

DEFINITSELF
SELFOBSERVERS
SELSURFACE NONE
SELFDATA
SELFCLS NONE
SELSURFACETYPE 0
DEFCHANGEDSELF EVENT
NOTIFY THE OBSERVERS
FOROBSERVER INSELFOBSERVERS
OBSERVERUPDATEEVENT SELF
DEFADDOBSERVERSELF OBSERVER
REGISTER AN OBSERVER
SELFOBSERVERSAPPENDOBSERVER
DEFSETSURFACESELF SURFACE
SELSURFACE SURFACE
DEFDUMPSVMLIGHTFILESELF FILE
DATA NPARRAYSELFDATA
X DATA 02
Y DATA 2
DUMPSVMLIGHTFILEX Y FILE
CLASS CONTROLLER
DEFINITSELF MODEL
SELFMODEL MODEL
SELFKERNEL TKINTVAR
SELSURFACETYPE TKINTVAR
WHETHER OR NOT A MODEL HAS BEEN FITTED
SELFFITTED FALSE
DEFFITSELF
PRINTFIT THE MODEL
TRAIN NPARRAYSELFMODELDATA
X TRAIN 02
Y TRAIN 2
52 EXAMPLES BASED ON REAL WORLD DATASETS 751
```

SCIKITLEARN USER GUIDE RELEASE 0213  
C FLOATSELFCOMPLEXITYGET  
GAMMA FLOATSELFGAMMAGET  
COEF0 FLOATSELFCOEF0GET  
DEGREE INTSELFDEGREEGET  
KERNELMAP 0 LINEAR 1 RBF 2 POLY  
IFLENNPUNIQUEY 1  
CLF SVMONECLASSSVMKERNELKERNELMAPSELFKERNELGET  
GAMMAGAMMA COEF0COEF0 DEGREEDEGREE  
CLFFITX  
ELSE  
CLF SVMSVCKERNELKERNELMAPSELFKERNELGET CC  
GAMMAGAMMA COEF0COEF0 DEGREEDEGREE  
CLFFITX Y  
IFHASATTRCLF SCORE  
PRINTACCURACY CLFSCOREX Y 100  
X1 X2 Z SELFDECISIONSURFACECLF  
SELFMODELCLF CLF  
SELFMODELSETSURFACEX1 X2 Z  
SELFMODELSURFACETYPE SELFSURFACETYPEGET  
SELFFITTED TRUE  
SELFMODELCHANGEDSURFACE  
DEFDECISIONSURFACESELF CLS  
DELTA 1  
X NPARANGEXMIN XMAX DELTA DELTA  
Y NPARANGEYMIN YMAX DELTA DELTA  
X1 X2 NPMESHGRIDX Y  
Z CLSDECISIONFUNCTIONNPCX1RAVEL X2RAVEL  
Z ZRESHAPEX1SHAPE  
RETURNX1 X2 Z  
DEFCLEARDATASELF  
SELFMODELDATA  
SELFFITTED FALSE  
SELFMODELCHANGEDCLEAR  
DEFADDEXAMPLESELF X Y LABEL  
SELFMODELDATAAPPENDX Y LABEL  
SELFMODELCHANGEDEXAMPLEADDED  
UPDATE DECISION SURFACE IF ALREADY FITTED  
SELFREFIT  
DEFREFITSELF  
REFIT THE MODEL IF ALREADY FITTED  
IFSELFFITTED  
SELFFIT  
CLASS VIEW  
TEST DOCSTRING  
DEFINITSELF ROOT CONTROLLER  
F FIGURE  
AX FADDSUBPLOT111  
AXSETXTICKS  
AXSETYTIMES  
AXSETXLIMXMIN XMAX  
AXSETYLIMYMIN YMAX  
752 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
CANVAS FIGURECANVASTKAGGF MASTERROOT  
CANVASSHOW  
CANVASGETTKWIDGETPACKSIDETKTOP FILLTKBOTH EXPAND1  
CANVASTKCANVASPACKSIDETKTOP FILLTKBOTH EXPAND1  
CANVASMPLCONNECTBUTTONNPRESSEVENT SELFONCLICK  
TOOLBAR NAVIGATIONTOOLBAR2TKAGGCANVAS ROOT  
TOOLBARUPDATE  
SELFCONTROLLBAR CONTROLLBARROOT CONTROLLER  
SELFF F  
SELFAX AX  
SELFCANVAS CANVAS  
SELFCONTROLLER CONTROLLER  
SELFCONTOURS  
SELFCLABELS NONE  
SELFPLOTKERNELS  
DEFPLOTKERNELSSELF  
SELFAXTEXT50 60 LINEAR UT V  
SELFAXTEXT20 60 RRB EXP GAMMA UV 2  
SELFAXTEXT10 60 RPOLY GAMMA UT V RD  
DEFONCLICKSELF EVENT  
IFEVENTXDATA ANDEVENTYDATA  
IFEVENTBUTTON 1  
SELFCONTROLLERADDEXAMPLEEVENTXDATA EVENTYDATA 1  
ELIFEVENTBUTTON 3  
SELFCONTROLLERADDEXAMPLEEVENTXDATA EVENTYDATA 1  
DEFUPDATEEXAMPLESELF MODEL IDX  
X Y L MODELDATAIDX  
IFL 1  
COLOR W  
ELIFL 1  
COLOR K  
SELFAXPLOTX Y SO COLOR SCALEX00 SCALEY00  
DEFUPDATESELF EVENT MODEL  
IFEVENT EXAMPLESLOADED  
FORIINRANGELENMODELDATA  
SELFUPDATEEXAMPLEMODEL I  
IFEVENT EXAMPLEADDED  
SELFUPDATEEXAMPLEMODEL 1  
IFEVENT CLEAR  
SELFAXCLEAR  
SELFAXSETXTICKS  
SELFAXSETYTICKS  
SELFCONTOURS  
SELFCLABELS NONE  
SELFPLOTKERNELS  
IFEVENT SURFACE  
SELFREMOVESURFACE  
SELFPLOTSUPPORTVECTORSMODELCLFSUPPORTVECTORS  
SELFPLOTDECISIONSURFACEMODELSURFACE MODELSURFACETYPE  
SELFCANVASDRAW  
52 EXAMPLES BASED ON REAL WORLD DATASETS 753

SCIKITLEARN USER GUIDE RELEASE 0213  
DEFREMOVESURFACESELF  
REMOVE OLD DECISION SURFACE  
IFLENSSELFCONTOURS 0  
FORCONTOUR INSELFCONTOURS  
IFISINSTANCECONTOUR CONTOURSET  
FORLINESET INCONTOURCOLLECTIONS  
LINESETREMOVE  
ELSE  
CONTOURREMOVE  
SELFCONTOURS  
DEFPLOTSUPPORTVECTORSSSELF SUPPORTVECTORS  
PLOT THE SUPPORT VECTORS BY PLACING CIRCLES OVER THE  
CORRESPONDING DATA POINTS AND ADDS THE CIRCLE COLLECTION  
TO THE CONTOURS LIST  
CS SELFAXSCATTERSUPPORTVECTORS 0 SUPPORTVECTORS 1  
S80 EDGECOLORSK FACECOLORSNONE  
SELFCONTOURSAPPENDCS  
DEFPLOTDECISIONSURFACESELF SURFACE TYPE  
X1 X2 Z SURFACE  
IFTYPE 0  
LEVELS 10 00 10  
LINESTYLES DASHED SOLID DASHED  
COLORS K  
SELFCONTOURSAPPENDSELFAXCONTOURX1 X2 Z LEVELS  
COLORSCOLORS  
LINESTYLESLINESTYLES  
ELIFTYPE 1  
SELFCONTOURSAPPENDSELFAXCONTOURFX1 X2 Z 10  
CMAPMATPLOTLIBCBONE  
ORIGINLOWER ALPHA085  
SELFCONTOURSAPPENDSELFAXCONTOURX1 X2 Z 00 COLORSK  
LINESTYLESSOLID  
ELSE  
RAISEVALUEERRORSURFACE TYPE UNKNOWN  
CLASS CONTROLLBAR  
DEFINITSELF ROOT CONTROLLER  
FM TKFRAMEROOT  
KERNELGROUP TKFRAMEFM  
TKRADIOBUTTONKERNELGROUP TEXTLINEAR VARIABLECONTROLLERKERNEL  
VALUE0 COMMANDCONTROLLERREFITPACKANCHORTKW  
TKRADIOBUTTONKERNELGROUP TEXTRBF VARIABLECONTROLLERKERNEL  
VALUE1 COMMANDCONTROLLERREFITPACKANCHORTKW  
TKRADIOBUTTONKERNELGROUP TEXTPOLY VARIABLECONTROLLERKERNEL  
VALUE2 COMMANDCONTROLLERREFITPACKANCHORTKW  
KERNELGROUPPACKSIDETKLEFT  
VALBOX TKFRAMEFM  
CONTROLLERCOMPLEXITY TKSTRINGVAR  
CONTROLLERCOMPLEXITYSET10  
C TKFRAMEVALBOX  
TKLABELC TEXTC ANCHORE WIDTH7PACKSIDETKLEFT  
TKENTRYC WIDTH6 TEXTVARIABLECONTROLLERCOMPLEXITYPACK  
SIDETKLEFT  
754 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213  
CPACK  
CONTROLLERGAMMA TKSTRINGVAR  
CONTROLLERGAMMASET001  
G TKFRAMEVALBOX  
TKLABELG TEXTGAMMA ANCHORE WIDTH7PACKSIDETKLEFT  
TKENTRYG WIDTH6 TEXTVARIABLECONTROLLERGAMMAPACKSIDETKLEFT  
GPACK  
CONTROLLERDEGREE TKSTRINGVAR  
CONTROLLERDEGREESET3  
D TKFRAMEVALBOX  
TKLABELD TEXTDEGREE ANCHORE WIDTH7PACKSIDETKLEFT  
TKENTRYD WIDTH6 TEXTVARIABLECONTROLLERDEGREEPACKSIDETKLEFT  
DPACK  
CONTROLLERCOEF0 TKSTRINGVAR  
CONTROLLERCOEF0SET0  
R TKFRAMEVALBOX  
TKLABELR TEXTCOEF0 ANCHORE WIDTH7PACKSIDETKLEFT  
TKENTRYR WIDTH6 TEXTVARIABLECONTROLLERCOEF0PACKSIDETKLEFT  
RPACK  
VALBOXPACKSIDETKLEFT  
CMAPGROUP TKFRAMEFM  
TKRADIOBUTTONCMAPGROUP TEXTHYPERPLANES  
VARIABLECONTROLLERSURFACETYPE VALUE0  
COMMANDCONTROLLERREFITPACKANCHORTKW  
TKRADIOBUTTONCMAPGROUP TEXTSURFACE  
VARIABLECONTROLLERSURFACETYPE VALUE1  
COMMANDCONTROLLERREFITPACKANCHORTKW  
CMAPGROUPPACKSIDETKLEFT  
TRAINBUTTON TKBUTTONFM TEXTFIT WIDTH5  
COMMANDCONTROLLERFIT  
TRAINBUTTONPACK  
FMPACKSIDETKLEFT  
TKBUTTONFM TEXTCLEAR WIDTH5  
COMMANDCONTROLLERCLEARDATAPACKSIDETKLEFT  
DEFGETPARSER  
FROM OPTPARSE IMPORT OPTIONPARSER  
OP OPTIONPARSER  
OPADDOPTIONOUTPUT  
ACTIONSTORE TYPESTR DESTOUTPUT  
HELPPATH WHERE TO DUMP DATA  
RETURNOP  
DEFMAINARGV  
OP GETPARSER  
OPTS ARGS OPPARSEARGSARGV1  
ROOT TKTK  
MODEL MODEL  
CONTROLLER CONTROLLERMODEL  
ROOTWMTITLESCIKITLEARN LIBSVM GUI  
52 EXAMPLES BASED ON REAL WORLD DATASETS 755

SCIKITLEARN USER GUIDE RELEASE 0213

VIEW VIEWROOT CONTROLLER

MODELADDOBSERVERVIEW

TKMAINLOOP

IFOPTSOUTPUT

MODELDUMPSVMLIGHTFILEOPTSOUTPUT

IFNAME MAIN

MAINSYSARGV

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0000 SECONDS

NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)

5210 PREDICTION LATENCY

THIS IS AN EXAMPLE SHOWING THE PREDICTION LATENCY OF VARIOUS SCIKITLEARN ESTIMATORS

THE GOAL IS TO MEASURE THE LATENCY ONE CAN EXPECT WHEN DOING PREDICTIONS EITHER IN BULK OR ATOMIC IE ONE BY ONE  
MODE

THE PLOTS REPRESENT THE DISTRIBUTION OF THE PREDICTION LATENCY AS A BOXPLOT

•

756 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

- 
- 

52 EXAMPLES BASED ON REAL WORLD DATASETS 757

SCIKITLEARN USER GUIDE RELEASE 0213

```
•
OUT
BENCHMARKING SGDREGRESSORALPHA001 L1RATIO025 PENALTYELASTICNET TOL00001
BENCHMARKING RANDOMFORESTREGRESSORNESTIMATORS100
BENCHMARKING SVR
BENCHMARKING WITH 100 FEATURES
BENCHMARKING WITH 250 FEATURES
BENCHMARKING WITH 500 FEATURES
EXAMPLE RUN IN 960S
AUTHORS EUSTACHE DIEMERT EUSTACHEDIEMERTFR
LICENSE BSD 3 CLAUSE
FROM COLLECTIONS IMPORT DEFAULTDICT
IMPORT TIME
IMPORT GC
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARNPREPROCESSING IMPORT STANDARDSCALER
FROM SKLEARNMODELSELECTION IMPORT TRAINTESTSPLIT
FROM SKLEARNDATASETSSAMPLESGENERATOR IMPORT MAKEREGRESSION
FROM SKLEARNENSEMBLEFOREST IMPORT RANDOMFORESTREGRESSOR
FROM SKLEARNLINEARMODELRIDGE IMPORT RIDGE
FROM SKLEARNLINEARMODELSTOCHASTICGRADIENT IMPORT SGDREGRESSOR
FROM SKLEARN SVMCLASSES IMPORT SVR
FROM SKLEARNUTILS IMPORT SHUFFLE
758 CHAPTER 5 EXAMPLES
```

SCIKITLEARN USER GUIDE RELEASE 0213

DEFNOTINSPHINX

HACK TO DETECT WHETHER WE ARE RUNNING BY THE SPHINX BUILDER

RETURNFILE INGLOBALS

DEFATOMICBENCHMARKESTIMATORESTIMATOR XTEST VERBOSEFALSE

MEASURE RUNTIME PREDICTION OF EACH INSTANCE

NINSTANCES XTESTSHAPE0

RUNTIMES NPZEROSNINSTANCES DYPENPFLOAT

FORIINRANGENINSTANCES

INSTANCE XTESTI

START TIMETIME

ESTIMATORPREDICTINSTANCE

RUNTIMESI TIMETIME START

IFVERBOSE

PRINTATOMICBENCHMARK RUNTIMES MINRUNTIMES NPPERCENTILE

RUNTIMES 50 MAXRUNTIMES

RETURNRUNTIMES

DEFBULKBENCHMARKESTIMATORESTIMATOR XTEST NBULKREPEATS VERBOSE

MEASURE RUNTIME PREDICTION OF THE WHOLE INPUT

NINSTANCES XTESTSHAPE0

RUNTIMES NPZEROSNBULKREPEATS DYPENPFLOAT

FORIINRANGENBULKREPEATS

START TIMETIME

ESTIMATORPREDICTXTEST

RUNTIMESI TIMETIME START

RUNTIMES NPARRAYLISTMAP LAMBDA X FLOATNINSTANCES RUNTIMES

IFVERBOSE

PRINTBULKBENCHMARK RUNTIMES MINRUNTIMES NPPERCENTILE

RUNTIMES 50 MAXRUNTIMES

RETURNRUNTIMES

DEFBENCHMARKESTIMATORESTIMATOR XTEST NBULKREPEATS30 VERBOSEFALSE

MEASURE RUNTIMES OF PREDICTION IN BOTH ATOMIC AND BULK MODE

PARAMETERS

ESTIMATOR ALREADY TRAINED ESTIMATOR SUPPORTING PREDICT

XTEST TEST INPUT

NBULKREPEATS HOW MANY TIMES TO REPEAT WHEN EVALUATING BULK MODE

RETURNS

ATOMICRUNTIMES BULKRUNTIMES A PAIR OF NPARRAY WHICH CONTAIN THE

RUNTIMES IN SECONDS

ATOMICRUNTIMES ATOMICBENCHMARKESTIMATORESTIMATOR XTEST VERBOSE

BULKRUNTIMES BULKBENCHMARKESTIMATORESTIMATOR XTEST NBULKREPEATS

VERBOSE

RETURNATOMICRUNTIMES BULKRUNTIMES

52 EXAMPLES BASED ON REAL WORLD DATASETS 759

SCIKITLEARN USER GUIDE RELEASE 0213  
 DEFGENERATEDDATASETNTTRAIN NTEST NFEATURES NOISE01 VERBOSEFALSE  
 GENERATE A REGRESSION DATASET WITH THE GIVEN PARAMETERS  
 IFVERBOSE  
 PRINTGENERATING DATASET  
 X Y COEF MAKEREGRESSIONNSAMPLESNTTRAIN NTEST  
 NFEATURESNFEATURES NOISENOISE COEFTRUE  
 RANDOMSEED 13  
 XTRAIN XTEST YTRAIN YTEST TRAINTESTSPLIT  
 X Y TRAINSIZENTRAIN TESTSIZENTEST RANDOMSTATERANDOMSEED  
 XTRAIN YTRAIN SHUFFLEXTRAIN YTRAIN RANDOMSTATERANDOMSEED  
 XSCALER STANDARDSCALER  
 XTRAIN XSCALERFITTRANSFORMXTRAIN  
 XTEST XSCALERTRANSFORMXTEST  
 YSCALER STANDARDSCALER  
 YTRAIN YSCALERFITTRANSFORMYTRAIN NONE 0  
 YTEST YSCALERTRANSFORMYTEST NONE 0  
 GCCOLLECT  
 IFVERBOSE  
 PRINTOK  
 RETURNXTRAIN YTRAIN XTEST YTEST  
 DEFBOXPLOTRUNTIMESRUNTIMES PREDTYPE CONFIGURATION  
  
 PLOT A NEW FIGURE WITH BOXPLOTS OF PREDICTION RUNTIMES  
 PARAMETERS  
  
 RUNTIMES LIST OF NPARRAY OF LATENCIES IN MICROSECONDS  
 CLSNAMES LIST OF ESTIMATOR CLASS NAMES THAT GENERATED THE RUNTIMES  
 PREDTYPE BULK OR ATOMIC  
  
 FIG AX1 PLTSUBPLOTSFIGSIZE10 6  
 BP PLTBOXPLOTRUNTIMES  
 CLSINFOS SND S ESTIMATORCONFNAME  
 ESTIMATORCONFCOMPLEXITYCOMPUTER  
 ESTIMATORCONFINSTANCE  
 ESTIMATORCONFCOMPLEXITYLABEL FOR  
 ESTIMATORCONF INCONFIGURATIONESTIMATORS  
 PLTSETPAX1 XTICKLABELSCLSINFOS  
 PLTSETPBPBOXES COLORBLACK  
 PLTSETPBPWHISKERS COLORBLACK  
 PLTSETPBPFLIERS COLORRED MARKER  
 AX1YAXISGRIDTRUE LINESTYLE WHICHMAJOR COLORLIGHTGREY  
 ALPHA05  
 AX1SETAXISBELOWTRUE  
 AX1SETTITLEPREDICTION TIME PER INSTANCE SDFEATS  
 760 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
PREDTYPECAPITALIZE  
CONFIGURATIONNFEATURES  
AX1SETYLABELPREDICTION TIME US  
PLTSHOW  
DEFBENCHMARKCONFIGURATION  
RUN THE WHOLE BENCHMARK  
XTRAIN YTRAIN XTEST YTEST GENERATEDDATASET  
CONFIGURATIONNTRAIN CONFIGURATIONNTEST  
CONFIGURATIONNFEATURES  
STATS  
FORESTIMATORCONF INCONFIGURATIONESTIMATORS  
PRINTBENCHMARKING ESTIMATORCONFINSTANCE  
ESTIMATORCONFINSTANCEFITXTRAIN YTRAIN  
GCCOLLECT  
A B BENCHMARKESTIMATORESTIMATORCONFINSTANCE XTEST  
STATSESTIMATORCONFNAME ATOMIC A BULK B  
CLSNAMES ESTIMATORCONFNAME FORESTIMATORCONF INCONFIGURATION  
ESTIMATORS  
RUNTIMES 1E6 STATSCLFNAMEATOMIC FORCLFNAME INCLSNAMES  
BOXPLOTRUNTIMESRUNTIMES ATOMIC CONFIGURATION  
RUNTIMES 1E6 STATSCLFNAMEBULK FORCLFNAME INCLSNAMES  
BOXPLOTRUNTIMESRUNTIMES BULK D CONFIGURATIONNTEST  
CONFIGURATION  
DEFNFEATUREINFLUENCEESTIMATORS NTRAIN NTEST NFEATURES PERCENTILE

ESTIMATE INFLUENCE OF THE NUMBER OF FEATURES ON PREDICTION TIME  
PARAMETERS

ESTIMATORS DICT OF NAME STR ESTIMATOR TO BENCHMARK  
NTRAIN NBER OF TRAINING INSTANCES INT  
NTEST NBER OF TESTING INSTANCES INT  
NFEATURES LIST OF FEATURES SPACE DIMENSIONALITY TO TEST INT  
PERCENTILE PERCENTILE AT WHICH TO MEASURE THE SPEED INT 0100  
RETURNS

PERCENTILES DICTESTIMATORNAME  
DICTNFEATURES PERCENTILEPERFINUS

PERCENTILES DEFAULTDICTDEFAULTDICT  
FORNINNFEATURES  
PRINTBENCHMARKING WITH DFEATURES N  
XTRAIN YTRAIN XTEST YTEST GENERATEDDATASETNTRAIN NTEST N  
FORCLSNAME ESTIMATOR INESTIMATORSITEMS  
ESTIMATORFITXTRAIN YTRAIN  
GCCOLLECT  
RUNTIMES BULKBENCHMARKESTIMATORESTIMATOR XTEST 30 FALSE  
52 EXAMPLES BASED ON REAL WORLD DATASETS 761

SCIKITLEARN USER GUIDE RELEASE 0213  
PERCENTILESCLSNAME 1E6 NPPERCENTILERUNTIMES  
PERCENTILE  
RETURNPERCENTILES  
DEFPLOTNFEATURESINFLUENCEPERCENTILES PERCENTILE  
FIG AX1 PLTSUBPLOTSFIGSIZE10 6  
COLORS R G B  
FORI CLSNAME INENUMERATEPERCENTILESKEYS  
X NPARRAYSORTEDN FORNINPERCENTILESCLSNAMEKEYS  
Y NPARRAYPERCENTILESCLSNAME FORNINX  
PLTPLOTX Y COLORCOLORSI  
AX1YAXISGRIDTRUE LINESTYLE WHICHMAJOR COLORLIGHTGREY  
ALPHA05  
AX1SETAXISBELOWTRUE  
AX1SETTITLEEVOLUTION OF PREDICTION TIME WITH FEATURES  
AX1SETXLABELFEATURES  
AX1SETYLABELPREDICTION TIME AT DILE US PERCENTILE  
PLTSHOW  
DEFBENCHMARKTHROUGHPUTSCONFIGURATION DURATIONSECS01  
BENCHMARK THROUGHPUT FOR DIFFERENT ESTIMATORS  
XTRAIN YTRAIN XTEST YTEST GENERATEDDATASET  
CONFIGURATIONNTRAIN CONFIGURATIONNTEST  
CONFIGURATIONNFEATURES  
THROUGHPUTS DICT  
FORESTIMATORCONFIG INCONFIGURATIONESTIMATORS  
ESTIMATORCONFIGINSTANCEFITXTRAIN YTRAIN  
STARTTIME TIMETIME  
NPREDICTIONS 0  
WHILETIMETIME STARTTIME DURATIONSECS  
ESTIMATORCONFIGINSTANCEPREDICTXTEST0  
NPREDICTIONS 1  
THROUGHPUTSESTIMATORCONFIGNAME NPREDICTIONS DURATIONSECS  
RETURNTHROUGHPUTS  
DEFPLOTBENCHMARKTHROUGHPUTTHROUGHPUTS CONFIGURATION  
FIG AX PLTSUBPLOTSFIGSIZE10 6  
COLORS R G B  
CLSINFOS SND S ESTIMATORCONFNAME  
ESTIMATORCONFCOMPLEXITYCOMPUTER  
ESTIMATORCONFINSTANCE  
ESTIMATORCONFCOMPLEXITYLABEL FOR  
ESTIMATORCONF INCONFIGURATIONESTIMATORS  
CLSVALUES THROUGHPUTSESTIMATORCONFNAME FORESTIMATORCONF IN  
CONFIGURATIONESTIMATORS  
PLTBARRANGELENTHROUGHPUTS CLSVALUES WIDTH05 COLORCOLORS  
AXSETXTICKSNPLINSPACE025 LENTHROUGHPUTS 075 LENTHROUGHPUTS  
AXSETXTICKLABELSCLSINFOS FONTSIZE10  
YMAX MAXCLSVALUES 12  
AXSETYLIM0 YMAX  
AXSETYLABELTHROUGHPUT PREDICTIONSSEC  
AXSETTITLEPREDICTION THROUGHPUT FOR DIFFERENT ESTIMATORS D  
FEATURES CONFIGURATIONNFEATURES  
PLTSHOW  
762 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213

MAIN CODE  
STARTTIME TIMETIME

BENCHMARK BULKATOMIC PREDICTION SPEED FOR VARIOUS REGRESSORS  
CONFIGURATION  
NTRAIN INT1E3  
NTEST INT1E2  
NFEATURES INT1E2  
ESTIMATORS  
NAME LINEAR MODEL  
INSTANCE SGDREGRESSORPENALTYELASTICNET ALPHA001  
L1RATIO025 FITINTERCEPTTRUE  
TOL1E4  
COMPLEXITYLABEL NONZERO COEFFICIENTS  
COMPLEXITYCOMPUTER LAMBDACLF NPCOUNTNONZEROCLFCOEF  
NAME RANDOMFOREST  
INSTANCE RANDOMFORESTREGRESSORNESTIMATORS100  
COMPLEXITYLABEL ESTIMATORS  
COMPLEXITYCOMPUTER LAMBDACLF CLFNESTIMATORS  
NAME SVR  
INSTANCE SVRKERNELRBF  
COMPLEXITYLABEL SUPPORT VECTORS  
COMPLEXITYCOMPUTER LAMBDACLF LENCLFSUPPORTVECTORS

BENCHMARKCONFIGURATION  
BENCHMARK NFEATURES INFLUENCE ON PREDICTION SPEED  
PERCENTILE 90  
PERCENTILES NFEATUREINFLUENCERIDGE RIDGE  
CONFIGURATIONNTRAIN  
CONFIGURATIONNTEST  
100 250 500 PERCENTILE  
PLOTNFEATURESINFLUENCEPERCENTILES PERCENTILE  
BENCHMARK THROUGHPUT  
THROUGHPUTS BENCHMARKTHROUGHPUTSCONFIGURATION  
PLOTBENCHMARKTHROUGHPUTTHROUGHPUTS CONFIGURATION  
STOPTIME TIMETIME  
PRINTEXAMPLE RUN IN 2FS STOPTIME STARTTIME  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 9597 SECONDS  
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE  
5211 OUTOF CORE CLASSIFICATION OF TEXT DOCUMENTS  
THIS IS AN EXAMPLE SHOWING HOW SCIKITLEARN CAN BE USED FOR CLASSIFICATION USING AN OUTOF CORE APPROACH LEARNING  
FROM DATA THAT DOESN'T FIT INTO MAIN MEMORY WE MAKE USE OF AN ONLINE CLASSIFIER IE ONE THAT SUPPORTS THE PARTIALFIT  
METHOD THAT WILL BE FED WITH BATCHES OF EXAMPLES TO GUARANTEE THAT THE FEATURES SPACE REMAINS THE SAME OVER TIME  
52 EXAMPLES BASED ON REAL WORLD DATASETS 763

SCIKITLEARN USER GUIDE RELEASE 0213

WE LEVERAGE A HASHINGVECTORIZER THAT WILL PROJECT EACH EXAMPLE INTO THE SAME FEATURE SPACE THIS IS ESPECIALLY USEFUL IN THE CASE OF TEXT CLASSIFICATION WHERE NEW FEATURES WORDS MAY APPEAR IN EACH BATCH

THE DATASET USED IN THIS EXAMPLE IS REUTERS21578 AS PROVIDED BY THE UCI ML REPOSITORY IT WILL BE AUTOMATICALLY DOWNLOADED AND UNCOMPRESSED ON FIRST RUN

THE PLOT REPRESENTS THE LEARNING CURVE OF THE CLASSIFIER THE EVOLUTION OF CLASSIFICATION ACCURACY OVER THE COURSE OF THE MINIBATCHES ACCURACY IS MEASURED ON THE FIRST 1000 SAMPLES HELD OUT AS A VALIDATION SET

TO LIMIT THE MEMORY CONSUMPTION WE QUEUE EXAMPLES UP TO A FIXED AMOUNT BEFORE FEEDING THEM TO THE LEARNER

AUTHORS EUSTACHE DIEMERT EUSTACHEDIEMERTFR  
FEDERICOV HTTPSGITHUBCOMFEDERICOV  
LICENSE BSD 3 CLAUSE

FROM GLOB IMPORT GLOB  
IMPORT ITERTOOLS  
IMPORT OSPATH  
IMPORT RE  
IMPORT TARFILE  
IMPORT TIME  
IMPORT SYS  
IMPORT NUMPY AS NP  
IMPORT MATPLOTLIBPYPLOT AS PLT  
FROM MATPLOTLIB IMPORT RCPARAMS  
FROM HTMLPARSER IMPORT HTMLPARSER  
FROM URLLIBREQUEST IMPORT URLRETRIEVE  
FROM SKLEARNDATASETS IMPORT GETDATAHOME  
FROM SKLEARNFEATUREEXTRACTIONTEXT IMPORT HASHINGVECTORIZER  
FROM SKLEARNLINEARMODEL IMPORT SGDCLASSIFIER  
FROM SKLEARNLINEARMODEL IMPORT PASSIVEAGGRESSIVECLASSIFIER  
FROM SKLEARNLINEARMODEL IMPORT PERCEPTRON  
FROM SKLEARNNAIVEBAYES IMPORT MULTINOMIALNB

DEFNOTINSPHINX  
HACK TO DETECT WHETHER WE ARE RUNNING BY THE SPHINX BUILDER

RETURNFILE INGLOBS

REUTERS DATASET RELATED ROUTINES

CLASS REUTERSPARSER HTMLPARSER

UTILITY CLASS TO PARSE A SGML FILE AND YIELD DOCUMENTS ONE AT A TIME

DEFINITSELF ENCODINGLATIN1  
HTMLPARSERINITSELF  
SELFRESET  
SELFENCODING ENCODING  
DEFHANDLESTARTTAGSELF TAG ATTRS  
METHOD START TAG  
GETATTRSELF METHOD LAMBDA NONEATTRS  
DEFHANDLEENDTAGSELF TAG  
METHOD END TAG  
GETATTRSELF METHOD LAMBDA NONE

764 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
DEFRESETSELF  
SELFINTITLE 0  
SELFINBODY 0  
SELFINTOPICS 0  
SELFINTOPICD 0  
SELFTITLE  
SELFBODY  
SELFTOPICS  
SELFTOPICD  
DEFPARSESELF FD  
SELFDOCS  
FORCHUNKINFD  
SELFFEEDCHUNKDECODESELFENCODING  
FORDOCINSELFDOCS  
YIELDDOC  
SELFDOCS  
SELFFCLOSE  
DEFHANDLEDATASELF DATA  
IFSELFINBODY  
SELFBODY DATA  
ELIFSELFINTITLE  
SELFTITLE DATA  
ELIFSELFINTOPICD  
SELFTOPICD DATA  
DEFSTARTREUTERSSELF ATTRIBUTES  
PASS  
DEFENDREUTERSSELF  
SELFBODY RESUBRS R SELFBODY  
SELFDOCSAPPENDTITLE SELFTITLE  
BODY SELFBODY  
TOPICS SELFTOPICS  
SELFRESET  
DEFSTARTTITLESELF ATTRIBUTES  
SELFINTITLE 1  
DEFENDTITLESELF  
SELFINTITLE 0  
DEFSTARTBODYSELF ATTRIBUTES  
SELFINBODY 1  
DEFENDBODYSELF  
SELFINBODY 0  
DEFSTARTTOPICSSELF ATTRIBUTES  
SELFINTOPICS 1  
DEFENDTOPICSSELF  
SELFINTOPICS 0  
DEFSTARTDSELF ATTRIBUTES  
SELFINTOPICD 1  
52 EXAMPLES BASED ON REAL WORLD DATASETS 765

SCIKITLEARN USER GUIDE RELEASE 0213  
DEFENDDSELF  
SELFINTOPICD 0  
SELFTOPICSAPPENDSELFTOPICD  
SELFTOPICD  
DEFSTREAMREUTERSDOCUMENTSDATAPATHNONE  
ITERATE OVER DOCUMENTS OF THE REUTERS DATASET  
THE REUTERS ARCHIVE WILL AUTOMATICALLY BE DOWNLOADED AND UNCOMPRESSED IF  
THE DATAPATH DIRECTORY DOES NOT EXIST  
DOCUMENTS ARE REPRESENTED AS DICTIONARIES WITH BODY STR  
TITLE STR TOPICS LISTSTR KEYS  
  
DOWNLOADURL HTTPARCHIVEICSDUCIEDUMLMACHINELEARNINGDATABASES  
REUTERS21578MLDREUTERS21578TARGZ  
ARCHIVEFILENAME REUTERS21578TARGZ  
IFDATAPATH ISNONE  
DATAPATH OSPATHJOINGETDATAHOME REUTERS  
IF NOTOSPATHEXISTSDATAPATH  
DOWNLOAD THE DATASET  
PRINTDOWNLOADING DATASET ONCE AND FOR ALL INTO S  
DATAPATH  
OSMKDIRDATAPATH  
DEFPROGRESSBLOCKNUM BS SIZE  
TOTALSZMB 2FMB SIZE 1E6  
CURRENTSZMB 2FMB BLOCKNUM BS 1E6  
IFNOTINSPHINX  
SYSSTDOUTWRITE  
RDOWNLOADED SS CURRENTSZMB TOTALSZMB  
ARCHIVEPATH OSPATHJOINDATAPATH ARCHIVEFILENAME  
URLRETRIEVEDOWNLOADURL FILENAMEARCHIVEPATH  
REPORTHOOKPROGRESS  
IFNOTINSPHINX  
SYSSTDOUTWRITE R  
PRINTUNTARRING REUTERS DATASET  
TARFILEOPENARCHIVEPATH RGZEXTRACTALLDATAPATH  
PRINTDONE  
PARSER REUTERSPARSER  
FORFILENAME INGLOBOSPATHJOINDATAPATH SGM  
FORDOCINPARSERPARSEOPENFILENAME RB  
YELDDOC  
MAIN  
CREATE THE VECTORIZER AND LIMIT THE NUMBER OF FEATURES TO A REASONABLE MAXIMUM  
VECTORIZER HASHINGVECTORIZERDECODEERRORIGNORE NFEATURES2 18  
ALTERNATESIGNFALSE  
766 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
ITERATOR OVER PARSED REUTERS SGML FILES  
DATASTREAM STREAMREUTERSDOCUMENTS  
WE LEARN A BINARY CLASSIFICATION BETWEEN THE ACQ CLASS AND ALL THE OTHERS  
ACQ WAS CHOSEN AS IT IS MORE OR LESS EVENLY DISTRIBUTED IN THE REUTERS  
FILES FOR OTHER DATASETS ONE SHOULD TAKE CARE OF CREATING A TEST SET WITH  
A REALISTIC PORTION OF POSITIVE INSTANCES  
ALLCLASSES NPARRAY0 1  
POSITIVECLASS ACQ  
HERE ARE SOME CLASSIFIERS THAT SUPPORT THE PARTIALFIT METHOD  
PARTIALFITCLASSIFIERS  
SGD SGDCLASSIFIERMAXITER5 TOL1E3  
PERCEPTRON PERCEPTRONTOL1E3  
NB MULTINOMIAL MULTINOMIALNBALPHA001  
PASSIVEAGGRESSIVE PASSIVEAGGRESSIVECLASSIFIERTOL1E3  
  
DEFGETMINIBATCHDOCITER SIZE POSCLASSPOSITIVECLASS  
EXTRACT A MINIBATCH OF EXAMPLES RETURN A TUPLE XTEXT Y  
NOTE SIZE IS BEFORE EXCLUDING INVALID DOCS WITH NO TOPICS ASSIGNED  
  
DATA TITLE NNBODYFORMAT DOC POSCLASS INDOCTOPICS  
FORDOCINITERTOOLSISLICEDOCITER SIZE  
IFDOCTOPICS  
IF NOTLENDATA  
RETURNNPASARRAY DTYPEINT NPASARRAY DTYPEINT  
XTEXT Y ZIP DATA  
RETURNXTEXT NPASARRAYY DTYPEINT  
DEFITERMINIBATCHESDOCITER MINIBATCHSIZE  
GENERATOR OF MINIBATCHES  
XTEXT Y GETMINIBATCHDOCITER MINIBATCHSIZE  
WHILELENXTEXT  
YIELDXTEXT Y  
XTEXT Y GETMINIBATCHDOCITER MINIBATCHSIZE  
TEST DATA STATISTICS  
TESTSTATS NTEST 0 NTESTPOS 0  
FIRST WE HOLD OUT A NUMBER OF EXAMPLES TO ESTIMATE ACCURACY  
NTESTDOCUMENTS 1000  
TICK TIMETIME  
XTESTTEXT YTEST GETMINIBATCHDATASTREAM 1000  
PARSINGTIME TIMETIME TICK  
TICK TIMETIME  
XTEST VECTORIZERTRANSFORMXTESTTEXT  
VECTORIZINGTIME TIMETIME TICK  
TESTSTATSNTEST LENYTEST  
TESTSTATSNTESTPOS SUMYTEST  
PRINTTEST SET IS DDOCUMENTS DPOSITIVE LENYTEST SUMYTEST  
52 EXAMPLES BASED ON REAL WORLD DATASETS 767

SCIKITLEARN USER GUIDE RELEASE 0213  
DEFPROGRESSCLSNAME STATS  
REPORT PROGRESS INFORMATION RETURN A STRING  
DURATION TIMETIME STATST0  
S 20SCCLASSIFIER T CLSNAME  
S NTRAIN6D TRAIN DOCS NTRAINPOS6D POSITIVE STATS  
S NTEST6D TEST DOCS NTESTPOS6D POSITIVE TESTSTATS  
S ACCURACY ACCURACY3F STATS  
S IN 2FS 5DDOCSS DURATION STATSNTRAIN DURATION  
RETURNS  
CLSSTATS  
FORCLSNAME INPARTIALFITCLASSIFIERS  
STATS NTRAIN 0 NTRAINPOS 0  
ACCURACY 00 ACCURACYHISTORY 0 0 TO TIMETIME  
RUNTIMEHISTORY 0 0 TOTALFITTIME 00  
CLSSTATSCLSNAME STATS  
GETMINIBATCHDATASTREAM NTESTDOCUMENTS  
DISCARD TEST SET  
WE WILL FEED THE CLASSIFIER WITH MINIBATCHES OF 1000 DOCUMENTS THIS MEANS  
WE HAVE AT MOST 1000 DOCS IN MEMORY AT ANY TIME THE SMALLER THE DOCUMENT  
BATCH THE BIGGER THE RELATIVE OVERHEAD OF THE PARTIAL FIT METHODS  
MINIBATCHSIZE 1000  
CREATE THE DATASTREAM THAT PARSES REUTERS SGML FILES AND ITERATES ON  
DOCUMENTS AS A STREAM  
MINIBATCHITERATORS ITERMINIBATCHESDATASTREAM MINIBATCHSIZE  
TOTALVECTTIME 00  
MAIN LOOP ITERATE ON MINIBATCHES OF EXAMPLES  
FORI XTRAINTEXT YTRAIN INENUMERATEMINIBATCHITERATORS  
TICK TIMETIME  
XTRAIN VECTORIZERTRANSFORMXTRAINTEXT  
TOTALVECTTIME TIMETIME TICK  
FORCLSNAME CLS INPARTIALFITCLASSIFIERSITEMS  
TICK TIMETIME  
UPDATE ESTIMATOR WITH EXAMPLES IN THE CURRENT MINIBATCH  
CLSPARTIALFITXTRAIN YTRAIN CLASSESALLCLASSES  
ACCUMULATE TEST ACCURACY STATS  
CLSSTATSCLSNAMETOTALFITTIME TIMETIME TICK  
CLSSTATSCLSNAMEENTRAIN XTRAINSHAPE0  
CLSSTATSCLSNAMEENTRAINPOS SUMYTRAIN  
TICK TIMETIME  
CLSSTATSCLSNAMEACCURACY CLSSCOREXTEST YTEST  
CLSSTATSCLSNAMEPREDICTIONTIME TIMETIME TICK  
ACCHISTORY CLSSTATSCLSNAMEACCURACY  
CLSSTATSCLSNAMEENTRAIN  
CLSSTATSCLSNAMEACCURACYHISTORYAPPENDACCHISTORY  
RUNHISTORY CLSSTATSCLSNAMEACCURACY  
TOTALVECTTIME CLSSTATSCLSNAMETOTALFITTIME  
768 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
CLSSTATSCLSNAMERUNTIMEHISTORYAPPENDRUNHISTORY  
IFI 3 0  
PRINTPROGRESSCLSNAME CLSSTATSCLSNAME  
IFI 3 0  
PRINTN  
OUT  
DOWNLOADING DATASET ONCE AND FOR ALL INTO HOMECIRCLECISCIKITLEARNDATAAREUTERS  
UNTARRING REUTERS DATASET  
DONE  
TEST SET IS 966 DOCUMENTS 88 POSITIVE  
SGD CLASSIFIER 994 TRAIN DOCS 121 POSITIVE 966  
↳TEST DOCS 88 POSITIVE ACCURACY 0907 IN 074S 1342 DOCSS  
PERCEPTRON CLASSIFIER 994 TRAIN DOCS 121 POSITIVE 966  
↳TEST DOCS 88 POSITIVE ACCURACY 0922 IN 074S 1338 DOCSS  
NB MULTINOMIAL CLASSIFIER 994 TRAIN DOCS 121 POSITIVE 966  
↳TEST DOCS 88 POSITIVE ACCURACY 0909 IN 075S 1327 DOCSS  
PASSIVEAGGRESSIVE CLASSIFIER 994 TRAIN DOCS 121 POSITIVE 966  
↳TEST DOCS 88 POSITIVE ACCURACY 0944 IN 075S 1323 DOCSS  
SGD CLASSIFIER 3924 TRAIN DOCS 496 POSITIVE 966  
↳TEST DOCS 88 POSITIVE ACCURACY 0957 IN 220S 1784 DOCSS  
PERCEPTRON CLASSIFIER 3924 TRAIN DOCS 496 POSITIVE 966  
↳TEST DOCS 88 POSITIVE ACCURACY 0952 IN 220S 1782 DOCSS  
NB MULTINOMIAL CLASSIFIER 3924 TRAIN DOCS 496 POSITIVE 966  
↳TEST DOCS 88 POSITIVE ACCURACY 0917 IN 221S 1778 DOCSS  
PASSIVEAGGRESSIVE CLASSIFIER 3924 TRAIN DOCS 496 POSITIVE 966  
↳TEST DOCS 88 POSITIVE ACCURACY 0967 IN 221S 1776 DOCSS  
SGD CLASSIFIER 6836 TRAIN DOCS 793 POSITIVE 966  
↳TEST DOCS 88 POSITIVE ACCURACY 0967 IN 382S 1791 DOCSS  
PERCEPTRON CLASSIFIER 6836 TRAIN DOCS 793 POSITIVE 966  
↳TEST DOCS 88 POSITIVE ACCURACY 0950 IN 382S 1790 DOCSS  
NB MULTINOMIAL CLASSIFIER 6836 TRAIN DOCS 793 POSITIVE 966  
↳TEST DOCS 88 POSITIVE ACCURACY 0921 IN 382S 1787 DOCSS  
PASSIVEAGGRESSIVE CLASSIFIER 6836 TRAIN DOCS 793 POSITIVE 966  
↳TEST DOCS 88 POSITIVE ACCURACY 0967 IN 383S 1786 DOCSS  
SGD CLASSIFIER 9597 TRAIN DOCS 1139 POSITIVE 966  
↳TEST DOCS 88 POSITIVE ACCURACY 0963 IN 539S 1779 DOCSS  
PERCEPTRON CLASSIFIER 9597 TRAIN DOCS 1139 POSITIVE 966  
↳TEST DOCS 88 POSITIVE ACCURACY 0955 IN 540S 1778 DOCSS  
NB MULTINOMIAL CLASSIFIER 9597 TRAIN DOCS 1139 POSITIVE 966  
↳TEST DOCS 88 POSITIVE ACCURACY 0930 IN 540S 1776 DOCSS  
PASSIVEAGGRESSIVE CLASSIFIER 9597 TRAIN DOCS 1139 POSITIVE 966  
↳TEST DOCS 88 POSITIVE ACCURACY 0965 IN 541S 1775 DOCSS  
SGD CLASSIFIER 12513 TRAIN DOCS 1558 POSITIVE 966  
↳TEST DOCS 88 POSITIVE ACCURACY 0969 IN 702S 1783 DOCSS  
PERCEPTRON CLASSIFIER 12513 TRAIN DOCS 1558 POSITIVE 966  
↳TEST DOCS 88 POSITIVE ACCURACY 0843 IN 702S 1783 DOCSS  
NB MULTINOMIAL CLASSIFIER 12513 TRAIN DOCS 1558 POSITIVE 966  
↳TEST DOCS 88 POSITIVE ACCURACY 0936 IN 702S 1781 DOCSS  
52 EXAMPLES BASED ON REAL WORLD DATASETS 769

SCIKITLEARN USER GUIDE RELEASE 0213

PASSIVEAGGRESSIVE CLASSIFIER 12513 TRAIN DOCS 1558 POSITIVE 966  
'→TEST DOCS 88 POSITIVE ACCURACY 0965 IN 703S 1781 DOCSS

SGD CLASSIFIER 14935 TRAIN DOCS 1863 POSITIVE 966  
'→TEST DOCS 88 POSITIVE ACCURACY 0966 IN 852S 1752 DOCSS

PERCEPTRON CLASSIFIER 14935 TRAIN DOCS 1863 POSITIVE 966  
'→TEST DOCS 88 POSITIVE ACCURACY 0958 IN 853S 1751 DOCSS

NB MULTINOMIAL CLASSIFIER 14935 TRAIN DOCS 1863 POSITIVE 966  
'→TEST DOCS 88 POSITIVE ACCURACY 0938 IN 853S 1750 DOCSS

PASSIVEAGGRESSIVE CLASSIFIER 14935 TRAIN DOCS 1863 POSITIVE 966  
'→TEST DOCS 88 POSITIVE ACCURACY 0955 IN 854S 1749 DOCSS

SGD CLASSIFIER 17417 TRAIN DOCS 2171 POSITIVE 966  
'→TEST DOCS 88 POSITIVE ACCURACY 0960 IN 998S 1745 DOCSS

PERCEPTRON CLASSIFIER 17417 TRAIN DOCS 2171 POSITIVE 966  
'→TEST DOCS 88 POSITIVE ACCURACY 0964 IN 998S 1744 DOCSS

NB MULTINOMIAL CLASSIFIER 17417 TRAIN DOCS 2171 POSITIVE 966  
'→TEST DOCS 88 POSITIVE ACCURACY 0943 IN 999S 1743 DOCSS

PASSIVEAGGRESSIVE CLASSIFIER 17417 TRAIN DOCS 2171 POSITIVE 966  
'→TEST DOCS 88 POSITIVE ACCURACY 0960 IN 999S 1743 DOCSS

PLOT RESULTS

DEFPLOTACCURACYX Y XLEGEND

PLOT ACCURACY AS A FUNCTION OF X

X NPARRAYX

Y NPARRAYY

PLTTITLECLASSIFICATION ACCURACY AS A FUNCTION OF S XLEGEND

PLTXLABEL S XLEGEND

PLTYLABELACCURACY

PLTGRIDTRUE

PLTPLOTX Y

RCPARAMSLEGENDFONTSIZE 10

CLS NAMES LISTSORTEDCLSSTATSKEYS

PLOT ACCURACY EVOLUTION

PLTFigure

FOR STATS INSORTEDCLSSTATSITEMS

PLOT ACCURACY EVOLUTION WITH EXAMPLES

ACCURACY NEXAMPLES ZIP STATSACCURACYHISTORY

PLOTACCURACYNEXAMPLES ACCURACY TRAINING EXAMPLES

AX PLTGCA

AXSETYLIM08 1

PLTLEGENDCLS NAMES LOCBEST

PLTFigure

FOR STATS INSORTEDCLSSTATSITEMS

PLOT ACCURACY EVOLUTION WITH RUNTIME

ACCURACY RUNTIME ZIP STATS RUNTIMEHISTORY

PLOTACCURACYRUNTIME ACCURACY RUNTIME S

AX PLTGCA

AXSETYLIM08 1

770 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213  
PLTLEGENDCLS NAMES LOCBEST  
PLOT FITTING TIMES  
PLTFigure  
FIG PLTGCF  
CLSRUNTIME STATSTOTALFITTIME  
FORCLSNAME STATS INSORTEDCLSSTATSITEMS  
CLSRUNTIMEAPPENDTOTALVECTTIME  
CLS NAMESAPPENDVECTORIZATION  
BARCOLORS B G R C M Y  
AX PLTSUBPLOT111  
RECTANGLES PLTBARRANGELENCLS NAMES CLSRUNTIME WIDTH05  
COLORBARCOLORS  
AXSETXTICKSNPLINSPACE0 LENCLS NAMES 1 LENCLS NAMES  
AXSETXTICKLABELSCLS NAMES FONTSIZE10  
YMAX MAXCLSRUNTIME 12  
AXSETYLIM0 YMAX  
AXSETYLABELRUNTIME S  
AXSETTITLETRAINING TIMES  
DEFAULTLABELRECTANGLES  
ATTACH SOME TEXT VI AUTOLABEL ON RECTANGLES  
FORRECTINRECTANGLES  
HEIGHT RECTGETHEIGHT  
AXTEXTRECTGETX RECTGETWIDTH 2  
105HEIGHT 4F HEIGHT  
HACENTER VABOTTOM  
PLTSETPPLXTICKS1 ROTATION30  
AUTOLABELRECTANGLES  
PLTTIGHTLAYOUT  
PLTSHOW  
PLOT PREDICTION TIMES  
PLTFigure  
CLSRUNTIME  
CLS NAMES LISTSORTEDCLSSTATSKEYS  
FORCLSNAME STATS INSORTEDCLSSTATSITEMS  
CLSRUNTIMEAPPENDSTATSPREDICTIONTIME  
CLSRUNTIMEAPPENDPARSINGTIME  
CLS NAMESAPPENDREADPARSE NFEATEXTR  
CLSRUNTIMEAPPENDVECTORIZINGTIME  
CLS NAMESAPPENDHASHING NVECT  
AX PLTSUBPLOT111  
RECTANGLES PLTBARRANGELENCLS NAMES CLSRUNTIME WIDTH05  
COLORBARCOLORS  
AXSETXTICKSNPLINSPACE0 LENCLS NAMES 1 LENCLS NAMES  
AXSETXTICKLABELSCLS NAMES FONTSIZE8  
PLTSETPPLXTICKS1 ROTATION30  
YMAX MAXCLSRUNTIME 12  
AXSETYLIM0 YMAX  
52 EXAMPLES BASED ON REAL WORLD DATASETS 771

SCIKITLEARN USER GUIDE RELEASE 0213  
AXSETYLABELRUNTIME S  
AXSETTITLEPREDICTION TIMES DINSTANCES NTESTDOCUMENTS  
AUTOLABELRECTANGLES  
PLTTIGHTLAYOUT  
PLTSHOW  
•  
772 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

- 

52 EXAMPLES BASED ON REAL WORLD DATASETS 773



SCIKITLEARN USER GUIDE RELEASE 0213

•

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 12190 SECONDS

53 BICLUSTERING

EXAMPLES CONCERNING THE SKLEARNCLUSTERBICLUSTER MODULE

NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)

531 A DEMO OF THE SPECTRAL COCLUSTERING ALGORITHM

THIS EXAMPLE DEMONSTRATES HOW TO GENERATE A DATASET AND BICLUSTER IT USING THE SPECTRAL COCLUSTERING ALGORITHM

THE DATASET IS GENERATED USING THE MAKEBICLUSTERS FUNCTION WHICH CREATES A MATRIX OF SMALL VALUES AND IM

PLANTS BICLUSTER WITH LARGE VALUES THE ROWS AND COLUMNS ARE THEN SHUFFLED AND PASSED TO THE SPECTRAL COCLUSTERING ALGORITHM REARRANGING THE SHUFFLED MATRIX TO MAKE BICLUSTERS CONTIGUOUS SHOWS HOW ACCURATELY THE ALGORITHM FOUND THE BICLUSTERS

53 BICLUSTERING 775





SCIKITLEARN USER GUIDE RELEASE 0213

•

OUT  
CONSENSUS SCORE 1000  
PRINTDOC  
AUTHOR KEMAL EREN KEMALKEMALERENCOM  
LICENSE BSD 3 CLAUSE  
IMPORT NUMPY AS NP  
FROM MATPLOTLIB IMPORT PYPLLOTASPLT  
FROM SKLEARNDATASETS IMPORT MAKEBICCLUSTERS  
FROM SKLEARNDATASETS IMPORT SAMPLESGENERATOR ASSG  
FROM SKLEARNCLUSTERBICCLUSTER IMPORT SPECTRALCOCLUSTERING  
FROM SKLEARNMETRICS IMPORT CONSENSUSSCORE  
DATA ROWS COLUMNS MAKEBICCLUSTERS  
SHAPE300 300 NCLUSTERS5 NOISE5  
SHUFFLEFALSE RANDOMSTATE0  
778 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213  
PLTMATSHOWDATA CMAPPLTCMBLUES  
PLTTITLEORIGINAL DATASET  
DATA ROWIDX COLIDX SGSHUFFLEDATA RANDOMSTATE0  
PLTMATSHOWDATA CMAPPLTCMBLUES  
PLTTITLESHUFFLED DATASET  
MODEL SPECTRALCOCLUSTERINGNCLUSTERS55 RANDOMSTATE0  
MODELFITDATA  
SCORE CONSENSUSSCOREMODELBICLUSTERS  
ROWS ROWIDX COLUMNS COLIDX  
PRINTCONSENSUS SCORE 3FFORMATSCORE  
FITDATA DATANPARGSORTMODELROWLABELS  
FITDATA FITDATA NPARGSORTMODELCOLUMNLABELS  
PLTMATSHOWFITDATA CMAPPLTCMBLUES  
PLTTITLEAFTER BICLUSTERING REARRANGED TO SHOW BICLUSTERS  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0061 SECONDS  
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE  
532 A DEMO OF THE SPECTRAL BICLUSTERING ALGORITHM  
THIS EXAMPLE DEMONSTRATES HOW TO GENERATE A CHECKERBOARD DATASET AND BICLUSTER IT USING THE SPECTRAL BICLUSTERING ALGORITHM  
THE DATA IS GENERATED WITH THE MAKECHECKERBOARD FUNCTION THEN SHUFFLED AND PASSED TO THE SPECTRAL BICLUSTERING ALGORITHM THE ROWS AND COLUMNS OF THE SHUFFLED MATRIX ARE REARRANGED TO SHOW THE BICLUSTERS FOUND BY THE ALGORITHM  
THE OUTER PRODUCT OF THE ROW AND COLUMN LABEL VECTORS SHOWS A REPRESENTATION OF THE CHECKERBOARD STRUCTURE  
53 BICLUSTERING 779







SCIKITLEARN USER GUIDE RELEASE 0213

•

```
OUT
CONSENSUS SCORE 10
PRINTDOC
AUTHOR KEMAL EREN KEMALKEMALERENCOM
LICENSE BSD 3 CLAUSE
IMPORT NUMPY AS NP
FROM MATPLOTLIB IMPORT PYPLASPLT
FROM SKLEARNDATASETS IMPORT MAKECHECKERBOARD
FROM SKLEARNDATASETS IMPORT SAMPLESGENERATOR ASSG
FROM SKLEARNCLUSTERBICLUSTER IMPORT SPECTRALBICLUSTERING
FROM SKLEARNMETRICS IMPORT CONSENSUSSCORE
NCLUSTERS 4 3
DATA ROWS COLUMNS MAKECHECKERBOARD
SHAPE300 300 NCLUSTERSNCLUSTERS NOISE10
SHUFFLEFALSE RANDOMSTATE0
53 BICLUSTERING 783
```

SCIKITLEARN USER GUIDE RELEASE 0213  
PLTMATSHOWDATA CMAPPLTCMBLUES  
PLTTITLEORIGINAL DATASET  
DATA ROWIDX COLIDX SGSHUFFLEDATA RANDOMSTATE0  
PLTMATSHOWDATA CMAPPLTCMBLUES  
PLTTITLESHUFFLED DATASET  
MODEL SPECTRALBICLUSTERINGNCLUSTERSNCLUSTERS METHODLOG  
RANDOMSTATE0  
MODELFITDATA  
SCORE CONSENSUSSCOREMODELBICLUSTERS  
ROWS ROWIDX COLUMNS COLIDX  
PRINTCONSENSUS SCORE 1FFORMATSCORE  
FITDATA DATANPARGSORTMODELROWLABELS  
FITDATA FITDATA NPARGSORTMODELCOLUMNLABELS  
PLTMATSHOWFITDATA CMAPPLTCMBLUES  
PLTTITLEAFTER BICLUSTERING REARRANGED TO SHOW BICLUSTERS  
PLTMATSHOWNPOUTERNPSORTMODELROWLABELS 1  
NPSORTMODELCOLUMNLABELS 1  
CMAPPLTCMBLUES  
PLTTITLECHECKERBOARD STRUCTURE OF REARRANGED DATA  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0447 SECONDS  
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE  
533 BICLUSTERING DOCUMENTS WITH THE SPECTRAL COCLUSTERING ALGORITHM  
THIS EXAMPLE DEMONSTRATES THE SPECTRAL COCLUSTERING ALGORITHM ON THE TWENTY NEWSGROUPS DATASET THE 'COMPOSMS  
WINDOWSMISC' CATEGORY IS EXCLUDED BECAUSE IT CONTAINS MANY POSTS CONTAINING NOTHING BUT DATA  
THE TFIDF VECTORIZED POSTS FORM A WORD FREQUENCY MATRIX WHICH IS THEN BICLUSTERED USING DHILLON'S SPECTRAL CO  
CLUSTERING ALGORITHM THE RESULTING DOCUMENTWORD BICLUSTERS INDICATE SUBSETS WORDS USED MORE OFTEN IN THOSE SUBSETS  
DOCUMENTS  
FOR A FEW OF THE BEST BICLUSTERS ITS MOST COMMON DOCUMENT CATEGORIES AND ITS TEN MOST IMPORTANT WORDS GET PRINTED  
THE BEST BICLUSTERS ARE DETERMINED BY THEIR NORMALIZED CUT THE BEST WORDS ARE DETERMINED BY COMPARING THEIR SUMS  
INSIDE AND OUTSIDE THE BICLUSTER  
FOR COMPARISON THE DOCUMENTS ARE ALSO CLUSTERED USING MINIBATCHKMEANS THE DOCUMENT CLUSTERS DERIVED FROM THE  
BICLUSTERS ACHIEVE A BETTER VMEASURE THAN CLUSTERS FOUND BY MINIBATCHKMEANS  
OUT  
VECTORIZING  
COCLUSTERING  
DONE IN 410S VMEASURE 04435  
MINIBATCHKMEANS  
DONE IN 636S VMEASURE 03344  
784 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
BEST BICLUSTERS

BICLUSTER 0 1957 DOCUMENTS 4363 WORDS  
CATEGORIES 23 TALKPOLITICSGUNS 18 TALKPOLITICSMISC 17 SCIMED  
WORDS GUN GUNS GEB BANKS GORDON CLINTON PITT CDT SURRENDER VEAL  
BICLUSTER 1 1263 DOCUMENTS 3551 WORDS  
CATEGORIES 27 SOCRELIGIONCHRISTIAN 25 TALKPOLITICSMIDEAST 24 ALTATHEISM  
WORDS GOD JESUS CHRISTIANS SIN OBJECTIVE KENT BELIEF CHRIST FAITH  
↪MORAL  
BICLUSTER 2 2212 DOCUMENTS 2774 WORDS  
CATEGORIES 18 COMPSYSMACHARDWARE 17 COMPSYSIBMPCHARDWARE 15 COMP  
↪GRAPHICS  
WORDS VOLTAGE BOARD DSP STEREO RECEIVER PACKAGES SHIPPING CIRCUIT  
↪PACKAGE COMPRESSION  
BICLUSTER 3 1774 DOCUMENTS 2629 WORDS  
CATEGORIES 27 RECMOTORCYCLES 23 RECAUTOS 13 MISCFORSALE  
WORDS BIKE CAR DOD ENGINE MOTORCYCLE RIDE HONDA BIKES HELMET BMW  
BICLUSTER 4 200 DOCUMENTS 1167 WORDS  
CATEGORIES 81 TALKPOLITICSMIDEAST 10 ALTATHEISM 8 SOCRELIGIONCHRISTIAN  
WORDS TURKISH ARMENIA ARMENIAN ARMENIANS TURKS PETCH SERA ZUMA ARGIC  
↪GVG47  
FROM COLLECTIONS IMPORT DEFAULTDICT  
IMPORT OPERATOR  
FROM TIME IMPORT TIME  
IMPORT NUMPY AS NP  
FROM SKLEARNCLUSTERBICLUSTER IMPORT SPECTRALCOCLUSTERING  
FROM SKLEARNCLUSTER IMPORT MINIBATCHKMEANS  
FROM SKLEARNDATASETSTWENTYNEWSGROUPS IMPORT FETCH20NEWSGROUPS  
FROM SKLEARNFEATUREEXTRACTIONTEXT IMPORT TFIDFVECTORIZER  
FROM SKLEARNMETRICSCUSTER IMPORT VMEASURESCORE  
PRINTDOC  
DEFNUMBERNORMALIZERTOKENS  
MAP ALL NUMERIC TOKENS TO A PLACEHOLDER  
FOR MANY APPLICATIONS TOKENS THAT BEGIN WITH A NUMBER ARE NOT DIRECTLY  
USEFUL BUT THE FACT THAT SUCH A TOKEN EXISTS CAN BE RELEVANT BY APPLYING  
THIS FORM OF DIMENSIONALITY REDUCTION SOME METHODS MAY PERFORM BETTER  
  
RETURNNUMBER IFTOKEN0ISDIGIT ELSETOKENFORTOKENINTOKENS  
CLASS NUMBERNORMALIZINGVECTORIZER TFIDFVECTORIZER  
53 BICLUSTERING 785

SCIKITLEARN USER GUIDE RELEASE 0213  
DEFBUILDTOKENIZERSELF  
TOKENIZE SUPERBUILDTOKENIZER  
RETURN LAMBDA DOC LISTNUMBERNORMALIZERTOKENIZEDOC  
EXCLUDE COMPOSMSWINDOWSMISC  
CATEGORIES ALTATHEISM COMPGRAPHICS  
COMPSYSIBMPCHARDWARE COMPSYSMACHARDWARE  
COMPWINDOWSX MISCFORSALE RECAUTOS  
RECMOTORCYCLES RECSPORTBASEBALL  
RECSPORTHOCKEY SCICRYPT SCIELECTRONICS  
SCIMED SCISPACE SOCRELIGIONCHRISTIAN  
TALKPOLITICSGUNS TALKPOLITICSMIDEAST  
TALKPOLITICSMISC TALKRELIGIONMISC  
NEWSGROUPS FETCH20NEWSGROUPSCATEGORIESCATEGORIES  
YTRUE NEWSGROUPSTARGET  
VECTORIZER NUMBERNORMALIZINGVECTORIZERSTOPWORDSENGLISH MINDF5  
COCLUSTER SPECTRALCOCLUSTERINGNCLUSTERSLENCATEGORIES  
SVDMETHODARPACK RANDOMSTATE0  
KMEANS MINIBATCHKMEANSNCLUSTERSLENCATEGORIES BATCHSIZE20000  
RANDOMSTATE0  
PRINTVECTORIZING  
X VECTORIZERFITTRANSFORMNEWSGROUPSDATA  
PRINTCOCLUSTERING  
STARTTIME TIME  
COCLUSTERFITX  
YCOCLUSTER COCLUSTERROWLABELS  
PRINTDONE IN 2FS VMEASURE 4FFORMAT  
TIME STARTTIME  
VMEASURESCOREYCOCLUSTER YTRUE  
PRINTMINIBATCHKMEANS  
STARTTIME TIME  
YKMEANS KMEANSFITPREDICTX  
PRINTDONE IN 2FS VMEASURE 4FFORMAT  
TIME STARTTIME  
VMEASURESCOREYKMEANS YTRUE  
FEATURENAMES VECTORIZERGETFEATURENAMES  
DOCUMENTNAMES LISTNEWSGROUPSTARGETNAMESI FORIINNEWSGROUPSTARGET  
DEFBICLUSTERNCUTI  
ROWS COLS COCLUSTERGETINDICESI  
IF NOTNPANYROWS ANDNPANYCOLS  
IMPORT SYS  
RETURNSYSFLOATINFOMAX  
ROWCOMPLEMENT NPNONZERONPLOGICALNOTCOCLUSTERROWSIO  
COLCOMPLEMENT NPNONZERONPLOGICALNOTCOCLUSTERCOLUMNSIO  
NOTE THE FOLLOWING IS IDENTICAL TO XROWS NPNEWAXIS  
COLSSUM BUT MUCH FASTER IN SCIPY 016  
WEIGHT XROWS COLSSUM  
CUT XROWCOMPLEMENT COLSSUM  
XROWS COLCOMPLEMENTSUM  
RETURN CUT WEIGHT  
786 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213  
DEFMOSTCOMMOND  
ITEMS OF A DEFAULTDICTINT WITH THE HIGHEST VALUES  
LIKE COUNTERMOSTCOMMON IN PYTHON 27

RETURNSORTEDDDITEMS KEYOPERATORITEMGETTER1 REVERSETRUE  
BICLUSTERNCUTS LISTBICLUSTERNCUTI  
FORIINRANGELENNEWSGROUPSTARGETNAMES  
BESTIDX NPARGSORTBICLUSTERNCUTS5  
PRINT  
PRINTBEST BICLUSTERS  
PRINT  
FORIDX CLUSTER INENUMERATEBESTIDX  
NROWS NCOLS COCLUSTERGETSHAPECLUSTER  
CLUSTERDOCS CLUSTERWORDS COCLUSTERGETINDICESCLUSTER  
IF NOTLENCLUSTERDOCS OR NOTLENCLUSTERWORDS  
CONTINUE  
CATEGORIES  
COUNTER DEFAULTDICTINT  
FORIINCLUSTERDOCS  
COUNTERDOCUMENTNAMESI 1  
CATSTRING JOINOF FORMATFLOATC NROWS 100 NAME  
FORNAME C INMOSTCOMMONCOUNTER3  
WORDS  
OUTOFCLUSTERDOCS COCLUSTERROWLABELS CLUSTER  
OUTOFCLUSTERDOCS NPWHEREOUTOFCLUSTERDOCS0  
WORDCOL X CLUSTERWORDS  
WORDSCORES NPARRAYWORDCOLCLUSTERDOCS SUMAXIS0  
WORDCOLOUTOFCLUSTERDOCS SUMAXIS0  
WORDSCORES WORDSCORESRAVEL  
IMPORTANTWORDS LISTFEATURENAMESCLUSTERWORDSI  
FORIINWORDSCORESARGSORT111  
PRINTBICLUSTER DOCUMENTS WORDSFORMAT  
IDX NROWS NCOLS  
PRINTCATEGORIES FORMATCATSTRING  
PRINTWORDS NFORMAT JOINIMPORTANTWORDS  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 13291 SECONDS  
54 CALIBRATION  
EXAMPLES ILLUSTRATING THE CALIBRATION OF PREDICTED PROBABILITIES OF CLASSIFIERS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
54 CALIBRATION 787

SCIKITLEARN USER GUIDE RELEASE 0213

541 COMPARISON OF CALIBRATION OF CLASSIFIERS

WELL CALIBRATED CLASSIFIERS ARE PROBABILISTIC CLASSIFIERS FOR WHICH THE OUTPUT OF THE PREDICTPROBA METHOD CAN BE DIRECTLY INTERPRETED AS A CONFIDENCE LEVEL FOR INSTANCE A WELL CALIBRATED BINARY CLASSIFIER SHOULD CLASSIFY THE SAMPLES SUCH THAT AMONG THE SAMPLES TO WHICH IT GAVE A PREDICTPROBA VALUE CLOSE TO 0.8 APPROX 80 ACTUALLY BELONG TO THE POSITIVE CLASS

LOGISTICREGRESSION RETURNS WELL CALIBRATED PREDICTIONS AS IT DIRECTLY OPTIMIZES LOGLOSS IN CONTRAST THE OTHER METHODS RETURN BIASED PROBABILITIES WITH DIFFERENT BIASES PER METHOD

- GAUSSIANNAIVEBAYES TENDS TO PUSH PROBABILITIES TO 0 OR 1 NOTE THE COUNTS IN THE HISTOGRAMS THIS IS MAINLY BECAUSE IT MAKES THE ASSUMPTION THAT FEATURES ARE CONDITIONALLY INDEPENDENT GIVEN THE CLASS WHICH IS NOT THE CASE IN THIS DATASET WHICH CONTAINS 2 REDUNDANT FEATURES
- RANDOMFORESTCLASSIFIER SHOWS THE OPPOSITE BEHAVIOR THE HISTOGRAMS SHOW PEAKS AT APPROX 0.2 AND 0.9 PROBABILITY WHILE PROBABILITIES CLOSE TO 0 OR 1 ARE VERY RARE AN EXPLANATION FOR THIS IS GIVEN BY NICULESCUMIZIL AND CARUANA1 “METHODS SUCH AS BAGGING AND RANDOM FORESTS THAT AVERAGE PREDICTIONS FROM A BASE SET OF MODELS CAN HAVE DIFFICULTY MAKING PREDICTIONS NEAR 0 AND 1 BECAUSE VARIANCE IN THE UNDERLYING BASE MODELS WILL BIAS PREDICTIONS THAT SHOULD BE NEAR ZERO OR ONE AWAY FROM THESE VALUES BECAUSE PREDICTIONS ARE RESTRICTED TO THE INTERVAL [0, 1] ERRORS CAUSED BY VARIANCE TEND TO BE ONE SIDED NEAR ZERO AND ONE FOR EXAMPLE IF A MODEL SHOULD PREDICT 0 FOR A CASE THE ONLY WAY BAGGING CAN ACHIEVE THIS IS IF ALL BAGGED TREES PREDICT ZERO IF WE ADD NOISE TO THE TREES THAT BAGGING IS AVERAGING OVER THIS NOISE WILL CAUSE SOME TREES TO PREDICT VALUES LARGER THAN 0 FOR THIS CASE THUS MOVING THE AVERAGE PREDICTION OF THE BAGGED ENSEMBLE AWAY FROM 0 WE OBSERVE THIS EFFECT MOST STRONGLY WITH RANDOM FORESTS BECAUSE THE BASELEVEL TREES TRAINED WITH RANDOM FORESTS HAVE RELATIVELY HIGH VARIANCE DUE TO FEATURE SUBSETTING” AS A RESULT THE CALIBRATION CURVE SHOWS A CHARACTERISTIC SIGMOID SHAPE INDICATING THAT THE CLASSIFIER COULD TRUST ITS “INTUITION” MORE AND RETURN PROBABILITIES CLOSER TO 0 OR 1 TYPICALLY
- SUPPORT VECTOR CLASSIFICATION SVC SHOWS AN EVEN MORE SIGMOID CURVE AS THE RANDOMFORESTCLASSIFIER WHICH IS TYPICAL FOR MAXIMUMMARGIN METHODS COMPARE NICULESCUMIZIL AND CARUANA1 WHICH FOCUS ON HARD SAMPLES THAT ARE CLOSE TO THE DECISION BOUNDARY THE SUPPORT VECTORS

REFERENCES

1. PREDICTING GOOD PROBABILITIES WITH SUPERVISED LEARNING A NICULESCUMIZIL R CARUANA ICML 2005

788 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
PRINTDOC  
AUTHOR JAN HENDRIK METZEN JHMINFORMATIKUNIBREMENDE  
LICENSE BSD STYLE  
IMPORT NUMPY AS NP  
NPRANDOMSEED0  
IMPORT MATPLOTLIBPYPLOT AS PLT  
FROM SKLEARN IMPORT DATASETS  
FROM SKLEARNNAIVEBAYES IMPORT GAUSSIANNB  
FROM SKLEARNLINEARMODEL IMPORT LOGISTICREGRESSION  
FROM SKLEARNENSEMBLE IMPORT RANDOMFORESTCLASSIFIER  
FROM SKLEARN SVM IMPORT LINEARSVC  
54 CALIBRATION 789

SCIKITLEARN USER GUIDE RELEASE 0213  
FROM SKLEARNCALIBRATION IMPORT CALIBRATIONCURVE  
X Y DATASETSMAKECLASSIFICATIONNSAMPLES100000 NFEATURES20  
NINFORMATIVE2 NREDUNDANT2  
TRAINSAMPLES 100 SAMPLES USED FOR TRAINING THE MODELS  
XTRAIN XTRAINSAMPLES  
XTEST XTRAINSAMPLES  
YTRAIN YTRAINSAMPLES  
YTEST YTRAINSAMPLES  
CREATE CLASSIFIERS  
LR LOGISTICREGRESSIONSOLVERLBFGS  
GNB GAUSSIANNB  
SVC LINEARSVCC10  
RFC RANDOMFORESTCLASSIFIERNESTIMATORS100

PLOT CALIBRATION PLOTS  
PLTFIGUREFIGSIZE10 10  
AX1 PLTSUBPLOT2GRID3 1 0 0 ROWSPAN2  
AX2 PLTSUBPLOT2GRID3 1 2 0  
AX1PLOT0 1 0 1 K LABELPERFECTLY CALIBRATED  
FORCLF NAME INLR LOGISTIC  
GNB NAIVE BAYES  
SVC SUPPORT VECTOR CLASSIFICATION  
RFC RANDOM FOREST  
CLFFITXTRAIN YTRAIN  
IFHASATTRCLF PREDICTPROBA  
PROBPOS CLFPREDICTPROBAXTEST 1  
ELSE USE DECISION FUNCTION  
PROBPOS CLFDECISIONFUNCTIONXTEST  
PROBPOS  
PROBPOS PROBPOSMIN PROBPOSMAX PROBPOSMIN  
FRACTIONOFPOSITIVES MEANPREDICTEDVALUE  
CALIBRATIONCURVEYTEST PROBPOS NBINS10  
AX1PLOTMEANPREDICTEDVALUE FRACTIONOFPOSITIVES S  
LABELS NAME  
AX2HISTPROBPOS RANGE0 1 BINS10 LABELNAME  
HISTTYPESTEP LW2  
AX1SETYLABELFRACTION OF POSITIVES  
AX1SETYLIM005 105  
AX1LEGENDLOCLOWER RIGHT  
AX1SETTITLECALIBRATION PLOTS RELIABILITY CURVE  
AX2SETXLABELMEAN PREDICTED VALUE  
AX2SETYLABELCOUNT  
AX2LEGENDLOCUPPER CENTER NCOL2  
PLTTIGHTLAYOUT  
PLTSHOW  
790 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 1028 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
542 PROBABILITY CALIBRATION CURVES  
WHEN PERFORMING CLASSIFICATION ONE OFTEN WANTS TO PREDICT NOT ONLY THE CLASS LABEL BUT ALSO THE ASSOCIATED PROBABILITY  
THIS PROBABILITY GIVES SOME KIND OF CONFIDENCE ON THE PREDICTION THIS EXAMPLE DEMONSTRATES HOW TO DISPLAY HOW WELL  
CALIBRATED THE PREDICTED PROBABILITIES ARE AND HOW TO CALIBRATE AN UNCALIBRATED CLASSIFIER  
THE EXPERIMENT IS PERFORMED ON AN ARTIFICIAL DATASET FOR BINARY CLASSIFICATION WITH 100000 SAMPLES 1000 OF THEM ARE  
USED FOR MODEL FITTING WITH 20 FEATURES OF THE 20 FEATURES ONLY 2 ARE INFORMATIVE AND 10 ARE REDUNDANT THE FIRST  
FIGURE SHOWS THE ESTIMATED PROBABILITIES OBTAINED WITH LOGISTIC REGRESSION GAUSSIAN NAIVE BAYES AND GAUSSIAN NAIVE  
BAYES WITH BOTH ISOTONIC CALIBRATION AND SIGMOID CALIBRATION THE CALIBRATION PERFORMANCE IS EVALUATED WITH BRIER  
SCORE REPORTED IN THE LEGEND THE SMALLER THE BETTER ONE CAN OBSERVE HERE THAT LOGISTIC REGRESSION IS WELL CALIBRATED  
WHILE RAW GAUSSIAN NAIVE BAYES PERFORMS VERY BADLY THIS IS BECAUSE OF THE REDUNDANT FEATURES WHICH VIOLATE THE  
ASSUMPTION OF FEATUREINDEPENDENCE AND RESULT IN AN OVERLY CONFIDENT CLASSIFIER WHICH IS INDICATED BY THE TYPICAL  
TRANPOSEDSIGMOID CURVE  
CALIBRATION OF THE PROBABILITIES OF GAUSSIAN NAIVE BAYES WITH ISOTONIC REGRESSION CAN FIX THIS ISSUE AS CAN BE SEEN FROM  
THE NEARLY DIAGONAL CALIBRATION CURVE SIGMOID CALIBRATION ALSO IMPROVES THE BRIER SCORE SLIGHTLY ALBEIT NOT AS STRONGLY  
AS THE NONPARAMETRIC ISOTONIC REGRESSION THIS CAN BE ATTRIBUTED TO THE FACT THAT WE HAVE PLENTY OF CALIBRATION DATA SUCH  
THAT THE GREATER FLEXIBILITY OF THE NONPARAMETRIC MODEL CAN BE EXPLOITED  
THE SECOND FIGURE SHOWS THE CALIBRATION CURVE OF A LINEAR SUPPORTVECTOR CLASSIFIER LINEARSVC LINEARSVC SHOWS  
THE OPPOSITE BEHAVIOR AS GAUSSIAN NAIVE BAYES THE CALIBRATION CURVE HAS A SIGMOID CURVE WHICH IS TYPICAL FOR AN  
UNDERCONFIDENT CLASSIFIER IN THE CASE OF LINEARSVC THIS IS CAUSED BY THE MARGIN PROPERTY OF THE HINGE LOSS WHICH  
LETS THE MODEL FOCUS ON HARD SAMPLES THAT ARE CLOSE TO THE DECISION BOUNDARY THE SUPPORT VECTORS  
BOTH KINDS OF CALIBRATION CAN FIX THIS ISSUE AND YIELD NEARLY IDENTICAL RESULTS THIS SHOWS THAT SIGMOID CALIBRATION CAN  
DEAL WITH SITUATIONS WHERE THE CALIBRATION CURVE OF THE BASE CLASSIFIER IS SIGMOID EG FOR LINEARSVC BUT NOT WHERE IT  
IS TRANPOSEDSIGMOID EG GAUSSIAN NAIVE BAYES  
54 CALIBRATION 791



SCIKITLEARN USER GUIDE RELEASE 0213

- OUT  
LOGISTIC  
BRIER 0099  
PRECISION 0872  
RECALL 0851  
F1 0862  
NAIVE BAYES  
BRIER 0118  
PRECISION 0857  
RECALL 0876  
F1 0867  
NAIVE BAYES ISOTONIC  
BRIER 0098  
PRECISION 0883  
54 CALIBRATION 793

SCIKITLEARN USER GUIDE RELEASE 0213  
RECALL 0836  
F1 0859  
NAIVE BAYES SIGMOID  
BRIER 0109  
PRECISION 0861  
RECALL 0871  
F1 0866  
LOGISTIC  
BRIER 0099  
PRECISION 0872  
RECALL 0851  
F1 0862  
SVC  
BRIER 0163  
PRECISION 0872  
RECALL 0852  
F1 0862  
SVC ISOTONIC  
BRIER 0100  
PRECISION 0853  
RECALL 0878  
F1 0865  
SVC SIGMOID  
BRIER 0099  
PRECISION 0874  
RECALL 0849  
F1 0861  
PRINTDOC  
AUTHOR ALEXANDRE GRAMFORT ALEXANDREGAMFORTTELECOMPARISTECHFR  
JAN HENDRIK METZEN JHMINFORMATIKUNIBREMENDE  
LICENSE BSD STYLE  
IMPORT MATPLOTLIBPYPLOT AS PLT  
FROM SKLEARN IMPORT DATASETS  
FROM SKLEARNNAIVEBAYES IMPORT GAUSSIANNB  
FROM SKLEARN SVM IMPORT LINEARSVC  
FROM SKLEARN LINEAR MODEL IMPORT LOGISTIC REGRESSION  
FROM SKLEARN METRICS IMPORT BRIER SCORE LOSS PRECISION SCORE RECALL SCORE  
F1 SCORE  
FROM SKLEARN CALIBRATION IMPORT CALIBRATED CLASSIFIER CV CALIBRATION CURVE  
FROM SKLEARN MODEL SELECTION IMPORT TRAIN TESTS SPLIT  
CREATE DATASET OF CLASSIFICATION TASK WITH MANY REDUNDANT AND FEW  
INFORMATIVE FEATURES  
794 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213  
X Y DATASETSMAKECLASSIFICATIONNSAMPLES100000 NFEATURES20  
NINFORMATIVE2 NREDUNDANT10  
RANDOMSTATE42  
XTRAIN XTEST YTRAIN YTEST TRAINTESTSPLITX Y TESTSIZE099  
RANDOMSTATE42  
DEFPLOTCALIBRATIONCURVEEST NAME FIGINDEX  
PLOT CALIBRATION CURVE FOR EST WO AND WITH CALIBRATION  
CALIBRATED WITH ISOTONIC CALIBRATION  
ISOTONIC CALIBRATEDCLASSIFIERCVEST CV2 METHODISOTONIC  
CALIBRATED WITH SIGMOID CALIBRATION  
SIGMOID CALIBRATEDCLASSIFIERCVEST CV2 METHODSIGMOID  
LOGISTIC REGRESSION WITH NO CALIBRATION AS BASELINE  
LR LOGISTICREGRESSIONC1 SOLVERLBFGS  
FIG PLTFIGUREFIGINDEX FIGSIZE10 10  
AX1 PLTSUBPLOT2GRID3 1 0 0 ROWSPAN2  
AX2 PLTSUBPLOT2GRID3 1 2 0  
AX1PLOT0 1 0 1 K LABELPERFECTLY CALIBRATED  
FORCLF NAME INLR LOGISTIC  
EST NAME  
ISOTONIC NAME ISOTONIC  
SIGMOID NAME SIGMOID  
CLFFITXTRAIN YTRAIN  
YPRED CLFPREDICTXTEST  
IFHASATTRCLF PREDICTPROBA  
PROBPOS CLFPREDICTPROBAXTEST 1  
ELSE USE DECISION FUNCTION  
PROBPOS CLFDECISIONFUNCTIONXTEST  
PROBPOS  
PROBPOS PROBPOSMIN PROBPOS MAX PROBPOS MIN  
CLFScore BRIERSCORELOSSYTEST PROBPOS POSLABELYMAX  
PRINTS NAME  
PRINTTBRIER13F CLFScore  
PRINTTPRECISSION 13F PRECISIONSCOREYTEST YPRED  
PRINTTRECALL 13F RECALLSCOREYTEST YPRED  
PRINTTF113FN F1SCOREYTEST YPRED  
FRACTIONOFPOSITIVES MEANPREDICTEDVALUE  
CALIBRATIONCURVEYTEST PROBPOS NBINS10  
AX1PLOTMEANPREDICTEDVALUE FRACTIONOFPOSITIVES S  
LABELS13F NAME CLFScore  
AX2HISTPROBPOS RANGE0 1 BINS10 LABELNAME  
HISTTYPESTEP LW2  
AX1SETYLABELFRACTION OF POSITIVES  
AX1SETYLIM005 105  
AX1LEGENDLOCLOWER RIGHT  
AX1SETTITLECALIBRATION PLOTS RELIABILITY CURVE  
54 CALIBRATION 795

SCIKITLEARN USER GUIDE RELEASE 0213  
AX2SETXLABELMEAN PREDICTED VALUE  
AX2SETYLABELCOUNT  
AX2LEGENDLOCUPPER CENTER NCOL2  
PLTTIGHTLAYOUT  
PLOT CALIBRATION CURVE FOR GAUSSIAN NAIVE BAYES  
PLOT CALIBRATION CURVE GAUSSIAN NB NAIVE BAYES 1  
PLOT CALIBRATION CURVE FOR LINEAR SVC  
PLOT CALIBRATION CURVE LINEAR SVC MAXITER 10000 SVC 2  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 1993 SECONDS  
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE  
543 PROBABILITY CALIBRATION OF CLASSIFIERS  
WHEN PERFORMING CLASSIFICATION YOU OFTEN WANT TO PREDICT NOT ONLY THE CLASS LABEL BUT ALSO THE ASSOCIATED PROBABILITY  
THIS PROBABILITY GIVES YOU SOME KIND OF CONFIDENCE ON THE PREDICTION HOWEVER NOT ALL CLASSIFIERS PROVIDE WELL  
CALIBRATED PROBABILITIES SOME BEING OVERCONFIDENT WHILE OTHERS BEING UNDERCONFIDENT THUS A SEPARATE CALIBRATION  
OF PREDICTED PROBABILITIES IS OFTEN DESIRABLE AS A POSTPROCESSING THIS EXAMPLE ILLUSTRATES TWO DIFFERENT METHODS FOR  
THIS CALIBRATION AND EVALUATES THE QUALITY OF THE RETURNED PROBABILITIES USING BRIER'S SCORE SEE [HTTPS://EN.WIKIPEDIA.ORG](https://en.wikipedia.org/wiki/Brier_score)  
WIKI/BRIERSCORE  
COMPARED ARE THE ESTIMATED PROBABILITY USING A GAUSSIAN NAIVE BAYES CLASSIFIER WITHOUT CALIBRATION WITH A SIGMOID  
CALIBRATION AND WITH A NONPARAMETRIC ISOTONIC CALIBRATION ONE CAN OBSERVE THAT ONLY THE NONPARAMETRIC MODEL IS  
ABLE TO PROVIDE A PROBABILITY CALIBRATION THAT RETURNS PROBABILITIES CLOSE TO THE EXPECTED 0.5 FOR MOST OF THE SAMPLES  
BELONGING TO THE MIDDLE CLUSTER WITH HETEROGENEOUS LABELS THIS RESULTS IN A SIGNIFICANTLY IMPROVED BRIER SCORE  
796 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213

•

OUT  
BRIER SCORES THE SMALLER THE BETTER  
NO CALIBRATION 0104  
WITH ISOTONIC CALIBRATION 0084  
WITH SIGMOID CALIBRATION 0109  
PRINTDOC  
AUTHOR MATHIEU BLONDEL MATHIEUMBLONDELORG  
ALEXANDRE GRAMFORT ALEXANDREGRAMFORTTELECOMPARISTECHFR  
BALAZS KEGL BALAZSKEGLGMAILCOM  
JAN HENDRIK METZEN JHMINFORMATIKUNIBREMENDE  
LICENSE BSD STYLE  
IMPORT NUMPY AS NP  
IMPORT MATPLOTLIBPYPLOT AS PLT  
FROM MATPLOTLIB IMPORT CM  
FROM SKLEARNDATASETS IMPORT MAKEBLOBS  
FROM SKLEARNNAIVEBAYES IMPORT GAUSSIANNB  
FROM SKLEARNMETRICS IMPORT BRIERSCORELOSS  
798 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
FROM SKLEARNCALIBRATION IMPORT CALIBRATEDCLASSIFIERCV
FROM SKLEARNMODELSELECTION IMPORT TRAINTESTSPLIT
NSAMPLES 50000
NBINS 3 USE 3 BINS FOR CALIBRATIONCURVE AS WE HAVE 3 CLUSTERS HERE
GENERATE 3 BLOBS WITH 2 CLASSES WHERE THE SECOND BLOB CONTAINS
HALF POSITIVE SAMPLES AND HALF NEGATIVE SAMPLES PROBABILITY IN THIS
BLOB IS THEREFORE 0.5
CENTERS 5 5 0 0 5 5
X Y MAKEBLOBSNSAMPLESNSAMPLES NFEATURES2 CLUSTERSTD10
CENTERSCENTERS SHUFFLEFALSE RANDOMSTATE42
YNSAMPLES 2 0
YNSAMPLES 2 1
SAMPLEWEIGHT NPRANDOMRANDOMSTATE42RANDYSHAPE0
SPLIT TRAIN TEST FOR CALIBRATION
XTRAIN XTEST YTRAIN YTEST SWTRAIN SWTEST
TRAINTESTSPLITX Y SAMPLEWEIGHT TESTSIZE0.9 RANDOMSTATE42
GAUSSIAN NAIVEBAYES WITH NO CALIBRATION
CLF GAUSSIANNB
CLFFITXTRAIN YTRAIN GAUSSIANNB ITSELF DOES NOT SUPPORT SAMPLEWEIGHTS
PROBPOSCLF CLFPREDICTPROBAXTEST 1
GAUSSIAN NAIVEBAYES WITH ISOTONIC CALIBRATION
CLFISOTONIC CALIBRATEDCLASSIFIERCVCLF CV2 METHODISOTONIC
CLFISOTONICFITXTRAIN YTRAIN SWTRAIN
PROBPOSISOTONIC CLFISOTONICPREDICTPROBAXTEST 1
GAUSSIAN NAIVEBAYES WITH SIGMOID CALIBRATION
CLFSIGMOID CALIBRATEDCLASSIFIERCVCLF CV2 METHODSIGMOID
CLFSIGMOIDFITXTRAIN YTRAIN SWTRAIN
PROBPOSSIGMOID CLFSIGMOIDPREDICTPROBAXTEST 1
PRINTBRIER SCORES THE SMALLER THE BETTER
CLFScore BRIERSCORELOSSYTEST PROBPOSCLF SWTEST
PRINTNO CALIBRATION 13F CLFScore
CLFISOTONICSCORE BRIERSCORELOSSYTEST PROBPOSISOTONIC SWTEST
PRINTWITH ISOTONIC CALIBRATION 13F CLFISOTONICSCORE
CLFSIGMOIDSCORE BRIERSCORELOSSYTEST PROBPOSSIGMOID SWTEST
PRINTWITH SIGMOID CALIBRATION 13F CLFSIGMOIDSCORE

PLOT THE DATA AND THE PREDICTED PROBABILITIES
PLTFigure
YUNIQUE NPUNIQUEY
COLORS CMRAINBOWNPLINSPACE0.0 10 YUNIQUESIZE
FORTHISY COLOR INZIPYUNIQUE COLORS
THISX XTRAINYTRAIN THISY
THISSW SWTRAINYTRAIN THISY
PLTSCATTERTHISX 0 THISX 1 STHISSW 50
CCOLORNPNEWAXIS
54 CALIBRATION 799
```

SCIKITLEARN USER GUIDE RELEASE 0213

ALPHA05 EDGECOLORK

LABELCLASS S THISY

PLTLEGENDLOCBEST

PLTTITLEDATA

PLTFigure

ORDER NPLEXSORTPROBPOSCLF

PLTPLOTPROBPOSCLFORDER R LABELNO CALIBRATION 13F CLFScore

PLTPLOTPROBPOSISOTONICORDER G LINEWIDTH3

LABELISOTONIC CALIBRATION 13F CLFISOTONICScore

PLTPLOTPROBPOSSIGMOIDORDER B LINEWIDTH3

LABELSIGMOID CALIBRATION 13F CLFSIGMOIDScore

PLTPLOTNPLINSPACE0 YTESTSIZE 5112

YTESTORDERRESHAPE25 1MEAN1

K LINEWIDTH3 LABELREMPirical

PLTYLIM005 105

PLTXLABELINSTANCES SORTED ACCORDING TO PREDICTED PROBABILITY

UNCALIBRATED GNB

PLTYLABELPY1

PLTLEGENDLOCUPPER LEFT

PLTTITLEGAUSSIAN NAIVE BAYES PROBABILITIES

PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0108 SECONDS

NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)

544 PROBABILITY CALIBRATION FOR 3CLASS CLASSIFICATION

THIS EXAMPLE ILLUSTRATES HOW SIGMOID CALIBRATION CHANGES PREDICTED PROBABILITIES FOR A 3CLASS CLASSIFICATION PROBLEM

ILLUSTRATED IS THE STANDARD 2SIMPLEX WHERE THE THREE CORNERS CORRESPOND TO THE THREE CLASSES ARROWS POINT FROM THE

PROBABILITY VECTORS PREDICTED BY AN UNCALIBRATED CLASSIFIER TO THE PROBABILITY VECTORS PREDICTED BY THE SAME CLASSIFIER

AFTER SIGMOID CALIBRATION ON A HOLDOUT VALIDATION SET COLORS INDICATE THE TRUE CLASS OF AN INSTANCE RED CLASS 1 GREEN

CLASS 2 BLUE CLASS 3

THE BASE CLASSIFIER IS A RANDOM FOREST CLASSIFIER WITH 25 BASE ESTIMATORS TREES IF THIS CLASSIFIER IS TRAINED ON ALL 800

TRAINING DATAPOINTS IT IS OVERLY CONFIDENT IN ITS PREDICTIONS AND THUS INCURS A LARGE LOGLOSS CALIBRATING AN IDENTICAL

CLASSIFIER WHICH WAS TRAINED ON 600 DATAPOINTS WITH METHOD'SIGMOID' ON THE REMAINING 200 DATAPOINTS REDUCES THE

CONFIDENCE OF THE PREDICTIONS IE MOVES THE PROBABILITY VECTORS FROM THE EDGES OF THE SIMPLEX TOWARDS THE CENTER

THIS CALIBRATION RESULTS IN A LOWER LOGLOSS NOTE THAT AN ALTERNATIVE WOULD HAVE BEEN TO INCREASE THE NUMBER OF BASE

ESTIMATORS WHICH WOULD HAVE RESULTED IN A SIMILAR DECREASE IN LOGLOSS

800 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213

- 

OUT  
LOGLOSS OF  
UNCALIBRATED CLASSIFIER TRAINED ON 800 DATAPOINTS 1280  
CLASSIFIER TRAINED ON 600 DATAPOINTS AND CALIBRATED ON 200 DATAPOINT 0534  
PRINTDOC  
AUTHOR JAN HENDRIK METZEN JHMFORMATIKUNIBREMENDE  
LICENSE BSD STYLE  
IMPORT MATPLOTLIBPYPLOT AS PLT  
IMPORT NUMPY AS NP  
FROM SKLEARNDATASETS IMPORT MAKEBLOBS  
FROM SKLEARNENSEMBLE IMPORT RANDOMFORESTCLASSIFIER  
FROM SKLEARNCALIBRATION IMPORT CALIBRATEDCLASSIFIERCV  
FROM SKLEARNMETRICS IMPORT LOGLOSS  
NPRANDOMSEED0  
802 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213  
GENERATE DATA  
X Y MAKEBLOBSNSAMPLES1000 NFEATURES2 RANDOMSTATE42  
CLUSTERSTD50  
XTRAIN YTRAIN X600 Y600  
XVALID YVALID X600800 Y600800  
XTRAININVALID YTRAININVALID X800 Y800  
XTEST YTEST X800 Y800  
TRAIN UNCALIBRATED RANDOM FOREST CLASSIFIER ON WHOLE TRAIN AND VALIDATION  
DATA AND EVALUATE ON TEST DATA  
CLF RANDOMFORESTCLASSIFIERNESTIMATORS25  
CLFFITXTRAININVALID YTRAININVALID  
CLFPROBS CLFPREDICTPROBAXTEST  
SCORE LOGLOSSYTEST CLFPROBS  
TRAIN RANDOM FOREST CLASSIFIER CALIBRATE ON VALIDATION DATA AND EVALUATE  
ON TEST DATA  
CLF RANDOMFORESTCLASSIFIERNESTIMATORS25  
CLFFITXTRAIN YTRAIN  
CLFPROBS CLFPREDICTPROBAXTEST  
SIGCLF CALIBRATEDCLASSIFIERCVCLF METHODSIGMOID CVPREFIT  
SIGCLFFITXVALID YVALID  
SIGCLFPROBS SIGCLFPREDICTPROBAXTEST  
SIGSCORE LOGLOSSYTEST SIGCLFPROBS  
PLOT CHANGES IN PREDICTED PROBABILITIES VIA ARROWS  
PLTFigure  
COLORS R G B  
FORIIRANGECLFPROBSSHAPE0  
PLTARROWCLFPROBSI 0 CLFPROBSI 1  
SIGCLFPROBSI 0 CLFPROBSI 0  
SIGCLFPROBSI 1 CLFPROBSI 1  
COLORCOLORSYTESTI HEADWIDTH1E2  
PLOT PERFECT PREDICTIONS  
PLTPLOT10 00 RO MS20 LABELCLASS 1  
PLTPLOT00 10 GO MS20 LABELCLASS 2  
PLTPLOT00 00 BO MS20 LABELCLASS 3  
PLOT BOUNDARIES OF UNIT SIMPLEX  
PLTPLOT00 10 00 00 00 00 10 00 K LABELSIMPLEX  
ANNOTATE POINTS ON THE SIMPLEX  
PLTANNOTATERFRAC13 FRAC13 FRAC13  
XY103 103 XYTEXT103 23 XYCOORDSDATA  
ARROWPROPSDICTFACECOLORBLACK SHRINK005  
HORIZONTALALIGNMENTCENTER VERTICALALIGNMENTCENTER  
PLTPLOT103 103 KO MS5  
PLTANNOTATERFRAC12 0 FRAC12  
XY5 0 XYTEXT5 1 XYCOORDSDATA  
ARROWPROPSDICTFACECOLORBLACK SHRINK005  
HORIZONTALALIGNMENTCENTER VERTICALALIGNMENTCENTER  
PLTANNOTATER0 FRAC12 FRAC12  
XY0 5 XYTEXT1 5 XYCOORDSDATA  
ARROWPROPSDICTFACECOLORBLACK SHRINK005  
HORIZONTALALIGNMENTCENTER VERTICALALIGNMENTCENTER  
PLTANNOTATERFRAC12 FRAC12 0  
54 CALIBRATION 803

SCIKITLEARN USER GUIDE RELEASE 0213  
XY5 5 XYTEXT6 6 XYCOORDSDATA  
ARROWPROPSDICTFACECOLORBLACK SHRINK005  
HORIZONTALALIGNMENTCENTER VERTICALALIGNMENTCENTER  
PLTANNOTATER0 0 1  
XY0 0 XYTEXT1 1 XYCOORDSDATA  
ARROWPROPSDICTFACECOLORBLACK SHRINK005  
HORIZONTALALIGNMENTCENTER VERTICALALIGNMENTCENTER  
PLTANNOTATER1 0 0  
XY1 0 XYTEXT1 1 XYCOORDSDATA  
ARROWPROPSDICTFACECOLORBLACK SHRINK005  
HORIZONTALALIGNMENTCENTER VERTICALALIGNMENTCENTER  
PLTANNOTATER0 1 0  
XY0 1 XYTEXT1 1 XYCOORDSDATA  
ARROWPROPSDICTFACECOLORBLACK SHRINK005  
HORIZONTALALIGNMENTCENTER VERTICALALIGNMENTCENTER  
ADD GRID  
PLTGRIDFALSE  
FORXIN00 01 02 03 04 05 06 07 08 09 10  
PLTPLOT0 X X 0 K ALPHA02  
PLTPLOT0 0 1X2 X X 1X2 K ALPHA02  
PLTPLOTX X 1X2 0 0 1X2 K ALPHA02  
PLTTITLECHANGE OF PREDICTED PROBABILITIES AFTER SIGMOID CALIBRATION  
PLTXLABELPROBABILITY CLASS 1  
PLTYLABELPROBABILITY CLASS 2  
PLTXLIM005 105  
PLTYLIM005 105  
PLTLEGENDLOCBEST  
PRINTLOGLOSS OF  
PRINTUNCALIBRATED CLASSIFIER TRAINED ON 800 DATAPOINTS 3F  
SCORE  
PRINTCLASSIFIER TRAINED ON 600 DATAPOINTS AND CALIBRATED ON  
200 DATAPOINT 3F SIGSCORE  
ILLUSTRATE CALIBRATOR  
PLTFigure  
GENERATE GRID OVER 2SIMPLEX  
P1D NPLINSPACE0 1 20  
P0 P1 NPMESHGRIDP1D P1D  
P2 1 P0 P1  
P NPCP0RAVEL P1RAVEL P2RAVEL  
P PP 2 0  
CALIBRATEDCLASSIFIER SIGCLFCALIBRATEDCLASSIFIERS0  
PREDICTION NPVSTACKCALIBRATORPREDICTTHISP  
FORCALIBRATOR THISP IN  
ZIPCALIBRATEDCLASSIFIERCALIBRATORS PTT  
PREDICTION PREDICTIONSUMAXIS1 NONE  
PLOT MODIFICATIONS OF CALIBRATOR  
FORIINRANGEPREDICTIONSHAPE0  
PLTARROWPI 0 PI 1  
PREDICTIONI 0 PI 0 PREDICTIONI 1 PI 1  
HEADWIDTH1E2 COLORCOLORSNPARGMAXPI  
PLOT BOUNDARIES OF UNIT SIMPLEX  
PLTPLOT00 10 00 00 00 00 10 00 K LABELSIMPLEX  
804 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
PLTGRIDFALSE  
FORXIN00 01 02 03 04 05 06 07 08 09 10  
PLTPLOT0 X X 0 K ALPHA02  
PLTPLOT0 0 1X2 X X 1X2 K ALPHA02  
PLTPLOTX X 1X2 0 0 1X2 K ALPHA02  
PLTTITLEILLUSTRATION OF SIGMOID CALIBRATOR  
PLTXLABELPROBABILITY CLASS 1  
PLTYLABELPROBABILITY CLASS 2  
PLTXLIM005 105  
PLTYLIM005 105  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0318 SECONDS  
55 CLASSIFICATION  
GENERAL EXAMPLES ABOUT CLASSIFICATION ALGORITHMS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
551 RECOGNIZING HANDWRITTEN DIGITS  
AN EXAMPLE SHOWING HOW THE SCIKITLEARN CAN BE USED TO RECOGNIZE IMAGES OF HANDWRITTEN DIGITS  
THIS EXAMPLE IS COMMENTED IN THE TUTORIAL SECTION OF THE USER MANUAL  
55 CLASSIFICATION 805

SCIKITLEARN USER GUIDE RELEASE 0213  
 OUT  
 CLASSIFICATION REPORT FOR CLASSIFIER SVCGAMMA0001  
 PRECISION RECALL F1SCORE SUPPORT  
 0 100 099 099 88  
 1 099 097 098 91  
 2 099 099 099 86  
 3 098 087 092 91  
 4 099 096 097 92  
 5 095 097 096 91  
 6 099 099 099 91  
 7 096 099 097 89  
 8 094 100 097 88  
 9 093 098 095 92  
 ACCURACY 097 899  
 MACRO AVG 097 097 097 899  
 WEIGHTED AVG 097 097 097 899  
 CONFUSION MATRIX  
 87 0 0 0 1 0 0 0 0 0  
 0 88 1 0 0 0 0 0 1 1  
 0 0 85 1 0 0 0 0 0 0  
 0 0 0 79 0 3 0 4 5 0  
 806 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
0 0 0 0 88 0 0 0 0 4
0 0 0 0 0 88 1 0 0 2
0 1 0 0 0 0 90 0 0 0
0 0 0 0 0 1 0 88 0 0
0 0 0 0 0 0 0 0 88 0
0 0 0 1 0 1 0 0 0 90
PRINTDOC
AUTHOR GAELEVAROQUAUX GAELE DOT VAROQUAUX AT NORMALESUP DOT ORG
LICENSE BSD 3 CLAUSE
STANDARD SCIENTIFIC PYTHON IMPORTS
IMPORT MATPLOTLIBPYPLOT AS PLT
IMPORT DATASETS CLASSIFIERS AND PERFORMANCE METRICS
FROM SKLEARN IMPORT DATASETS SVM METRICS
THE DIGITS DATASET
DIGITS DATASETSLOADDIGITS
THE DATA THAT WE ARE INTERESTED IN IS MADE OF 8X8 IMAGES OF DIGITS LETS
HAVE A LOOK AT THE FIRST 4 IMAGES STORED IN THE IMAGES ATTRIBUTE OF THE
DATASET IF WE WERE WORKING FROM IMAGE FILES WE COULD LOAD THEM USING
MATPLOTLIBPYPLOTIMREAD NOTE THAT EACH IMAGE MUST HAVE THE SAME SIZE FOR THESE
IMAGES WE KNOW WHICH DIGIT THEY REPRESENT IT IS GIVEN IN THE TARGET OF
THE DATASET
IMAGESANDLABELS LISTZIPDIGITSIMAGES DIGITSTARGET
FORINDEX IMAGE LABEL INENUMERATEIMAGESANDLABELS4
PLTSUBPLOT2 4 INDEX 1
PLTAXISOFF
PLTIMSHOWIMAGE CMAPPLTCMGRAYR INTERPOLATIONNEAREST
PLTTITLETRAINING I LABEL
TO APPLY A CLASSIFIER ON THIS DATA WE NEED TO FLATTEN THE IMAGE TO
TURN THE DATA IN A SAMPLES FEATURE MATRIX
NSAMPLES LENDIGITSIMAGES
DATA DIGITSIMAGESRESHAPENSAMPLES 1
CREATE A CLASSIFIER A SUPPORT VECTOR CLASSIFIER
CLASSIFIER SVMSCVCGAMMA0001
WE LEARN THE DIGITS ON THE FIRST HALF OF THE DIGITS
CLASSIFIERFITDATANSAMPLES 2 DIGITSTARGETNSAMPLES 2
NOW PREDICT THE VALUE OF THE DIGIT ON THE SECOND HALF
EXPECTED DIGITSTARGETNSAMPLES 2
PREDICTED CLASSIFIERPREDICTDATANSAMPLES 2
PRINTCLASSIFICATION REPORT FOR CLASSIFIER SNSN
CLASSIFIER METRICSClassificationReportEXPECTED PREDICTED
PRINTCONFUSION MATRIX NS METRICSCONFUSIONMATRIXEXPECTED PREDICTED
55 CLASSIFICATION 807
```

SCIKITLEARN USER GUIDE RELEASE 0213  
IMAGESANDPREDICTIONS LISTZIPDIGITSIMAGESNSAMPLES 2 PREDICTED  
FORINDEX IMAGE PREDICTION INENUMERATEIMAGESANDPREDICTIONS4  
PLTSUBPLOT2 4 INDEX 5  
PLTAXISOFF  
PLTIMSHOWIMAGE CMAPPLTCMGRAYR INTERPOLATIONNEAREST  
PLTTITLEPREDICTION I PREDICTION  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0237 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
552 NORMAL AND SHRINKAGE LINEAR DISCRIMINANT ANALYSIS FOR CLASSIFICATION  
SHOWS HOW SHRINKAGE IMPROVES CLASSIFICATION  
IMPORT NUMPY AS NP  
IMPORT MATPLOTLIBPYPLOT AS PLT  
FROM SKLEARNDATASETS IMPORT MAKEBLOBS  
FROM SKLEARNDISCRIMINANTANALYSIS IMPORT LINEARDISCRIMINANTANALYSIS  
808 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
NTRAIN 20 SAMPLES FOR TRAINING  
NTEST 200 SAMPLES FOR TESTING  
NAVERAGES 50 HOW OFTEN TO REPEAT CLASSIFICATION  
NFEATURESMAX 75 MAXIMUM NUMBER OF FEATURES  
STEP 4 STEP SIZE FOR THE CALCULATION  
DEFGENERATEDATANSAMPLES NFEATURES  
GENERATE RANDOM BLOBISH DATA WITH NOISY FEATURES  
THIS RETURNS AN ARRAY OF INPUT DATA WITH SHAPE NSAMPLES NFEATURES  
AND AN ARRAY OF NSAMPLES TARGET LABELS  
ONLY ONE FEATURE CONTAINS DISCRIMINATIVE INFORMATION THE OTHER FEATURES  
CONTAIN ONLY NOISE

X Y MAKEBLOBSNSAMPLESNSAMPLES NFEATURES1 CENTERS2 2  
ADD NONDISCRIMINATIVE FEATURES  
IFNFEATURES 1  
X NPHSTACKX NPRANDOMRANDNNSAMPLES NFEATURES 1  
RETURNX Y  
ACCCLF1 ACCCLF2  
NFEATURESRANGE RANGE1 NFEATURESMAX 1 STEP  
FORNFEATURES INNFEATURESRANGE  
SCORECLF1 SCORECLF2 0 0  
FORINRANGENAVERAGES  
X Y GENERATEDATANTRAIN NFEATURES  
CLF1 LINEARDISCRIMINANTANALYSIS SOLVERLSQR SHRINKAGEAUTOFITX Y  
CLF2 LINEARDISCRIMINANTANALYSIS SOLVERLSQR SHRINKAGENONEFITX Y  
X Y GENERATEDATANTEST NFEATURES  
SCORECLF1 CLF1SCOREX Y  
SCORECLF2 CLF2SCOREX Y  
ACCCLF1APPENDSCORECLF1 NVERAGES  
ACCCLF2APPENDSCORECLF2 NVERAGES  
FEATURESSAMPLESRATIO NPARRAYNFEATURESRANGE NTRAIN  
PLTPLOTFEATURESSAMPLESRATIO ACCCLF1 LINEWIDTH2  
LABELLINEAR DISCRIMINANT ANALYSIS WITH SHRINKAGE COLORNAVY  
PLTPLOTFEATURESSAMPLESRATIO ACCCLF2 LINEWIDTH2  
LABELLINEAR DISCRIMINANT ANALYSIS COLORGOLD  
PLTXLABELNFEATURES NSAMPLES  
PLTYLABELCLASSIFICATION ACCURACY  
PLTLEGENDLOC1 PROPSIZE 12  
PLTSUPTITLELINEAR DISCRIMINANT ANALYSIS VS  
SHRINKAGE LINEAR DISCRIMINANT ANALYSIS 1 DISCRIMINATIVE FEATURE  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 6160 SECONDS  
55 CLASSIFICATION 809

SCIKITLEARN USER GUIDE RELEASE 0213

[NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)

553 PLOT CLASSIFICATION PROBABILITY

PLOT THE CLASSIFICATION PROBABILITY FOR DIFFERENT CLASSIFIERS WE USE A 3 CLASS DATASET AND WE CLASSIFY IT WITH A SUP  
PORT VECTOR CLASSIFIER L1 AND L2 PENALIZED LOGISTIC REGRESSION WITH EITHER A ONEVSREST OR MULTINOMIAL SETTING AND  
GAUSSIAN PROCESS CLASSIFICATION

LINEAR SVC IS NOT A PROBABILISTIC CLASSIFIER BY DEFAULT BUT IT HAS A BUILTIN CALIBRATION OPTION ENABLED IN THIS EXAMPLE  
PROBABILITYTRUE

THE LOGISTIC REGRESSION WITH ONEVSREST IS NOT A MULTICLASS CLASSIFIER OUT OF THE BOX AS A RESULT IT HAS MORE TROUBLE  
IN SEPARATING CLASS 2 AND 3 THAN THE OTHER ESTIMATORS

810 CHAPTER 5 EXAMPLES





```
SCIKITLEARN USER GUIDE RELEASE 0213
OUT
ACCURACY TRAIN FOR L1 LOGISTIC 833
ACCURACY TRAIN FOR L2 LOGISTIC MULTINOMIAL 827
ACCURACY TRAIN FOR L2 LOGISTIC OVR 793
ACCURACY TRAIN FOR LINEAR SVC 820
ACCURACY TRAIN FOR GPC 827
PRINTDOC
AUTHOR ALEXANDRE GRAMFORT ALEXANDREGRAMFORTINRIA.FR
LICENSE BSD 3 CLAUSE
IMPORT MATPLOTLIB.PY.PLOT AS PLT
IMPORT NUMPY AS NP
FROM SKLEARN.METRICS IMPORT ACCURACY_SCORE
FROM SKLEARN.LINEAR_MODEL IMPORT LOGISTIC_REGRESSION
FROM SKLEARN.SVM IMPORT SVC
FROM SKLEARN.GAUSSIAN_PROCESS IMPORT GAUSSIAN_PROCESS_CLASSIFIER
FROM SKLEARN.GAUSSIAN_PROCESS KERNELS IMPORT RBF
FROM SKLEARN IMPORT DATASETS
IRIS = DATASETS.LOAD_IRIS
X = IRIS.DATA[0:2] WE ONLY TAKE THE FIRST TWO FEATURES FOR VISUALIZATION
Y = IRIS.TARGET
N_FEATURES = X.SHAPE[1]
C = 10
KERNEL = 10 * RBF(10 * 10) FOR GPC
CREATE DIFFERENT CLASSIFIERS
CLASSIFIERS
L1 = LOGISTIC_REGRESSION(C = PENALTY_L1)
SOLVER_SAGA
MULTICLASS_MULTINOMIAL
MAX_ITER = 10000
L2 = LOGISTIC_MULTINOMIAL_REGRESSION(C = PENALTY_L2)
SOLVER_SAGA
MULTICLASS_MULTINOMIAL
MAX_ITER = 10000
L2 = LOGISTIC_OVR_REGRESSION(C = PENALTY_L2)
SOLVER_SAGA
MULTICLASS_OVR
MAX_ITER = 10000
LINEAR_SVC = SVC(KERNEL = 'linear', CC_PROBABILITY = True)
RANDOM_STATE = 0
GPC = GAUSSIAN_PROCESS_CLASSIFIER(KERNEL)

N_CLASSIFIERS = LEN(CLASSIFIERS)
812 CHAPTER 5 EXAMPLES
```

SCIKITLEARN USER GUIDE RELEASE 0213  
PLTFIGUREFIGSIZE3 2 NCLASSIFIERS 2  
PLTSUBPLOTSADJUSTBOTTOM2 TOP95  
XX NPLINSPACE3 9 100  
YY NPLINSPACE1 5 100T  
XX YY NPMESHGRIDXX YY  
XFULL NPCXXRAVEL YYRAVEL  
FORINDEX NAME CLASSIFIER INENUMERATECLASSIFIERSITEMS  
CLASSIFIERFITX Y  
YPRED CLASSIFIERPREDICTX  
ACCURACY ACCURACYSOREY YPRED  
PRINTACCURACY TRAIN FOR S01F NAME ACCURACY 100  
VIEW PROBABILITIES  
PROBAS CLASSIFIERPREDICTPROBAXFULL  
NCLASSES NPUNIQUEYPREDSIZE  
FORINRANGENCLASSES  
PLTSUBPLOTNCLASSIFIERS NCLASSES INDEX NCLASSES K 1  
PLTTITLECLASS D K  
IFK 0  
PLTYLABELNAME  
IMSHOWHANDLE PLTIMSHOWPROBAS KRESHAPE100 100  
EXTENT3 9 1 5 ORIGINLOWER  
PLXTTICKS  
PLTYTICKS  
IDX YPRED K  
IFIDXANY  
PLTSCATTERXIDX 0 XIDX 1 MARKERO CW EDGECOLORK  
AX PLTAXES015 004 07 005  
PLTTITLEPROBABILITY  
PLTCOLORBARIMSHOWHANDLE CAXAX ORIENTATIONHORIZONTAL  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 1409 SECONDS  
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE  
554 CLASSIFIER COMPARISON  
A COMPARISON OF A SEVERAL CLASSIFIERS IN SCIKITLEARN ON SYNTHETIC DATASETS THE POINT OF THIS EXAMPLE IS TO ILLUSTRATE  
THE NATURE OF DECISION BOUNDARIES OF DIFFERENT CLASSIFIERS THIS SHOULD BE TAKEN WITH A GRAIN OF SALT AS THE INTUITION  
CONVEYED BY THESE EXAMPLES DOES NOT NECESSARILY CARRY OVER TO REAL DATASETS  
PARTICULARLY IN HIGHDIMENSIONAL SPACES DATA CAN MORE EASILY BE SEPARATED LINEARLY AND THE SIMPLICITY OF CLASSIFIERS  
SUCH AS NAIVE BAYES AND LINEAR SVMs MIGHT LEAD TO BETTER GENERALIZATION THAN IS ACHIEVED BY OTHER CLASSIFIERS  
THE PLOTS SHOW TRAINING POINTS IN SOLID COLORS AND TESTING POINTS SEMITRANSSPARENT THE LOWER RIGHT SHOWS THE CLASSIFI  
CATION ACCURACY ON THE TEST SET  
55 CLASSIFICATION 813

SCIKITLEARN USER GUIDE RELEASE 0213  
PRINTDOC  
CODE SOURCE GAËL VAROQUAUX  
ANDREAS MÜLLER  
MODIFIED FOR DOCUMENTATION BY JAUQUES GROBLER  
LICENSE BSD 3 CLAUSE  
IMPORT NUMPY AS NP  
IMPORT MATPLOTLIBPYPLOT AS PLT  
FROM MATPLOTLIBCOLORS IMPORT LISTEDCOLORMAP  
FROM SKLEARNMODELSELECTION IMPORT TRAINTESTSPLIT  
FROM SKLEARNPREPROCESSING IMPORT STANDARDSCALER  
FROM SKLEARNDATASETS IMPORT MAKEMOONS MAKECIRCLES MAKECLASSIFICATION  
FROM SKLEARNNEURALNETWORK IMPORT MLPCLASSIFIER  
FROM SKLEARNNEIGHBORS IMPORT KNEIGHBORSCLASSIFIER  
FROM SKLEARN SVM IMPORT SVC  
FROM SKLEARN GAUSSIANPROCESS IMPORT GAUSSIANPROCESSCLASSIFIER  
FROM SKLEARN GAUSSIANPROCESS KERNELS IMPORT RBF  
FROM SKLEARN TREE IMPORT DECISIONTREECLASSIFIER  
FROM SKLEARN ENSEMBLE IMPORT RANDOMFORESTCLASSIFIER ADABOOSTCLASSIFIER  
FROM SKLEARN NAIVEBAYES IMPORT GAUSSIANNB  
FROM SKLEARN DISCRIMINANTANALYSIS IMPORT QUADRATICDISCRIMINANTANALYSIS  
H 02 STEP SIZE IN THE MESH  
NAMES NEAREST NEIGHBORS LINEAR SVM RBF SVM GAUSSIAN PROCESS  
DECISION TREE RANDOM FOREST NEURAL NET ADABOOST  
NAIVE BAYES QDA  
CLASSIFIERS  
KNEIGHBORSCLASSIFIER3  
SVCKERNELLINEAR C0025  
SVCGAMMA2 C1  
GAUSSIANPROCESSCLASSIFIER10 RBF10  
DECISIONTREECLASSIFIERMAXDEPTH5  
RANDOMFORESTCLASSIFIERMAXDEPTH5 NESTIMATORS10 MAXFEATURES1  
MLPCLASSIFIERALPHA1 MAXITER1000  
ADABOOSTCLASSIFIER  
GAUSSIANNB  
QUADRATICDISCRIMINANTANALYSIS  
X Y MAKECLASSIFICATIONNFEATURES2 NREDUNDANT0 NINFORMATIVE2  
814 CHAPTER 5 EXAMPLES

```

SCIKITLEARN USER GUIDE RELEASE 0213
RANDOMSTATE1 NCLUSTERSPERCLASS1
RNG NPRANDOMRANDOMSTATE2
X 2RNGUNIFORMSIZEXSHAPE
LINEARLYSEPARABLE X Y
DATASETS MAKEMOONSNOISE03 RANDOMSTATE0
MAKECIRCLESNOISE02 FACTOR05 RANDOMSTATE1
LINEARLYSEPARABLE

FIGURE PLTFIGUREFIGSIZE27 9
I 1
  ITERATE OVER DATASETS
  FORDSCNT DS INENUMERATEDDATASETS
  PREPROCESS DATASET SPLIT INTO TRAINING AND TEST PART
  X Y DS
  X STANDARDSCALERFITTRANSFORMX
  XTRAIN XTEST YTRAIN YTEST
  TRAINTESTSPLITX Y TESTSIZE4 RANDOMSTATE42
  XMIN XMAX X 0MIN 5 X 0MAX 5
  YMIN YMAX X 1MIN 5 X 1MAX 5
  XX YY NPMESHGRIDNPARANGEXMIN XMAX H
  NPARANGEYMIN YMAX H
  JUST PLOT THE DATASET FIRST
  CM PLTCMRDBU
  CMBRIGHT LISTEDCOLORMAPFF0000 0000FF
  AX PLTSUBPLOTLENDATASETS LENCLASSIFIERS 1 I
  IFDSCNT 0
  AXSETTITLEINPUT DATA
  PLOT THE TRAINING POINTS
  AXSCATTERXTRAIN 0 XTRAIN 1 CYTRAIN CMAPCMBRIGHT
  EDGECOLORSK
  PLOT THE TESTING POINTS
  AXSCATTERXTEST 0 XTEST 1 CYTEST CMAPCMBRIGHT ALPHA06
  EDGECOLORSK
  AXSETXLMXXMIN XXMAX
  AXSETYLIMYYMIN YYMAX
  AXSETXTICKS
  AXSETYTICKS
I 1
  ITERATE OVER CLASSIFIERS
  FORNAME CLF INZIPNAMES CLASSIFIERS
  AX PLTSUBPLOTLENDATASETS LENCLASSIFIERS 1 I
  CLFFITXTRAIN YTRAIN
  SCORE CLFSCOREXTEST YTEST
  PLOT THE DECISION BOUNDARY FOR THAT WE WILL ASSIGN A COLOR TO EACH
  POINT IN THE MESH XMIN XMAXXYMIN YMAX
  IFHASATTRCLF DECISIONFUNCTION
  Z CLFDECISIONFUNCTIONNPCXXRAVEL YYRAVEL
  ELSE
  Z CLFPREDICTPROBANPCXXRAVEL YYRAVEL 1
  PUT THE RESULT INTO A COLOR PLOT
  Z ZRESHAPEXXSHAPE
55 CLASSIFICATION 815
```

SCIKITLEARN USER GUIDE RELEASE 0213  
AXCONTOURFXX YY Z CMAPCM ALPHA8  
PLOT THE TRAINING POINTS  
AXSCATTERXTRAIN 0 XTRAIN 1 CYTRAIN CMAPCMBRIGHT  
EDGECOLORSK  
PLOT THE TESTING POINTS  
AXSCATTERXTEST 0 XTEST 1 CYTEST CMAPCMBRIGHT  
EDGECOLORSK ALPHA06  
AXSETXLIMXXMIN XXMAX  
AXSETYLIMYYMIN YYMAX  
AXSETXTICKS  
AXSETYTICKS  
IFDSCNT 0  
AXSETTITLENAME  
AXTEXTXXMAX 3 YYMIN 3 2F SCORELSTRIPO  
SIZE15 HORIZONTALALIGNMENTRIGHT  
I 1  
PLTTIGHTLAYOUT  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 5178 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
555 LINEAR AND QUADRATIC DISCRIMINANT ANALYSIS WITH COVARIANCE ELLIPSOID  
THIS EXAMPLE PLOTS THE COVARIANCE ELLIPSOIDS OF EACH CLASS AND DECISION BOUNDARY LEARNED BY LDA AND QDA THE  
ELLIPSOIDS DISPLAY THE DOUBLE STANDARD DEVIATION FOR EACH CLASS WITH LDA THE STANDARD DEVIATION IS THE SAME FOR ALL  
THE CLASSES WHILE EACH CLASS HAS ITS OWN STANDARD DEVIATION WITH QDA  
816 CHAPTER 5 EXAMPLES

```

SCIKITLEARN USER GUIDE RELEASE 0213
PRINTDOC
FROM SCIPY IMPORT LINALG
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
IMPORT MATPLOTLIB AS MPL
FROM MATPLOTLIB IMPORT COLORS
FROM SKLEARNDISCRIMINANTANALYSIS IMPORT LINEARDISCRIMINANTANALYSIS
FROM SKLEARNDISCRIMINANTANALYSIS IMPORT QUADRATICDISCRIMINANTANALYSIS

COLORMAP
CMAP COLORSLINEARSEGMENTEDCOLORMAP
REDBLUECLASSES
RED 0 1 1 1 07 07
GREEN 0 07 07 1 07 07
BLUE 0 07 07 1 1 1
PLTCMREGISTERCMAPCMAPCMAP

GENERATE DATASETS
55 CLASSIFICATION 817

```

```
SCIKITLEARN USER GUIDE RELEASE 0213
DEFDATASETFIXEDCOV
GENERATE 2 GAUSSIANS SAMPLES WITH THE SAME COVARIANCE MATRIX
N DIM 300 2
NPRANDOMSEED0
C NPARRAY0 023 083 23
X NPRNPDOTNPRANDOMRANDNN DIM C
NPDOTNPRANDOMRANDNN DIM C NPARRAY1 1
Y NPHSTACKNPZEROSN NPONESN
RETURNX Y
DEFDATASETCOV
GENERATE 2 GAUSSIANS SAMPLES WITH DIFFERENT COVARIANCE MATRICES
N DIM 300 2
NPRANDOMSEED0
C NPARRAY0 1 25 7 2
X NPRNPDOTNPRANDOMRANDNN DIM C
NPDOTNPRANDOMRANDNN DIM CT NPARRAY1 4
Y NPHSTACKNPZEROSN NPONESN
RETURNX Y
```

```
PLOT FUNCTIONS
DEFPLOTDATAALDA X Y YPRED FIGINDEX
SPLOT PLTSUBPLOT2 2 FIGINDEX
IFFIGINDEX 1
PLTTITLELINEAR DISCRIMINANT ANALYSIS
PLTYLABELDATA WITH NFIXED COVARIANCE
ELIFFIGINDEX 2
PLTTITLEQUADRATIC DISCRIMINANT ANALYSIS
ELIFFIGINDEX 3
PLTYLABELDATA WITH NVARYING COVARIANCES
TP Y YPRED TRUE POSITIVE
TP0 TP1 TPY 0 TPY 1
X0 X1 XY 0 XY 1
X0TP X0FP X0TP0 X0TP0
X1TP X1FP X1TP1 X1TP1
CLASS 0 DOTS
PLTSCATTERX0TP 0 X0TP 1 MARKER COLORRED
PLTSCATTERX0FP 0 X0FP 1 MARKERX
S20 COLOR990000 DARK RED
CLASS 1 DOTS
PLTSCATTERX1TP 0 X1TP 1 MARKER COLORBLUE
PLTSCATTERX1FP 0 X1FP 1 MARKERX
S20 COLOR000099 DARK BLUE
CLASS 0 AND 1 AREAS
NX NY 200 100
XMIN XMAX PLTXLIM
YMIN YMAX PLTYLIM
XX YY NPMESHGRIDNPLINSPACEXMIN XMAX NX
NPLINSPACEYMIN YMAX NY
Z LDAPREDICTPROBANPCXXRAVEL YYRAVEL
Z Z 1RESHAPEXXSHAPE
818 CHAPTER 5 EXAMPLES
```



SCIKITLEARN USER GUIDE RELEASE 0213  
PLTPCOLORMESHXX YY Z CMAPREDBLUECLASSES  
NORMCOLORSNORMALIZE0 1 ZORDER0  
PLTCONTOURXX YY Z 05 LINEWIDTHS2 COLORSWHITE  
MEANS  
PLTPLOTLDAMEANS00 LDAMEANS01  
COLORYELLOW MARKERSIZE15 MARKEREDGECOLORGREY  
PLTPLOTLDAMEANS10 LDAMEANS11  
COLORYELLOW MARKERSIZE15 MARKEREDGECOLORGREY  
RETURNSPLOT  
DEFPLOTELLIPSESLOT MEAN COV COLOR  
V W LINALGEIGHCOV  
U W0 LINALGNORMW0  
ANGLE NPARCTANU1 U0  
ANGLE 180 ANGLE NPPI CONVERT TO DEGREES  
FILLED GAUSSIAN AT 2 STANDARD DEVIATION  
ELL MPLPATCHESELLIPSEMEAN 2 V005 2V105  
180 ANGLE FACECOLORCOLOR  
EDGECOLORBLACK LINEWIDTH2  
ELLSETCLIPBOXSPLOTBBOX  
ELLSETALPHA02  
SPLOTADDARTISTELL  
SPLOTSETXTICKS  
SPLOTSETYTICKS  
DEFPLOTLDACOV LDA SPLOT  
PLOTELLIPSESLOT LDAMEANS0 LDACOVARIANCE RED  
PLOTELLIPSESLOT LDAMEANS1 LDACOVARIANCE BLUE  
DEFPLOTQDACOVQDA SPLOT  
PLOTELLIPSESLOT QDAMEANS0 QDACOVARIANCE0 RED  
PLOTELLIPSESLOT QDAMEANS1 QDACOVARIANCE1 BLUE  
PLTFIGUREFIGSIZE10 8 FACECOLORWHITE  
PLTSUPTITLELINEAR DISCRIMINANT ANALYSIS VS QUADRATIC DISCRIMINANT ANALYSIS  
Y098 FONTSIZE15  
FORI X Y INENUMERATEDDATASETFIXEDCOV DATASETCOV  
LINEAR DISCRIMINANT ANALYSIS  
LDA LINEARDISCRIMINANTANALYSIS SOLVERSVD STORECOVARIANCETRUE  
YPRED LDAFITX YPREDICTX  
SPLOT PLOTDATA LDA X Y YPRED FIGINDEX2 I 1  
PLOTLDACOV LDA SPLOT  
PLTAXISTIGHT  
QUADRATIC DISCRIMINANT ANALYSIS  
QDA QUADRATICDISCRIMINANTANALYSIS STORECOVARIANCETRUE  
YPRED QDAFITX YPREDICTX  
SPLOT PLOTDATA QDA X Y YPRED FIGINDEX2 I 2  
PLOTQDACOVQDA SPLOT  
PLTAXISTIGHT  
PLTTIGHTLAYOUT  
PLTSUBPLOTSADJUSTTOP092  
55 CLASSIFICATION 819

SCIKITLEARN USER GUIDE RELEASE 0213  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0231 SECONDS  
56 CLUSTERING  
EXAMPLES CONCERNING THE SKLEARNCLUSTER MODULE  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
561 FEATURE AGGLOMERATION  
THESE IMAGES HOW SIMILAR FEATURES ARE MERGED TOGETHER USING FEATURE AGGLOMERATION  
PRINTDOC  
CODE SOURCE GAËL VAROQUAUX  
MODIFIED FOR DOCUMENTATION BY JAQUES GROBLER  
LICENSE BSD 3 CLAUSE  
IMPORT NUMPY AS NP  
IMPORT MATPLOTLIBPYPLOT AS PLT  
FROM SKLEARN IMPORT DATASETS CLUSTER  
FROM SKLEARNFEATUREEXTRACTIONIMAGE IMPORT GRIDTOGRAPH  
DIGITS DATASETSLOADDIGITS  
IMAGES DIGITSIMAGES  
X NPRESHAPEIMAGES LENIMAGES 1  
820 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
CONNECTIVITY GRIDTOGRAPH IMAGES0SHAPE  
AGGLO CLUSTERFEATUREAGGLOMERATIONCONNECTIVITYCONNECTIVITY  
NCLUSTERS32  
AGGLOFITX  
XREDUCED AGGLOTRANSFORMX  
XRESTORED AGGLOINVERSETRANSFORMXREDUCED  
IMAGESRESTORED NPRESHAPEXRESTORED IMAGESSHAPE  
PLTFigure1 FIGSIZE4 35  
PLTCLF  
PLTSUBPLOTSADJUSTLEFT01 RIGHT99 BOTTOM01 TOP91  
FORIIRANGE4  
PLTSUBPLOT3 4 | 1  
PLTIMSHOWIMAGESI CMAPPLTCMGRAY VMAX16 INTERPOLATIONNEAREST  
PLXTICKS  
PLTYTICKS  
IFI 1  
PLTTITLEORIGINAL DATA  
PLTSUBPLOT3 4 4 | 1  
PLTIMSHOWIMAGESRESTOREDI CMAPPLTCMGRAY VMAX16  
INTERPOLATIONNEAREST  
IFI 1  
PLTTITLEAGGLOMERATED DATA  
PLXTICKS  
PLTYTICKS  
PLTSUBPLOT3 4 10  
PLTIMSHOWNPRESHAPEAGGLOLABELS IMAGES0SHAPE  
INTERPOLATIONNEAREST CMAPPLTCMNIPYSPECTRAL  
PLXTICKS  
PLTYTICKS  
PLTTITLELABELS  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0174 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
562 A DEMO OF THE MEANSHIFT CLUSTERING ALGORITHM  
REFERENCE  
DORIN COMANICIU AND PETER MEER “MEAN SHIFT A ROBUST APPROACH TOWARD FEATURE SPACE ANALYSIS” IEEE TRANSACTIONS  
ON PATTERN ANALYSIS AND MACHINE INTELLIGENCE 2002 PP 603619  
56 CLUSTERING 821

```
SCIKITLEARN USER GUIDE RELEASE 0213
OUT
NUMBER OF ESTIMATED CLUSTERS 3
PRINTDOC
IMPORT NUMPY AS NP
FROM SKLEARNCLUSTER IMPORT MEANSHIFT ESTIMATEBANDWIDTH
FROM SKLEARNDATASETSSAMPLESGENERATOR IMPORT MAKEBLOBS

GENERATE SAMPLE DATA
CENTERS 1 1 1 1 1 1
X MAKEBLOBSNSAMPLES10000 CENTERSCENTERS CLUSTERSTD06

COMPUTE CLUSTERING WITH MEANSHIFT
THE FOLLOWING BANDWIDTH CAN BE AUTOMATICALLY DETECTED USING
BANDWIDTH ESTIMATEBANDWIDTHX QUANTILE02 NSAMPLES500
822 CHAPTER 5 EXAMPLES
```

```
SCIKITLEARN USER GUIDE RELEASE 0213
MS MEANSHIFTBANDWIDTHBANDWIDTH BINSEEDINGTRUE
MSFITX
LABELS MSLABELS
CLUSTERCENTERS MSCLUSTERCENTERS
LABELSUNIQUE NPUNIQUELABELS
NCLUSTERS LENLABELSUNIQUE
PRINTNUMBER OF ESTIMATED CLUSTERS D NCLUSTERS

PLOT RESULT
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM ITERTOOLS IMPORT CYCLE
PLTFigure1
PLTCLF
COLORS CYCLEBGRCMYKBGRCMYKBGRCMYKBGRCMYK
FORK COLINZIPRANGENCLUSTERS COLORS
MYMEMBERS LABELS K
CLUSTERCENTER CLUSTERCENTERSK
PLTPLOTXMYMEMBERS 0 XMYMEMBERS 1 COL
PLTPLOTCLUSTERCENTER0 CLUSTERCENTER1 O MARKERFACECOLORCOL
MARKEREDGECOLORK MARKERSIZE14
PLTTITLEESTIMATED NUMBER OF CLUSTERS D NCLUSTERS
PLTSHOW
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0397 SECONDS
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE
563 DEMONSTRATION OF KMEANS ASSUMPTIONS
THIS EXAMPLE IS MEANT TO ILLUSTRATE SITUATIONS WHERE KMEANS WILL PRODUCE UNINTUITIVE AND POSSIBLY UNEXPECTED CLUSTERS
IN THE FIRST THREE PLOTS THE INPUT DATA DOES NOT CONFORM TO SOME IMPLICIT ASSUMPTION THAT KMEANS MAKES AND UNDESIRABLE
CLUSTERS ARE PRODUCED AS A RESULT IN THE LAST PLOT KMEANS RETURNS INTUITIVE CLUSTERS DESPITE UNEVENLY SIZED BLOBS
56 CLUSTERING 823
```

SCIKITLEARN USER GUIDE RELEASE 0213  
PRINTDOC  
AUTHOR PHIL ROTH MRPHILROTHGMAILCOM  
LICENSE BSD 3 CLAUSE  
IMPORT NUMPY AS NP  
IMPORT MATPLOTLIBPYPLOT AS PLT  
FROM SKLEARNCLUSTER IMPORT KMEANS  
FROM SKLEARNDATASETS IMPORT MAKEBLOBS  
PLTFIGUREFIGSIZE12 12  
NSAMPLES 1500  
RANDOMSTATE 170  
824 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
X Y MAKEBLOBSNSAMPLESNSAMPLES RANDOMSTATERRANDOMSTATE  
INCORRECT NUMBER OF CLUSTERS  
YPRED KMEANSNCLUSTERS2 RANDOMSTATERRANDOMSTATEFITPREDICTX  
PLTSUBPLOT221  
PLTSCATTERX 0 X 1 CYPRED  
PLTTITLEINCORRECT NUMBER OF BLOBS  
ANISOTROPICLY DISTRIBUTED DATA  
TRANSFORMATION 060834549 063667341 040887718 085253229  
XANISO NPDOTX TRANSFORMATION  
YPRED KMEANSNCLUSTERS3 RANDOMSTATERRANDOMSTATEFITPREDICTXANISO  
PLTSUBPLOT222  
PLTSCATTERXANISO 0 XANISO 1 CYPRED  
PLTTITLEANISOTROPICLY DISTRIBUTED BLOBS  
DIFFERENT VARIANCE  
XVARIED YVARIED MAKEBLOBSNSAMPLESNSAMPLES  
CLUSTERSTD10 25 05  
RANDOMSTATERRANDOMSTATE  
YPRED KMEANSNCLUSTERS3 RANDOMSTATERRANDOMSTATEFITPREDICTXVARIED  
PLTSUBPLOT223  
PLTSCATTERXVARIED 0 XVARIED 1 CYPRED  
PLTTITLEUNEQUAL VARIANCE  
UNEVENLY SIZED BLOBS  
XFILTERED NPVSTACKXY 0500 XY 1100 XY 210  
YPRED KMEANSNCLUSTERS3  
RANDOMSTATERRANDOMSTATEFITPREDICTXFILTERED  
PLTSUBPLOT224  
PLTSCATTERXFILTERED 0 XFILTERED 1 CYPRED  
PLTTITLEUNEVENLY SIZED BLOBS  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0163 SECONDS  
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE  
564 ONLINE LEARNING OF A DICTIONARY OF PARTS OF FACES  
THIS EXAMPLE USES A LARGE DATASET OF FACES TO LEARN A SET OF 20 X 20 IMAGES PATCHES THAT CONSTITUTE FACES  
FROM THE PROGRAMMING STANDPOINT IT IS INTERESTING BECAUSE IT SHOWS HOW TO USE THE ONLINE API OF THE SCIKITLEARN  
TO PROCESS A VERY LARGE DATASET BY CHUNKS THE WAY WE PROCEED IS THAT WE LOAD AN IMAGE AT A TIME AND EXTRACT  
RANDOMLY 50 PATCHES FROM THIS IMAGE ONCE WE HAVE ACCUMULATED 500 OF THESE PATCHES USING 10 IMAGES WE RUN THE  
PARTIALFIT METHOD OF THE ONLINE KMEANS OBJECT MINIBATCHKMEANS  
THE VERBOSE SETTING ON THE MINIBATCHKMEANS ENABLES US TO SEE THAT SOME CLUSTERS ARE REASSIGNED DURING THE SUCCESSIVE  
CALLS TO PARTIALFIT THIS IS BECAUSE THE NUMBER OF PATCHES THAT THEY REPRESENT HAS BECOME TOO LOW AND IT IS BETTER TO  
CHOOSE A RANDOM NEW CLUSTER  
56 CLUSTERING 825

SCIKITLEARN USER GUIDE RELEASE 0213  
OUT  
LEARNING THE DICTIONARY  
PARTIAL FIT OF 100 OUT OF 2400  
PARTIAL FIT OF 200 OUT OF 2400  
MINIBATCHKMEANS REASSIGNING 16 CLUSTER CENTERS  
PARTIAL FIT OF 300 OUT OF 2400  
PARTIAL FIT OF 400 OUT OF 2400  
PARTIAL FIT OF 500 OUT OF 2400  
PARTIAL FIT OF 600 OUT OF 2400  
PARTIAL FIT OF 700 OUT OF 2400  
PARTIAL FIT OF 800 OUT OF 2400  
PARTIAL FIT OF 900 OUT OF 2400  
PARTIAL FIT OF 1000 OUT OF 2400  
PARTIAL FIT OF 1100 OUT OF 2400  
PARTIAL FIT OF 1200 OUT OF 2400  
PARTIAL FIT OF 1300 OUT OF 2400  
PARTIAL FIT OF 1400 OUT OF 2400  
PARTIAL FIT OF 1500 OUT OF 2400  
PARTIAL FIT OF 1600 OUT OF 2400  
PARTIAL FIT OF 1700 OUT OF 2400  
PARTIAL FIT OF 1800 OUT OF 2400  
PARTIAL FIT OF 1900 OUT OF 2400  
PARTIAL FIT OF 2000 OUT OF 2400  
PARTIAL FIT OF 2100 OUT OF 2400  
PARTIAL FIT OF 2200 OUT OF 2400  
PARTIAL FIT OF 2300 OUT OF 2400  
PARTIAL FIT OF 2400 OUT OF 2400  
DONE IN 506S  
826 CHAPTER 5 EXAMPLES



```
SCIKITLEARN USER GUIDE RELEASE 0213
PRINTDOC
IMPORT TIME
IMPORT MATPLOTLIBPYPLOT AS PLT
IMPORT NUMPY AS NP
FROM SKLEARN IMPORT DATASETS
FROM SKLEARNCLUSTER IMPORT MINIBATCHKMEANS
FROM SKLEARNFEATUREEXTRACTIONIMAGE IMPORT EXTRACTPATCHES2D
FACES DATASETSFETCHOLIVETTIFACES

LEARN THE DICTIONARY OF IMAGES
PRINTLEARNING THE DICTIONARY
RNG NPRANDOMRANDOMSTATE0
KMEANS MINIBATCHKMEANSNCLUSTERS81 RANDOMSTATERNG VERBOSETRUE
PATCHSIZE 20 20
BUFFER
T0 TIMETIME
THE ONLINE LEARNING PART CYCLE OVER THE WHOLE DATASET 6 TIMES
INDEX 0
FORINRANGE6
FORIMGINFACESIMAGES
DATA EXTRACTPATCHES2DIMG PATCHSIZE MAXPATCHES50
RANDOMSTATERNG
DATA NPRESHAPEDATA LENDATA 1
BUFFERAPPENDDATA
INDEX 1
IFINDEX 10 0
DATA NPCONCATENATEBUFFER AXIS0
DATA NPMEANDATA AXIS0
DATA NPSTDDATA AXIS0
KMEANSPARTIALFITDATA
BUFFER
IFINDEX 100 0
PRINTPARTIAL FIT OF 4IOUT OFI
INDEX 6 LENFACESIMAGES
DT TIMETIME T0
PRINTDONE IN 2FS DT

PLOT THE RESULTS
PLTFIGUREFIGSIZE42 4
FORI PATCH INENUMERATEKMEANSCLUSTERCENTERS
PLTSUBPLOT9 9 I 1
PLTIMSHOWPATCHRESHAPEPATCHSIZE CMAPPLTCMGRAY
INTERPOLATIONNEAREST
PLXTICKS
56 CLUSTERING 827
```

SCIKITLEARN USER GUIDE RELEASE 0213  
PLTYTICKS  
PLTSUPTITLEPATCHES OF FACES NTRAIN TIME 1FS ONDPATCHES  
DT 8LENFACESIMAGES FONTSIZE16  
PLTSUBPLOTSADJUST008 002 092 085 008 023  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 6148 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
565 VECTOR QUANTIZATION EXAMPLE  
FACE A 1024 X 768 SIZE IMAGE OF A RACCOON FACE IS USED HERE TO ILLUSTRATE HOW KMEANS IS USED FOR VECTOR QUANTIZATION

- 
- 

828 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

- 
- 

```
PRINTDOC
CODE SOURCE GAËL VAROQUAUX
MODIFIED FOR DOCUMENTATION BY JAQUES GROBLER
LICENSE BSD 3 CLAUSE
IMPORT NUMPY AS NP
IMPORT SCIPY AS SP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARN IMPORT CLUSTER
TRY SCIPY_016 HAVE FACE IN MISC
FROM SCIPYMISC IMPORT FACE
FACE_FACEGRAYTRUE
EXCEPTIMPORTERROR
FACE_SPFACEGRAYTRUE
NCLUSTERS_5
NPRANDOMSEED0
X_FACERESHAPE1_1 WE NEED AN NSAMPLE NFEATURE ARRAY
KMEANS_CLUSTERKMEANSNCLUSTERSNCLUSTERS_NINIT4
KMEANSFITX
VALUES_KMEANSCLUSTERCENTERSSQUEEZE
LABELS_KMEANSLABELS
56 CLUSTERING 829
```

SCIKITLEARN USER GUIDE RELEASE 0213  
CREATE AN ARRAY FROM LABELS AND VALUES  
FACECOMPRESSED NPCHOOSLABELS VALUES  
FACECOMPRESSED SHAPE FACESHAPE  
VMIN FACEMIN  
VMAX FACEMAX  
ORIGINAL FACE  
PLTFigure1 FIGSIZE3 22  
PLTIMSHOWFACE CMAPPLTCMGRAY VMINVMIN VMAX256  
COMPRESSED FACE  
PLTFigure2 FIGSIZE3 22  
PLTIMSHOWFACECOMPRESSED CMAPPLTCMGRAY VMINVMIN VMAXVMAX  
EQUAL BINS FACE  
REGULARVALUES NPLINSPACE0 256 NCLUSTERS 1  
REGULARLABELS NPSEARCHSORTEDREGULARVALUES FACE 1  
REGULARVALUES 5 REGULARVALUES1 REGULARVALUES1 MEAN  
REGULARFACE NPCHOOSEREGULARLABELSRAVEL REGULARVALUES MODECLIP  
REGULARFACESHAPE FACESHAPE  
PLTFigure3 FIGSIZE3 22  
PLTIMSHOWREGULARFACE CMAPPLTCMGRAY VMINVMIN VMAXVMAX  
HISTOGRAM  
PLTFigure4 FIGSIZE3 22  
PLTCLF  
PLTAXES01 01 98 98  
PLTHISTX BINS256 COLOR5 EDGECOLOR5  
PLTYTICKS  
PLTXTICKSREGULARVALUES  
VALUES NPSORTVALUES  
FORCENTER1 CENTER2 INZIPVALUES1 VALUES1  
PLTAXVLINE5 CENTER1 CENTER2 COLORB  
FORCENTER1 CENTER2 INZIPREGULARVALUES1 REGULARVALUES1  
PLTAXVLINE5 CENTER1 CENTER2 COLORB LINESTYLE  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 3752 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
566 AGGLOMERATIVE CLUSTERING WITH AND WITHOUT STRUCTURE  
THIS EXAMPLE SHOWS THE EFFECT OF IMPOSING A CONNECTIVITY GRAPH TO CAPTURE LOCAL STRUCTURE IN THE DATA THE GRAPH IS  
SIMPLY THE GRAPH OF 20 NEAREST NEIGHBORS  
TWO CONSEQUENCES OF IMPOSING A CONNECTIVITY CAN BE SEEN FIRST CLUSTERING WITH A CONNECTIVITY MATRIX IS MUCH FASTER  
SECOND WHEN USING A CONNECTIVITY MATRIX SINGLE AVERAGE AND COMPLETE LINKAGE ARE UNSTABLE AND TEND TO CREATE A FEW  
CLUSTERS THAT GROW VERY QUICKLY INDEED AVERAGE AND COMPLETE LINKAGE FIGHT THIS PERCOLATION BEHAVIOR BY CONSIDERING ALL  
THE DISTANCES BETWEEN TWO CLUSTERS WHEN MERGING THEM WHILE SINGLE LINKAGE EXAGGERATES THE BEHAVIOUR BY CONSIDERING  
830 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

ONLY THE SHORTEST DISTANCE BETWEEN CLUSTERS THE CONNECTIVITY GRAPH BREAKS THIS MECHANISM FOR AVERAGE AND COMPLETE LINKAGE MAKING THEM RESEMBLE THE MORE BRITTLE SINGLE LINKAGE THIS EFFECT IS MORE PRONOUNCED FOR VERY SPARSE GRAPHS TRY DECREASING THE NUMBER OF NEIGHBORS IN KNEIGHBORSGRAPH AND WITH COMPLETE LINKAGE IN PARTICULAR HAVING A VERY SMALL NUMBER OF NEIGHBORS IN THE GRAPH IMPOSES A GEOMETRY THAT IS CLOSE TO THAT OF SINGLE LINKAGE WHICH IS WELL KNOWN TO HAVE THIS PERCOLATION INSTABILITY

- 
- 
- 

56 CLUSTERING 831

SCIKITLEARN USER GUIDE RELEASE 0213

```
•
AUTHORS GAELE VAROQUAUX NELLE VAROQUAUX
LICENSE BSD 3 CLAUSE
IMPORT TIME
IMPORT MATPLOTLIBPYPLOT AS PLT
IMPORT NUMPY AS NP
FROM SKLEARNCLUSTER IMPORT AGGLOMERATIVECLUSTERING
FROM SKLEARNNEIGHBORS IMPORT KNEIGHBORSGRAPH
GENERATE SAMPLE DATA
NSAMPLES 1500
NPRANDOMSEED0
T 15 NPPI1 3NPRANDOMRAND1 NSAMPLES
X TNPCOST
Y TNPSINT
X NPCONCATENATEX Y
X 7 NPRANDOMRANDN2 NSAMPLES
X XT
CREATE A GRAPH CAPTURING LOCAL CONNECTIVITY LARGER NUMBER OF NEIGHBORS
WILL GIVE MORE HOMOGENEOUS CLUSTERS TO THE COST OF COMPUTATION
TIME A VERY LARGE NUMBER OF NEIGHBORS GIVES MORE EVENLY DISTRIBUTED
CLUSTER SIZES BUT MAY NOT IMPOSE THE LOCAL MANIFOLD STRUCTURE OF
THE DATA
KNNGRAPH KNEIGHBORSGRAPHX 30 INCLUDESELFFALSE
FORCONNECTIVITY INNONE KNNGRAPH
FORNCLUSTERS IN30 3
PLTFIGUREFIGSIZE10 4
FORINDEX LINKAGE INENUMERATEAVERAGE
COMPLETE
WARD
SINGLE
PLTSUBPLOT1 4 INDEX 1
MODEL AGGLOMERATIVECLUSTERINGLINKAGELINKAGE
CONNECTIVITYCONNECTIVITY
NCLUSTERSNCLUSTERS
T0 TIMETIME
832 CHAPTER 5 EXAMPLES
```

SCIKITLEARN USER GUIDE RELEASE 0213  
MODELFITX  
ELAPSEDTIME TIMETIME T0  
PLTSCATTERX 0 X 1 CMODELLABELS  
CMAPPLTCMNIPYSPECTRAL  
PLTTITLELINKAGE SNTIME2FS LINKAGE ELAPSEDTIME  
FONTDICTDICTVERTICALALIGNMENTTOP  
PLTAXISEQUAL  
PLTAXISOFF  
PLTSUBPLOTSADJUSTBOTTOM0 TOP89 WSPACE0  
LEFT0 RIGHT1  
PLTSUPTITLENCLUSTER I CONNECTIVITY R  
NCLUSTERS CONNECTIVITY IS NOTNONE SIZE17  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 1743 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
567 DEMO OF AFFINITY PROPAGATION CLUSTERING ALGORITHM  
REFERENCE BRENDAN J FREY AND DELBERT DUECK “CLUSTERING BY PASSING MESSAGES BETWEEN DATA POINTS” SCIENCE FEB  
2007  
56 CLUSTERING 833

SCIKITLEARN USER GUIDE RELEASE 0213  
OUT  
ESTIMATED NUMBER OF CLUSTERS 3  
HOMOGENEITY 0872  
COMPLETENESS 0872  
VMEASURE 0872  
ADJUSTED RAND INDEX 0912  
ADJUSTED MUTUAL INFORMATION 0871  
SILHOUETTE COEFFICIENT 0753  
PRINTDOC  
FROM SKLEARNCLUSTER IMPORT AFFINITYPROPAGATION  
FROM SKLEARN IMPORT METRICS  
FROM SKLEARNDATASETSSAMPLESGENERATOR IMPORT MAKEBLOBS  
  
GENERATE SAMPLE DATA  
CENTERS 1 1 1 1 1 1  
X LABELSTRUE MAKEBLOBSNSAMPLES300 CENTERSCENTERS CLUSTERSTD05  
RANDOMSTATE0  
834 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213

```
COMPUTE AFFINITY PROPAGATION
AF AFFINITYPROPAGATIONPREFERENCE50FITX
CLUSTERCENTERSINDICES AFCLUSTERCENTERSINDICES
LABELS AFLABELS
NCLUSTERS LENCLUSTERCENTERSINDICES
PRINTESTIMATED NUMBER OF CLUSTERS D NCLUSTERS
PRINTHOMOGENEITY 03F METRICSHOMOGENEITYSCORELABELSTRUE LABELS
PRINTCOMPLETENESS 03F METRICSCOMPLETENESSSCORELABELSTRUE LABELS
PRINTVMEASURE 03F METRICSVMEASURESCORELABELSTRUE LABELS
PRINTADJUSTED RAND INDEX 03F
METRICSadJUSTEDRANDSCORELABELSTRUE LABELS
PRINTADJUSTED MUTUAL INFORMATION 03F
METRICSadJUSTEDMUTUALINFOSCORELABELSTRUE LABELS
AVERAGEMETHODARITHMETIC
PRINTSILHOUETTE COEFFICIENT 03F
METRICSSILHOUETTESCOREX LABELS METRICSQEUCLIDEAN
```

```
PLOT RESULT
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM ITERTOOLS IMPORT CYCLE
PLTCLOSEALL
PLTFigure1
PLTCLF
COLORS CYCLEBGRCMYKBGRCMYKBGRCMYKBGRCMYK
FORK COLINZIPRANGENCLUSTERS COLORS
CLASSMEMBERS LABELS K
CLUSTERCENTER XCLUSTERCENTERSINDICESK
PLTPLOTXCLASSMEMBERS 0 XCLASSMEMBERS 1 COL
PLTPLOTCLUSTERCENTER0 CLUSTERCENTER1 O MARKERFACECOLORCOL
MARKEREDGECOLORK MARKERSIZE14
FORXINXCLASSMEMBERS
PLTPLOTCLUSTERCENTER0 X0 CLUSTERCENTER1 X1 COL
PLTTITLEESTIMATED NUMBER OF CLUSTERS D NCLUSTERS
PLTSHOW
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0569 SECONDS
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE
568 SEGMENTING THE PICTURE OF GREEK COINS IN REGIONS
THIS EXAMPLE USES SPECTRAL CLUSTERING ON A GRAPH CREATED FROM VOXELTOVOXEL DIFFERENCE ON AN IMAGE TO BREAK THIS
IMAGE INTO MULTIPLE PARTLYHOMOGENEOUS REGIONS
THIS PROCEDURE SPECTRAL CLUSTERING ON AN IMAGE IS AN EFFICIENT APPROXIMATE SOLUTION FOR FINDING NORMALIZED GRAPH CUTS
THERE ARE TWO OPTIONS TO ASSIGN LABELS
56 CLUSTERING 835
```

```
SCIKITLEARN USER GUIDE RELEASE 0213
• WITH ‘KMEANS’ SPECTRAL CLUSTERING WILL CLUSTER SAMPLES IN THE EMBEDDING SPACE USING A KMEANS ALGORITHM
• WHEREAS ‘DISCRETE’ WILL ITERATIVELY SEARCH FOR THE CLOSEST PARTITION SPACE TO THE EMBEDDING SPACE
PRINTDOC
AUTHOR GAEL VAROQUAUX GAELVAROQUAUXNORMALESUPORG BRIAN CHEUNG
LICENSE BSD 3 CLAUSE
IMPORT TIME
IMPORT NUMPY AS NP
FROM DISTUTILSVERSION IMPORT LOOSEVERSION
FROM SCIPYNDIMAGEFILTERS IMPORT GAUSSIANFILTER
IMPORT MATPLOTLIBPYPLOT AS PLT
IMPORT SKIMAGE
FROM SKIMAGEDATA IMPORT COINS
FROM SKIMAGETRANSFORM IMPORT RESCALE
FROM SKLEARNFEATUREEXTRACTION IMPORT IMAGE
FROM SKLEARNCLUSTER IMPORT SPECTRALCLUSTERING
THESE WERE INTRODUCED IN SKIMAGE014
IFLOOSEVERSIONSKIMAGEVERSION 014
RESCALEPARAMS ANTIALIASING FALSE MULTICHANNEL FALSE
ELSE
RESCALEPARAMS
LOAD THE COINS AS A NUMPY ARRAY
ORIGCOINS COINS
RESIZE IT TO 20 OF THE ORIGINAL SIZE TO SPEED UP THE PROCESSING
APPLYING A GAUSSIAN FILTER FOR SMOOTHING PRIOR TO DOWNSCALING
REDUCES ALIASING ARTIFACTS
SMOOTHENEDCOINS GAUSSIANFILTERORIGCOINS SIGMA2
RESCALEDCOINS RESCALESMOOTHENEDCOINS 02 MODEREFLECT
RESCALEPARAMS
CONVERT THE IMAGE INTO A GRAPH WITH THE VALUE OF THE GRADIENT ON THE
EDGES
GRAPH IMAGEIMGTOGRAPHRESCALEDCOINS
TAKE A DECREASING FUNCTION OF THE GRADIENT AN EXPONENTIAL
THE SMALLER BETA IS THE MORE INDEPENDENT THE SEGMENTATION IS OF THE
ACTUAL IMAGE FOR BETA1 THE SEGMENTATION IS CLOSE TO A VORONOI
BETA 10
EPS 1E6
GRAPHDATA NPEXPBETA GRAPHDATA GRAPHDATASTD EPS
APPLY SPECTRAL CLUSTERING THIS STEP GOES MUCH FASTER IF YOU HAVE PYAMG
INSTALLED
NREGIONS 25
VISUALIZE THE RESULTING REGIONS
FORASSIGNLABELS INKMEANS DISCRETIZE
TO TIMETIME
LABELS SPECTRALCLUSTERINGGRAPH NCLUSTERSNREGIONS
ASSIGNLABELSASSIGNLABELS RANDOMSTATE42
836 CHAPTER 5 EXAMPLES
```

SCIKITLEARN USER GUIDE RELEASE 0213  
T1 TIMETIME  
LABELS LABELSRESHAPERESCALEDCOINSSHAPE  
PLTFIGUREFIGSIZE5 5  
PLTIMSHOWRESCALEDCOINS CMAPPLTCMGRAY  
FORLINRANGENREGIONS  
PLTCONTOURLABELS L  
COLORSPLTCMNIPYSPECTRALL FLOATNREGIONS  
PLTXTICKS  
PLTYTICKS  
TITLE SPECTRAL CLUSTERING S2FS ASSIGNLABELS T1 TO  
PRINTTITLE  
PLTTITLETITLE  
PLTSHOW  
•  
56 CLUSTERING 837

SCIKITLEARN USER GUIDE RELEASE 0213

•

OUT

SPECTRAL CLUSTERING KMEANS 515S

SPECTRAL CLUSTERING DISCRETIZE 598S

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 11961 SECONDS

NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)

569 KMEANS CLUSTERING

THE PLOTS DISPLAY FIRSTLY WHAT A KMEANS ALGORITHM WOULD YIELD USING THREE CLUSTERS IT IS THEN SHOWN WHAT THE EFFECT OF A BAD INITIALIZATION IS ON THE CLASSIFICATION PROCESS BY SETTING NINIT TO ONLY 1 DEFAULT IS 10 THE AMOUNT OF TIMES THAT THE ALGORITHM WILL BE RUN WITH DIFFERENT CENTROID SEEDS IS REDUCED THE NEXT PLOT DISPLAYS WHAT USING EIGHT CLUSTERS WOULD DELIVER AND FINALLY THE GROUND TRUTH

838 CHAPTER 5 EXAMPLES

- 
-

SCIKITLEARN USER GUIDE RELEASE 0213

- 
- 

PRINTDOC

CODE SOURCE GAËL VAROQUAUX  
MODIFIED FOR DOCUMENTATION BY JAQUES GROBLER  
LICENSE BSD 3 CLAUSE

IMPORT NUMPY AS NP

IMPORT MATPLOTLIBPYPLOT AS PLT

THOUGH THE FOLLOWING IMPORT IS NOT DIRECTLY BEING USED IT IS REQUIRED  
FOR 3D PROJECTION TO WORK

FROM MPLTOOLKITSMPLOT3D IMPORT AXES3D

FROM SKLEARNCLUSTER IMPORT KMEANS

FROM SKLEARN IMPORT DATASETS

NPRANDOMSEED5

840 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
IRIS DATASETSLOADIRIS  
X IRISDATA  
Y IRISTARGET  
ESTIMATORS KMEANSIRIS8 KMEANSNCLUSTERS8  
KMEANSIRIS3 KMEANSNCLUSTERS3  
KMEANSIRISBADINIT KMEANSNCLUSTERS3 NINIT1  
INITRANDOM  
FIGNUM 1  
TITLES 8 CLUSTERS 3 CLUSTERS 3 CLUSTERS BAD INITIALIZATION  
FORNAME EST INESTIMATORS  
FIG PLTFIGUREFIGNUM FIGSIZE4 3  
AX AXES3DFIG RECT0 0 95 1 ELEV48 AZIM134  
ESTFITX  
LABELS ESTLABELS  
AXSCATTERX 3 X 0 X 2  
CLABELSASTYPENPFLOAT EDGECOLORK  
AXWXAXISSETTICKLABELS  
AXWYAXISSETTICKLABELS  
AXWZAXISSETTICKLABELS  
AXSETXLABELPETAL WIDTH  
AXSETYLABELSEPAL LENGTH  
AXSETZLABELPETAL LENGTH  
AXSETTITLETITLESFIGNUM 1  
AXDIST 12  
FIGNUM FIGNUM 1  
PLOT THE GROUND TRUTH  
FIG PLTFIGUREFIGNUM FIGSIZE4 3  
AX AXES3DFIG RECT0 0 95 1 ELEV48 AZIM134  
FORNAME LABEL INSETOSA 0  
VERSICOLOUR 1  
VIRGINICA 2  
AXTEXT3DXY LABEL 3MEAN  
XY LABEL 0MEAN  
XY LABEL 2MEAN 2 NAME  
HORIZONTALALIGNMENTCENTER  
BBOXDICTALPHA2 EDGECOLORW FACECOLORW  
REORDER THE LABELS TO HAVE COLORS MATCHING THE CLUSTER RESULTS  
Y NPCHOOSEY 1 2 0ASTYPENPFLOAT  
AXSCATTERX 3 X 0 X 2 CY EDGECOLORK  
AXWXAXISSETTICKLABELS  
AXWYAXISSETTICKLABELS  
AXWZAXISSETTICKLABELS  
AXSETXLABELPETAL WIDTH  
AXSETYLABELSEPAL LENGTH  
AXSETZLABELPETAL LENGTH  
AXSETTITLEGROUND TRUTH  
AXDIST 12  
FIGSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0241 SECONDS  
56 CLUSTERING 841

SCIKITLEARN USER GUIDE RELEASE 0213

NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)

5610 VARIOUS AGGLOMERATIVE CLUSTERING ON A 2D EMBEDDING OF DIGITS

AN ILLUSTRATION OF VARIOUS LINKAGE OPTION FOR AGGLOMERATIVE CLUSTERING ON A 2D EMBEDDING OF THE DIGITS DATASET

THE GOAL OF THIS EXAMPLE IS TO SHOW INTUITIVELY HOW THE METRICS BEHAVE AND NOT TO FIND GOOD CLUSTERS FOR THE DIGITS

THIS IS WHY THE EXAMPLE WORKS ON A 2D EMBEDDING

WHAT THIS EXAMPLE SHOWS US IS THE BEHAVIOR “RICH GETTING RICHER” OF AGGLOMERATIVE CLUSTERING THAT TENDS TO CREATE UNEVEN CLUSTER SIZES THIS BEHAVIOR IS PRONOUNCED FOR THE AVERAGE LINKAGE STRATEGY THAT ENDS UP WITH A COUPLE OF SINGLETON CLUSTERS WHILE IN THE CASE OF SINGLE LINKAGE WE GET A SINGLE CENTRAL CLUSTER WITH ALL OTHER CLUSTERS BEING DRAWN FROM NOISE POINTS AROUND THE FRINGES

- 

842 CHAPTER 5 EXAMPLES



- 
-

SCIKITLEARN USER GUIDE RELEASE 0213

- OUT  
COMPUTING EMBEDDING  
DONE  
WARD 046S  
AVERAGE 036S  
COMPLETE 038S  
SINGLE 017S  
AUTHORS GAELEVAROQUAUX  
LICENSE BSD 3 CLAUSE C INRIA 2014  
PRINTDOC  
FROM TIME IMPORT TIME  
IMPORT NUMPY AS NP  
FROM SCIPY IMPORT NDIMAGE  
FROM MATPLOTLIB IMPORT PYPLASPLT  
FROM SKLEARN IMPORT MANIFOLD DATASETS  
DIGITS DATASETSLOADDIGITSNCLASS10  
X DIGITSDATA  
Y DIGITSTARGET  
NSAMPLES NFEATURES XSHAPE  
844 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
NPRANDOMSEED0
DEFNUDGEIMAGESX Y
HAVING A LARGER DATASET SHOWS MORE CLEARLY THE BEHAVIOR OF THE
METHODS BUT WE MULTIPLY THE SIZE OF THE DATASET ONLY BY 2 AS THE
COST OF THE HIERARCHICAL CLUSTERING METHODS ARE STRONGLY
SUPERLINEAR IN NSAMPLES
SHIFT LAMBDA X NDIMAGESHIFTXRESHAPE8 8
3NPRANDOMNORMALSIZE2
MODECONSTANT
RAVEL
X NPMCONCATENATEX NPAPPLYALONGAXISSHIFT 1 X
Y NPMCONCATENATEY Y AXIS0
RETURNX Y
X Y NUDGEIMAGESX Y

VISUALIZE THE CLUSTERING
DEFPLOTCLUSTERINGXRED LABELS TITLENONE
XMIN XMAX NPMINXRED AXIS0 NPMAXXRED AXIS0
XRED XRED XMIN XMAX XMIN
PLTFIGUREFIGSIZE6 4
FORIINRANGEXREDSHAPE0
PLTTTEXTXREDI 0 XREDI 1 STRYI
COLORPLTCMNIPYSPECTRALLABELSI 10
FONTDICTWEIGHT BOLD SIZE 9
PLTXTICKS
PLTTYTICKS
IFTITLEIS NOTNONE
PLTTITLETITLE SIZE17
PLTAXISOFF
PLTTIGHTLAYOUTRECT0 003 1 095

2D EMBEDDING OF THE DIGITS DATASET
PRINTCOMPUTING EMBEDDING
XRED MANIFOLDSPECTRALEMBEDDINGNCOMPONENTS2FITTRANSFORMX
PRINTDONE
FROM SKLEARNCLUSTER IMPORT AGGLOMERATIVECLUSTERING
FORLINKAGE INWARD AVERAGE COMPLETE SINGLE
CLUSTERING AGGLOMERATIVECLUSTERINGLINKAGELINKAGE NCLUSTERS10
TO TIME
CLUSTERINGFITXRED
PRINTST2FS LINKAGE TIME TO
PLOTCLUSTERINGXRED CLUSTERINGLABELS SLINKAGE LINKAGE
PLTSHOW
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 26873 SECONDS
56 CLUSTERING 845
```

SCIKITLEARN USER GUIDE RELEASE 0213

NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)

5611 SPECTRAL CLUSTERING FOR IMAGE SEGMENTATION

IN THIS EXAMPLE AN IMAGE WITH CONNECTED CIRCLES IS GENERATED AND SPECTRAL CLUSTERING IS USED TO SEPARATE THE CIRCLES IN THESE SETTINGS THE SPECTRAL CLUSTERING APPROACH SOLVES THE PROBLEM KNOW AS ‘NORMALIZED GRAPH CUTS’ THE IMAGE IS SEEN AS A GRAPH OF CONNECTED VOXELS AND THE SPECTRAL CLUSTERING ALGORITHM AMOUNTS TO CHOOSING GRAPH CUTS DEFINING REGIONS WHILE MINIMIZING THE RATIO OF THE GRADIENT ALONG THE CUT AND THE VOLUME OF THE REGION AS THE ALGORITHM TRIES TO BALANCE THE VOLUME IE BALANCE THE REGION SIZES IF WE TAKE CIRCLES WITH DIFFERENT SIZES THE SEGMENTATION FAILS

IN ADDITION AS THERE IS NO USEFUL INFORMATION IN THE INTENSITY OF THE IMAGE OR ITS GRADIENT WE CHOOSE TO PERFORM THE SPECTRAL CLUSTERING ON A GRAPH THAT IS ONLY WEAKLY INFORMED BY THE GRADIENT THIS IS CLOSE TO PERFORMING A V ORONOI PARTITION OF THE GRAPH

IN ADDITION WE USE THE MASK OF THE OBJECTS TO RESTRICT THE GRAPH TO THE OUTLINE OF THE OBJECTS IN THIS EXAMPLE WE ARE INTERESTED IN SEPARATING THE OBJECTS ONE FROM THE OTHER AND NOT FROM THE BACKGROUND

•





SCIKITLEARN USER GUIDE RELEASE 0213

```
•
PRINTDOC
AUTHORS EMMANUELLE GOUILLART EMMANUELLEGOUILLARTNORMALESUPORG
Gael VAROQUAUX GaelVAROQUAUXNORMALESUPORG
LICENSE BSD 3 CLAUSE
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARNFEATUREEXTRACTION IMPORT IMAGE
FROM SKLEARNCLUSTER IMPORT SPECTRALCLUSTERING
L 100
X Y NPINDICESL L
CENTER1 28 24
CENTER2 40 50
CENTER3 67 58
CENTER4 24 70
RADIUS1 RADIUS2 RADIUS3 RADIUS4 16 14 15 14
CIRCLE1 X CENTER10 2 Y CENTER11 2 RADIUS1 2
CIRCLE2 X CENTER20 2 Y CENTER21 2 RADIUS2 2
CIRCLE3 X CENTER30 2 Y CENTER31 2 RADIUS3 2
CIRCLE4 X CENTER40 2 Y CENTER41 2 RADIUS4 2
56 CLUSTERING 849
```

SCIKITLEARN USER GUIDE RELEASE 0213

```
4 CIRCLES
IMG CIRCLE1 CIRCLE2 CIRCLE3 CIRCLE4
WE USE A MASK THAT LIMITS TO THE FOREGROUND THE PROBLEM THAT WE ARE
INTERESTED IN HERE IS NOT SEPARATING THE OBJECTS FROM THE BACKGROUND
BUT SEPARATING THEM ONE FROM THE OTHER
MASK IMGASTYPEBOOL
IMG IMGASTYPEFLOAT
IMG 1 02 NPRANDOMRANDN IMGSHAPE
CONVERT THE IMAGE INTO A GRAPH WITH THE VALUE OF THE GRADIENT ON THE
EDGES
GRAPH IMAGEIMGTOGRAPHIMG MASKMASK
TAKE A DECREASING FUNCTION OF THE GRADIENT WE TAKE IT WEAKLY
DEPENDENT FROM THE GRADIENT THE SEGMENTATION IS CLOSE TO A VORONOI
GRAPHDATA NPEXPGRAPHDATA GRAPHDATASTD
FORCE THE SOLVER TO BE ARPACK SINCE AMG IS NUMERICALLY
UNSTABLE ON THIS EXAMPLE
LABELS SPECTRALCLUSTERINGGRAPH NCLUSTERS4 EIGENSOLVERARPACK
LABELIM NPFULLMASKSHAPE 1
LABELIMMASK LABELS
PLTMATSHOWIMG
PLTMATSHOWLABELIM
```

```
2 CIRCLES
IMG CIRCLE1 CIRCLE2
MASK IMGASTYPEBOOL
IMG IMGASTYPEFLOAT
IMG 1 02 NPRANDOMRANDN IMGSHAPE
GRAPH IMAGEIMGTOGRAPHIMG MASKMASK
GRAPHDATA NPEXPGRAPHDATA GRAPHDATASTD
LABELS SPECTRALCLUSTERINGGRAPH NCLUSTERS2 EIGENSOLVERARPACK
LABELIM NPFULLMASKSHAPE 1
LABELIMMASK LABELS
PLTMATSHOWIMG
PLTMATSHOWLABELIM
PLTSHOW
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0675 SECONDS
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE
850 CHAPTER 5 EXAMPLES
```



SCIKITLEARN USER GUIDE RELEASE 0213  
5612 A DEMO OF STRUCTURED WARD HIERARCHICAL CLUSTERING ON AN IMAGE OF COINS  
COMPUTE THE SEGMENTATION OF A 2D IMAGE WITH WARD HIERARCHICAL CLUSTERING THE CLUSTERING IS SPATIALLY CONSTRAINED IN  
ORDER FOR EACH SEGMENTED REGION TO BE IN ONE PIECE  
OUT  
COMPUTE STRUCTURED HIERARCHICAL CLUSTERING  
ELAPSED TIME 02273859977722168  
NUMBER OF PIXELS 4697  
NUMBER OF CLUSTERS 27  
AUTHOR VINCENT MICHEL 2010  
ALEXANDRE GRAMFORT 2011  
LICENSE BSD 3 CLAUSE  
PRINTDOC  
IMPORT TIME AS TIME  
56 CLUSTERING 851

```
SCIKITLEARN USER GUIDE RELEASE 0213
IMPORT NUMPY AS NP
FROM DISTUTILSVERSION IMPORT LOOSEVERSION
FROM SCIPYNDIMAGEFILTERS IMPORT GAUSSIANFILTER
IMPORT MATPLOTLIBPYPLOT AS PLT
IMPORT SKIMAGE
FROM SKIMAGEDATA IMPORT COINS
FROM SKIMAGETRANSFORM IMPORT RESCALE
FROM SKLEARNFEATUREEXTRACTIONIMAGE IMPORT GRIDTOGRAPH
FROM SKLEARNCLUSTER IMPORT AGGLOMERATIVECLUSTERING
THESE WERE INTRODUCED IN SKIMAGE014
IFLOOSEVERSIONSKIMAGEVERSION 014
RESCALEPARAMS ANTIALIASING FALSE MULTICHANNEL FALSE
ELSE
RESCALEPARAMS

GENERATE DATA
ORIGCOINS COINS
RESIZE IT TO 20 OF THE ORIGINAL SIZE TO SPEED UP THE PROCESSING
APPLYING A GAUSSIAN FILTER FOR SMOOTHING PRIOR TO DOWNSCALING
REDUCES ALIASING ARTIFACTS
SMOOTHENEDCOINS GAUSSIANFILTERORIGCOINS SIGMA2
RESCALEDCOINS RESCALESMOOTHENEDCOINS 02 MODEREFLECT
RESCALEPARAMS
X NPRESHAPEXRESCALEDCOINS 1 1

DEFINE THE STRUCTURE A OF THE DATA PIXELS CONNECTED TO THEIR NEIGHBORS
CONNECTIVITY GRIDTOGRAPH RESCALEDCOINSSHAPE

COMPUTE CLUSTERING
PRINTCOMPUTE STRUCTURED HIERARCHICAL CLUSTERING
ST TIMETIME
NCLUSTERS 27 NUMBER OF REGIONS
WARD AGGLOMERATIVECLUSTERINGNCLUSTERSNCLUSTERS LINKAGEWARD
CONNECTIVITYCONNECTIVITY
WARDFITX
LABEL NPRESHAPEWARDLABELS RESCALEDCOINSSHAPE
PRINTELAPSED TIME TIMETIME ST
PRINTNUMBER OF PIXELS LABELSIZE
PRINTNUMBER OF CLUSTERS NPUNIQUELABELSIZE

PLOT THE RESULTS ON AN IMAGE
PLTFIGUREFIGSIZE5 5
PLTIMSHOWRESCALEDCOINS CMAPPLTCMGRAY
FORLINRANGENCLUSTERS
PLTCONTOURLABEL L
COLORSPLTCMNIPYSPECTRALL FLOATNCLUSTERS
PLXTICKS
PLTYTICKS
852 CHAPTER 5 EXAMPLES
```

SCIKITLEARN USER GUIDE RELEASE 0213  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0897 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
5613 DEMO OF DBSCAN CLUSTERING ALGORITHM  
FINDS CORE SAMPLES OF HIGH DENSITY AND EXPANDS CLUSTERS FROM THEM  
OUT  
ESTIMATED NUMBER OF CLUSTERS 3  
ESTIMATED NUMBER OF NOISE POINTS 18  
HOMOGENEITY 0953  
COMPLETENESS 0883  
VMEASURE 0917  
ADJUSTED RAND INDEX 0952  
ADJUSTED MUTUAL INFORMATION 0916  
SILHOUETTE COEFFICIENT 0626  
56 CLUSTERING 853

```
SCIKITLEARN USER GUIDE RELEASE 0213
PRINTDOC
IMPORT NUMPY AS NP
FROM SKLEARNCLUSTER IMPORT DBSCAN
FROM SKLEARN IMPORT METRICS
FROM SKLEARNDATASETSAMPLESGENERATOR IMPORT MAKEBLOBS
FROM SKLEARNPREPROCESSING IMPORT STANDARDSCALER

GENERATE SAMPLE DATA
CENTERS 1 1 1 1 1 1
X LABELSTRUE MAKEBLOBSNSAMPLES750 CENTERSCENTERS CLUSTERSTD04
RANDOMSTATE0
X STANDARDSCALERFITTRANSFORMX

COMPUTE DBSCAN
DB DBSCANEPS03 MINSAMPLES10FITX
CORESAMPLESMASK NPZEROSLIKEDBLABELS DTYPEBOOL
CORESAMPLESMASKDBCORESAMPLEINDICES TRUE
LABELS DBLABELS
NUMBER OF CLUSTERS IN LABELS IGNORING NOISE IF PRESENT
NCLUSTERS LENSETLABELS 1 IF1INLABELSELSE0
NNOISE LISTLABELSCOUNT1
PRINTESTIMATED NUMBER OF CLUSTERS D NCLUSTERS
PRINTESTIMATED NUMBER OF NOISE POINTS D NNOISE
PRINTHOMOGENEITY 03F METRICSHOMOGENEITYSCORELABELSTRUE LABELS
PRINTCOMPLETENESS 03F METRICSCOMPLETENESSSCORELABELSTRUE LABELS
PRINTVMEASURE 03F METRICSVMEASURESCORELABELSTRUE LABELS
PRINTADJUSTED RAND INDEX 03F
METRICSadJUSTEDRANDSCORELABELSTRUE LABELS
PRINTADJUSTED MUTUAL INFORMATION 03F
METRICSadJUSTEDMUTUALINFOSCORELABELSTRUE LABELS
AVERAGEMETHODARITHMETIC
PRINTSILHOUETTE COEFFICIENT 03F
METRICSSILHOUETTESCOREX LABELS

PLOT RESULT
IMPORT MATPLOTLIBPYPLOT AS PLT
BLACK REMOVED AND IS USED FOR NOISE INSTEAD
UNIQUELABELS SETLABELS
COLORS PLTCMSPECTRALEACH
FOREACHINNPLINSPACE0 1 LENUNIQUELABELS
FORK COLINZIPUNIQUELABELS COLORS
IFK 1
BLACK USED FOR NOISE
COL 0 0 0 1
854 CHAPTER 5 EXAMPLES
```

SCIKITLEARN USER GUIDE RELEASE 0213  
CLASSMEMBERMASK LABELS K  
XY XCLASSMEMBERMASK CORESAMPLESMASK  
PLTPLOTXY 0 XY 1 O MARKERFACECOLORTUPLECOL  
MARKEREDGECOLORK MARKERSIZE14  
XY XCLASSMEMBERMASK CORESAMPLESMASK  
PLTPLOTXY 0 XY 1 O MARKERFACECOLORTUPLECOL  
MARKEREDGECOLORK MARKERSIZE6  
PLTTITLEESTIMATED NUMBER OF CLUSTERS D NCLUSTERS  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0043 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
5614 COLOR QUANTIZATION USING KMEANS  
PERFORMS A PIXELWISE VECTOR QUANTIZATION VQ OF AN IMAGE OF THE SUMMER PALACE CHINA REDUCING THE NUMBER OF  
COLORS REQUIRED TO SHOW THE IMAGE FROM 96615 UNIQUE COLORS TO 64 WHILE PRESERVING THE OVERALL APPEARANCE QUALITY  
IN THIS EXAMPLE PIXELS ARE REPRESENTED IN A 3DSPACE AND KMEANS IS USED TO FIND 64 COLOR CLUSTERS IN THE IMAGE  
PROCESSING LITERATURE THE CODEBOOK OBTAINED FROM KMEANS THE CLUSTER CENTERS IS CALLED THE COLOR PALETTE USING A  
SINGLE BYTE UP TO 256 COLORS CAN BE ADDRESSED WHEREAS AN RGB ENCODING REQUIRES 3 BYTES PER PIXEL THE GIF FILE  
FORMAT FOR EXAMPLE USES SUCH A PALETTE  
FOR COMPARISON A QUANTIZED IMAGE USING A RANDOM CODEBOOK COLORS PICKED UP RANDOMLY IS ALSO SHOWN  
56 CLUSTERING 855





SCIKITLEARN USER GUIDE RELEASE 0213

- OUT  
FITTING MODEL ON A SMALL SUBSAMPLE OF THE DATA  
DONE IN 0264S  
PREDICTING COLOR INDICES ON THE FULL IMAGE KMEANS  
DONE IN 0216S  
PREDICTING COLOR INDICES ON THE FULL IMAGE RANDOM  
DONE IN 0310S  
AUTHORS ROBERT LAYTON ROBERTLAYTONGMAILCOM  
OLIVIER GRISEL OLIVIERGRISELENSTAORG  
MATHIEU BLONDEL MATHIEUMBLONDELORG  
  
LICENSE BSD 3 CLAUSE  
PRINTDOC  
IMPORT NUMPY AS NP  
IMPORT MATPLOTLIBPYPLOT AS PLT  
FROM SKLEARNCLUSTER IMPORT KMEANS  
FROM SKLEARNMETRICS IMPORT PAIRWISEDISTANCESARGMIN  
FROM SKLEARNDATASETS IMPORT LOADSAMPLEIMAGE  
FROM SKLEARNUTILS IMPORT SHUFFLE  
858 CHAPTER 5 EXAMPLES



```
SCIKITLEARN USER GUIDE RELEASE 0213
FROM TIME IMPORT TIME
NCOLORS 64
LOAD THE SUMMER PALACE PHOTO
CHINA LOADSAMPLEIMAGECHINA.JPG
CONVERT TO FLOATS INSTEAD OF THE DEFAULT 8 BITS INTEGER CODING DIVIDING BY
255 IS IMPORTANT SO THAT PLTIMSHOW BEHAVES WORKS WELL ON FLOAT DATA NEED TO
BE IN THE RANGE 01
CHINA NPARRAYCHINA DTYPE NPFLOAT64 255
LOAD IMAGE AND TRANSFORM TO A 2D NUMPY ARRAY
W H D ORIGINALSHAPE TUPLECHINASHAPE
ASSERTD 3
IMAGEARRAY NPRESHAPECHINA W H D
PRINTFITTING MODEL ON A SMALL SUBSAMPLE OF THE DATA
TO TIME
IMAGEARRAYSAMPLE SHUFFLEIMAGEARRAY RANDOMSTATE01000
KMEANS KMEANSNCLUSTERSNCOLORS RANDOMSTATE0FITIMAGEARRAYSAMPLE
PRINTDONE IN 03FS TIME TO
GET LABELS FOR ALL POINTS
PRINTPREDICTING COLOR INDICES ON THE FULL IMAGE KMEANS
TO TIME
LABELS KMEANSPREDICTIMAGEARRAY
PRINTDONE IN 03FS TIME TO
CODEBOOKRANDOM SHUFFLEIMAGEARRAY RANDOMSTATE0NCOLORS
PRINTPREDICTING COLOR INDICES ON THE FULL IMAGE RANDOM
TO TIME
LABELSRANDOM PAIRWISEDISTANCESARGMINCODEBOOKRANDOM
IMAGEARRAY
AXIS0
PRINTDONE IN 03FS TIME TO
DEFRECREATEIMAGECODEBOOK LABELS W H
RECREATE THE COMPRESSED IMAGE FROM THE CODE BOOK LABELS
D CODEBOOKSHAPE1
IMAGE NPZEROSW H D
LABELIDX 0
FORIINRANGEW
FORJINRANGEH
IMAGEIJ CODEBOOKLABELSLABELIDX
LABELIDX 1
RETURNIMAGE
DISPLAY ALL RESULTS ALONGSIDE ORIGINAL IMAGE
PLTFigure1
PLTCLF
PLTAXISOFF
PLTTITLEORIGINAL IMAGE 96615 COLORS
PLTIMSHOWCHINA
PLTFigure2
56 CLUSTERING 859
```

SCIKITLEARN USER GUIDE RELEASE 0213

PLTCLF  
PLTAXISOFF  
PLTTITLEQUANTIZED IMAGE 64 COLORS KMEANS  
PLTIMSHOWRECREATEIMAGEKMEANSCLUSTERCENTERS LABELS W H  
PLTFigure3  
PLTCLF  
PLTAXISOFF  
PLTTITLEQUANTIZED IMAGE 64 COLORS RANDOM  
PLTIMSHOWRECREATEIMAGECODEBOOKRANDOM LABELSRANDOM W H  
PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 1400 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
5615 HIERARCHICAL CLUSTERING STRUCTURED VS UNSTRUCTURED WARD  
EXAMPLE BUILDS A SWISS ROLL DATASET AND RUNS HIERARCHICAL CLUSTERING ON THEIR POSITION  
FOR MORE INFORMATION SEE HIERARCHICAL CLUSTERING  
IN A FIRST STEP THE HIERARCHICAL CLUSTERING IS PERFORMED WITHOUT CONNECTIVITY CONSTRAINTS ON THE STRUCTURE AND IS SOLELY  
BASED ON DISTANCE WHEREAS IN A SECOND STEP THE CLUSTERING IS RESTRICTED TO THE KNEAREST NEIGHBORS GRAPH IT’S A  
HIERARCHICAL CLUSTERING WITH STRUCTURE PRIOR  
SOME OF THE CLUSTERS LEARNED WITHOUT CONNECTIVITY CONSTRAINTS DO NOT RESPECT THE STRUCTURE OF THE SWISS ROLL AND EXTEND  
ACROSS DIFFERENT FOLDS OF THE MANIFOLDS ON THE OPPOSITE WHEN OPPOSING CONNECTIVITY CONSTRAINTS THE CLUSTERS FORM A  
NICE PARCELLATION OF THE SWISS ROLL  
860 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213

•

OUT  
COMPUTE UNSTRUCTURED HIERARCHICAL CLUSTERING  
ELAPSED TIME 008S  
NUMBER OF POINTS 1500  
COMPUTE STRUCTURED HIERARCHICAL CLUSTERING  
ELAPSED TIME 007S  
NUMBER OF POINTS 1500  
AUTHORS VINCENT MICHEL 2010  
ALEXANDRE GRAMFORT 2010  
GAEL VAROQUAUX 2010  
LICENSE BSD 3 CLAUSE  
PRINTDOC  
IMPORT TIME AS TIME  
IMPORT NUMPY AS NP  
IMPORT MATPLOTLIBPYPLOT AS PLT  
IMPORT MPLTOOLKITSMPLLOT3DAXES3D AS P3  
FROM SKLEARNCLUSTER IMPORT AGGLOMERATIVECLUSTERING  
FROM SKLEARNDATASETSSAMPLESGENERATOR IMPORT MAKESWISSROLL  
862 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

GENERATE DATA SWISS ROLL DATASET  
NSAMPLES 1500  
NOISE 005  
X MAKESWISSROLLNSAMPLES NOISE  
MAKE IT THINNER  
X 1 5

COMPUTE CLUSTERING  
PRINTCOMPUTE UNSTRUCTURED HIERARCHICAL CLUSTERING  
ST TIMETIME  
WARD AGGLOMERATIVECLUSTERINGNCLUSTERS6 LINKAGEWARDFITX  
ELAPSEDTIME TIMETIME ST  
LABEL WARDLABELS  
PRINTELAPSED TIME 2FS ELAPSEDTIME  
PRINTNUMBER OF POINTS I LABELSIZE

PLOT RESULT  
FIG PLTFigure  
AX P3AXES3DFIG  
AXVIEWINIT7 80  
FORLINNPUNIQUELABEL  
AXSCATTERXLABEL L 0 XLABEL L 1 XLABEL L 2  
COLORPLTCMJETNPFLOATL NPMAXLABEL 1  
S20 EDGECOLORK  
PLTTITLEWITHOUT CONNECTIVITY CONSTRAINTS TIME 2FS ELAPSEDTIME

DEFINE THE STRUCTURE A OF THE DATA HERE A 10 NEAREST NEIGHBORS  
FROM SKLEARNNEIGHBORS IMPORT KNEIGHBORSGRAPH  
CONNECTIVITY KNEIGHBORSGRAPHX NNEIGHBORS10 INCLUDESELFFALSE

COMPUTE CLUSTERING  
PRINTCOMPUTE STRUCTURED HIERARCHICAL CLUSTERING  
ST TIMETIME  
WARD AGGLOMERATIVECLUSTERINGNCLUSTERS6 CONNECTIVITYCONNECTIVITY  
LINKAGEWARDFITX  
ELAPSEDTIME TIMETIME ST  
LABEL WARDLABELS  
PRINTELAPSED TIME 2FS ELAPSEDTIME  
PRINTNUMBER OF POINTS I LABELSIZE

PLOT RESULT  
FIG PLTFigure  
AX P3AXES3DFIG  
AXVIEWINIT7 80  
FORLINNPUNIQUELABEL  
AXSCATTERXLABEL L 0 XLABEL L 1 XLABEL L 2  
COLORPLTCMJETFLOATL NPMAXLABEL 1  
S20 EDGECOLORK  
PLTTITLEWITH CONNECTIVITY CONSTRAINTS TIME 2FS ELAPSEDTIME  
56 CLUSTERING 863

SCIKITLEARN USER GUIDE RELEASE 0213

PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0202 SECONDS

NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)

5616 AGGLOMERATIVE CLUSTERING WITH DIFFERENT METRICS

DEMONSTRATES THE EFFECT OF DIFFERENT METRICS ON THE HIERARCHICAL CLUSTERING

THE EXAMPLE IS ENGINEERED TO SHOW THE EFFECT OF THE CHOICE OF DIFFERENT METRICS IT IS APPLIED TO WAVEFORMS WHICH CAN BE SEEN AS HIGHDIMENSIONAL VECTOR INDEED THE DIFFERENCE BETWEEN METRICS IS USUALLY MORE PRONOUNCED IN HIGH DIMENSION IN PARTICULAR FOR EUCLIDEAN AND CITYBLOCK

WE GENERATE DATA FROM THREE GROUPS OF WAVEFORMS TWO OF THE WAVEFORMS WAVEFORM 1 AND WAVEFORM 2 ARE PROPORTIONAL ONE TO THE OTHER THE COSINE DISTANCE IS INVARIANT TO A SCALING OF THE DATA AS A RESULT IT CANNOT DISTINGUISH THESE TWO WAVEFORMS THUS EVEN WITH NO NOISE CLUSTERING USING THIS DISTANCE WILL NOT SEPARATE OUT WAVEFORM 1 AND 2 WE ADD OBSERVATION NOISE TO THESE WAVEFORMS WE GENERATE VERY SPARSE NOISE ONLY 6 OF THE TIME POINTS CONTAIN NOISE AS A RESULT THE L1 NORM OF THIS NOISE IE "CITYBLOCK" DISTANCE IS MUCH SMALLER THAN IT'S L2 NORM "EUCLIDEAN" DISTANCE THIS CAN BE SEEN ON THE INTERCLASS DISTANCE MATRICES THE VALUES ON THE DIAGONAL THAT CHARACTERIZE THE SPREAD OF THE CLASS ARE MUCH BIGGER FOR THE EUCLIDEAN DISTANCE THAN FOR THE CITYBLOCK DISTANCE

WHEN WE APPLY CLUSTERING TO THE DATA WE FIND THAT THE CLUSTERING REFLECTS WHAT WAS IN THE DISTANCE MATRICES INDEED FOR THE EUCLIDEAN DISTANCE THE CLASSES ARE ILLSEPARATED BECAUSE OF THE NOISE AND THUS THE CLUSTERING DOES NOT SEPARATE THE WAVEFORMS FOR THE CITYBLOCK DISTANCE THE SEPARATION IS GOOD AND THE WAVEFORM CLASSES ARE RECOVERED FINALLY THE COSINE DISTANCE DOES NOT SEPARATE AT ALL WAVEFORM 1 AND 2 THUS THE CLUSTERING PUTS THEM IN THE SAME CLUSTER

864 CHAPTER 5 EXAMPLES















SCIKITLEARN USER GUIDE RELEASE 0213

```
•
AUTHOR GAELEVAROQUAUX
LICENSE BSD 3CLAUSE OR CC0
IMPORT MATPLOTLIBPYPLOT AS PLT
IMPORT NUMPY AS NP
FROM SKLEARNCLUSTER IMPORT AGGLOMERATIVECLUSTERING
FROM SKLEARNMETRICS IMPORT PAIRWISEDISTANCES
NPRANDOMSEED0
GENERATE WAVEFORM DATA
NFEATURES 2000
T NPPI NPLinspace0 1 NFEATURES
DEFSQRX
RETURNNPSIGNNPCOSX
X LIST
Y LIST
FORI PHI A INENUMERATE5 15 5 6 3 2
FORINRANGE30
PHASENOISE 01 NPRANDOMNORMAL
AMPLITUDENOISE 04 NPRANDOMNORMAL
ADDITIONALNOISE 1 2 NPRANDOMRANDNFEATURES
MAKE THE NOISE SPARSE
56 CLUSTERING 871
```

SCIKITLEARN USER GUIDE RELEASE 0213  
ADDITIONALNOISENPABSADDITIONALNOISE 997 0  
XAPPEND12 A AMPLITUDENOISE  
SQR6T PHI PHASENOISE  
ADDITIONALNOISE  
YAPPENDI  
X NPARRAYX  
Y NPARRAYY  
NCLUSTERS 3  
LABELS WAVEFORM 1 WAVEFORM 2 WAVEFORM 3  
PLOT THE GROUNDTRUTH LABELLING  
PLTFigure  
PLTAXES0 0 1 1  
FORL C N INZIPRANGENCLUSTERS RGB  
LABELS  
LINES PLTPLOTXY LT CC ALPHA5  
LINES0SETLABELN  
PLTLEGENDLOCBEST  
PLTAXISTIGHT  
PLTAXISOFF  
PLTSUPTITLEGROUND TRUTH SIZE20  
PLOT THE DISTANCES  
FORINDEX METRIC INENUMERATECOSINE EUCLIDEAN CITYBLOCK  
AVGDIST NPZEROSNCLUSTERS NCLUSTERS  
PLTFigureFIGSIZE5 45  
FORIINRANGENCLUSTERS  
FORJINRANGENCLUSTERS  
AVGDIST I J PAIRWISEDISTANCESXY I XY J  
METRICMETRICMEAN  
AVGDIST AVGDISTMAX  
FORIINRANGENCLUSTERS  
FORJINRANGENCLUSTERS  
PLTTEXT I J 53F AVGDIST I J  
VERTICALALIGNMENTCENTER  
HORIZONTALALIGNMENTCENTER  
PLTIMSHOWAVGDIST INTERPOLATIONNEAREST CMAPPLTCMGNUPLOT2  
VMINO  
PLXTICKSRANGENCLUSTERS LABELS ROTATION45  
PLTYTICKSRANGENCLUSTERS LABELS  
PLTCOLORBAR  
PLTSUPTITLEINTERCLASS SDISTANCES METRIC SIZE18  
PLTTIGHTLAYOUT  
PLOT CLUSTERING RESULTS  
FORINDEX METRIC INENUMERATECOSINE EUCLIDEAN CITYBLOCK  
MODEL AGGLOMERATIVECLUSTERINGNCLUSTERSNCLUSTERS  
LINKAGEAVERAGE AFFINITYMETRIC  
MODELFITX  
872 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
PLTFigure
plt.axes([0, 0, 1, 1])
for l in zip(nparange(model_n_clusters, 1, 10), range(1, 10)):
    plt.plot(x, model_labels, 'lt', c=alpha)
plt.axis('tight')
plt.axis('off')
plt.suptitle('Agglomerative Clustering Affinity S Metric Size 20')
plt.show()

TOTAL RUNNING TIME OF THE SCRIPT: 0 MINUTES 0584 SECONDS
NOTE: CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE

5617 INDUCTIVE CLUSTERING
CLUSTERING CAN BE EXPENSIVE ESPECIALLY WHEN OUR DATASET CONTAINS MILLIONS OF DATAPOINTS. MANY CLUSTERING ALGORITHMS
ARE NOT INDUCTIVE AND SO CANNOT BE DIRECTLY APPLIED TO NEW DATA SAMPLES WITHOUT RECOMPUTING THE CLUSTERING WHICH
MAY BE INTRACTABLE. INSTEAD, WE CAN USE CLUSTERING TO THEN LEARN AN INDUCTIVE MODEL WITH A CLASSIFIER WHICH HAS SEVERAL
BENEFITS:
• IT ALLOWS THE CLUSTERS TO SCALE AND APPLY TO NEW DATA
• UNLIKE REFITTING THE CLUSTERS TO NEW SAMPLES, IT MAKES SURE THE LABELLING PROCEDURE IS CONSISTENT OVER TIME
• IT ALLOWS US TO USE THE INFERENCE CAPABILITIES OF THE CLASSIFIER TO DESCRIBE OR EXPLAIN THE CLUSTERS
THIS EXAMPLE ILLUSTRATES A GENERIC IMPLEMENTATION OF A METAESTIMATOR WHICH EXTENDS CLUSTERING BY INDUCING A CLASSIFIER
FROM THE CLUSTER LABELS.
AUTHORS: CHIRAG NAGPAL
CHRISTOS ARIDAS
PRINTDOC
import numpy as np
import matplotlib.pyplot as plt
from sklearn.base import BaseEstimator, Clone
from sklearn.cluster import AgglomerativeClustering
from sklearn.datasets import make_blobs
from sklearn.ensemble import RandomForestClassifier

56 CLUSTERING 873
```

```
SCIKITLEARN USER GUIDE RELEASE 0213
FROM SKLEARNUTILSMETAESTIMATORS IMPORT IFDELEGATEHASMETHOD
NSAMPLES 5000
RANDOMSTATE 42
CLASS INDUCTIVECLUSTERER BASEESTIMATOR
DEFINITSELF CLUSTERER CLASSIFIER
SELFCLUSTERER CLUSTERER
SELFCLASSIFIER CLASSIFIER
DEFFITSELF X YNONE
SELFCLUSTERER CLONESELFCLUSTERER
SELFCLASSIFIER CLONESELFCLASSIFIER
Y SELFCLUSTERERFITPREDICTX
SELFCLASSIFIERFITX Y
RETURNSELF
IFDELEGATEHASMETHOD DELEGATECLASSIFIER
DEFPREDICTSELF X
RETURNSELFCLASSIFIERPREDICTX
IFDELEGATEHASMETHOD DELEGATECLASSIFIER
DEFDECISIONFUNCTIONSELF X
RETURNSELFCLASSIFIERDECISIONFUNCTIONX
DEFPLOTSCATTERX COLOR ALPHA05
RETURNPLTSCATTERX 0
X 1
CCOLOR
ALPHAALPHA
EDGECOLORK
GENERATE SOME TRAINING DATA FROM CLUSTERING
X Y MAKEBLOBSNSAMPLESNSAMPLES
CLUSTERSTD10 10 05
CENTERS5 5 0 0 5 5
RANDOMSTATERANDOMSTATE
TRAIN A CLUSTERING ALGORITHM ON THE TRAINING DATA AND GET THE CLUSTER LABELS
CLUSTERER AGGLOMERATIVECLUSTERINGNCLUSERS3
CLUSTERLABELS CLUSTERERFITPREDICTX
PLTFIGUREFIGSIZE12 4
PLTSUBPLOT131
PLOTSCATTERX CLUSTERLABELS
PLTTITLEWARD LINKAGE
GENERATE NEW SAMPLES AND PLOT THEM ALONG WITH THE ORIGINAL DATASET
XNEW YNEW MAKEBLOBSNSAMPLES10
CENTERS7 1 2 4 3 6
RANDOMSTATERANDOMSTATE
874 CHAPTER 5 EXAMPLES
```



SCIKITLEARN USER GUIDE RELEASE 0213

PLTSUBPLOT132

PLOTSCATTERX CLUSTERLABELS

PLOTSCATTERXNEW BLACK 1

PLTTITLEUNKNOWN INSTANCES

DECLARE THE INDUCTIVE LEARNING MODEL THAT IT WILL BE USED TO

PREDICT CLUSTER MEMBERSHIP FOR UNKNOWN INSTANCES

CLASSIFIER RANDOMFORESTCLASSIFIERRANDOMSTATERANDOMSTATE

INDUCTIVELEARNER INDUCTIVECLUSTERERCLUSTERER CLASSIFIERFITX

PROBABLECLUSTERS INDUCTIVELEARNERPREDICTXNEW

PLTSUBPLOT133

PLOTSCATTERX CLUSTERLABELS

PLOTSCATTERXNEW PROBABLECLUSTERS

PLOTTING DECISION REGIONS

XMIN XMAX X 0MIN 1 X 0MAX 1

YMIN YMAX X 1MIN 1 X 1MAX 1

XX YY NPMESHGRIDNPARANGEXMIN XMAX 01

NPARANGEYMIN YMAX 01

Z INDUCTIVELEARNERPREDICTNPCXXRAVEL YYRAVEL

Z ZRESHAPEXXSHAPE

PLTCONTOURFXX YY Z ALPHA04

PLTTITLECLASSIFY UNKNOWN INSTANCES

PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 1436 SECONDS

NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)

5618 DEMO OF OPTICS CLUSTERING ALGORITHM

FINDS CORE SAMPLES OF HIGH DENSITY AND EXPANDS CLUSTERS FROM THEM THIS EXAMPLE USES DATA THAT IS GENERATED SO

THAT THE CLUSTERS HAVE DIFFERENT DENSITIES THE SKLEARNCLUSTEROPTICS IS FIRST USED WITH ITS XI CLUSTER DETEC

TION METHOD AND THEN SETTING SPECIFIC THRESHOLDS ON THE REACHABILITY WHICH CORRESPONDS TO SKLEARNCLUSTER

DBSCAN WE CAN SEE THAT THE DIFFERENT CLUSTERS OF OPTICS'S XI METHOD CAN BE RECOVERED WITH DIFFERENT CHOICES OF

THRESHOLDS IN DBSCAN

56 CLUSTERING 875

```
SCIKITLEARN USER GUIDE RELEASE 0213
AUTHORS SHANE GRIGSBY REFUGEROCKTALUSCOM
ADRIN JALALI ADRINJALALIGMAILCOM
LICENSE BSD 3 CLAUSE
FROM SKLEARNCLUSTER IMPORT OPTICS CLUSTEROPTICSDBCAN
IMPORT MATPLOTLIBGRIDSPEC AS GRIDSPEC
IMPORT MATPLOTLIBPYPLOT AS PLT
IMPORT NUMPY AS NP
GENERATE SAMPLE DATA
NPRANDOMSEED0
NPOINTSPERCLUSTER 250
C1 5 2 8 NPRANDOMRANDNNPOINTSPERCLUSTER 2
C2 4 1 1 NPRANDOMRANDNNPOINTSPERCLUSTER 2
C3 1 2 2 NPRANDOMRANDNNPOINTSPERCLUSTER 2
C4 2 3 3 NPRANDOMRANDNNPOINTSPERCLUSTER 2
C5 3 2 16 NPRANDOMRANDNNPOINTSPERCLUSTER 2
C6 5 6 2 NPRANDOMRANDNNPOINTSPERCLUSTER 2
X NPVSTACKC1 C2 C3 C4 C5 C6
CLUST OPTICSMINSAMPLES50 XI05 MINCLUSTERSIZE05
RUN THE FIT
CLUSTFITX
876 CHAPTER 5 EXAMPLES
```

SCIKITLEARN USER GUIDE RELEASE 0213  
LABELS050 CLUSTEROPTICSDBSCANREACHABILITYCLUSTREACHABILITY  
COREDISTANCESCLUSTCOREDISTANCES  
ORDERINGCLUSTORDERING EPS05  
LABELS200 CLUSTEROPTICSDBSCANREACHABILITYCLUSTREACHABILITY  
COREDISTANCESCLUSTCOREDISTANCES  
ORDERINGCLUSTORDERING EPS2  
SPACE NPARANGELENX  
REACHABILITY CLUSTREACHABILITYCLUSTORDERING  
LABELS CLUSTLABELSCLUSTORDERING  
PLTFIGUREFIGSIZE10 7  
G GRIDSPECGRIDSPEC2 3  
AX1 PLTSUBPLOTG0  
AX2 PLTSUBPLOTG1 0  
AX3 PLTSUBPLOTG1 1  
AX4 PLTSUBPLOTG1 2  
REACHABILITY PLOT  
COLORS G R B Y C  
FORCLASS COLOR INZIPRANGE0 5 COLORS  
XK SPACELABELS KCLASS  
RK REACHABILITYLABELS KCLASS  
AX1PLOTXK RK COLOR ALPHA03  
AX1PLOTSPACECLABELS 1 REACHABILITYLABELS 1 K ALPHA03  
AX1PLOTSPACE NPFULLLIKESPACE 2 DTYPEFLOAT K ALPHA05  
AX1PLOTSPACE NPFULLLIKESPACE 05 DTYPEFLOAT K ALPHA05  
AX1SETYLABELREACHABILITY EPSILON DISTANCE  
AX1SETTITLEREACHABILITY PLOT  
OPTICS  
COLORS G R B Y C  
FORCLASS COLOR INZIPRANGE0 5 COLORS  
XK XCLUSTLABELS KCLASS  
AX2PLOTXK 0 XK 1 COLOR ALPHA03  
AX2PLOTXCLUSTLABELS 1 0 XCLUSTLABELS 1 1 K ALPHA01  
AX2SETTITLEAUTOMATIC CLUSTERING NOPTICS  
DBSCAN AT 05  
COLORS G GREENYELLOW OLIVE R B C  
FORCLASS COLOR INZIPRANGE0 6 COLORS  
XK XLABELS050 KCLASS  
AX3PLOTXK 0 XK 1 COLOR ALPHA03 MARKER  
AX3PLOTXLABELS050 1 0 XLABELS050 1 1 K ALPHA01  
AX3SETTITLECLUSTERING AT 05 EPSILON CUT NDBSCAN  
DBSCAN AT 2  
COLORS G M Y C  
FORCLASS COLOR INZIPRANGE0 4 COLORS  
XK XLABELS200 KCLASS  
AX4PLOTXK 0 XK 1 COLOR ALPHA03  
AX4PLOTXLABELS200 1 0 XLABELS200 1 1 K ALPHA01  
AX4SETTITLECLUSTERING AT 20 EPSILON CUT NDBSCAN  
PLTTIGHTLAYOUT  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0935 SECONDS  
56 CLUSTERING 877

SCIKITLEARN USER GUIDE RELEASE 0213

NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)

5619 COMPARE BIRCH AND MINIBATCHKMEANS

THIS EXAMPLE COMPARES THE TIMING OF BIRCH WITH AND WITHOUT THE GLOBAL CLUSTERING STEP AND MINIBATCHKMEANS ON A SYNTHETIC DATASET HAVING 100000 SAMPLES AND 2 FEATURES GENERATED USING MAKEBLOBS

IFNCLUSTERS IS SET TO NONE THE DATA IS REDUCED FROM 100000 SAMPLES TO A SET OF 158 CLUSTERS THIS CAN BE VIEWED AS A PREPROCESSING STEP BEFORE THE FINAL GLOBAL CLUSTERING STEP THAT FURTHER REDUCES THESE 158 CLUSTERS TO 100 CLUSTERS OUT

BIRCH WITHOUT GLOBAL CLUSTERING AS THE FINAL STEP TOOK 247 SECONDS

NCLUSTERS 158

BIRCH WITH GLOBAL CLUSTERING AS THE FINAL STEP TOOK 290 SECONDS

NCLUSTERS 100

TIME TAKEN TO RUN MINIBATCHKMEANS 360 SECONDS

AUTHORS MANOJ KUMAR MANOJKUMARSIVARAJ334GMAILCOM

ALEXANDRE GRAMFORT ALEXANDREGRAMFORTTELECOMPARISTECHFR

LICENSE BSD 3 CLAUSE

PRINTDOC

FROM ITERTOOLS IMPORT CYCLE

FROM TIME IMPORT TIME

IMPORT NUMPY AS NP

IMPORT MATPLOTLIBPYPLOT AS PLT

IMPORT MATPLOTLIBCOLORS AS COLORS

FROM SKLEARNCLUSTER IMPORT BIRCH MINIBATCHKMEANS

FROM SKLEARNDATASETSAMPLESGENERATOR IMPORT MAKEBLOBS

GENERATE CENTERS FOR THE BLOBS SO THAT IT FORMS A 10 X 10 GRID

878 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
XX NPLinspace22 22 10  
YY NPLinspace22 22 10  
XX YY NPMESHGRIDXX YY  
NCENTRES NPHSTACKNPRAVELXX NPNEWAXIS  
NPRAVELYY NPNEWAXIS  
GENERATE BLOBS TO DO A COMPARISON BETWEEN MINIBATCHKMEANS AND BIRCH  
X Y MAKEBLOBSNSAMPLES100000 CENTERSNCENTRES RANDOMSTATE0  
USE ALL COLORS THAT MATPLOTLIB PROVIDES BY DEFAULT  
COLORS CYCLECOLORSCNAMESKEYS  
FIG PLTFIGUREFIGSIZE12 4  
FIGSUBPLOTSADJUSTLEFT004 RIGHT098 BOTTOM01 TOP09  
COMPUTE CLUSTERING WITH BIRCH WITH AND WITHOUT THE FINAL CLUSTERING STEP  
AND PLOT  
BIRCHMODELS BIRCHTHRESHOLD17 NCLUSTERSNONE  
BIRCHTHRESHOLD17 NCLUSTERS100  
FINALSTEP WITHOUT GLOBAL CLUSTERING WITH GLOBAL CLUSTERING  
FORIND BIRCHMODEL INFO INENUMERATEZIPBIRCHMODELS FINALSTEP  
T TIME  
BIRCHMODELFITX  
TIME TIME T  
PRINTBIRCH SAS THE FINAL STEP TOOK 02FSECONDS  
INFO TIME T  
PLOT RESULT  
LABELS BIRCHMODELLABELS  
CENTROIDS BIRCHMODELSUBCLUSTERCENTERS  
NCLUSTERS NPUNIQUELABELSSIZE  
PRINTNCLUSTERS D NCLUSTERS  
AX FIGADDSUBPLOT1 3 IND 1  
FORTHISCENTROID K COL INZIPCENTROIDS RANGENCLUSTERS COLORS  
MASK LABELS K  
AXSCATTERXMASK 0 XMASK 1  
CW EDGECOLORCOL MARKER ALPHA05  
IFBIRCHMODELNCLUSTERS ISNONE  
AXSCATTERTHISCENTROID0 THISCENTROID1 MARKER  
CK S25  
AXSETYLIM25 25  
AXSETXLIM25 25  
AXSETAUTOSCALEYONFALSE  
AXSETTITLEBIRCH S INFO  
COMPUTE CLUSTERING WITH MINIBATCHKMEANS  
MBK MINIBATCHKMEANSINITKMEANS NCLUSTERS100 BATCHSIZE100  
NINIT10 MAXNOIMPROVEMENT10 VERBOSE0  
RANDOMSTATE0  
TO TIME  
MBKFITX  
TMINIBATCH TIME TO  
PRINTTIME TAKEN TO RUN MINIBATCHKMEANS 02FSECONDS TMINIBATCH  
MBKMEANSLABELSUNIQUE NPUNIQUEMBKLABELS  
AX FIGADDSUBPLOT1 3 3  
56 CLUSTERING 879

SCIKITLEARN USER GUIDE RELEASE 0213  
FORTHISCENTROID K COL INZIPMBKCLUSTERCENTERS  
RANGENCLUSTERS COLORS  
MASK MBKLABELS K  
AXSCATTERXMASK 0 XMASK 1 MARKER  
CW EDGECOLORCOL ALPHA05  
AXSCATTERTHISCENTROID0 THISCENTROID1 MARKER  
CK S25  
AXSETXLIM25 25  
AXSETYLIM25 25  
AXSETTITLEMINIBATCHKMEANS  
AXSETAUTOSCALEYONFALSE  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 10943 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
5620 EMPIRICAL EVALUATION OF THE IMPACT OF KMEANS INITIALIZATION  
EVALUATE THE ABILITY OF KMEANS INITIALIZATIONS STRATEGIES TO MAKE THE ALGORITHM CONVERGENCE ROBUST AS MEASURED BY  
THE RELATIVE STANDARD DEVIATION OF THE INERTIA OF THE CLUSTERING IE THE SUM OF SQUARED DISTANCES TO THE NEAREST CLUSTER  
CENTER  
THE FIRST PLOT SHOWS THE BEST INERTIA REACHED FOR EACH COMBINATION OF THE MODEL KMEANS ORMINIBATCHKMEANS  
AND THE INIT METHOD INITRANDOM ORINITKMEANS FOR INCREASING VALUES OF THE NINIT PARAMETER  
THAT CONTROLS THE NUMBER OF INITIALIZATIONS  
THE SECOND PLOT DEMONSTRATE ONE SINGLE RUN OF THE MINIBATCHKMEANS ESTIMATOR USING A INITRANDOM AND  
NINIT1 THIS RUN LEADS TO A BAD CONVERGENCE LOCAL OPTIMUM WITH ESTIMATED CENTERS STUCK BETWEEN GROUND TRUTH  
CLUSTERS  
THE DATASET USED FOR EVALUATION IS A 2D GRID OF ISOTROPIC GAUSSIAN CLUSTERS WIDELY SPACED  
880 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213

•  
OUT  
EVALUATION OF KMEANS WITH KMEANS INIT  
EVALUATION OF KMEANS WITH RANDOM INIT  
EVALUATION OF MINIBATCHKMEANS WITH KMEANS INIT  
EVALUATION OF MINIBATCHKMEANS WITH RANDOM INIT  
PRINTDOC  
AUTHOR OLIVIER GRISEL OLIVIERGRISELENSTAORG  
LICENSE BSD 3 CLAUSE  
IMPORT NUMPY AS NP  
IMPORT MATPLOTLIBPYPLOT AS PLT  
IMPORT MATPLOTLIBCM AS CM  
FROM SKLEARNUTILS IMPORT SHUFFLE  
FROM SKLEARNUTILS IMPORT CHECKRANDOMSTATE  
FROM SKLEARNCLUSTER IMPORT MINIBATCHKMEANS  
FROM SKLEARNCLUSTER IMPORT KMEANS  
RANDOMSTATE NPRANDOMRANDOMSTATE0  
882 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213  
NUMBER OF RUN WITH RANDOMLY GENERATED DATASET FOR EACH STRATEGY SO AS  
TO BE ABLE TO COMPUTE AN ESTIMATE OF THE STANDARD DEVIATION  
NRUNS 5  
KMEANS MODELS CAN DO SEVERAL RANDOM INITES SO AS TO BE ABLE TO TRADE  
CPU TIME FOR CONVERGENCE ROBUSTNESS  
NINITRANGE NPARRAY1 5 10 15 20  
DATASETS GENERATION PARAMETERS  
NSAMPLESPERCENTER 100  
GRIDSIZES 3  
SCALE 01  
NCLUSTERS GRIDSIZES 2  
DEFMAKEDATARANDOMSTATE NSAMPLESPERCENTER GRIDSIZES SCALE  
RANDOMSTATE CHECKRANDOMSTATERANDOMSTATE  
CENTERS NPARRAYIJ  
FORIINRANGEGRIDSIZES  
FORJINRANGEGRIDSIZES  
NCLUSTERSTRUE NFEATURES CENTERSSSHAPE  
NOISE RANDOMSTATENORMAL  
SCALESIZE SIZESAMPLESPERCENTER CENTERSSSHAPE1  
X NPCONCATENATEC NOISE FORCINCENTERS  
Y NPCONCATENATEI NSAMPLESPERCENTER  
FORIINRANGENCLUSTERSTRUE  
RETURNSSHUFFLEX Y RANDOMSTATERANDOMSTATE  
PART 1 QUANTITATIVE EVALUATION OF VARIOUS INIT METHODS  
PLTFigure  
PLOTS  
LEGENDS  
CASES  
KMEANS KMEANS  
KMEANS RANDOM  
MINIBATCHKMEANS KMEANS MAXNOIMPROVEMENT 3  
MINIBATCHKMEANS RANDOM MAXNOIMPROVEMENT 3 INITSIZE 500  
FORFACTORY INIT PARAMS INCASES  
PRINTEVALUATION OF SWITHSINIT FACTORYNAME INIT  
INERTIA NPEMPTYLENNINITRANGE NRUNS  
FORRUNIDINRANGENRUNS  
X Y MAKEDATARUNID NSAMPLESPERCENTER GRIDSIZES SCALE  
FORI NINIT INENUMERATENINITRANGE  
KM FACTORYNCLUSTERSNCLUSTERS INITINIT RANDOMSTATERUNID  
NINITNINIT PARAMSFITX  
INERTIAI RUNID KMINERTIA  
P PLTERRORBARNINITRANGE INERTIAI MEANAXIS1 INERTIASTDAXIS1  
PLOTSAPPENDP0  
LEGENDSAPPEND SWITHSINIT FACTORYNAME INIT  
56 CLUSTERING 883

SCIKITLEARN USER GUIDE RELEASE 0213

PLTXLABELNINIT

PLTYLABELINERTIA

PLTLEGENDPLOTS LEGENDS

PLTTITLEMEAN INERTIA FOR VARIOUS KMEANS INIT ACROSS DRUNS NRUNS

PART 2 QUALITATIVE VISUAL INSPECTION OF THE CONVERGENCE

X Y MAKEDATARANDOMSTATE NSAMPLESPERCENTER GRIDSIZE SCALE

KM MINIBATCHKMEANSNCLUSTERSNCLUSTERS INITRANDOM NINIT1

RANDOMSTATERANDOMSTATEFITX

PLTFigure

FORKINRANGENCLUSTERS

MYMEMBERS KMLABELS K

COLOR CMNIPYSPECTRALFLOATK NCLUSTERS 1

PLTPLOTXMYMEMBERS 0 XMYMEMBERS 1 O MARKER CCOLOR

CLUSTERCENTER KMCLUSTERCENTERSK

PLTPLOTCLUSTERCENTER0 CLUSTERCENTER1 O

MARKERFACECOLORCOLOR MARKEREDGECELOK MARKERSIZE6

PLTTITLEEXAMPLE CLUSTER ALLOCATION WITH A SINGLE RANDOM INIT N

WITH MINIBATCHKMEANS

PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 2695 SECONDS

NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE

5621 ADJUSTMENT FOR CHANCE IN CLUSTERING PERFORMANCE EVALUATION

THE FOLLOWING PLOTS DEMONSTRATE THE IMPACT OF THE NUMBER OF CLUSTERS AND NUMBER OF SAMPLES ON VARIOUS CLUSTERING PERFORMANCE EVALUATION METRICS

NONADJUSTED MEASURES SUCH AS THE VMEASURE SHOW A DEPENDENCY BETWEEN THE NUMBER OF CLUSTERS AND THE NUMBER OF SAMPLES THE MEAN VMEASURE OF RANDOM LABELING INCREASES SIGNIFICANTLY AS THE NUMBER OF CLUSTERS IS CLOSER TO THE TOTAL NUMBER OF SAMPLES USED TO COMPUTE THE MEASURE

ADJUSTED FOR CHANCE MEASURE SUCH AS ARI DISPLAY SOME RANDOM VARIATIONS CENTERED AROUND A MEAN SCORE OF 00 FOR ANY NUMBER OF SAMPLES AND CLUSTERS

ONLY ADJUSTED MEASURES CAN HENCE SAFELY BE USED AS A CONSENSUS INDEX TO EVALUATE THE AVERAGE STABILITY OF CLUSTERING ALGORITHMS FOR A GIVEN VALUE OF K ON VARIOUS OVERLAPPING SUBSAMPLES OF THE DATASET

884 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213

•

OUT

COMPUTING ADJUSTEDRANDSCORE FOR 10 VALUES OF NCLUSTERS AND NSAMPLES100

DONE IN 0026S

COMPUTING VMEASURESCORE FOR 10 VALUES OF NCLUSTERS AND NSAMPLES100

DONE IN 0040S

COMPUTING AMIScore FOR 10 VALUES OF NCLUSTERS AND NSAMPLES100

DONE IN 0349S

COMPUTING MUTUALINFOScore FOR 10 VALUES OF NCLUSTERS AND NSAMPLES100

DONE IN 0033S

COMPUTING ADJUSTEDRANDSCORE FOR 10 VALUES OF NCLUSTERS AND NSAMPLES1000

DONE IN 0041S

COMPUTING VMEASURESCORE FOR 10 VALUES OF NCLUSTERS AND NSAMPLES1000

DONE IN 0053S

COMPUTING AMIScore FOR 10 VALUES OF NCLUSTERS AND NSAMPLES1000

DONE IN 0208S

COMPUTING MUTUALINFOScore FOR 10 VALUES OF NCLUSTERS AND NSAMPLES1000

DONE IN 0043S

PRINTDOC

AUTHOR OLIVIER GRISEL OLIVIERGRISELENSTAORG

886 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
LICENSE BSD 3 CLAUSE  
IMPORT NUMPY AS NP  
IMPORT MATPLOTLIBPYPLOT AS PLT  
FROM TIME IMPORT TIME  
FROM SKLEARN IMPORT METRICS  
DEFUNIFORMLABELINGSSCORESSCOREFUNC NSAMPLES NCLUSTERSRANGE  
FIXEDNCLASSESNONE NRUNS5 SEED42  
COMPUTE SCORE FOR 2 RANDOM UNIFORM CLUSTER LABELINGS  
BOTH RANDOM LABELINGS HAVE THE SAME NUMBER OF CLUSTERS FOR EACH VALUE  
POSSIBLE VALUE IN NCLUSTERSRANGE  
WHEN FIXEDNCLASSES IS NOT NONE THE FIRST LABELING IS CONSIDERED A GROUND  
TRUTH CLASS ASSIGNMENT WITH FIXED NUMBER OF CLASSES

RANDOMLABELS NPRANDOMRANDOMSTATESEEDRANDINT  
SCORES NPZEROSLENNCLUSTERSRANGE NRUNS  
IFFIXEDNCLASSES IS NOTNONE  
LABELSA RANDOMLABELSLOW0 HIGHFIXEDNCLASSES SIZENSAMPLES  
FORI KINENUMERATENCLUSTERSRANGE  
FORJINRANGENRUNS  
IFFIXEDNCLASSES ISNONE  
LABELSA RANDOMLABELSLOW0 HIGHK SIZENSAMPLES  
LABELSB RANDOMLABELSLOW0 HIGHK SIZENSAMPLES  
SCORESI J SCOREFUNCLABELSA LABELSB  
RETURNSCORES  
DEFAMISCOREU V  
RETURNMETRICSADJUSTEDMUTUALINFOSCOREU V  
AVERAGEMETHODARITHMETIC  
SCOREFUNCS  
METRICSADJUSTEDRANDSCORE  
METRICSVMEASURESCORE  
AMISCORE  
METRICSMUTUALINFOSCORE

2 INDEPENDENT RANDOM CLUSTERINGS WITH EQUAL CLUSTER NUMBER  
NSAMPLES 100  
NCLUSTERSRANGE NPLinspace2 NSAMPLES 10ASTYPENPINT  
PLTFigure1  
PLOTS  
NAMES  
FORSCOREFUNC INSCOREFUNCS  
PRINTCOMPUTING SFORDVALUES OF NCLUSTERS AND NSAMPLES D  
SCOREFUNCNAME LENNCLUSTERSRANGE NSAMPLES  
T0 TIME  
SCORES UNIFORMLABELINGSSCORESSCOREFUNC NSAMPLES NCLUSTERSRANGE  
56 CLUSTERING 887

SCIKITLEARN USER GUIDE RELEASE 0213  
PRINTDONE IN 03FS TIME TO  
PLOTSAPPENDPLTERRORBAR  
NCLUSTERSRANGE NPMEDIANScores AXIS1 SCORESSTDAXIS10  
NAMESAPPENDSCOREFUNCNAME  
PLTTITLECLUSTERING MEASURES FOR 2 RANDOM UNIFORM LABELINGS N  
WITH EQUAL NUMBER OF CLUSTERS  
PLTXLABELNUMBER OF CLUSTERS NUMBER OF SAMPLES IS FIXED TO D NSAMPLES  
PLTYLABELSCORE VALUE  
PLTLEGENDPLOTS NAMES  
PLTYLIMBOTTOM005 TOP105  
RANDOM LABELING WITH VARYING NCLUSTERS AGAINST GROUND CLASS LABELS  
WITH FIXED NUMBER OF CLUSTERS  
NSAMPLES 1000  
NCLUSTERSRANGE NPLINSPACE2 100 10ASTYPENPINT  
NCLASSES 10  
PLTFigure2  
PLOTS  
NAMES  
FORSCOREFUNC INSCOREFUNCS  
PRINTCOMPUTING SFORDVALUES OF NCLUSTERS AND NSAMPLES D  
SCOREFUNCNAME LENNCLUSTERSRANGE NSAMPLES  
TO TIME  
SCORES UNIFORMLABELINGSSCORESSCOREFUNC NSAMPLES NCLUSTERSRANGE  
FIXEDNCLASSESNCLASSES  
PRINTDONE IN 03FS TIME TO  
PLOTSAPPENDPLTERRORBAR  
NCLUSTERSRANGE SCORESMEANAXIS1 SCORESSTDAXIS10  
NAMESAPPENDSCOREFUNCNAME  
PLTTITLECLUSTERING MEASURES FOR RANDOM UNIFORM LABELING N  
AGAINST REFERENCE ASSIGNMENT WITH DCLASSES NCLASSES  
PLTXLABELNUMBER OF CLUSTERS NUMBER OF SAMPLES IS FIXED TO D NSAMPLES  
PLTYLABELSCORE VALUE  
PLTYLIMBOTTOM005 TOP105  
PLTLEGENDPLOTS NAMES  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0835 SECONDS  
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE  
5622 FEATURE AGGLOMERATION VS UNIVARIATE SELECTION  
THIS EXAMPLE COMPARES 2 DIMENSIONALITY REDUCTION STRATEGIES  
• UNIVARIATE FEATURE SELECTION WITH ANOVA  
• FEATURE AGGLOMERATION WITH WARD HIERARCHICAL CLUSTERING  
888 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
BOTH METHODS ARE COMPARED IN A REGRESSION PROBLEM USING A BAYESIANRIDGE AS SUPERVISED ESTIMATOR  
OUT

MEMORY CALLING SKLEARNCLUSTERHIERARCHICALWARDTREE  
WARDTREEARRAY0451933 0675318

0275706 1085711  
1600X1600 SPARSE MATRIX OF TYPE CLASS NUMPYINT64  
WITH 7840 STORED ELEMENTS IN COORDINATE FORMAT NCLUSTERSNONE RETURN  
'→DISTANCEFALSE  
WARDTREE 01S 00MIN

MEMORY CALLING SKLEARNCLUSTERHIERARCHICALWARDTREE  
WARDTREEARRAY 0905206 0161245

0849835 1091621  
1600X1600 SPARSE MATRIX OF TYPE CLASS NUMPYINT64  
WITH 7840 STORED ELEMENTS IN COORDINATE FORMAT NCLUSTERSNONE RETURN  
'→DISTANCEFALSE  
WARDTREE 01S 00MIN

MEMORY CALLING SKLEARNCLUSTERHIERARCHICALWARDTREE  
WARDTREEARRAY 0905206 0675318

0849835 1085711  
1600X1600 SPARSE MATRIX OF TYPE CLASS NUMPYINT64  
WITH 7840 STORED ELEMENTS IN COORDINATE FORMAT NCLUSTERSNONE RETURN  
'→DISTANCEFALSE  
WARDTREE 01S 00MIN

MEMORY CALLING SKLEARNFEATURESELECTIONUNIVARIATESELECTIONFREGRESSION  
FREGRESSIONARRAY0451933 0275706

0675318 1085711  
ARRAY 25267703 25026711  
FREGRESSION 00S 00MIN

MEMORY CALLING SKLEARNFEATURESELECTIONUNIVARIATESELECTIONFREGRESSION  
FREGRESSIONARRAY 0905206 0849835

56 CLUSTERING 889

```
SCIKITLEARN USER GUIDE RELEASE 0213
0161245 1091621
ARRAY 27447268 112638768
FREGRESSION 00S 00MIN

MEMORY CALLING SKLEARNFEATURESELECTIONUNIVARIATESELECTIONFREGRESSION
FREGRESSIONARRAY 0905206 0849835

0675318 1085711
ARRAY27447268 25026711
FREGRESSION 00S 00MIN
AUTHOR ALEXANDRE GRAMFORT ALEXANDREGGRAMFORTINRIA.FR
LICENSE BSD 3 CLAUSE
PRINTDOC
IMPORT SHUTIL
IMPORT TEMPFILE
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPY.PLOT AS PLT
FROM SCIPY IMPORT LINALG NDIMAGE
FROM JOBLIB IMPORT MEMORY
FROM SKLEARNFEATUREEXTRACTIONIMAGE IMPORT GRIDTOGRAPH
FROM SKLEARN IMPORT FEATURESELECTION
FROM SKLEARNCLUSTER IMPORT FEATUREAGGLOMERATION
FROM SKLEARNLINEARMODEL IMPORT BAYESIANRIDGE
FROM SKLEARNPIPELINE IMPORT PIPELINE
FROM SKLEARNMODELSELECTION IMPORT GRIDSEARCHCV
FROM SKLEARNMODELSELECTION IMPORT KFOLD

GENERATE DATA
NSAMPLES 200
SIZE 40 IMAGE SIZE
ROISIZE 15
SNR 5
NPRANDOMSEED0
MASK NPONESIZE SIZE DTYPE.NP.BOOL
COEF NPZEROSIZE SIZE
COEFROISIZE 0ROISIZE 1
COEFROISIZE ROISIZE 1
X NPRANDOMRANDNNSAMPLES SIZE 2
FORXINX SMOOTH DATA
X NDIMAGEGAUSSIANFILTERXRESHAPE.SIZE SIZE SIGMA10RAVEL
X XMEANAXIS0
X XSTDAXIS0
Y NPDOTX COEFRAVEL
NOISE NPRANDOMRANDNYSHAPE0
890 CHAPTER 5 EXAMPLES
```



```
SCIKITLEARN USER GUIDE RELEASE 0213
NOISECOEF LINALGNORMY 2 NPEXPSNR 20 LINALGNORMNOISE 2
Y NOISECOEF NOISE ADD NOISE

COMPUTE THE COEFS OF A BAYESIAN RIDGE WITH GRIDSEARCH
CV KFOLD2 CROSSVALIDATION GENERATOR FOR MODEL SELECTION
RIDGE BAYESIANRIDGE
CACHEDIR TEMPFILEMKDTEMP
MEM MEMORYLOCATIONCACHEDIR VERBOSE1
WARD AGGLOMERATION FOLLOWED BY BAYESIANRIDGE
CONNECTIVITY GRIDTOGRAPHNXSIZE NYSIZE
WARD FEATUREAGGLOMERATIONNNCLUSTERS10 CONNECTIVITYCONNECTIVITY
MEMORYMEM
CLF PIPELINEWARD WARD RIDGE RIDGE
SELECT THE OPTIMAL NUMBER OF PARCELS WITH GRID SEARCH
CLF GRIDSEARCHCVCLF WARDNCLUSTERS 10 20 30 NJOBS1 CVCV
CLFFITX Y SET THE BEST PARAMETERS
COEF CLFBESTESTIMATORSTEPS11COEF
COEF CLFBESTESTIMATORSTEPS01INVERSETRANSFORMCOEF
COEFAGGLOMERATION COEFRESHAPESIZE SIZE
ANOVA UNIVARIATE FEATURE SELECTION FOLLOWED BY BAYESIANRIDGE
FREGRESSION MEMCACHEFEATURESELECTIONFREGRESSION CACHING FUNCTION
ANOVA FEATURESELECTIONSELECTPERCENTILEFREGRESSION
CLF PIPELINEANOVA ANOVA RIDGE RIDGE
SELECT THE OPTIMAL PERCENTAGE OF FEATURES WITH GRID SEARCH
CLF GRIDSEARCHCVCLF ANOVAPERCENTILE 5 10 20 CVCV
CLFFITX Y SET THE BEST PARAMETERS
COEF CLFBESTESTIMATORSTEPS11COEF
COEF CLFBESTESTIMATORSTEPS01INVERSETRANSFORMCOEFRESHAPE1 1
COEFSELECTION COEFRESHAPESIZE SIZE

INVERSE THE TRANSFORMATION TO PLOT THE RESULTS ON AN IMAGE
PLTCLOSEALL
PLTFIGUREFIGSIZE73 27
PLTSUBPLOT1 3 1
PLTIMSHOWCOEF INTERPOLATIONNEAREST CMAPPLTCMRDBUR
PLTTITLETRUE WEIGHTS
PLTSUBPLOT1 3 2
PLTIMSHOWCOEFSELECTION INTERPOLATIONNEAREST CMAPPLTCMRDBUR
PLTTITLEFEATURE SELECTION
PLTSUBPLOT1 3 3
PLTIMSHOWCOEFAGGLOMERATION INTERPOLATIONNEAREST CMAPPLTCMRDBUR
PLTTITLEFEATURE AGGLOMERATION
PLTSUBPLOTSADJUST004 00 098 094 016 026
PLTSHOW
ATTEMPT TO REMOVE THE TEMPORARY CACHEDIR BUT DONT WORRY IF IT FAILS
SHUTILRMTREECACHEDIR IGNOREERRORSTRUE
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0623 SECONDS
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE
56 CLUSTERING 891
```

```
SCIKITLEARN USER GUIDE RELEASE 0213
5623 COMPARISON OF THE KMEANS AND MINIBATCHKMEANS CLUSTERING ALGORITHMS
WE WANT TO COMPARE THE PERFORMANCE OF THE MINIBATCHKMEANS AND KMEANS THE MINIBATCHKMEANS IS FASTER BUT
GIVES SLIGHTLY DIFFERENT RESULTS SEE MINI BATCH KMEANS
WE WILL CLUSTER A SET OF DATA FIRST WITH KMEANS AND THEN WITH MINIBATCHKMEANS AND PLOT THE RESULTS WE WILL ALSO
PLOT THE POINTS THAT ARE LABELLED DIFFERENTLY BETWEEN THE TWO ALGORITHMS
PRINTDOC
IMPORT TIME
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARNCLUSTER IMPORT MINIBATCHKMEANS KMEANS
FROM SKLEARNMETRICSPAIRWISE IMPORT PAIRWISEDISTANCESARGMIN
FROM SKLEARNDATASETSAMPLESGENERATOR IMPORT MAKEBLOBS

GENERATE SAMPLE DATA
NPRANDOMSEED0
BATCHSIZE 45
CENTERS 1 1 1 1 1 1
NCLUSTERS LENCENTERS
X LABELSTRUE MAKEBLOBSNSAMPLES3000 CENTERSCENTERS CLUSTERSTD07

COMPUTE CLUSTERING WITH MEANS
KMEANS KMEANSINITKMEANS NCLUSTERS3 NINIT10
T0 TIMETIME
KMEANSFITX
TBATCH TIMETIME T0

COMPUTE CLUSTERING WITH MINIBATCHKMEANS
MBK MINIBATCHKMEANSINITKMEANS NCLUSTERS3 BATCHSIZEBATCHSIZE
NINIT10 MAXNOIMPROVEMENT10 VERBOSE0
T0 TIMETIME
892 CHAPTER 5 EXAMPLES
```

SCIKITLEARN USER GUIDE RELEASE 0213  
MBKFITX  
TMINIBATCH TIMETIME TO

PLOT RESULT  
FIG PLTFIGUREFIGSIZE8 3  
FIGSUBPLOTSADJUSTLEFT002 RIGHT098 BOTTOM005 TOP09  
COLORS 4EACC5 FF9C34 4E9A06  
WE WANT TO HAVE THE SAME COLORS FOR THE SAME CLUSTER FROM THE  
MINIBATCHKMEANS AND THE KMEANS ALGORITHM LETS PAIR THE CLUSTER CENTERS PER  
CLOSEST ONE  
KMEANSCLUSTERCENTERS NPSORTKMEANSCLUSTERCENTERS AXIS0  
MBKMEANSCLUSTERCENTERS NPSORTMBKCLUSTERCENTERS AXIS0  
KMEANSLABELS PAIRWISEDISTANCESARGMINX KMEANSCLUSTERCENTERS  
MBKMEANSLABELS PAIRWISEDISTANCESARGMINX MBKMEANSCLUSTERCENTERS  
ORDER PAIRWISEDISTANCESARGMINKMEANSCLUSTERCENTERS  
MBKMEANSCLUSTERCENTERS  
KMEANS  
AX FIGADDSUBPLOT1 3 1  
FORK COLINZIPRANGENCLUSTERS COLORS  
MYMEMBERS KMEANSLABELS K  
CLUSTERCENTER KMEANSCLUSTERCENTERSK  
AXPLOTXMYMEMBERS 0 XMYMEMBERS 1 W  
MARKERFACECOLORCOL MARKER  
AXPLOTCLUSTERCENTER0 CLUSTERCENTER1 O MARKERFACECOLORCOL  
MARKEREDGECOLORK MARKERSIZE6  
AXSETTITLEKMEANS  
AXSETXTICKS  
AXSETYTIMES  
PLTTEXT35 18 TRAIN TIME 2FSNINERTIA F  
TBATCH KMEANSINERTIA  
MINIBATCHKMEANS  
AX FIGADDSUBPLOT1 3 2  
FORK COLINZIPRANGENCLUSTERS COLORS  
MYMEMBERS MBKMEANSLABELS ORDERK  
CLUSTERCENTER MBKMEANSCLUSTERCENTERSORDERK  
AXPLOTXMYMEMBERS 0 XMYMEMBERS 1 W  
MARKERFACECOLORCOL MARKER  
AXPLOTCLUSTERCENTER0 CLUSTERCENTER1 O MARKERFACECOLORCOL  
MARKEREDGECOLORK MARKERSIZE6  
AXSETTITLEMINIBATCHKMEANS  
AXSETXTICKS  
AXSETYTIMES  
PLTTEXT35 18 TRAIN TIME 2FSNINERTIA F  
TMINIBATCH MBKINERTIA  
INITIALISE THE DIFFERENT ARRAY TO ALL FALSE  
DIFFERENT MBKMEANSLABELS 4  
AX FIGADDSUBPLOT1 3 3  
FORKINRANGENCLUSTERS  
DIFFERENT KMEANSLABELS K MBKMEANSLABELS ORDERK  
IDENTIC NPLOGICALNOTDIFFERENT  
56 CLUSTERING 893

SCIKITLEARN USER GUIDE RELEASE 0213  
AXPLOTXIDENTIC 0 XIDENTIC 1 W  
MARKERFACECOLORBBBBBB MARKER  
AXPLOTXDIFFERENT 0 XDIFFERENT 1 W  
MARKERFACECOLORM MARKER  
AXSETTITLEDIFFERENCE  
AXSETXTICKS  
AXSETYTICKS  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0127 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
5624 A DEMO OF KMEANS CLUSTERING ON THE HANDWRITTEN DIGITS DATA  
IN THIS EXAMPLE WE COMPARE THE VARIOUS INITIALIZATION STRATEGIES FOR KMEANS IN TERMS OF RUNTIME AND QUALITY OF THE RESULTS  
AS THE GROUND TRUTH IS KNOWN HERE WE ALSO APPLY DIFFERENT CLUSTER QUALITY METRICS TO JUDGE THE GOODNESS OF FIT OF THE CLUSTER LABELS TO THE GROUND TRUTH  
CLUSTER QUALITY METRICS EVALUATED SEE CLUSTERING PERFORMANCE EVALUATION FOR DEFINITIONS AND DISCUSSIONS OF THE METRICS  
SHORTHAND FULL NAME  
HOMO HOMOGENEITY SCORE  
COMPL COMPLETENESS SCORE  
VMEAS V MEASURE  
ARI ADJUSTED RAND INDEX  
AMI ADJUSTED MUTUAL INFORMATION  
SILHOUETTE SILHOUETTE COEFFICIENT  
894 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
OUT
NDIGITS 10 NSAMPLES 1797 NFEATURES 64

INIT TIME INERTIA HOMO COMPL VMEAS ARI AMI SILHOUETTE
KMEANS 037S 69432 0602 0650 0625 0465 0621 0146
RANDOM 027S 69694 0669 0710 0689 0553 0686 0147
PCABASED 003S 70804 0671 0698 0684 0561 0681 0118

PRINTDOC
FROM TIME IMPORT TIME
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARN IMPORT METRICS
FROM SKLEARNCLUSTER IMPORT KMEANS
FROM SKLEARNDATASETS IMPORT LOADDIGITS
FROM SKLEARNDECOMPOSITION IMPORT PCA
FROM SKLEARNPREPROCESSING IMPORT SCALE
56 CLUSTERING 895
```

SCIKITLEARN USER GUIDE RELEASE 0213  
NPRANDOMSEED42  
DIGITS LOADDIGITS  
DATA SCALEDIGITS  
DATA SCALEDIGITS  
NSAMPLES NFEATURES DATASHAPE  
NDIGITS LENNPUNIQUEDIGITSTARGET  
LABELS DIGITSTARGET  
SAMPLESIZE 300  
PRINTNDIGITS DTNSAMPLES DTNFEATURES D  
NDIGITS NSAMPLES NFEATURES  
PRINT82  
PRINTINITTTTTIMETINERTIATHOMOTCOMPLTVMEASTARITAMITSILHOUETTE  
DEFBENCHKMEANSESTIMATOR NAME DATA  
TO TIME  
ESTIMATORFITDATA  
PRINT9ST2FSTIT3FT3FT3FT3FT3F  
NAME TIME TO ESTIMATORINERTIA  
METRICSHOMOGENEITYSCORELABELS ESTIMATORLABELS  
METRICSCOMPLETENESSSCORELABELS ESTIMATORLABELS  
METRICSVMEASURESCORELABELS ESTIMATORLABELS  
METRICSDJUSTEDRANDSCORELABELS ESTIMATORLABELS  
METRICSDJUSTEDMUTUALINFOSCORELABELS ESTIMATORLABELS  
AVERAGEMETHODARITHMETIC  
METRICSSILHOUETTESCOREDATA ESTIMATORLABELS  
METRICEUCLIDEAN  
SAMPLESIZESAMPLESIZE  
BENCHKMEANSKMEANSINITKMEANS NCLUSTERSNDIGITS NINIT10  
NAMEKMEANS DATADATA  
BENCHKMEANSKMEANSINITRANDOM NCLUSTERSNDIGITS NINIT10  
NAMERANDOM DATADATA  
IN THIS CASE THE SEEDING OF THE CENTERS IS DETERMINISTIC HENCE WE RUN THE  
KMEANS ALGORITHM ONLY ONCE WITH NINIT1  
PCA PCANCOMPONENTSNDIGITSFITDATA  
BENCHKMEANSKMEANSINITPCACOMPONENTS NCLUSTERSNDIGITS NINIT1  
NAMEPCABASED  
DATADATA  
PRINT82  
  
VISUALIZE THE RESULTS ON PCAREduced DATA  
REDUCEDDATA PCANCOMPONENTS2FITTRANSFORMDATA  
KMEANS KMEANSINITKMEANS NCLUSTERSNDIGITS NINIT10  
KMEANSFITREDUCEDDATA  
STEP SIZE OF THE MESH DECREASE TO INCREASE THE QUALITY OF THE VQ  
H 02 POINT IN THE MESH XMIN XMAXXYMIN YMAX  
896 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

PLOT THE DECISION BOUNDARY FOR THAT WE WILL ASSIGN A COLOR TO EACH  
XMIN XMAX REDUCEDDATA 0MIN 1 REDUCEDDATA 0MAX 1  
YMIN YMAX REDUCEDDATA 1MIN 1 REDUCEDDATA 1MAX 1  
XX YY NPMESHGRIDNPARANGEXMIN XMAX H NPARANGHEYMIN YMAX H  
OBTAIN LABELS FOR EACH POINT IN MESH USE LAST TRAINED MODEL  
Z KMEANSPREDICTNPCXXRAVEL YYRAVEL  
PUT THE RESULT INTO A COLOR PLOT  
Z ZRESHAPEXXSHAPE  
PLTFigure1  
PLTCLF  
PLTImSHOWZ INTERPOLATIONNEAREST  
EXTENTXXMIN XXMAX YYMIN YYMAX  
CMAPPLTCMPAIED  
ASPECTAUTO ORIGINLOWER  
PLTPLOTREDUCEDDATA 0 REDUCEDDATA 1 K MARKERSIZE2  
PLOT THE CENTROIDS AS A WHITE X  
CENTROIDS KMEANSCLUSTERCENTERS  
PLTSCATTERCENTROIDS 0 CENTROIDS 1  
MARKERX S169 LINEWIDTHS3  
COLORW ZORDER10  
PLTTITLEKMEANS CLUSTERING ON THE DIGITS DATASET PCAREduced DATA N  
CENTROIDS ARE MARKED WITH WHITE CROSS  
PLTXLIMXMIN XMAX  
PLTYLIMYMIN YMAX  
PLXTTICKS  
PLTYTICKS  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 1230 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
5625 COMPARING DIFFERENT HIERARCHICAL LINKAGE METHODS ON TOY DATASETS  
THIS EXAMPLE SHOWS CHARACTERISTICS OF DIFFERENT LINKAGE METHODS FOR HIERARCHICAL CLUSTERING ON DATASETS THAT ARE “INTERESTING” BUT STILL IN 2D  
THE MAIN OBSERVATIONS TO MAKE ARE

- SINGLE LINKAGE IS FAST AND CAN PERFORM WELL ON NONGLOBULAR DATA BUT IT PERFORMS POORLY IN THE PRESENCE OF NOISE
- AVERAGE AND COMPLETE LINKAGE PERFORM WELL ON CLEANLY SEPARATED GLOBULAR CLUSTERS BUT HAVE MIXED RESULTS OTHERWISE
- WARD IS THE MOST EFFECTIVE METHOD FOR NOISY DATA

WHILE THESE EXAMPLES GIVE SOME INTUITION ABOUT THE ALGORITHMS THIS INTUITION MIGHT NOT APPLY TO VERY HIGH DIMENSIONAL DATA

PRINTDOC  
IMPORT TIME  
56 CLUSTERING 897

SCIKITLEARN USER GUIDE RELEASE 0213

IMPORT WARNINGS

IMPORT NUMPY AS NP

IMPORT MATPLOTLIBPYPLOT AS PLT

FROM SKLEARN IMPORT CLUSTER DATASETS

FROM SKLEARNPREPROCESSING IMPORT STANDARDSCALER

FROM ITERTOOLS IMPORT CYCLE ISLICE

NPRANDOMSEED0

GENERATE DATASETS WE CHOOSE THE SIZE BIG ENOUGH TO SEE THE SCALABILITY OF THE ALGORITHMS BUT NOT TOO BIG TO AVOID TOO LONG RUNNING TIMES

NSAMPLES 1500

NOISYCIRCLES DATASETSMAKECIRCLESNSAMPLESNSAMPLES FACTOR5

NOISE05

NOISYMOONS DATASETSMAKEMOONSNSAMPLESNSAMPLES NOISE05

BLOBS DATASETSMAKEBLOBSNSAMPLESNSAMPLES RANDOMSTATE8

NOSTRUCTURE NPRANDOMRANDNSAMPLES 2 NONE

ANISOTROPICLY DISTRIBUTED DATA

RANDOMSTATE 170

X Y DATASETSMAKEBLOBSNSAMPLESNSAMPLES RANDOMSTATERANDOMSTATE

TRANSFORMATION 06 06 04 08

XANISO NPDOTX TRANSFORMATION

ANISO XANISO Y

BLOBS WITH VARIED VARIANCES

VARIED DATASETSMAKEBLOBSNSAMPLESNSAMPLES

CLUSTERSTD10 25 05

RANDOMSTATERANDOMSTATE

RUN THE CLUSTERING AND PLOT

SET UP CLUSTER PARAMETERS

PLTFIGUREFIGSIZE9 13 2 145

PLTSUBPLOTSADJUSTLEFT02 RIGHT98 BOTTOM001 TOP96 WSPACE05

HSPACE01

PLOTNUM 1

DEFAULTBASE NNEIGHBORS 10

NCLUSTERS 3

DATASETS

NOISYCIRCLES NCLUSTERS 2

NOISYMOONS NCLUSTERS 2

VARIED NNEIGHBORS 2

ANISO NNEIGHBORS 2

BLOBS

NOSTRUCTURE

FORIDATASET DATASET ALGOPARAMS INENUMERATEDDATASETS

UPDATE PARAMETERS WITH DATASETSPECIFIC VALUES

PARAMS DEFAULTBASECOPY

PARAMSUPDATEALGOPARAMS

898 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213  
X Y DATASET  
NORMALIZE DATASET FOR EASIER PARAMETER SELECTION  
X STANDARDSCALERFITTRANSFORMX  
  
CREATE CLUSTER OBJECTS  
  
WARD CLUSTERAGGLOMERATIVECLUSTERING  
NCLUSTERSPARAMSNCLUSTERS LINKAGEWARD  
COMPLETE CLUSTERAGGLOMERATIVECLUSTERING  
NCLUSTERSPARAMSNCLUSTERS LINKAGECOMPLETE  
AVERAGE CLUSTERAGGLOMERATIVECLUSTERING  
NCLUSTERSPARAMSNCLUSTERS LINKAGEAVERAGE  
SINGLE CLUSTERAGGLOMERATIVECLUSTERING  
NCLUSTERSPARAMSNCLUSTERS LINKAGESINGLE  
CLUSTERINGALGORITHMS  
SINGLE LINKAGE SINGLE  
AVERAGE LINKAGE AVERAGE  
COMPLETE LINKAGE COMPLETE  
WARD LINKAGE WARD  
  
FORNAME ALGORITHM INCLUSTERINGALGORITHMS  
T0 TIMETIME  
CATCH WARNINGS RELATED TO KNEIGHBORSGRAPH  
WITHWARNINGSCATCHWARNINGS  
WARNINGSFILTERWARNINGS  
IGNORE  
MESSAGE THE NUMBER OF CONNECTED COMPONENTS OF THE  
CONNECTIVITY MATRIX IS 0912  
1 COMPLETING IT TO AVOID STOPPING THE TREE EARLY  
CATEGORYUSERWARNING  
ALGORITHMFITX  
T1 TIMETIME  
IFHASATTRALGORITHM LABELS  
YPRED ALGORITHM LABELSASTYPENPINT  
ELSE  
YPRED ALGORITHM PREDICTX  
PLTSUBPLOTLENDATASETS LENCLUSTERINGALGORITHMS PLOTNUM  
IFIDATASET 0  
PLTTITLENAME SIZE18  
COLORS NPARRAYLISTISLICECYCLE377EB8 FF7F00 4DAF4A  
F781BF A65628 984EA3  
999999 E41A1C DEDE00  
INTMAXYPRED 1  
PLTSCATTERX 0 X 1 S10 COLORCOLORSYYPRED  
PLTXLIM25 25  
PLTYLIM25 25  
PLXTTICKS  
PLTYTICKS  
PLTTEXT99 01 2FS T1 TOLSTRIPO  
56 CLUSTERING 899

SCIKITLEARN USER GUIDE RELEASE 0213  
TRANSFORMPLTGCATRANSAXES SIZE15  
HORIZONTALALIGNMENTRIGHT  
PLOTNUM 1  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 1577 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
900 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
5626 SELECTING THE NUMBER OF CLUSTERS WITH SILHOUETTE ANALYSIS ON KMEANS CLUSTERING

SILHOUETTE ANALYSIS CAN BE USED TO STUDY THE SEPARATION DISTANCE BETWEEN THE RESULTING CLUSTERS THE SILHOUETTE PLOT DISPLAYS A MEASURE OF HOW CLOSE EACH POINT IN ONE CLUSTER IS TO POINTS IN THE NEIGHBORING CLUSTERS AND THUS PROVIDES A WAY TO ASSESS PARAMETERS LIKE NUMBER OF CLUSTERS VISUALLY THIS MEASURE HAS A RANGE OF 1 1  
SILHOUETTE COEFFICIENTS AS THESE VALUES ARE REFERRED TO AS NEAR 1 INDICATE THAT THE SAMPLE IS FAR AWAY FROM THE NEIGHBORING CLUSTERS A VALUE OF 0 INDICATES THAT THE SAMPLE IS ON OR VERY CLOSE TO THE DECISION BOUNDARY BETWEEN TWO NEIGHBORING CLUSTERS AND NEGATIVE VALUES INDICATE THAT THOSE SAMPLES MIGHT HAVE BEEN ASSIGNED TO THE WRONG CLUSTER  
IN THIS EXAMPLE THE SILHOUETTE ANALYSIS IS USED TO CHOOSE AN OPTIMAL VALUE FOR NCLUSTERS THE SILHOUETTE PLOT SHOWS THAT THENCLUSTERS VALUE OF 3 5 AND 6 ARE A BAD PICK FOR THE GIVEN DATA DUE TO THE PRESENCE OF CLUSTERS WITH BELOW AVERAGE SILHOUETTE SCORES AND ALSO DUE TO WIDE FLUCTUATIONS IN THE SIZE OF THE SILHOUETTE PLOTS SILHOUETTE ANALYSIS IS MORE AMBIVALENT IN DECIDING BETWEEN 2 AND 4  
ALSO FROM THE THICKNESS OF THE SILHOUETTE PLOT THE CLUSTER SIZE CAN BE VISUALIZED THE SILHOUETTE PLOT FOR CLUSTER 0 WHEN NCLUSTERS IS EQUAL TO 2 IS BIGGER IN SIZE OWING TO THE GROUPING OF THE 3 SUB CLUSTERS INTO ONE BIG CLUSTER HOWEVER WHEN THENCLUSTERS IS EQUAL TO 4 ALL THE PLOTS ARE MORE OR LESS OF SIMILAR THICKNESS AND HENCE ARE OF SIMILAR SIZES AS CAN BE ALSO VERIFIED FROM THE LABELLED SCATTER PLOT ON THE RIGHT

- 
- 

56 CLUSTERING 901

SCIKITLEARN USER GUIDE RELEASE 0213

- 
- 
- 

OUT

FOR NCLUSTERS 2 THE AVERAGE SILHOUETTESCORE IS 07049787496083262  
FOR NCLUSTERS 3 THE AVERAGE SILHOUETTESCORE IS 05882004012129721  
FOR NCLUSTERS 4 THE AVERAGE SILHOUETTESCORE IS 06505186632729437  
FOR NCLUSTERS 5 THE AVERAGE SILHOUETTESCORE IS 056376469026194  
FOR NCLUSTERS 6 THE AVERAGE SILHOUETTESCORE IS 04504666294372765  
902 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
FROM SKLEARNDATASETS IMPORT MAKEBLOBS
FROM SKLEARNCLUSTER IMPORT KMEANS
FROM SKLEARNMETRICS IMPORT SILHOUETTESAMPLES SILHOUETTESCORE
IMPORT MATPLOTLIBPYPLOT AS PLT
IMPORT MATPLOTLIBCM AS CM
IMPORT NUMPY AS NP
PRINTDOC
GENERATING THE SAMPLE DATA FROM MAKEBLOBS
THIS PARTICULAR SETTING HAS ONE DISTINCT CLUSTER AND 3 CLUSTERS PLACED CLOSE
TOGETHER
X Y MAKEBLOBSNSAMPLES500
NFEATURES2
CENTERS4
CLUSTERSTD1
CENTERBOX100 100
SHUFFLETRUE
RANDOMSTATE1 FOR REPRODUCIBILITY
RANGENCLUSTERS 2 3 4 5 6
FORNCLUSTERS INRANGENCLUSTERS
CREATE A SUBPLOT WITH 1 ROW AND 2 COLUMNS
FIG AX1 AX2 PLTSUBPLOTS1 2
FIGSETSIZEINCHES18 7
THE 1ST SUBPLOT IS THE SILHOUETTE PLOT
THE SILHOUETTE COEFFICIENT CAN RANGE FROM 1 1 BUT IN THIS EXAMPLE ALL
LIE WITHIN 01 1
AX1SETXLIM01 1
THE NCLUSTERS1 10 IS FOR INSERTING BLANK SPACE BETWEEN SILHOUETTE
PLOTS OF INDIVIDUAL CLUSTERS TO DEMARCATATE THEM CLEARLY
AX1SETYLIM0 LENX NCLUSTERS 1 10
INITIALIZE THE CLUSTERER WITH NCLUSTERS VALUE AND A RANDOM GENERATOR
SEED OF 10 FOR REPRODUCIBILITY
CLUSTERER KMEANSNCLUSTERSNCLUSTERS RANDOMSTATE10
CLUSTERLABELS CLUSTERERFITPREDICTX
THE SILHOUETTESCORE GIVES THE AVERAGE VALUE FOR ALL THE SAMPLES
THIS GIVES A PERSPECTIVE INTO THE DENSITY AND SEPARATION OF THE FORMED
CLUSTERS
SILHOUETTEAVG SILHOUETTESCOREX CLUSTERLABELS
PRINTFOR NCLUSTERS NCLUSTERS
THE AVERAGE SILHOUETTESCORE IS SILHOUETTEAVG
COMPUTE THE SILHOUETTE SCORES FOR EACH SAMPLE
SAMPLESILHOUETTEVALUES SILHOUETTESAMPLESX CLUSTERLABELS
YLOWER 10
FORIINRANGENCLUSTERS
AGGREGATE THE SILHOUETTE SCORES FOR SAMPLES BELONGING TO
CLUSTER I AND SORT THEM
56 CLUSTERING 903
```

SCIKITLEARN USER GUIDE RELEASE 0213  
ITHCLUSTERSILHOUETTEVALUES  
SAMPLESILHOUETTEVALUESCLUSTERLABELS I  
ITHCLUSTERSILHOUETTEVALUESSORT  
SIZECLUSTERI ITHCLUSTERSILHOUETTEVALUESSHAPE0  
YUPPER YLOWER SIZECLUSTERI  
COLOR CMNIPYSPECTRALFLOATI NCLUSTERS  
AX1FILLBETWEENXNPARANGEYLOWER YUPPER  
0 ITHCLUSTERSILHOUETTEVALUES  
FACECOLORCOLOR EDGECOLORCOLOR ALPHA07  
LABEL THE SILHOUETTE PLOTS WITH THEIR CLUSTER NUMBERS AT THE MIDDLE  
AX1TEXT005 YLOWER 05 SIZECLUSTERI STRI  
COMPUTE THE NEW YLOWER FOR NEXT PLOT  
YLOWER YUPPER 10 10 FOR THE 0 SAMPLES  
AX1SETTITLETHE SILHOUETTE PLOT FOR THE VARIOUS CLUSTERS  
AX1SETXLABELTHE SILHOUETTE COEFFICIENT VALUES  
AX1SETYLABELCLUSTER LABEL  
THE VERTICAL LINE FOR AVERAGE SILHOUETTE SCORE OF ALL THE VALUES  
AX1AXVLINEXSILHOUETTEAVG COLORRED LINESTYLE  
AX1SETYTICKS CLEAR THE YAXIS LABELS TICKS  
AX1SETXTICKS01 0 02 04 06 08 1  
2ND PLOT SHOWING THE ACTUAL CLUSTERS FORMED  
COLORS CMNIPYSPECTRALCLUSTERLABELSASTYPEFLOAT NCLUSTERS  
AX2SCATTERX 0 X 1 MARKER S30 LW0 ALPHA07  
CCOLORS EDGECOLORK  
LABELING THE CLUSTERS  
CENTERS CLUSTERERCLUSTERCENTERS  
DRAW WHITE CIRCLES AT CLUSTER CENTERS  
AX2SCATTERCENTERS 0 CENTERS 1 MARKERO  
CWHITE ALPHA1 S200 EDGECOLORK  
FORI CINENUMERATECENTERS  
AX2SCATTERC0 C1 MARKER D I ALPHA1  
S50 EDGECOLORK  
AX2SETTITLETHE VISUALIZATION OF THE CLUSTERED DATA  
AX2SETXLABELFEATURE SPACE FOR THE 1ST FEATURE  
AX2SETYLABELFEATURE SPACE FOR THE 2ND FEATURE  
PLTSUPTITLESILHOUETTE ANALYSIS FOR KMEANS CLUSTERING ON SAMPLE DATA  
WITH NCLUSTERS D NCLUSTERS  
FONTSIZE14 FONTWEIGHTBOLD  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0629 SECONDS  
904 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE
5627 COMPARING DIFFERENT CLUSTERING ALGORITHMS ON TOY DATASETS
THIS EXAMPLE SHOWS CHARACTERISTICS OF DIFFERENT CLUSTERING ALGORITHMS ON DATASETS THAT ARE “INTERESTING” BUT STILL IN 2D
WITH THE EXCEPTION OF THE LAST DATASET THE PARAMETERS OF EACH OF THESE DATASETALGORITHM PAIRS HAS BEEN TUNED TO PRODUCE
GOOD CLUSTERING RESULTS SOME ALGORITHMS ARE MORE SENSITIVE TO PARAMETER VALUES THAN OTHERS
THE LAST DATASET IS AN EXAMPLE OF A ‘NULL’ SITUATION FOR CLUSTERING THE DATA IS HOMOGENEOUS AND THERE IS NO GOOD
CLUSTERING FOR THIS EXAMPLE THE NULL DATASET USES THE SAME PARAMETERS AS THE DATASET IN THE ROW ABOVE IT WHICH
REPRESENTS A MISMATCH IN THE PARAMETER VALUES AND THE DATA STRUCTURE
WHILE THESE EXAMPLES GIVE SOME INTUITION ABOUT THE ALGORITHMS THIS INTUITION MIGHT NOT APPLY TO VERY HIGH DIMENSIONAL
DATA
PRINTDOC
IMPORT TIME
IMPORT WARNINGS
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARN IMPORT CLUSTER DATASETS MIXTURE
FROM SKLEARNNEIGHBORS IMPORT KNEIGHBORSGRAPH
FROM SKLEARNPREPROCESSING IMPORT STANDARDSCALER
FROM ITERTOOLS IMPORT CYCLE ISLICE
NPRANDOMSEED0
56 CLUSTERING 905
```

SCIKITLEARN USER GUIDE RELEASE 0213

GENERATE DATASETS WE CHOOSE THE SIZE BIG ENOUGH TO SEE THE SCALABILITY  
OF THE ALGORITHMS BUT NOT TOO BIG TO AVOID TOO LONG RUNNING TIMES

NSAMPLES 1500  
NOISYCIRCLES DATASETSMAKECIRCLESNSAMPLESNSAMPLES FACTOR5  
NOISE05  
NOISYMOONS DATASETSMAKEMOONSNSAMPLESNSAMPLES NOISE05  
BLOBS DATASETSMAKEBLOBSNSAMPLESNSAMPLES RANDOMSTATE8  
NOSTRUCTURE NPRANDOMRANDNSAMPLES 2 NONE  
ANISOTROPICLY DISTRIBUTED DATA  
RANDOMSTATE 170  
X Y DATASETSMAKEBLOBSNSAMPLESNSAMPLES RANDOMSTATERANDOMSTATE  
TRANSFORMATION 06 06 04 08  
XANISO NPDOTX TRANSFORMATION  
ANISO XANISO Y  
BLOBS WITH VARIED VARIANCES  
VARIED DATASETSMAKEBLOBSNSAMPLESNSAMPLES  
CLUSTERSTD10 25 05  
RANDOMSTATERANDOMSTATE

SET UP CLUSTER PARAMETERS

PLTFIGUREFIGSIZE9 2 3 125  
PLTSUBPLOTSADJUSTLEFT02 RIGHT98 BOTTOM001 TOP96 WSPACE05  
HSPACE01  
PLOTNUM 1  
DEFAULTBASE QUANTILE 3  
EPS 3  
DAMPING 9  
PREFERENCE 200  
NNEIGHBORS 10  
NCLUSTERS 3  
MINSAMPLES 20  
XI 005  
MINCLUSTERSIZE 01  
DATASETS  
NOISYCIRCLES DAMPING 77 PREFERENCE 240  
QUANTILE 2 NCLUSTERS 2  
MINSAMPLES 20 XI 025  
NOISYMOONS DAMPING 75 PREFERENCE 220 NCLUSTERS 2  
VARIED EPS 18 NNEIGHBORS 2  
MINSAMPLES 5 XI 0035 MINCLUSTERSIZE 2  
ANISO EPS 15 NNEIGHBORS 2  
MINSAMPLES 20 XI 01 MINCLUSTERSIZE 2  
BLOBS  
NOSTRUCTURE  
FORIDATASET DATASET ALGOPARAMS INENUMERATEDDATASETS  
UPDATE PARAMETERS WITH DATASETSPECIFIC VALUES  
PARAMS DEFAULTBASECOPY  
906 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213  
PARAMSUPDATEALGOPARAMS  
X Y DATASET  
NORMALIZE DATASET FOR EASIER PARAMETER SELECTION  
X STANDARDSCALERFITTRANSFORMX  
ESTIMATE BANDWIDTH FOR MEAN SHIFT  
BANDWIDTH CLUSTERESTIMATEBANDWIDTHX QUANTILEPARAMSQUNTILE  
CONNECTIVITY MATRIX FOR STRUCTURED WARD  
CONNECTIVITY KNEIGHBORSGRAPH  
X NNEIGHBORSPARAMSNNEIGHBORS INCLUDESELFFALSE  
MAKE CONNECTIVITY SYMMETRIC  
CONNECTIVITY 05 CONNECTIVITY CONNECTIVITYT

CREATE CLUSTER OBJECTS

MS CLUSTERMEANSHIFTBANDWIDTHBANDWIDTH BINSEEDINGTRUE  
TWOMEANS CLUSTERMINIBATCHKMEANSNCLUSTERSPARAMSNCLUSTERS  
WARD CLUSTERAGGLOMERATIVECLUSTERING  
NCLUSTERSPARAMSNCLUSTERS LINKAGEWARD  
CONNECTIVITYCONNECTIVITY  
SPECTRAL CLUSTERSPECTRALCLUSTERING  
NCLUSTERSPARAMSNCLUSTERS EIGENSOLVERARPACK  
AFFINITYNEARESTNEIGHBORS  
DBSCAN CLUSTERDBSCANEPSPARAMSEPS  
OPTICS CLUSTEROPTICSMINSAMPLESPARAMSMINSAMPLES  
XIPARAMSXI  
MINCLUSTERSIZEPARAMSMINCLUSTERSIZE  
AFFINITYPROPAGATION CLUSTERAFFINITYPROPAGATION  
DAMPINGPARAMSDAMPING PREFERENCEPARAMSPREFERENCE  
AVERAGELINKAGE CLUSTERAGGLOMERATIVECLUSTERING  
LINKAGEAVERAGE AFFINITYCITYBLOCK  
NCLUSTERSPARAMSNCLUSTERS CONNECTIVITYCONNECTIVITY  
BIRCH CLUSTERBIRCHNCLUSTERSPARAMSNCLUSTERS  
GMM MIXTUREGAUSSIANMIXTURE  
NCOMPONENTSPARAMSNCLUSTERS COVARIANCETYPEFULL  
CLUSTERINGALGORITHMS  
MINIBATCHKMEANS TWOMEANS  
AFFINITYPROPAGATION AFFINITYPROPAGATION  
MEANSHIFT MS  
SPECTRALCLUSTERING SPECTRAL  
WARD WARD  
AGGLOMERATIVECLUSTERING AVERAGELINKAGE  
DBSCAN DBSCAN  
OPTICS OPTICS  
BIRCH BIRCH  
GAUSSIANMIXTURE GMM

FORNAME ALGORITHM INCLUSTERINGALGORITHMS  
T0 TIMETIME  
CATCH WARNINGS RELATED TO KNEIGHBORSGRAPH  
WITHWARNINGSCATCHWARNINGS  
56 CLUSTERING 907

SCIKITLEARN USER GUIDE RELEASE 0213  
WARNINGSFILTERWARNINGS  
IGNORE  
MESSAGE THE NUMBER OF CONNECTED COMPONENTS OF THE  
CONNECTIVITY MATRIX IS 0912  
1 COMPLETING IT TO AVOID STOPPING THE TREE EARLY  
CATEGORYUSERWARNING  
WARNINGSFILTERWARNINGS  
IGNORE  
MESSAGEGRAPH IS NOT FULLY CONNECTED SPECTRAL EMBEDDING  
MAY NOT WORK AS EXPECTED  
CATEGORYUSERWARNING  
ALGORITHMFITX  
T1 TIMETIME  
IFHASATTRALGORITHM LABELS  
YPRED ALGORITHM LABELSASTYPENPINT  
ELSE  
YPRED ALGORITHM PREDICTX  
PLTSUBPLOTLENDATASETS LENCLUSTERINGALGORITHMS PLOTNUM  
IFIDATASET 0  
PLTTITLENAME SIZE18  
COLORS NPARRAYLISTISLICECYCLE377EB8 FF7F00 4DAF4A  
F781BF A65628 984EA3  
999999 E41A1C DEDE00  
INTMAXYPRED 1  
ADD BLACK COLOR FOR OUTLIERS IF ANY  
COLORS NPAPPENDCOLORS 000000  
PLTSCATTERX 0 X 1 S10 COLORCOLORSYMPRED  
PLTXLIM25 25  
PLTYLIM25 25  
PLXTTICKS  
PLTYTICKS  
PLTTTEXT99 01 2FS T1 T0LSTRIPO  
TRANSFORMPLTGCATTRANSAXES SIZE15  
HORIZONTALALIGNMENTRIGHT  
PLOTNUM 1  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 39530 SECONDS  
57 PIPELINES AND COMPOSITE ESTIMATORS  
EXAMPLES OF HOW TO COMPOSE TRANSFORMERS AND PIPELINES FROM OTHER ESTIMATORS SEE THE USER GUIDE  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
908 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

571 CONCATENATING MULTIPLE FEATURE EXTRACTION METHODS

IN MANY REALWORLD EXAMPLES THERE ARE MANY WAYS TO EXTRACT FEATURES FROM A DATASET OFTEN IT IS BENEFICIAL TO COMBINE SEVERAL METHODS TO OBTAIN GOOD PERFORMANCE THIS EXAMPLE SHOWS HOW TO USE FEATUREUNION TO COMBINE FEATURES OBTAINED BY PCA AND UNIVARIATE SELECTION

COMBINING FEATURES USING THIS TRANSFORMER HAS THE BENEFIT THAT IT ALLOWS CROSS VALIDATION AND GRID SEARCHES OVER THE WHOLE PROCESS

THE COMBINATION USED IN THIS EXAMPLE IS NOT PARTICULARLY HELPFUL ON THIS DATASET AND IS ONLY USED TO ILLUSTRATE THE USAGE OF FEATUREUNION

OUT

COMBINED SPACE HAS 3 FEATURES

FITTING 5 FOLDS FOR EACH OF 18 CANDIDATES TOTALLING 90 FITS

CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK1 SVMC01

CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK1 SVMC01 SCORE0

↪933 TOTAL 00S

CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK1 SVMC01

CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK1 SVMC01 SCORE0

↪933 TOTAL 00S

CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK1 SVMC01

CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK1 SVMC01 SCORE0

↪867 TOTAL 00S

CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK1 SVMC01

CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK1 SVMC01 SCORE0

↪933 TOTAL 00S

CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK1 SVMC01

CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK1 SVMC01 SCORE1

↪000 TOTAL 00S

CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK1 SVMC1

CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK1 SVMC1 SCORE0

↪900 TOTAL 00S

CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK1 SVMC1

CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK1 SVMC1 SCORE1

↪000 TOTAL 00S

CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK1 SVMC1

CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK1 SVMC1 SCORE0

↪867 TOTAL 00S

CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK1 SVMC1

CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK1 SVMC1 SCORE0

↪933 TOTAL 00S

CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK1 SVMC1

CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK1 SVMC1 SCORE1

↪000 TOTAL 00S

CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK1 SVMC10

CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK1 SVMC10 SCORE0

↪933 TOTAL 00S

CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK1 SVMC10

CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK1 SVMC10 SCORE1

↪000 TOTAL 00S

CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK1 SVMC10

CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK1 SVMC10 SCORE0

↪900 TOTAL 00S

CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK1 SVMC10

CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK1 SVMC10 SCORE0

↪933 TOTAL 00S

CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK1 SVMC10

57 PIPELINES AND COMPOSITE ESTIMATORS 909

SCIKITLEARN USER GUIDE RELEASE 0213

CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK1 SVMC10 SCORE1  
↪000 TOTAL 00S  
CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK2 SVMC01  
CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK2 SVMC01 SCORE0  
↪933 TOTAL 00S  
CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK2 SVMC01  
CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK2 SVMC01 SCORE0  
↪967 TOTAL 00S  
CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK2 SVMC01  
CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK2 SVMC01 SCORE0  
↪933 TOTAL 00S  
CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK2 SVMC01  
CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK2 SVMC01 SCORE0  
↪933 TOTAL 00S  
CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK2 SVMC01  
CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK2 SVMC01 SCORE1  
↪000 TOTAL 01S  
CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK2 SVMC1  
CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK2 SVMC1 SCORE0  
↪933 TOTAL 00S  
CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK2 SVMC1  
CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK2 SVMC1 SCORE0  
↪967 TOTAL 00S  
CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK2 SVMC1  
CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK2 SVMC1 SCORE0  
↪933 TOTAL 00S  
CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK2 SVMC1  
CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK2 SVMC1 SCORE0  
↪933 TOTAL 00S  
CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK2 SVMC1  
CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK2 SVMC1 SCORE1  
↪000 TOTAL 00S  
CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK2 SVMC10  
CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK2 SVMC10 SCORE0  
↪967 TOTAL 00S  
CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK2 SVMC10  
CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK2 SVMC10 SCORE0  
↪967 TOTAL 00S  
CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK2 SVMC10  
CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK2 SVMC10 SCORE0  
↪933 TOTAL 00S  
CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK2 SVMC10  
CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK2 SVMC10 SCORE0  
↪933 TOTAL 00S  
CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK2 SVMC10  
CV FEATURESPCANCOMPONENTS1 FEATURESUNIVSELECTK2 SVMC10 SCORE1  
↪000 TOTAL 00S  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK1 SVMC01  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK1 SVMC01 SCORE0  
↪933 TOTAL 00S  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK1 SVMC01  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK1 SVMC01 SCORE1  
↪000 TOTAL 00S  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK1 SVMC01  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK1 SVMC01 SCORE0  
↪867 TOTAL 00S  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK1 SVMC01  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK1 SVMC01 SCORE0  
↪933 TOTAL 00S  
910 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK1 SVMC01  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK1 SVMC01 SCORE1  
↪000 TOTAL 00S  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK1 SVMC1  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK1 SVMC1 SCORE0  
↪967 TOTAL 00S  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK1 SVMC1  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK1 SVMC1 SCORE1  
↪000 TOTAL 00S  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK1 SVMC1  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK1 SVMC1 SCORE0  
↪933 TOTAL 00S  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK1 SVMC1  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK1 SVMC1 SCORE0  
↪933 TOTAL 00S  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK1 SVMC1  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK1 SVMC1 SCORE1  
↪000 TOTAL 00S  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK1 SVMC10  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK1 SVMC10 SCORE0  
↪967 TOTAL 00S  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK1 SVMC10  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK1 SVMC10 SCORE0  
↪967 TOTAL 00S  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK1 SVMC10  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK1 SVMC10 SCORE0  
↪900 TOTAL 00S  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK1 SVMC10  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK1 SVMC10 SCORE0  
↪933 TOTAL 00S  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK1 SVMC10  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK1 SVMC10 SCORE1  
↪000 TOTAL 00S  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK2 SVMC01  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK2 SVMC01 SCORE0  
↪967 TOTAL 00S  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK2 SVMC01  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK2 SVMC01 SCORE1  
↪000 TOTAL 00S  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK2 SVMC01  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK2 SVMC01 SCORE0  
↪933 TOTAL 00S  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK2 SVMC01  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK2 SVMC01 SCORE0  
↪933 TOTAL 00S  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK2 SVMC01  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK2 SVMC01 SCORE1  
↪000 TOTAL 00S  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK2 SVMC1  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK2 SVMC1 SCORE0  
↪967 TOTAL 00S  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK2 SVMC1  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK2 SVMC1 SCORE1  
↪000 TOTAL 00S  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK2 SVMC1  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK2 SVMC1 SCORE0  
↪933 TOTAL 00S  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK2 SVMC1

SCIKITLEARN USER GUIDE RELEASE 0213

CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK2 SVMC1 SCORE0  
↪967 TOTAL 00S

CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK2 SVMC1  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK2 SVMC1 SCORE1  
↪000 TOTAL 00S

CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK2 SVMC10  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK2 SVMC10 SCORE0  
↪967 TOTAL 00S

CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK2 SVMC10  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK2 SVMC10 SCORE1  
↪000 TOTAL 00S

CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK2 SVMC10  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK2 SVMC10 SCORE0  
↪900 TOTAL 00S

CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK2 SVMC10  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK2 SVMC10 SCORE0  
↪933 TOTAL 00S

CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK2 SVMC10  
CV FEATURESPCANCOMPONENTS2 FEATURESUNIVSELECTK2 SVMC10 SCORE1  
↪000 TOTAL 00S

CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK1 SVMC01  
CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK1 SVMC01 SCORE0  
↪967 TOTAL 00S

CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK1 SVMC01  
CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK1 SVMC01 SCORE1  
↪000 TOTAL 00S

CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK1 SVMC01  
CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK1 SVMC01 SCORE0  
↪933 TOTAL 00S

CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK1 SVMC01  
CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK1 SVMC01 SCORE0  
↪967 TOTAL 00S

CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK1 SVMC01  
CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK1 SVMC01 SCORE1  
↪000 TOTAL 00S

CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK1 SVMC1  
CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK1 SVMC1 SCORE0  
↪967 TOTAL 00S

CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK1 SVMC1  
CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK1 SVMC1 SCORE1  
↪000 TOTAL 00S

CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK1 SVMC1  
CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK1 SVMC1 SCORE0  
↪933 TOTAL 00S

CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK1 SVMC1  
CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK1 SVMC1 SCORE0  
↪967 TOTAL 00S

CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK1 SVMC1  
CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK1 SVMC1 SCORE1  
↪000 TOTAL 00S

CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK1 SVMC10  
CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK1 SVMC10 SCORE1  
↪000 TOTAL 00S

CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK1 SVMC10  
CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK1 SVMC10 SCORE1  
↪000 TOTAL 00S

CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK1 SVMC10  
CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK1 SVMC10 SCORE0  
↪933 TOTAL 00S

912 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK1 SVMC10  
CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK1 SVMC10 SCORE0  
↪967 TOTAL 00S

CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK1 SVMC10  
CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK1 SVMC10 SCORE1  
↪000 TOTAL 00S

CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK2 SVMC01  
CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK2 SVMC01 SCORE0  
↪967 TOTAL 00S

CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK2 SVMC01  
CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK2 SVMC01 SCORE1  
↪000 TOTAL 00S

CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK2 SVMC01  
CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK2 SVMC01 SCORE0  
↪933 TOTAL 00S

CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK2 SVMC01  
CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK2 SVMC01 SCORE0  
↪967 TOTAL 00S

CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK2 SVMC01  
CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK2 SVMC01 SCORE1  
↪000 TOTAL 00S

CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK2 SVMC1  
CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK2 SVMC1 SCORE0  
↪967 TOTAL 00S

CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK2 SVMC1  
CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK2 SVMC1 SCORE1  
↪000 TOTAL 00S

CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK2 SVMC1  
CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK2 SVMC1 SCORE0  
↪967 TOTAL 00S

CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK2 SVMC1  
CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK2 SVMC1 SCORE0  
↪967 TOTAL 00S

CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK2 SVMC1  
CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK2 SVMC1 SCORE1  
↪000 TOTAL 00S

CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK2 SVMC10  
CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK2 SVMC10 SCORE1  
↪000 TOTAL 00S

CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK2 SVMC10  
CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK2 SVMC10 SCORE0  
↪900 TOTAL 00S

CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK2 SVMC10  
CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK2 SVMC10 SCORE0  
↪967 TOTAL 00S

CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK2 SVMC10  
CV FEATURESPCANCOMPONENTS3 FEATURESUNIVSELECTK2 SVMC10 SCORE1  
↪000 TOTAL 00S

PIPELINESTEPSFEATURES  
FEATUREUNIONTRANSFORMERLISTPCA PCANCOMPONENTS3  
UNIVSELECT  
SELECTKBESTK1  
SVM SVCC10 KERNELLINEAR  
57 PIPELINES AND COMPOSITE ESTIMATORS 913

SCIKITLEARN USER GUIDE RELEASE 0213  
AUTHOR ANDREAS MUELLER AMUELLERAISUNIBONNDE

```
LICENSE BSD 3 CLAUSE
FROM SKLEARNPIPELINE IMPORT PIPELINE FEATUREUNION
FROM SKLEARNMODELSELECTION IMPORT GRIDSEARCHCV
FROM SKLEARN SVM IMPORT SVC
FROM SKLEARN DATASETS IMPORT LOADIRIS
FROM SKLEARN DECOMPOSITION IMPORT PCA
FROM SKLEARN FEATURE SELECTION IMPORT SELECTK BEST
IRIS = LOADIRIS
X Y = IRIS DATA IRIS TARGET
THIS DATASET IS WAY TOO HIGH DIMENSIONAL BETTER DO PCA
PCA = PCANCOMPONENTS2
MAYBE SOME ORIGINAL FEATURES WERE GOOD TOO
SELECTION = SELECTK BEST K1
BUILD ESTIMATOR FROM PCA AND UNIVARIATE SELECTION
COMBINED FEATURES = FEATUREUNION PCA PCA UNIV SELECT SELECTION
USE COMBINED FEATURES TO TRANSFORM DATASET
X FEATURES = COMBINED FEATURES FIT X Y TRANSFORM X
PRINT COMBINED SPACE HAS X FEATURES SHAPE 1 FEATURES
SVM = SVCKERNEL LINEAR
DO GRID SEARCH OVER K NCOMPONENTS AND C
PIPELINE = PIPELINE FEATURES COMBINED FEATURES SVM SVM
PARAM GRID = DICT FEATURES PCANCOMPONENTS 1 2 3
FEATURES UNIV SELECT K1 2
SVM C01 1 10
GRID SEARCH = GRID SEARCH CV PIPELINE PARAM GRID PARAM GRID CV5 VERBOSE 10
GRID SEARCH FIT X Y
PRINT GRID SEARCH BEST ESTIMATOR
TOTAL RUNNING TIME OF THE SCRIPT = 0 MINUTES 0865 SECONDS
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE
572 PIPELINING CHAINING A PCA AND A LOGISTIC REGRESSION
THE PCA DOES AN UNSUPERVISED DIMENSIONALITY REDUCTION WHILE THE LOGISTIC REGRESSION DOES THE PREDICTION
WE USE A GRID SEARCH CV TO SET THE DIMENSIONALITY OF THE PCA
914 CHAPTER 5 EXAMPLES
```



SCIKITLEARN USER GUIDE RELEASE 0213  
OUT  
BEST PARAMETER CV SCORE0917  
LOGISTICALPHA 001 PCANCOMPONENTS 64  
PRINTDOC  
CODE SOURCE GAËL VAROQUAUX  
MODIFIED FOR DOCUMENTATION BY JAQUES GROBLER  
LICENSE BSD 3 CLAUSE  
57 PIPELINES AND COMPOSITE ESTIMATORS 915

```
SCIKITLEARN USER GUIDE RELEASE 0213
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
IMPORT PANDAS AS PD
FROM SKLEARN IMPORT DATASETS
FROM SKLEARNDECOMPOSITION IMPORT PCA
FROM SKLEARNLINEARMODEL IMPORT SGDCLASSIFIER
FROM SKLEARNPIPELINE IMPORT PIPELINE
FROM SKLEARNMODELSELECTION IMPORT GRIDSEARCHCV
DEFINE A PIPELINE TO SEARCH FOR THE BEST COMBINATION OF PCA TRUNCATION
AND CLASSIFIER REGULARIZATION
LOGISTIC SGDCLASSIFIERLOSSLOG PENALTYL2 EARLYSTOPPINGTRUE
MAXITER10000 TOL1E5 RANDOMSTATE0
PCA PCA
PIPE PIPELINESTEPSPCA PCA LOGISTIC LOGISTIC
DIGITS DATASETSLOADDIGITS
XDIGITS DIGITSDATA
YDIGITS DIGITSTARGET
PARAMETERS OF PIPELINES CAN BE SET USING " SEPARATED PARAMETER NAMES
PARAMGRID
PCANCOMPONENTS 5 20 30 40 50 64
LOGISTICALPHA NPLOGSPACE4 4 5

SEARCH GRIDSEARCHCVPIPE PARAMGRID IIDFALSE CV5
SEARCHFITXDIGITS YDIGITS
PRINTBEST PARAMETER CV SCORE 03F SEARCHBESTSCORE
PRINTSEARCHBESTPARAMS
PLOT THE PCA SPECTRUM
PCAFITXDIGITS
FIG AX0 AX1 PLTSUBPLOTSNROWS2 SHAREXTRUE FIGSIZE6 6
AX0PLOTPCAEXPLAINEDVARIANCERATIO LINEWIDTH2
AX0SETYLABELPCA EXPLAINED VARIANCE
AX0AXVLINERSEARCHBESTESTIMATORNAMEDSTEPSPCANCOMPONENTS
LINESTYLE LABELNCOMPONENTS CHOSEN
AX0LEGENDPROPDICTIONSIZE12
FOR EACH NUMBER OF COMPONENTS FIND THE BEST CLASSIFIER RESULTS
RESULTS PDDATAFRAMESEARCHCVRESULTS
COMPONENTSCOL PARAMPCANCOMPONENTS
BESTCLFS RESULTSGROUPBYCOMPONENTSCOLAPPLY
LAMBDA GNLARGEST1 MEANTESTSCORE
BESTCLFSPLOTXCOMPONENTSCOL YMEANTESTSCORE YERRSTDTESTSCORE
LEGENDFALSE AXAX1
AX1SETYLABELCLASSIFICATION ACCURACY VAL
AX1SETXLABELNCOMPONENTS
PLTTIGHTLAYOUT
PLTSHOW
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 20748 SECONDS
916 CHAPTER 5 EXAMPLES
```

SCIKITLEARN USER GUIDE RELEASE 0213

NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE

573 COLUMN TRANSFORMER WITH MIXED TYPES

THIS EXAMPLE ILLUSTRATES HOW TO APPLY DIFFERENT PREPROCESSING AND FEATURE EXTRACTION PIPELINES TO DIFFERENT SUBSETS OF FEATURES USING SKLEARNCOMPOSECOLUMNTRANSFORMER THIS IS PARTICULARLY HANDY FOR THE CASE OF DATASETS THAT CONTAIN HETEROGENEOUS DATA TYPES SINCE WE MAY WANT TO SCALE THE NUMERIC FEATURES AND ONEHOT ENCODE THE CATEGORICAL ONES

IN THIS EXAMPLE THE NUMERIC DATA IS STANDARDSCALED AFTER MEANIMPUTATION WHILE THE CATEGORICAL DATA IS ONEHOT ENCODED AFTER IMPUTING MISSING VALUES WITH A NEW CATEGORY MISSING

FINALLY THE PREPROCESSING PIPELINE IS INTEGRATED IN A FULL PREDICTION PIPELINE USING SKLEARNPIPELINE PIPELINE TOGETHER WITH A SIMPLE CLASSIFICATION MODEL

AUTHOR PEDRO MORALES PARTMORALESGMAILCOM

LICENSE BSD 3 CLAUSE

```
import pandas as pd
import numpy as np
from sklearn.compose import ColumnTransformer
from sklearn.pipeline import Pipeline
from sklearn.impute import SimpleImputer
from sklearn.preprocessing import StandardScaler, OneHotEncoder
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split, GridSearchCV
np.random.seed(0)

# Read data from Titanic dataset
titanic_url = 'https://raw.githubusercontent.com/amuel/SciPy2017Sklearn091D371Notebooks/master/titanic3.csv'
data = pd.read_csv(titanic_url)

# We will train our classifier with the following features
numeric_features = ['age', 'fare']
categorical_features = ['embarked', 'sex', 'pclass']

# We create the preprocessing pipelines for both numeric and categorical data
numeric_transformer = Pipeline(steps=[
    ('imputer', SimpleImputer(strategy='median')),
    ('scaler', StandardScaler())
])

categorical_transformer = Pipeline(steps=[
    ('imputer', SimpleImputer(strategy='constant', fill_value='missing')),
    ('onehot', OneHotEncoder(handle_unknown='ignore'))
])

# 57 pipelines and composite estimators 917
```

SCIKITLEARN USER GUIDE RELEASE 0213  
PREPROCESSOR COLUMNTRANSFORMER  
TRANSFORMERS  
NUM NUMERICTRANSFORMER NUMERICFEATURES  
CAT CATEGORICALTRANSFORMER CATEGORICALFEATURES  
APPEND CLASSIFIER TO PREPROCESSING PIPELINE  
NOW WE HAVE A FULL PREDICTION PIPELINE  
CLF PIPELINESTEPSPREPROCESSOR PREPROCESSOR  
CLASSIFIER LOGISTICREGRESSIONSOLVERLBFGS  
X DATADROPSURVIVED AXIS1  
Y DATASURVIVED  
XTRAIN XTEST YTRAIN YTEST TRAINTESTSPLITX Y TESTSIZE02  
CLFFITXTRAIN YTRAIN  
PRINTMODEL SCORE 3F CLFSCOREXTEST YTEST  
OUT  
MODEL SCORE 0790  
USING THE PREDICTION PIPELINE IN A GRID SEARCH  
GRID SEARCH CAN ALSO BE PERFORMED ON THE DIFFERENT PREPROCESSING STEPS DEFINED IN THE  
COLUMNTRANSFORMER OBJECT TOGETHER WITH THE CLASSIFIER'S HYPERPARAMETERS AS PART OF THE PIPELINE  
WE WILL SEARCH FOR BOTH THE IMPUTER STRATEGY OF THE NUMERIC PREPROCESSING AND THE REGULARIZATION PARAMETER  
OF THE LOGISTIC REGRESSION USING SKLEARNMODELSELECTIONGRIDSEARCHCV  
PARAMGRID  
PREPROCESSORNUMIMPUTERSTRATEGY MEAN MEDIAN  
CLASSIFIERC 01 10 10 100  
  
GRIDSEARCH GRIDSEARCHCVCLF PARAMGRID CV10 IIDFALSE  
GRIDSEARCHFITXTRAIN YTRAIN  
PRINTBEST LOGISTIC REGRESSION FROM GRID SEARCH 3F  
GRIDSEARCHSCOREXTEST YTEST  
OUT  
BEST LOGISTIC REGRESSION FROM GRID SEARCH 0798  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 2103 SECONDS  
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE  
574 SELECTING DIMENSIONALITY REDUCTION WITH PIPELINE AND GRIDSEARCHCV  
THIS EXAMPLE CONSTRUCTS A PIPELINE THAT DOES DIMENSIONALITY REDUCTION FOLLOWED BY PREDICTION WITH A SUPPORT VECTOR  
CLASSIFIER IT DEMONSTRATES THE USE OF GRIDSEARCHCV ANDPIPELINE TO OPTIMIZE OVER DIFFERENT CLASSES OF ESTIMATORS  
918 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

IN A SINGLE CV RUN - UNSUPERVISED PCA ANDNMF DIMENSIONALITY REDUCTIONS ARE COMPARED TO UNIVARIATE FEATURE SELECTION DURING THE GRID SEARCH

ADDITIONALLY PIPELINE CAN BE INSTANTIATED WITH THE MEMORY ARGUMENT TO MEMOIZE THE TRANSFORMERS WITHIN THE PIPELINE AVOIDING TO FIT AGAIN THE SAME TRANSFORMERS OVER AND OVER

NOTE THAT THE USE OF MEMORY TO ENABLE CACHING BECOMES INTERESTING WHEN THE FITTING OF A TRANSFORMER IS COSTLY

### ILLUSTRATION OF PIPELINE AND GRIDSEARCHCV

THIS SECTION ILLUSTRATES THE USE OF A PIPELINE WITH `GRIDSEARCHCV`

AUTHORS ROBERT MCGIBBON JOEL NOTHMAN GUILLAUME LEMAITRE

```
import numpy as np
```

```
import matplotlib.pyplot as plt
```

```
FROM SKLEARNDATASETS IMPORT LOADDIGITS
```

```
FROM SKLEARNMODELSELECTION IMPORT GRIDSEARCHCV
```

```
FROM SKLEARNPIPELINE IMPORT PIPELINE
```

```
FROM SKLEARN.SVM IMPORT LINEARSVC
```

```
FROM SKLEARNDECOMPOSITION IMPORT PCA NMF
```

```
FROM SKLEARNFEATURESELECTION IMPORT SELECTKBEST CH12
```

PRINTDOC

PIPE PIPELINE

THE REDUCEDIM STAGE IS POPULATED BY THE PARAMGRID

REDUCEDIM PASSTHROUGH

```
CLASSIFY LINEARSVC DUAL FALSE MAXITER 10000
```

NFEATURESOPTIONS 2 4 8

COPTIONS 1 10 100 1000

PARAMGRID

REDUCEDIM PCAITERATEDPOWER7 NMF

REDUCEDIMNCOMPONENTS NFEATURESOPTIONS

## CLASSIFYC OPTIONS

REDUCEDIM SELECTKBESTCHI2

REDUCEDIMK NFEATURESOPTIONS

## CLASSIFYC OPTIONS

REDUCERLABELS PCA NMF KBESTCHI2

```
REDUCEREABLES TCA NMI RBESTCHIZ  
GRID GRIDSEARCHCVPIPE CV5 NIOBS1 PARAMGRIDPARAMGRID IIDFALSE
```

DIGITS LOADDIGITS

DIGITS\_LOADDIGITS  
GRIDFITDIGITS DATA DIGITSTARGET

MEANScores NPARRAYGRIDCVRESULTSMEANTESTSCORE

SCORES ARE IN THE ORDER OF PARAMGRID ITERATION WHICH IS ALPHABETICAL

MEANScores MEANScoresRESHAPELENCOptions 1 LENNFEATuRESOptions

SELECT SCORE FOR BEST C

MEANScores MEANScoresMAXAxis0

MEANScores MEANScoresMAXAxis30  
BAROffsets NPARANGELENNFEATuRESOptions

57 PIPELINES AND COMPOSITE ESTIMATORS 919

SCIKITLEARN USER GUIDE RELEASE 0213  
LENREDUCERLABELS 1 5  
PLTFigure  
COLORS BGRCMYK  
FORI LABEL REDUCERScores INENUMERATEZIPREDUCERLABELS MEANScores  
PLTBARBAROFFSETS 1 REDUCERScores LABELLABEL COLORCOLORSI  
PLTTITLECOMPARING FEATURE REDUCTION TECHNIQUES  
PLTXLABELREDUCED NUMBER OF FEATURES  
PLTXTICKSBAROFFSETS LENREDUCERLABELS 2 NFEATURESOPTIONS  
PLTYLABELDIGIT CLASSIFICATION ACCURACY  
PLTYLIM0 1  
PLTLEGENDLOCUPPER LEFT  
PLTSHOW  
CACHING TRANSFORMERS WITHIN A PIPELINE  
IT IS SOMETIMES WORTHWHILE STORING THE STATE OF A SPECIFIC TRANSFORMER SINCE IT COULD BE USED AGAIN USING A  
PIPELINE INGRIDSEARCHCV TRIGGERS SUCH SITUATIONS THEREFORE WE USE THE ARGUMENT MEMORY TO ENABLE  
CACHING  
920 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
WARNING NOTE THAT THIS EXAMPLE IS HOWEVER ONLY AN ILLUSTRATION SINCE FOR THIS SPECIFIC CASE FITTING
PCA IS NOT NECESSARILY SLOWER THAN LOADING THE CACHE HENCE USE THE MEMORY CONSTRUCTOR PARAMETER
WHEN THE FITTING OF A TRANSFORMER IS COSTLY
FROM TEMPFILE IMPORT MKDTEMP
FROM SHUTIL IMPORT RMTREE
FROM JOBLIB IMPORT MEMORY
    CREATE A TEMPORARY FOLDER TO STORE THE TRANSFORMERS OF THE PIPELINE
    CACHEDIR MKDTEMP
MEMORY MEMORYLOCATIONCACHEDIR VERBOSE10
CACHEDPIPE PIPELINEREDUCEDIM PCA
CLASSIFY LINEARSVC DUALFALSE MAXITER10000
MEMORYMEMORY
    THIS TIME A CACHED PIPELINE WILL BE USED WITHIN THE GRID SEARCH
GRID GRIDSEARCHCVCACHEDPIPE CV5 NJOBS1 PARAMGRIDPARAMGRID
IIDFALSE
DIGITS LOADDIGITS
GRIDFITDIGITS DATA DIGITSTARGET
    DELETE THE TEMPORARY CACHE BEFORE EXITING
RMTREECACHEDIR
OUT
```

```
MEMORY CALLING SKLEARNPIPELINEFITTRANSFORMONE
FITTRANSFORMONEPCAITERATEDPOWER7 NCOMPONENTS2 ARRAY0 0

0 0 ARRAY0 8 NONE MESSAGECLSNAMEPIPELINE
↵MESSAGEGENONE
FITTRANSFORMONE 00S 00MIN
```

```
MEMORY CALLING SKLEARNPIPELINEFITTRANSFORMONE
FITTRANSFORMONEPCAITERATEDPOWER7 NCOMPONENTS2 ARRAY0 0

0 0 ARRAY0 8 NONE MESSAGECLSNAMEPIPELINE
↵MESSAGEGENONE
FITTRANSFORMONE 00S 00MIN
```

```
MEMORY CALLING SKLEARNPIPELINEFITTRANSFORMONE
FITTRANSFORMONEPCAITERATEDPOWER7 NCOMPONENTS2 ARRAY0 0

0 0 ARRAY0 8 NONE MESSAGECLSNAMEPIPELINE
↵MESSAGEGENONE
FITTRANSFORMONE 00S 00MIN
```

```
MEMORY CALLING SKLEARNPIPELINEFITTRANSFORMONE
FITTRANSFORMONEPCAITERATEDPOWER7 NCOMPONENTS2 ARRAY0 0

0 0 ARRAY0 8 NONE MESSAGECLSNAMEPIPELINE
↵MESSAGEGENONE
FITTRANSFORMONE 00S 00MIN
```

```
MEMORY CALLING SKLEARNPIPELINEFITTRANSFORMONE
57 PIPELINES AND COMPOSITE ESTIMATORS 921
```

SCIKITLEARN USER GUIDE RELEASE 0213  
FITTRANSFORMONEPCAITERATEDPOWER7 NCOMPONENTS2 ARRAY0 0

0 0 ARRAY0 9 NONE MESSAGECLSNAMEPIPELINE  
↩→MESSAGENONE  
FITTRANSFORMONE 00S 00MIN

MEMORY CALLING SKLEARNPIPELINEFITTRANSFORMONE  
FITTRANSFORMONEPCAITERATEDPOWER7 NCOMPONENTS4 ARRAY0 0

0 0 ARRAY0 8 NONE MESSAGECLSNAMEPIPELINE  
↩→MESSAGENONE  
FITTRANSFORMONE 00S 00MIN

MEMORY CALLING SKLEARNPIPELINEFITTRANSFORMONE  
FITTRANSFORMONEPCAITERATEDPOWER7 NCOMPONENTS4 ARRAY0 0

0 0 ARRAY0 8 NONE MESSAGECLSNAMEPIPELINE  
↩→MESSAGENONE  
FITTRANSFORMONE 00S 00MIN

MEMORY CALLING SKLEARNPIPELINEFITTRANSFORMONE  
FITTRANSFORMONEPCAITERATEDPOWER7 NCOMPONENTS4 ARRAY0 0

0 0 ARRAY0 8 NONE MESSAGECLSNAMEPIPELINE  
↩→MESSAGENONE  
FITTRANSFORMONE 00S 00MIN

MEMORY CALLING SKLEARNPIPELINEFITTRANSFORMONE  
FITTRANSFORMONEPCAITERATEDPOWER7 NCOMPONENTS4 ARRAY0 0

0 0 ARRAY0 8 NONE MESSAGECLSNAMEPIPELINE  
↩→MESSAGENONE  
FITTRANSFORMONE 00S 00MIN

MEMORY CALLING SKLEARNPIPELINEFITTRANSFORMONE  
FITTRANSFORMONEPCAITERATEDPOWER7 NCOMPONENTS4 ARRAY0 0

0 0 ARRAY0 9 NONE MESSAGECLSNAMEPIPELINE  
↩→MESSAGENONE  
FITTRANSFORMONE 00S 00MIN

MEMORY CALLING SKLEARNPIPELINEFITTRANSFORMONE  
FITTRANSFORMONEPCAITERATEDPOWER7 NCOMPONENTS8 ARRAY0 0

0 0 ARRAY0 8 NONE MESSAGECLSNAMEPIPELINE  
↩→MESSAGENONE  
FITTRANSFORMONE 01S 00MIN

MEMORY CALLING SKLEARNPIPELINEFITTRANSFORMONE  
FITTRANSFORMONEPCAITERATEDPOWER7 NCOMPONENTS8 ARRAY0 0

0 0 ARRAY0 8 NONE MESSAGECLSNAMEPIPELINE  
↩→MESSAGENONE  
FITTRANSFORMONE 00S 00MIN

MEMORY CALLING SKLEARNPIPELINEFITTRANSFORMONE  
FITTRANSFORMONEPCAITERATEDPOWER7 NCOMPONENTS8 ARRAY0 0



SCIKITLEARN USER GUIDE RELEASE 0213  
0 0 ARRAY0 8 NONE MESSAGECLSNAMEPIPELINE  
↔MESSAGENONE  
FITTRANSFORMONE 00S 00MIN

MEMORY CALLING SKLEARNPIPELINEFITTRANSFORMONE  
FITTRANSFORMONEPCAITERATEDPOWER7 NCOMPONENTS8 ARRAY0 0

0 0 ARRAY0 8 NONE MESSAGECLSNAMEPIPELINE  
↔MESSAGENONE  
FITTRANSFORMONE 00S 00MIN

MEMORY CALLING SKLEARNPIPELINEFITTRANSFORMONE  
FITTRANSFORMONEPCAITERATEDPOWER7 NCOMPONENTS8 ARRAY0 0

0 0 ARRAY0 9 NONE MESSAGECLSNAMEPIPELINE  
↔MESSAGENONE  
FITTRANSFORMONE 00S 00MIN

MEMORY CALLING SKLEARNPIPELINEFITTRANSFORMONE  
FITTRANSFORMONENMFNCOMPONENTS2 ARRAY0 0

0 0 ARRAY0 8 NONE MESSAGECLSNAMEPIPELINE  
↔MESSAGENONE  
FITTRANSFORMONE 01S 00MIN

MEMORY CALLING SKLEARNPIPELINEFITTRANSFORMONE  
FITTRANSFORMONENMFNCOMPONENTS2 ARRAY0 0

0 0 ARRAY0 8 NONE MESSAGECLSNAMEPIPELINE  
↔MESSAGENONE  
FITTRANSFORMONE 00S 00MIN

MEMORY CALLING SKLEARNPIPELINEFITTRANSFORMONE  
FITTRANSFORMONENMFNCOMPONENTS2 ARRAY0 0

0 0 ARRAY0 8 NONE MESSAGECLSNAMEPIPELINE  
↔MESSAGENONE  
FITTRANSFORMONE 01S 00MIN

MEMORY CALLING SKLEARNPIPELINEFITTRANSFORMONE  
FITTRANSFORMONENMFNCOMPONENTS2 ARRAY0 0

0 0 ARRAY0 8 NONE MESSAGECLSNAMEPIPELINE  
↔MESSAGENONE  
FITTRANSFORMONE 00S 00MIN

MEMORY CALLING SKLEARNPIPELINEFITTRANSFORMONE  
FITTRANSFORMONENMFNCOMPONENTS2 ARRAY0 0

0 0 ARRAY0 9 NONE MESSAGECLSNAMEPIPELINE  
↔MESSAGENONE  
FITTRANSFORMONE 00S 00MIN

MEMORY CALLING SKLEARNPIPELINEFITTRANSFORMONE  
FITTRANSFORMONENMFNCOMPONENTS4 ARRAY0 0

0 0 ARRAY0 8 NONE MESSAGECLSNAMEPIPELINE  
↔MESSAGENONE

57 PIPELINES AND COMPOSITE ESTIMATORS 923

SCIKITLEARN USER GUIDE RELEASE 0213  
FITTRANSFORMONE 01S 00MIN

MEMORY CALLING SKLEARNPIPELINEFITTRANSFORMONE  
FITTRANSFORMONENMFNCOMPONENTS4 ARRAY0 0

0 0 ARRAY0 8 NONE MESSAGECLSNAMEPIPELINE  
↩️MESSAGENONE  
FITTRANSFORMONE 00S 00MIN

MEMORY CALLING SKLEARNPIPELINEFITTRANSFORMONE  
FITTRANSFORMONENMFNCOMPONENTS4 ARRAY0 0

0 0 ARRAY0 8 NONE MESSAGECLSNAMEPIPELINE  
↩️MESSAGENONE  
FITTRANSFORMONE 00S 00MIN

MEMORY CALLING SKLEARNPIPELINEFITTRANSFORMONE  
FITTRANSFORMONENMFNCOMPONENTS4 ARRAY0 0

0 0 ARRAY0 8 NONE MESSAGECLSNAMEPIPELINE  
↩️MESSAGENONE  
FITTRANSFORMONE 00S 00MIN

MEMORY CALLING SKLEARNPIPELINEFITTRANSFORMONE  
FITTRANSFORMONENMFNCOMPONENTS4 ARRAY0 0

0 0 ARRAY0 9 NONE MESSAGECLSNAMEPIPELINE  
↩️MESSAGENONE  
FITTRANSFORMONE 01S 00MIN

MEMORY CALLING SKLEARNPIPELINEFITTRANSFORMONE  
FITTRANSFORMONENMFNCOMPONENTS8 ARRAY0 0

0 0 ARRAY0 8 NONE MESSAGECLSNAMEPIPELINE  
↩️MESSAGENONE  
FITTRANSFORMONE 01S 00MIN

MEMORY CALLING SKLEARNPIPELINEFITTRANSFORMONE  
FITTRANSFORMONENMFNCOMPONENTS8 ARRAY0 0

0 0 ARRAY0 8 NONE MESSAGECLSNAMEPIPELINE  
↩️MESSAGENONE  
FITTRANSFORMONE 01S 00MIN

MEMORY CALLING SKLEARNPIPELINEFITTRANSFORMONE  
FITTRANSFORMONENMFNCOMPONENTS8 ARRAY0 0

0 0 ARRAY0 8 NONE MESSAGECLSNAMEPIPELINE  
↩️MESSAGENONE  
FITTRANSFORMONE 02S 00MIN

MEMORY CALLING SKLEARNPIPELINEFITTRANSFORMONE  
FITTRANSFORMONENMFNCOMPONENTS8 ARRAY0 0

0 0 ARRAY0 8 NONE MESSAGECLSNAMEPIPELINE  
↩️MESSAGENONE  
FITTRANSFORMONE 01S 00MIN

SCIKITLEARN USER GUIDE RELEASE 0213  
MEMORY CALLING SKLEARNPIPELINEFITTRANSFORMONE  
FITTRANSFORMONENMFCOMPONENTS8 ARRAY0 0

0 0 ARRAY0 9 NONE MESSAGECLNAMEPIPELINE  
→MESSAGENONE  
FITTRANSFORMONE 01S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
→SKLEARNPIPELINEFITTRANSFORMONE591ACAE8E60C9632AA990E4FA0962A67  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
→SKLEARNPIPELINEFITTRANSFORMONE323909961FCC2A4950DEFB581F08007F  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
→SKLEARNPIPELINEFITTRANSFORMONE3654E1D61587AEE4F9980C99541D55D3  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
→SKLEARNPIPELINEFITTRANSFORMONE835820E75BE3172B4E2E257028147BB2  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
→SKLEARNPIPELINEFITTRANSFORMONE90604CE8BB756E3FF1ACD4C1CE065B1A  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
→SKLEARNPIPELINEFITTRANSFORMONE3C66DF347F488DFE044ECF991CA2498B  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
→SKLEARNPIPELINEFITTRANSFORMONEBF4BEC5FE5245BF3C0D271E31BC2342F  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
→SKLEARNPIPELINEFITTRANSFORMONECE0F911DF93A8E38932C48F983E7FDE6  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
→SKLEARNPIPELINEFITTRANSFORMONE285EFBA3F9E5B1A35D6825E37A718019  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
→SKLEARNPIPELINEFITTRANSFORMONE74083D6E1ED743FA1E756A41BA467E6C  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
→SKLEARNPIPELINEFITTRANSFORMONEB72F745822A406F2191B6EBFD8CA9DF9  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
→SKLEARNPIPELINEFITTRANSFORMONE997E4CE30BA5143E90330DB664BF61C5  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
→SKLEARNPIPELINEFITTRANSFORMONE3140177F6CCBE72992772342347A166F  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
→SKLEARNPIPELINEFITTRANSFORMONEF843896AD42A20C4A50C0FD29512111A  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
→SKLEARNPIPELINEFITTRANSFORMONED91A2D083B8C7311C87EDD6592FCD446  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
→SKLEARNPIPELINEFITTRANSFORMONE50EB0F4152CD0F7DCE3B471ADD31D07D  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
→SKLEARNPIPELINEFITTRANSFORMONE3B28284B39C934CB519F2A31E03A951B  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
→SKLEARNPIPELINEFITTRANSFORMONEADED0EB932B9C692C801C6C56DC2BDA2  
57 PIPELINES AND COMPOSITE ESTIMATORS 925

SCIKITLEARN USER GUIDE RELEASE 0213

FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONECDF3D1530A8A2388BA80A16C66413409  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONEFF5666A0CB943711A4E4DFD5B5CD8587  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE536A5FBB5A6C8D69B4FB426F0B51D9EB  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE26CD3A8242ACB986F98CF38291EF5BAA  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE8B6094F421DAB9FC88FB54C2F32E16B4  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONED0E8598C62DCB5F059A5C2B541C23B09  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE34AFFB50A3D55C5302E1C7F97AEC9679  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONEB3635522E8F89BBE9E56AB7E7F4693FC  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE7389AD9BAE0AD1181B7DA9E4DD39449D  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONEAB59173E00AB4A7BDF7740E5C6D50123  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE42894538FD0AB2B4235E55EF53B9CD59  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE580596A10617556E39D5E069C4DE096A  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE591ACAE8E60C9632AA990E4FA0962A67  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE323909961FCC2A4950DEFB581F08007F  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE3654E1D61587AEE4F9980C99541D55D3  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE835820E75BE3172B4E2E257028147BB2  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE90604CE8BB756E3FF1ACD4C1CE065B1A  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE3C66DF347F488DFE044ECF991CA2498B  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONEBF4BEC5FE5245BF3C0D271E31BC2342F  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
926 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONECE0F911DF93A8E38932C48F983E7FDE6  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE285EFBA3F9E5B1A35D6825E37A718019  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE74083D6E1ED743FA1E756A41BA467E6C  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONEB72F745822A406F2191B6EBFD8CA9DF9  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE997E4CE30BA5143E90330DB664BF61C5  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE3140177F6CCBE72992772342347A166F  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONEF843896AD42A20C4A50C0FD29512111A  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONED91A2D083B8C7311C87EDD6592FCD446  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE50EB0F4152CD0F7DCE3B471ADD31D07D  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE3B28284B39C934CB519F2A31E03A951B  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONEADE0EB932B9C692C801C6C56DC2BDA2  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONECDF3D1530A8A2388BA80A16C66413409  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONEFF5666A0CB943711A4E4DFD5B5CD8587  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE536A5FBB5A6C8D69B4FB426F0B51D9EB  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE26CD3A8242ACB986F98CF38291EF5BAA  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE8B6094F421DAB9FC88FB54C2F32E16B4  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONED0E8598C62DCB5F059A5C2B541C23B09  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE34AFFB50A3D55C5302E1C7F97AEC9679  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONEB3635522E8F89BBE9E56AB7E7F4693FC  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE7389AD9BAE0AD1181B7DA9E4DD39449D  
57 PIPELINES AND COMPOSITE ESTIMATORS 927

SCIKITLEARN USER GUIDE RELEASE 0213

FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONEAB59173E00AB4A7BDF7740E5C6D50123  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE42894538FD0AB2B4235E55EF53B9CD59  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE580596A10617556E39D5E069C4DE096A  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE591ACAE8E60C9632AA990E4FA0962A67  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE323909961FCC2A4950DEFB581F08007F  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE3654E1D61587AE4F9980C99541D55D3  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE835820E75BE3172B4E2E257028147BB2  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE90604CE8BB756E3FF1ACD4C1CE065B1A  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE3C66DF347F488DFE044ECF991CA2498B  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONEBF4BEC5FE5245BF3C0D271E31BC2342F  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONECE0F911DF93A8E38932C48F983E7FDE6  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE285EFBA3F9E5B1A35D6825E37A718019  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE74083D6E1ED743FA1E756A41BA467E6C  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONEB72F745822A406F2191B6EBFD8CA9DF9  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE997E4CE30BA5143E90330DB664BF61C5  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE3140177F6CCBE72992772342347A166F  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONEF843896AD42A20C4A50C0FD29512111A  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONED91A2D083B8C7311C87EDD6592FCD446  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE50EB0F4152CD0F7DCE3B471ADD31D07D  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
928 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE3B28284B39C934CB519F2A31E03A951B  
FITTRANSFORMONE CACHE LOADED 00S 00MIN

MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONEADED0EB932B9C692C801C6C56DC2BDA2  
FITTRANSFORMONE CACHE LOADED 00S 00MIN

MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONECDF3D1530A8A2388BA80A16C66413409  
FITTRANSFORMONE CACHE LOADED 00S 00MIN

MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONEFF5666A0CB943711A4E4DFD5B5CD8587  
FITTRANSFORMONE CACHE LOADED 00S 00MIN

MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE536A5FBB5A6C8D69B4FB426F0B51D9EB  
FITTRANSFORMONE CACHE LOADED 00S 00MIN

MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE26CD3A8242ACB986F98CF38291EF5BAA  
FITTRANSFORMONE CACHE LOADED 00S 00MIN

MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE8B6094F421DAB9FC88FB54C2F32E16B4  
FITTRANSFORMONE CACHE LOADED 00S 00MIN

MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONED0E8598C62DCB5F059A5C2B541C23B09  
FITTRANSFORMONE CACHE LOADED 00S 00MIN

MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE34AFFB50A3D55C5302E1C7F97AEC9679  
FITTRANSFORMONE CACHE LOADED 00S 00MIN

MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONEB3635522E8F89BBE9E56AB7E7F4693FC  
FITTRANSFORMONE CACHE LOADED 00S 00MIN

MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE7389AD9BAE0AD1181B7DA9E4DD39449D  
FITTRANSFORMONE CACHE LOADED 00S 00MIN

MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONEAB59173E00AB4A7BDF7740E5C6D50123  
FITTRANSFORMONE CACHE LOADED 00S 00MIN

MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE42894538FD0AB2B4235E55EF53B9CD59  
FITTRANSFORMONE CACHE LOADED 00S 00MIN

MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE580596A10617556E39D5E069C4DE096A  
FITTRANSFORMONE CACHE LOADED 00S 00MIN

MEMORY CALLING SKLEARNPIPELINEFITTRANSFORMONE  
FITTRANSFORMONESELECTKBESTK2 SCOREFUNCFUNCTION CHI2 AT 0X7EFE30BB2268  
↪ARRAY0 0

0 0 ARRAY0 8 NONE MESSAGECLSNAMEPIPELINE  
↪MESSAGENONE  
FITTRANSFORMONE 00S 00MIN

MEMORY CALLING SKLEARNPIPELINEFITTRANSFORMONE  
FITTRANSFORMONESELECTKBESTK2 SCOREFUNCFUNCTION CHI2 AT 0X7EFE30BB2268  
↪ARRAY0 0

0 0 ARRAY0 8 NONE MESSAGECLSNAMEPIPELINE  
↪MESSAGENONE  
FITTRANSFORMONE 00S 00MIN  
57 PIPELINES AND COMPOSITE ESTIMATORS 929

SCIKITLEARN USER GUIDE RELEASE 0213

MEMORY CALLING SKLEARNPIPELINEFITTRANSFORMONE  
FITTRANSFORMONESELECTKBESTK2 SCOREFUNCFUNCTION CHI2 AT 0X7EFE30BB2268  
↩️ARRAY0 0

0 0 ARRAY0 8 NONE MESSAGECLSNAMEPIPELINE  
↩️MESSAGENONE  
FITTRANSFORMONE 00S 00MIN

MEMORY CALLING SKLEARNPIPELINEFITTRANSFORMONE  
FITTRANSFORMONESELECTKBESTK2 SCOREFUNCFUNCTION CHI2 AT 0X7EFE30BB2268  
↩️ARRAY0 0

0 0 ARRAY0 8 NONE MESSAGECLSNAMEPIPELINE  
↩️MESSAGENONE  
FITTRANSFORMONE 00S 00MIN

MEMORY CALLING SKLEARNPIPELINEFITTRANSFORMONE  
FITTRANSFORMONESELECTKBESTK2 SCOREFUNCFUNCTION CHI2 AT 0X7EFE30BB2268  
↩️ARRAY0 0

0 0 ARRAY0 9 NONE MESSAGECLSNAMEPIPELINE  
↩️MESSAGENONE  
FITTRANSFORMONE 00S 00MIN

MEMORY CALLING SKLEARNPIPELINEFITTRANSFORMONE  
FITTRANSFORMONESELECTKBESTK4 SCOREFUNCFUNCTION CHI2 AT 0X7EFE30BB2268  
↩️ARRAY0 0

0 0 ARRAY0 8 NONE MESSAGECLSNAMEPIPELINE  
↩️MESSAGENONE  
FITTRANSFORMONE 00S 00MIN

MEMORY CALLING SKLEARNPIPELINEFITTRANSFORMONE  
FITTRANSFORMONESELECTKBESTK4 SCOREFUNCFUNCTION CHI2 AT 0X7EFE30BB2268  
↩️ARRAY0 0

0 0 ARRAY0 8 NONE MESSAGECLSNAMEPIPELINE  
↩️MESSAGENONE  
FITTRANSFORMONE 00S 00MIN

MEMORY CALLING SKLEARNPIPELINEFITTRANSFORMONE  
FITTRANSFORMONESELECTKBESTK4 SCOREFUNCFUNCTION CHI2 AT 0X7EFE30BB2268  
↩️ARRAY0 0

0 0 ARRAY0 8 NONE MESSAGECLSNAMEPIPELINE  
↩️MESSAGENONE  
FITTRANSFORMONE 00S 00MIN

MEMORY CALLING SKLEARNPIPELINEFITTRANSFORMONE  
FITTRANSFORMONESELECTKBESTK4 SCOREFUNCFUNCTION CHI2 AT 0X7EFE30BB2268  
↩️ARRAY0 0

0 0 ARRAY0 8 NONE MESSAGECLSNAMEPIPELINE  
↩️MESSAGENONE  
FITTRANSFORMONE 00S 00MIN

MEMORY CALLING SKLEARNPIPELINEFITTRANSFORMONE  
930 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213  
FITTRANSFORMONESELECTKBESTK4 SCOREFUNCFUNCTION CHI2 AT 0X7EFE30BB2268  
↪ARRAY0 0

0 0 ARRAY0 9 NONE MESSAGECLSNAMEPIPELINE  
↪MESSAGENONE  
FITTRANSFORMONE 00S 00MIN

MEMORY CALLING SKLEARNPIPELINEFITTRANSFORMONE  
FITTRANSFORMONESELECTKBESTK8 SCOREFUNCFUNCTION CHI2 AT 0X7EFE30BB2268  
↪ARRAY0 0

0 0 ARRAY0 8 NONE MESSAGECLSNAMEPIPELINE  
↪MESSAGENONE  
FITTRANSFORMONE 00S 00MIN

MEMORY CALLING SKLEARNPIPELINEFITTRANSFORMONE  
FITTRANSFORMONESELECTKBESTK8 SCOREFUNCFUNCTION CHI2 AT 0X7EFE30BB2268  
↪ARRAY0 0

0 0 ARRAY0 8 NONE MESSAGECLSNAMEPIPELINE  
↪MESSAGENONE  
FITTRANSFORMONE 00S 00MIN

MEMORY CALLING SKLEARNPIPELINEFITTRANSFORMONE  
FITTRANSFORMONESELECTKBESTK8 SCOREFUNCFUNCTION CHI2 AT 0X7EFE30BB2268  
↪ARRAY0 0

0 0 ARRAY0 8 NONE MESSAGECLSNAMEPIPELINE  
↪MESSAGENONE  
FITTRANSFORMONE 00S 00MIN

MEMORY CALLING SKLEARNPIPELINEFITTRANSFORMONE  
FITTRANSFORMONESELECTKBESTK8 SCOREFUNCFUNCTION CHI2 AT 0X7EFE30BB2268  
↪ARRAY0 0

0 0 ARRAY0 8 NONE MESSAGECLSNAMEPIPELINE  
↪MESSAGENONE  
FITTRANSFORMONE 00S 00MIN

MEMORY CALLING SKLEARNPIPELINEFITTRANSFORMONE  
FITTRANSFORMONESELECTKBESTK8 SCOREFUNCFUNCTION CHI2 AT 0X7EFE30BB2268  
↪ARRAY0 0

0 0 ARRAY0 9 NONE MESSAGECLSNAMEPIPELINE  
↪MESSAGENONE  
FITTRANSFORMONE 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE0C252F8C2AF87FE4A595CA853EA983CD  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE5FEB279ACC70B729EDE7514B133D0172  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE6742BAD9E26DD16B831FC9ABD1883C9E  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE096380A7AD90487199A4A6E8EC8A5744  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
57 PIPELINES AND COMPOSITE ESTIMATORS 931

SCIKITLEARN USER GUIDE RELEASE 0213

MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONEE5107EF0D6468F91A5CD772A596BA876  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE2EFA56788B791BECEC2FCD015ECF751F  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONED27C6056FC9AB1F3FB4327FD17346312  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONED649C35672F6E2731789F5C6A9B03066  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE4906C4AFA619398FB77288E086D5DD82  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE7DC3C29273FD5C57E2913F82C9185294  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE63C1C1AF7115994B5CD17A1189691D5B  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE670746700924677F7A550FA4CAE5C9BB  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONEF9617292C7B763D23B3F6C9D96697C29  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE9A2BCB4639D1B192E5525C2BF0DCA317  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE87C7C41885D63769A203E26B244FD823  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE0C252F8C2AF87FE4A595CA853EA983CD  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE5FEB279ACC70B729EDE7514B133D0172  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE6742BAD9E26DD16B831FC9ABD1883C9E  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE096380A7AD90487199A4A6E8EC8A5744  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONEE5107EF0D6468F91A5CD772A596BA876  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE2EFA56788B791BECEC2FCD015ECF751F  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONED27C6056FC9AB1F3FB4327FD17346312  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONED649C35672F6E2731789F5C6A9B03066  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE4906C4AFA619398FB77288E086D5DD82  
932 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE7DC3C29273FD5C57E2913F82C9185294  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE63C1C1AF7115994B5CD17A1189691D5B  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE670746700924677F7A550FA4CAE5C9BB  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONEF9617292C7B763D23B3F6C9D96697C29  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE9A2BCB4639D1B192E5525C2BF0DCA317  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE87C7C41885D63769A203E26B244FD823  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE0C252F8C2AF87FE4A595CA853EA983CD  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE5FEB279ACC70B729EDE7514B133D0172  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE6742BAD9E26DD16B831FC9ABD1883C9E  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE096380A7AD90487199A4A6E8EC8A5744  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONEE5107EF0D6468F91A5CD772A596BA876  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE2EFA56788B791BECEC2FCD015ECF751F  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONED27C6056FC9AB1F3FB4327FD17346312  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE649C35672F6E2731789F5C6A9B03066  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE4906C4AFA619398FB77288E086D5DD82  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE7DC3C29273FD5C57E2913F82C9185294  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE63C1C1AF7115994B5CD17A1189691D5B  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONE670746700924677F7A550FA4CAE5C9BB  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
↪SKLEARNPIPELINEFITTRANSFORMONEF9617292C7B763D23B3F6C9D96697C29  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
57 PIPELINES AND COMPOSITE ESTIMATORS 933

SCIKITLEARN USER GUIDE RELEASE 0213  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
'→SKLEARNPIPELINEFITTRANSFORMONE9A2BCB4639D1B192E5525C2BF0DCA317  
FITTRANSFORMONE CACHE LOADED 00S 00MIN  
MEMORY00S 00MIN LOADING FITTRANSFORMONE FROM TMPTMPL5RIBDVUJOBLIB  
'→SKLEARNPIPELINEFITTRANSFORMONE87C7C41885D63769A203E26B244FD823  
FITTRANSFORMONE CACHE LOADED 00S 00MIN

MEMORY CALLING SKLEARNPIPELINEFITTRANSFORMONE  
FITTRANSFORMONEPCAITERATEDPOWER7 NCOMPONENTS8 ARRAY0 0

0 0 ARRAY0 8 NONE MESSAGECLSNAMEPIPELINE  
'→MESSAGENONE  
FITTRANSFORMONE 00S 00MIN  
THEPCA FITTING IS ONLY COMPUTED AT THE EVALUATION OF THE FIRST CONFIGURATION OF THE CPARAMETER OF THE LINEARSVC  
CLASSIFIER THE OTHER CONFIGURATIONS OF CWILL TRIGGER THE LOADING OF THE CACHED PCA ESTIMATOR DATA LEADING TO SAVE PRO  
CESSING TIME THEREFORE THE USE OF CACHING THE PIPELINE USING MEMORY IS HIGHLY BENEFICIAL WHEN FITTING A TRANSFORMER  
IS COSTLY  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 20695 SECONDS  
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE  
575 COLUMN TRANSFORMER WITH HETEROGENEOUS DATA SOURCES  
DATASETS CAN OFTEN CONTAIN COMPONENTS OF THAT REQUIRE DIFFERENT FEATURE EXTRACTION AND PROCESSING PIPELINES THIS  
SCENARIO MIGHT OCCUR WHEN  
1 YOUR DATASET CONSISTS OF HETEROGENEOUS DATA TYPES EG RASTER IMAGES AND TEXT CAPTIONS  
2 YOUR DATASET IS STORED IN A PANDAS DATAFRAME AND DIFFERENT COLUMNS REQUIRE DIFFERENT PROCESSING PIPELINES  
THIS EXAMPLE DEMONSTRATES HOW TO USE SKLEARNCOMPOSECOLUMNTRANSFORMER ON A DATASET CONTAINING DIF  
FERENT TYPES OF FEATURES WE USE THE 20NEWSGROUPS DATASET AND COMPUTE STANDARD BAGOFWORDS FEATURES FOR THE SUBJECT  
LINE AND BODY IN SEPARATE PIPELINES AS WELL AS AD HOC FEATURES ON THE BODY WE COMBINE THEM WITH WEIGHTS USING A  
COLUMNTRANSFORMER AND FINALLY TRAIN A CLASSIFIER ON THE COMBINED SET OF FEATURES  
THE CHOICE OF FEATURES IS NOT PARTICULARLY HELPFUL BUT SERVES TO ILLUSTRATE THE TECHNIQUE  
OUT  
PIPELINE STEP 1 OF 3 PROCESSING SUBJECTBODY TOTAL 01S  
PIPELINE STEP 2 OF 3 PROCESSING UNION TOTAL 03S  
PIPELINE STEP 3 OF 3 PROCESSING SVC TOTAL 03S  
PRECISION RECALL F1SCORE SUPPORT  
0 096 062 076 494  
1 025 084 039 76  
ACCURACY 065 570  
MACRO AVG 061 073 057 570  
WEIGHTED AVG 087 065 071 570  
934 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
AUTHOR MATT TERRY MATTTERRYGMAILCOM

LICENSE BSD 3 CLAUSE  
IMPORT NUMPY AS NP  
FROM SKLEARNBASE IMPORT BASEESTIMATOR TRANSFORMERMIXIN  
FROM SKLEARNDATASETS IMPORT FETCH20NEWSGROUPS  
FROM SKLEARNDATASETSTWENTYNEWSGROUPS IMPORT STRIPNEWSGROUPFOOTER  
FROM SKLEARNDATASETSTWENTYNEWSGROUPS IMPORT STRIPNEWSGROUPQUOTING  
FROM SKLEARNDECOMPOSITION IMPORT TRUNCATEDSVD  
FROM SKLEARNFEATUREEXTRACTION IMPORT DICTVECTORIZER  
FROM SKLEARNFEATUREEXTRACTIONTEXT IMPORT TFIDFVECTORIZER  
FROM SKLEARNMETRICS IMPORT CLASSIFICATIONREPORT  
FROM SKLEARNPIPELINE IMPORT PIPELINE  
FROM SKLEARNCOMPOSE IMPORT COLUMNTRANSFORMER  
FROM SKLEARN SVM IMPORT LINEARSVC  
CLASS TEXTSTATS BASEESTIMATOR TRANSFORMERMIXIN  
EXTRACT FEATURES FROM EACH DOCUMENT FOR DICTVECTORIZER  
DEFFITSELF X YNONE  
RETURNSELF  
DEFTRANSFORMSELF POSTS  
RETURNLENGTH LENTEXT  
NUMSENTENCES TEXTCOUNT  
FORTEXTINPOSTS  
CLASS SUBJECTBODYEXTRACTOR BASEESTIMATOR TRANSFORMERMIXIN  
EXTRACT THE SUBJECT BODY FROM A USENET POST IN A SINGLE PASS  
TAKES A SEQUENCE OF STRINGS AND PRODUCES A DICT OF SEQUENCES KEYS ARE  
SUBJECT AND BODY

DEFFITSELF X YNONE  
RETURNSELF  
DEFTRANSFORMSELF POSTS  
CONSTRUCT OBJECT DTYPE ARRAY WITH TWO COLUMNS  
FIRST COLUMN SUBJECT AND SECOND COLUMN BODY  
FEATURES NPEMPTYSHAPELENPOSTS 2 DTYPEOBJECT  
FORI TEXT INENUMERATEPOSTS  
HEADERS BOD TEXTPARTITION NN  
BOD STRIPNEWSGROUPFOOTERBOD  
BOD STRIPNEWSGROUPQUOTINGBOD  
FEATURESI 1 BOD  
PREFIX SUBJECT  
SUB  
FORLINEINHEADERSPLIT N  
IFLINESTARTSWITHPREFIX  
SUB LINELENPREFIX  
BREAK  
FEATURESI 0 SUB  
57 PIPELINES AND COMPOSITE ESTIMATORS 935

SCIKITLEARN USER GUIDE RELEASE 0213  
RETURNFEATURES  
PIPELINE PIPELINE  
EXTRACT THE SUBJECT BODY  
SUBJECTBODY SUBJECTBODYEXTRACTOR  
USE COLUMNTRANSFORMER TO COMBINE THE FEATURES FROM SUBJECT AND BODY  
UNION COLUMNTRANSFORMER

PULLING FEATURES FROM THE POSTS SUBJECT LINE FIRST COLUMN  
SUBJECT TFIDFVECTORIZERMINDF50 0  
PIPELINE FOR STANDARD BAGOFWORDS MODEL FOR BODY SECOND COLUMN  
BODYBOW PIPELINE  
TFIDF TFIDFVECTORIZER  
BEST TRUNCATEDSVDNCOMPONENTS50  
1  
PIPELINE FOR PULLING AD HOC FEATURES FROM POSTS BODY  
BODYSTATS PIPELINE  
STATS TEXTSTATS RETURNS A LIST OF DICTS  
VECT DICTVECTORIZER LIST OF DICTS FEATURE MATRIX  
1

WEIGHT COMPONENTS IN COLUMNTRANSFORMER  
TRANSFORMERWEIGHTS  
SUBJECT 08  
BODYBOW 05  
BODYSTATS 10

USE A SVC CLASSIFIER ON THE COMBINED FEATURES  
SVC LINEARSVC  
VERBOSETRUE  
LIMIT THE LIST OF CATEGORIES TO MAKE RUNNING THIS EXAMPLE FASTER  
CATEGORIES ALTATHEISM TALKRELIGIONMISC  
TRAIN FETCH20NEWSGROUPSRANDOMSTATE1  
SUBSETTRAIN  
CATEGORIESCATEGORIES

TEST FETCH20NEWSGROUPSRANDOMSTATE1  
SUBSETTEST  
CATEGORIESCATEGORIES

PIPELINEFITTRAINDATA TRAINTARGET  
Y PIPELINEPREDICTTESTDATA  
PRINTCLASSIFICATIONREPORTY TESTTARGET  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 1258 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
936 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

576 EFFECT OF TRANSFORMING THE TARGETS IN REGRESSION MODEL

IN THIS EXAMPLE WE GIVE AN OVERVIEW OF THE SKLEARNCOMPOSESTRANSFORMEDTARGETREGRESSOR TWO

EXAMPLES ILLUSTRATE THE BENEFIT OF TRANSFORMING THE TARGETS BEFORE LEARNING A LINEAR REGRESSION MODEL THE FIRST EXAMPLE

USES SYNTHETIC DATA WHILE THE SECOND EXAMPLE IS BASED ON THE BOSTON HOUSING DATA SET

AUTHOR GUILLAUME LEMAITRE GUILLAUMELEMAITREINRIA.FR

LICENSE BSD 3 CLAUSE

```
import numpy as np
import matplotlib
import matplotlib.pyplot as plt
from distutils.version import LooseVersion
print doc
SYNTHETIC EXAMPLE
from sklearn.datasets import make_regression
from sklearn.model_selection import train_test_split
from sklearn.linear_model import RidgeCV
from sklearn.compose import TransformedTargetRegressor
from sklearn.metrics import median_absolute_error, r2_score
# NORMED is being deprecated in favor of density in histograms
if LooseVersion(matplotlib.__version__) < 2.1:
    density_param = Density(True)
else:
    density_param = Normed(True)
A SYNTHETIC RANDOM REGRESSION PROBLEM IS GENERATED THE TARGETS Y ARE MODIFIED BY I TRANSLATING ALL TARGETS SUCH
THAT ALL ENTRIES ARE NONNEGATIVE AND II APPLYING AN EXPONENTIAL FUNCTION TO OBTAIN NONLINEAR TARGETS WHICH CANNOT BE
FITTED USING A SIMPLE LINEAR MODEL
THEREFORE A LOGARITHMIC np.log1p AND AN EXPONENTIAL FUNCTION np.expm1 WILL BE USED TO TRANSFORM THE TARGETS
BEFORE TRAINING A LINEAR REGRESSION MODEL AND USING IT FOR PREDICTION
X, y = make_regression(n_samples=10000, noise=100, random_state=0)
y = np.exp(y) * np.log1p(y)
y_trans = np.log1p(y)
THE FOLLOWING ILLUSTRATE THE PROBABILITY DENSITY FUNCTIONS OF THE TARGET BEFORE AND AFTER APPLYING THE LOGARITHMIC FUNC
TIONS
fig, (ax0, ax1) = plt.subplots(1, 2)
ax0.hist(y, bins=100, density_param=density_param)
ax0.set_xlim(0, 2000)
ax0.set_ylabel('PROBABILITY')
ax0.set_xlabel('TARGET')
ax0.set_title('TARGET DISTRIBUTION')
ax1.hist(y_trans, bins=100, density_param=density_param)
ax1.set_ylabel('PROBABILITY')
```

57 PIPELINES AND COMPOSITE ESTIMATORS 937

SCIKITLEARN USER GUIDE RELEASE 0213  
AX1SETXLABELTARGET  
AX1SETTITLETRANSFORMED TARGET DISTRIBUTION  
FSUPTITLESYNTHETIC DATA Y0035  
FTIGHTLAYOUTRECT005 005 095 095  
XTRAIN XTEST YTRAIN YTEST TRAINTESTSPLITX Y RANDOMSTATE0  
AT FIRST A LINEAR MODEL WILL BE APPLIED ON THE ORIGINAL TARGETS DUE TO THE NONLINEARITY THE MODEL TRAINED WILL NOT  
BE PRECISE DURING THE PREDICTION SUBSEQUENTLY A LOGARITHMIC FUNCTION IS USED TO LINEARIZE THE TARGETS ALLOWING BETTER  
PREDICTION EVEN WITH A SIMILAR LINEAR MODEL AS REPORTED BY THE MEDIAN ABSOLUTE ERROR MAE  
F AX0 AX1 PLTSUBPLOTS1 2 SHAREYTRUE  
REGR RIDGECV  
REGRFITXTRAIN YTRAIN  
YPRED REGRPREDICTXTEST  
AX0SCATTERYTEST YPRED  
AX0PLOT0 2000 0 2000 K  
AX0SETYLABELTARGET PREDICTED  
AX0SETXLABELTRUE TARGET  
AX0SETTITLERIDGE REGRESSION NWWITHOUT TARGET TRANSFORMATION  
AX0TEXT100 1750 RR2 2F MAE2F  
R2SCOREYTEST YPRED MEDIANABSOLUTEERRORYTEST YPRED  
AX0SETXLIM0 2000  
938 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213  
AX0SETYLIM0 2000  
REGRTRANS TRANSFORMEDTARGETREGRESSORREGRESSORRIDGECV  
FUNCPLOG1P  
INVERSEFUNCPPEXP1  
REGRTRANSFITXTRAIN YTRAIN  
YPRED REGRTRANSPREDICTXTEST  
AX1SCATTERYTEST YPRED  
AX1PLOT0 2000 0 2000 K  
AX1SETYLABELTARGET PREDICTED  
AX1SETXLABELTRUE TARGET  
AX1SETTITLERIDGE REGRESSION NWITH TARGET TRANSFORMATION  
AX1TEXT100 1750 RR2 2F MAE2F  
R2SCOREYTEST YPRED MEDIANABSOLUTEERRORYTEST YPRED  
AX1SETXLIM0 2000  
AX1SETYLIM0 2000  
FSUPTITLESYNTHETIC DATA Y0035  
FTIGHTLAYOUTRECT005 005 095 095  
57 PIPELINES AND COMPOSITE ESTIMATORS 939

SCIKITLEARN USER GUIDE RELEASE 0213  
REALWORLD DATA SET  
IN A SIMILAR MANNER THE BOSTON HOUSING DATA SET IS USED TO SHOW THE IMPACT OF TRANSFORMING THE TARGETS BEFORE LEARNING A  
MODEL IN THIS EXAMPLE THE TARGETS TO BE PREDICTED CORRESPONDS TO THE WEIGHTED DISTANCES TO THE FIVE BOSTON EMPLOYMENT  
CENTERS  
FROM SKLEARNDATASETS IMPORT LOADBOSTON  
FROM SKLEARNPREPROCESSING IMPORT QUANTILETRANSFORMER QUANTILETRANSFORM  
DATASET LOADBOSTON  
TARGET NPARRAYDATASETFEATURENAMES DIS  
X DATASETDATA NPLOGICALNOTTARGET  
Y DATASETDATA TARGETSQUEEZE  
YTRANS QUANTILETRANSFORMDATASETDATA TARGET  
NQUANTILES300  
OUTPUTDISTRIBUTIONNORMAL  
COPYTRUESQUEEZE  
ASKLEARNPREPROCESSINGQUANTILETRANSFORMER IS USED SUCH THAT THE TARGETS FOLLOWS A NORMAL DISTRI  
BUTION BEFORE APPLYING A SKLEARNLINEARMODELRIDGECV MODEL  
F AX0 AX1 PLTSUBPLOTS1 2  
AX0HISTY BINS100 DENSITYPARAM  
AX0SETYLABELPROBABILITY  
AX0SETXLABELTARGET  
AX0SETTITLETARGET DISTRIBUTION  
AX1HISTYTRANS BINS100 DENSITYPARAM  
AX1SETYLABELPROBABILITY  
AX1SETXLABELTARGET  
AX1SETTITLETRANSFORMED TARGET DISTRIBUTION  
FSUPTITLEBOSTON HOUSING DATA DISTANCE TO EMPLOYMENT CENTERS Y0035  
FTIGHTLAYOUTRECT005 005 095 095  
XTRAIN XTEST YTRAIN YTEST TRAINTESTSPLITX Y RANDOMSTATE1  
940 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

THE EFFECT OF THE TRANSFORMER IS WEAKER THAN ON THE SYNTHETIC DATA HOWEVER THE TRANSFORM INDUCES A DECREASE OF THE MAE

F AX0 AX1 PLTSUBPLOTS1 2 SHAREYTRUE

REGR RIDGECV

REGRFITXTRAIN YTRAIN

YPRED REGRPREDICTXTEST

AX0SCATTERYTEST YPRED

AX0PLOT0 10 0 10 K

AX0SETYLABELTARGET PREDICTED

AX0SETXLABELTRUE TARGET

AX0SETTITLERIDGE REGRESSION NWITHOUT TARGET TRANSFORMATION

AX0TEXT1 9 RR2 2F MAE2F

R2SCOREYTEST YPRED MEDIANABSOLUTEERRORYTEST YPRED

AX0SETXLIM0 10

AX0SETYLIM0 10

REGRTRANS TRANSFORMEDTARGETREGRESSOR

REGRESSORRIDGECV

TRANSFORMERQUANTILETRANSFORMERNQUANTILES300

OUTPUTDISTRIBUTIONNORMAL

REGRTRANSFITXTRAIN YTRAIN

YPRED REGRTRANSPREDICTXTEST

57 PIPELINES AND COMPOSITE ESTIMATORS 941

SCIKITLEARN USER GUIDE RELEASE 0213  
AX1SCATTERYTEST YPRED  
AX1PLOT0 10 0 10 K  
AX1SETYLABELTARGET PREDICTED  
AX1SETXLABELTRUE TARGET  
AX1SETTITLERIDGE REGRESSION NWITH TARGET TRANSFORMATION  
AX1TEXT1 9 RR2 2F MAE2F  
R2SCOREYTEST YPRED MEDIANABSOLUTEERRORYTEST YPRED  
AX1SETXLIM0 10  
AX1SETYLIM0 10  
FSUPTITLEBOSTON HOUSING DATA DISTANCE TO EMPLOYMENT CENTERS Y0035  
FTIGHTLAYOUTRECT005 005 095 095  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 1682 SECONDS  
58 COVARIANCE ESTIMATION  
EXAMPLES CONCERNING THE SKLEARNCOVARIANCE MODULE  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
942 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

581 LEDOITWOLF VS OAS ESTIMATION

THE USUAL COVARIANCE MAXIMUM LIKELIHOOD ESTIMATE CAN BE REGULARIZED USING SHRINKAGE LEDOIT AND WOLF PROPOSED A CLOSE FORMULA TO COMPUTE THE ASYMPTOTICALLY OPTIMAL SHRINKAGE PARAMETER MINIMIZING A MSE CRITERION YIELDING THE LEDOITWOLF COVARIANCE ESTIMATE

CHEN ET AL PROPOSED AN IMPROVEMENT OF THE LEDOITWOLF SHRINKAGE PARAMETER THE OAS COEFFICIENT WHOSE CONVERGENCE IS SIGNIFICANTLY BETTER UNDER THE ASSUMPTION THAT THE DATA ARE GAUSSIAN

THIS EXAMPLE INSPIRED FROM CHEN’S PUBLICATION 1 SHOWS A COMPARISON OF THE ESTIMATED MSE OF THE LW AND OAS METHODS USING GAUSSIAN DISTRIBUTED DATA

1 “SHRINKAGE ALGORITHMS FOR MMSE COVARIANCE ESTIMATION” CHEN ET AL IEEE TRANS ON SIGN PROC V OLUME 58 ISSUE 10 OCTOBER 2010

```
PRINTDOC
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SCIPYLINALG IMPORT TOEPLITZ CHOLESKY
FROM SKLEARNCOVARIANCE IMPORT LEDOITWOLF OAS
NPRANDOMSEED0
NFEATURES 100
SIMULATION COVARIANCE MATRIX AR1 PROCESS
R 01
REALCOV TOEPLITZR NPARANGENFEATURES
COLORINGMATRIX CHOLESKYREALCOV
NSAMPLESRANGE NPARANGE6 31 1
REPEAT 100
LWMSE NPZEROSNSAMPLESRANGESIZE REPEAT
OAMSE NPZEROSNSAMPLESRANGESIZE REPEAT
LWSHRINKAGE NPZEROSNSAMPLESRANGESIZE REPEAT
OASHRINKAGE NPZEROSNSAMPLESRANGESIZE REPEAT
FORI NSAMPLES INENUMERATENSAMPLESRANGE
FORJINRANGEREPEAT
X NPDOT
NPRANDOMNORMALSIZENSAMPLES NFEATURES COLORINGMATRIXT
LW LEDOITWOLFSTOREPRECISIONFALSE ASSUMECENTEREDTRUE
LWFITX
LWMSEI J LWERRORNORMREALCOV SCALINGFALSE
LWSHRINKAGEI J LWSHRINKAGE
OA OASSTOREPRECISIONFALSE ASSUMECENTEREDTRUE
OAFITX
OAMSEI J OAERRORNORMREALCOV SCALINGFALSE
OASHRINKAGEI J OASHRINKAGE
PLOT MSE
PLTSUBPLOT2 1 1
PLTERRORBARNSAMPLESRANGE LWMSEMEAN1 YERRRLWMSESTD1
LABELLEDOITWOLF COLORNAVY LW2
PLTERRORBARNSAMPLESRANGE OAMSEMEAN1 YERROAMSESTD1
LABELOAS COLORDARKORANGE LW2
PLTYLABELSQUARED ERROR
58 COVARIANCE ESTIMATION 943
```

SCIKITLEARN USER GUIDE RELEASE 0213  
PLTLEGENDLOCUPPER RIGHT  
PLTTITLECOMPARISON OF COVARIANCE ESTIMATORS  
PLTXLIM5 31  
PLOT SHRINKAGE COEFFICIENT  
PLTSUBPLOT2 1 2  
PLTERRORBARNSAMPLESRANGE LWSHRINKAGEMEAN1 YERRLWSHRINKAGESTD1  
LABELLEDOITWOLF COLORNNAVY LW2  
PLTERRORBARNSAMPLESRANGE OASHRINKAGEMEAN1 YERROASHRINKAGESTD1  
LABELOAS COLORDARKORANGE LW2  
PLTXLABELNSAMPLES  
PLTYLABELSHRINKAGE  
PLTLEGENDLOCLOWER RIGHT  
PLTYLIMPLTYLIM0 1 PLTYLIM1 PLTYLIM0 10  
PLTXLIM5 31  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 3818 SECONDS  
NOTE [CLICK HERE](#) TO DOWNLOAD THE FULL EXAMPLE CODE  
944 CHAPTER 5 EXAMPLES

582 SPARSE INVERSE COVARIANCE ESTIMATION

USING THE GRAPHICALASSO ESTIMATOR TO LEARN A COVARIANCE AND SPARSE PRECISION FROM A SMALL NUMBER OF SAMPLES TO ESTIMATE A PROBABILISTIC MODEL EG A GAUSSIAN MODEL ESTIMATING THE PRECISION MATRIX THAT IS THE INVERSE COVARIANCE MATRIX IS AS IMPORTANT AS ESTIMATING THE COVARIANCE MATRIX INDEED A GAUSSIAN MODEL IS PARAMETRIZED BY THE PRECISION MATRIX

TO BE IN FAVORABLE RECOVERY CONDITIONS WE SAMPLE THE DATA FROM A MODEL WITH A SPARSE INVERSE COVARIANCE MATRIX IN ADDITION WE ENSURE THAT THE DATA IS NOT TOO MUCH CORRELATED LIMITING THE LARGEST COEFFICIENT OF THE PRECISION MATRIX AND THAT THERE A NO SMALL COEFFICIENTS IN THE PRECISION MATRIX THAT CANNOT BE RECOVERED IN ADDITION WITH A SMALL NUMBER OF OBSERVATIONS IT IS EASIER TO RECOVER A CORRELATION MATRIX RATHER THAN A COVARIANCE THUS WE SCALE THE TIME SERIES HERE THE NUMBER OF SAMPLES IS SLIGHTLY LARGER THAN THE NUMBER OF DIMENSIONS THUS THE EMPIRICAL COVARIANCE IS STILL INVERTIBLE HOWEVER AS THE OBSERVATIONS ARE STRONGLY CORRELATED THE EMPIRICAL COVARIANCE MATRIX IS ILLCONDITIONED AND AS A RESULT ITS INVERSE -THE EMPIRICAL PRECISION MATRIX- IS VERY FAR FROM THE GROUND TRUTH

IF WE USE L2 SHRINKAGE AS WITH THE LEDOITWOLF ESTIMATOR AS THE NUMBER OF SAMPLES IS SMALL WE NEED TO SHRINK A LOT AS A RESULT THE LEDOITWOLF PRECISION IS FAIRLY CLOSE TO THE GROUND TRUTH PRECISION THAT IS NOT FAR FROM BEING DIAGONAL BUT THE OFFDIAGONAL STRUCTURE IS LOST

THE L1PENALIZED ESTIMATOR CAN RECOVER PART OF THIS OFFDIAGONAL STRUCTURE IT LEARNS A SPARSE PRECISION IT IS NOT ABLE TO RECOVER THE EXACT SPARSITY PATTERN IT DETECTS TOO MANY NONZERO COEFFICIENTS HOWEVER THE HIGHEST NONZERO COEFFICIENTS OF THE L1 ESTIMATED CORRESPOND TO THE NONZERO COEFFICIENTS IN THE GROUND TRUTH FINALLY THE COEFFICIENTS OF THE L1 PRECISION ESTIMATE ARE BIASED TOWARD ZERO BECAUSE OF THE PENALTY THEY ARE ALL SMALLER THAN THE CORRESPONDING GROUND TRUTH VALUE AS CAN BE SEEN ON THE FIGURE

NOTE THAT THE COLOR RANGE OF THE PRECISION MATRICES IS TWEAKED TO IMPROVE READABILITY OF THE FIGURE THE FULL RANGE OF VALUES OF THE EMPIRICAL PRECISION IS NOT DISPLAYED

THE ALPHA PARAMETER OF THE GRAPHICALASSO SETTING THE SPARSITY OF THE MODEL IS SET BY INTERNAL CROSSVALIDATION IN THE GRAPHICALASSOCV AS CAN BE SEEN ON FIGURE 2 THE GRID TO COMPUTE THE CROSSVALIDATION SCORE IS ITERATIVELY REFINED IN THE NEIGHBORHOOD OF THE MAXIMUM

SCIKITLEARN USER GUIDE RELEASE 0213

•

```
PRINTDOC
AUTHOR GAEL VAROQUAUX GAELVAROQUAUXINRIA.FR
LICENSE BSD 3 CLAUSE
COPYRIGHT INRIA
IMPORT NUMPY AS NP
FROM SCIPY IMPORT LINALG
FROM SKLEARNDATASETS IMPORT MAKESPARSESPDMATRIX
FROM SKLEARNCOVARIANCE IMPORT GRAPHICALLASSOCV LEDOITWOLF
IMPORT MATPLOTLIBPYPLOT AS PLT
```

```
GENERATE THE DATA
NSAMPLES 60
NFEATURES 20
PRNG NPRANDOMRANDOMSTATE1
PREC MAKESPARSESPDMATRIXNFEATURES ALPHA98
SMALLESTCOEF4
LARGESTCOEF7
RANDOMSTATEPRNG
COV LINALGINVPREC
D NPSQRTNPDIAGCOV
COV D
COV D NPNEWAXIS
PREC D
PREC D NPNEWAXIS
X PRNGMULTIVARIATENORMALNPZEROSNFEATURES COV SIZESAMPLES
X XMEANAXIS0
X XSTDAXIS0
```

```
ESTIMATE THE COVARIANCE
EMPCOV NPDOTXT X NSAMPLES
MODEL GRAPHICALLASSOCVCV5
MODELFITX
COV MODELCOVARIANCE
946 CHAPTER 5 EXAMPLES
```



```
SCIKITLEARN USER GUIDE RELEASE 0213
PREC MODELPRECISION
LWCOV LEDOITWOLFX
LWPREC LINALGINVLWCOV

PLOT THE RESULTS
PLTFIGUREFIGSIZE10 6
PLTSUBPLOTSADJUSTLEFT002 RIGHT098
PLOT THE COVARIANCES
COVS EMPIRICAL EMPCOV LEDOITWOLF LWCOV
GRAPHICALASSOCV COV TRUE COV
VMAX COVMAX
FORI NAME THISCOV INENUMERATECOVS
PLTSUBPLOT2 4 I 1
PLTIMSHOWTHISCOV INTERPOLATIONNEAREST VMINVMAX VMAXVMAX
CMAPPLTCMRDBUR
PLTXTICKS
PLTYTICKS
PLTTITLE SCOVARIANCE NAME
PLOT THE PRECISIONS
PRECS EMPIRICAL LINALGINVEMPCOV LEDOITWOLF LWPREC
GRAPHICALASSO PREC TRUE PREC
VMAX 9 PRECMAX
FORI NAME THISPREC INENUMERATEPRECS
AX PLTSUBPLOT2 4 I 5
PLTIMSHOWNPMAMASKEDEQUALTHISPREC 0
INTERPOLATIONNEAREST VMINVMAX VMAXVMAX
CMAPPLTCMRDBUR
PLTXTICKS
PLTYTICKS
PLTTITLE SPRECISION NAME
IFHASATTRAX SETFACECOLOR
AXSETFACECOLOR7
ELSE
AXSETAXISBGCOLOR7
PLOT THE MODEL SELECTION METRIC
PLTFIGUREFIGSIZE4 3
PLTAXES2 15 75 7
PLTPLOTMODELCVALPHAS NPMEANMODELGRIDSCORES AXIS1 0
PLTAXVLINEMODELALPHA COLOR5
PLTTITLEMODEL SELECTION
PLTYLABELCROSSVALIDATION SCORE
PLTXLABELALPHA
PLTSHOW
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0501 SECONDS
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE
58 COVARIANCE ESTIMATION 947
```

583 SHRINKAGE COVARIANCE ESTIMATION LEDOITWOLF VS OAS AND MAXLIKELIHOOD

WHEN WORKING WITH COVARIANCE ESTIMATION THE USUAL APPROACH IS TO USE A MAXIMUM LIKELIHOOD ESTIMATOR SUCH AS THE SKLEARN COVARIANCE EMPIRICAL COVARIANCE. IT IS UNBIASED IE IT CONVERGES TO THE TRUE POPULATION COVARIANCE WHEN GIVEN MANY OBSERVATIONS HOWEVER IT CAN ALSO BE BENEFICIAL TO REGULARIZE IT IN ORDER TO REDUCE ITS VARIANCE THIS IN TURN INTRODUCES SOME BIAS THIS EXAMPLE ILLUSTRATES THE SIMPLE REGULARIZATION USED IN SHRUNK COVARIANCE ESTIMATORS IN PARTICULAR IT FOCUSES ON HOW TO SET THE AMOUNT OF REGULARIZATION IE HOW TO CHOOSE THE BIAS VARIANCE TRADEOFF

HERE WE COMPARE 3 APPROACHES

- SETTING THE PARAMETER BY CROSSVALIDATING THE LIKELIHOOD ON THREE FOLDS ACCORDING TO A GRID OF POTENTIAL SHRINKAGE PARAMETERS

- A CLOSE FORMULA PROPOSED BY LEDOIT AND WOLF TO COMPUTE THE ASYMPTOTICALLY OPTIMAL REGULARIZATION PARAMETER MINIMIZING A MSE CRITERION YIELDING THE SKLEARN COVARIANCE LEDOITWOLF COVARIANCE ESTIMATE

- AN IMPROVEMENT OF THE LEDOITWOLF SHRINKAGE THE SKLEARN COVARIANCE OAS PROPOSED BY CHEN ET AL

ITS CONVERGENCE IS SIGNIFICANTLY BETTER UNDER THE ASSUMPTION THAT THE DATA ARE GAUSSIAN IN PARTICULAR FOR SMALL SAMPLES

TO QUANTIFY ESTIMATION ERROR WE PLOT THE LIKELIHOOD OF UNSEEN DATA FOR DIFFERENT VALUES OF THE SHRINKAGE PARAMETER WE ALSO SHOW THE CHOICES BY CROSSVALIDATION OR WITH THE LEDOITWOLF AND OAS ESTIMATES

NOTE THAT THE MAXIMUM LIKELIHOOD ESTIMATE CORRESPONDS TO NO SHRINKAGE AND THUS PERFORMS POORLY THE LEDOITWOLF ESTIMATE PERFORMS REALLY WELL AS IT IS CLOSE TO THE OPTIMAL AND IS COMPUTATIONAL NOT COSTLY IN THIS EXAMPLE THE OAS ESTIMATE IS A BIT FURTHER AWAY INTERESTINGLY BOTH APPROACHES OUTPERFORM CROSSVALIDATION WHICH IS SIGNIFICANTLY MOST COMPUTATIONALLY COSTLY

```
SCIKITLEARN USER GUIDE RELEASE 0213
PRINTDOC
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SCIPY IMPORT LINALG
FROM SKLEARNCOVARIANCE IMPORT LEDOITWOLF OAS SHRUNKCOVARIANCE
LOGLIKELIHOOD EMPIRICALCOVARIANCE
FROM SKLEARNMODELSELECTION IMPORT GRIDSEARCHCV

GENERATE SAMPLE DATA
NFEATURES NSAMPLES 40 20
NPRANDOMSEED42
BASEXTRAIN NPRANDOMNORMALSIZENSAMPLES NFEATURES
BASEXTEST NPRANDOMNORMALSIZENSAMPLES NFEATURES
COLOR SAMPLES
COLORINGMATRIX NPRANDOMNORMALSIZENFEATURES NFEATURES
XTRAIN NPDOTBASEXTRAIN COLORINGMATRIX
XTEST NPDOTBASEXTEST COLORINGMATRIX

COMPUTE THE LIKELIHOOD ON TEST DATA
58 COVARIANCE ESTIMATION 949
```

SCIKITLEARN USER GUIDE RELEASE 0213  
SPANNING A RANGE OF POSSIBLE SHRINKAGE COEFFICIENT VALUES  
SHRINKAGES NPLOGSPACE2 0 30  
NEGATIVELOGLIKS SHRUNKCOVARIANCESHRINKAGESFITXTRAINSCOREXTEST  
FORSINSHRINKAGES  
UNDER THE GROUNDTRUTH MODEL WHICH WE WOULD NOT HAVE ACCESS TO IN REAL  
SETTINGS  
REALCOV NPDOTCOLORINGMATRIXT COLORINGMATRIX  
EMPCOV EMPIRICALCOVARIANCEXTRAIN  
LOGLIKREAL LOGLIKELIHOODEMPCOV LINALGINVREALCOV

COMPARE DIFFERENT APPROACHES TO SETTING THE PARAMETER  
GRIDSEARCH FOR AN OPTIMAL SHRINKAGE COEFFICIENT  
TUNEDPARAMETERS SHRINKAGE SHRINKAGES  
CV GRIDSEARCHCVSHRUNKCOVARIANCE TUNEDPARAMETERS CV5  
CVFITXTRAIN  
LEDOITWOLF OPTIMAL SHRINKAGE COEFFICIENT ESTIMATE  
LW LEDOITWOLF  
LOGLIKW LWFITXTRAINSCOREXTEST  
OAS COEFFICIENT ESTIMATE  
OA OAS  
LOGLIKOA OAFITXTRAINSCOREXTEST

PLOT RESULTS  
FIG PLTFigure  
PLTTITLEREGULARIZED COVARIANCE LIKELIHOOD AND SHRINKAGE COEFFICIENT  
PLTXLABELREGULARIZATION PARAMETER SHRINKAGE COEFFICIENT  
PLTYLABELERROR NEGATIVE LOGLIKELIHOOD ON TEST DATA  
RANGE SHRINKAGE CURVE  
PLTLOGLOGSHRINKAGES NEGATIVELOGLIKS LABELNEGATIVE LOGLIKELIHOOD  
PLTPLOTPLTXLIM 2 LOGLIKREAL R  
LABELREAL COVARIANCE LIKELIHOOD  
ADJUST VIEW  
LIKMAX NPAMAXNEGATIVELOGLIKS  
LIKMIN NPAMINNEGATIVELOGLIKS  
YMIN LIKMIN 6 NPLOGPLTYLIM1 PLTYLIM0  
YMAX LIKMAX 10 NPLOGLIKMAX LIKMIN  
XMIN SHRINKAGES0  
XMAX SHRINKAGES1  
LW LIKELIHOOD  
PLTVLINESLWSHRINKAGE YMIN LOGLIKW COLORMAGENTA  
LINEWIDTH3 LABELLEDOITWOLF ESTIMATE  
OAS LIKELIHOOD  
PLTVLINESOASHRINKAGE YMIN LOGLIKOA COLORPURPLE  
LINEWIDTH3 LABELOAS ESTIMATE  
BEST CV ESTIMATOR LIKELIHOOD  
PLTVLINESCVBESTESTIMATORSHRINKAGE YMIN  
CVBESTESTIMATORSCOREXTEST COLORCYAN  
LINEWIDTH3 LABELCROSSVALIDATION BEST ESTIMATE  
PLTYLIMYMIN YMAX  
950 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

PLTXLIMXMIN XMAX

PLTLEGEND

PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0183 SECONDS

NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE

584 ROBUST COVARIANCE ESTIMATION AND MAHALANOBIS DISTANCES RELEVANCE

AN EXAMPLE TO SHOW COVARIANCE ESTIMATION WITH THE MAHALANOBIS DISTANCES ON GAUSSIAN DISTRIBUTED DATA

FOR GAUSSIAN DISTRIBUTED DATA THE DISTANCE OF AN OBSERVATION  $\mu$  TO THE MODE OF THE DISTRIBUTION CAN BE COMPUTED USING

ITS MAHALANOBIS DISTANCE  $\sqrt{(\mu - \mu_0)^T \Sigma^{-1} (\mu - \mu_0)}$  WHERE  $\mu_0$  AND  $\Sigma$  ARE THE LOCATION AND THE COVARIANCE OF

THE UNDERLYING GAUSSIAN DISTRIBUTION

IN PRACTICE  $\mu_0$  AND  $\Sigma$  ARE REPLACED BY SOME ESTIMATES THE USUAL COVARIANCE MAXIMUM LIKELIHOOD ESTIMATE IS VERY

SENSITIVE TO THE PRESENCE OF OUTLIERS IN THE DATA SET AND THEREFOR THE CORRESPONDING MAHALANOBIS DISTANCES ARE ONE

WOULD BETTER HAVE TO USE A ROBUST ESTIMATOR OF COVARIANCE TO GUARANTEE THAT THE ESTIMATION IS RESISTANT TO “ERRONEOUS”

OBSERVATIONS IN THE DATA SET AND THAT THE ASSOCIATED MAHALANOBIS DISTANCES ACCURATELY REFLECT THE TRUE ORGANISATION OF

THE OBSERVATIONS

THE MINIMUM COVARIANCE DETERMINANT ESTIMATOR IS A ROBUST HIGHBREAKDOWN POINT IE IT CAN BE USED TO ESTIMATE THE

COVARIANCE MATRIX OF HIGHLY CONTAMINATED DATASETS UP TO  $\sqrt[n]{n}$  SAMPLES –  $\sqrt[p]{p}$  FEATURES – 1

2 OUTLIERS ESTIMATOR OF COVARIANCE THE IDEA IS

TO FIND  $\sqrt[n]{n}$  SAMPLES  $\sqrt[p]{p}$  FEATURES 1

2 OBSERVATIONS WHOSE EMPIRICAL COVARIANCE HAS THE SMALLEST DETERMINANT YIELDING A “PURE” SUBSET

OF OBSERVATIONS FROM WHICH TO COMPUTE STANDARDS ESTIMATES OF LOCATION AND COVARIANCE

THE MINIMUM COVARIANCE DETERMINANT ESTIMATOR MCD HAS BEEN INTRODUCED BY PJROUSSEUW IN 1

THIS EXAMPLE ILLUSTRATES HOW THE MAHALANOBIS DISTANCES ARE AFFECTED BY OUTLYING DATA OBSERVATIONS DRAWN FROM A

CONTAMINATING DISTRIBUTION ARE NOT DISTINGUISHABLE FROM THE OBSERVATIONS COMING FROM THE REAL GAUSSIAN DISTRIBUTION

THAT ONE MAY WANT TO WORK WITH USING MCD BASED MAHALANOBIS DISTANCES THE TWO POPULATIONS BECOME DISTINGUISH

ABLE ASSOCIATED APPLICATIONS ARE OUTLIERS DETECTION OBSERVATIONS RANKING CLUSTERING FOR VISUALIZATION PURPOSE

THE CUBIC ROOT OF THE MAHALANOBIS DISTANCES ARE REPRESENTED IN THE BOXPLOT AS WILSON AND HILFERTY SUGGEST 2

1 P J ROUSSEUW LEAST MEDIAN OF SQUARES REGRESSION J AM STAT ASS 79871 1984

2 WILSON E B HILFERTY M M 1931 THE DISTRIBUTION OF CHISQUARE PROCEEDINGS OF THE NATIONAL

ACADEMY OF SCIENCES OF THE UNITED STATES OF AMERICA 17 684688

58 COVARIANCE ESTIMATION 951

```
SCIKITLEARN USER GUIDE RELEASE 0213
PRINTDOC
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARNCOVARIANCE IMPORT EMPIRICALCOVARIANCE MINCOVDET
NSAMPLES 125
NOUTLIERS 25
NFEATURES 2
GENERATE DATA
GENCOV NPEYENFEATURES
GENCOV0 0 2
X NPDOTNPRANDOMRANDNNSAMPLES NFEATURES GENCOV
ADD SOME OUTLIERS
OUTLIERSCOV NPEYENFEATURES
OUTLIERSCOVNPARANAGE1 NFEATURES NPARANGE1 NFEATURES 7
XNOUTLIERS NPDOTNPRANDOMRANDNNOUTLIERS NFEATURES OUTLIERSCOV
FIT A MINIMUM COVARIANCE DETERMINANT MCD ROBUST ESTIMATOR TO DATA
ROBUSTCOV MINCOVDETFITX
COMPARE ESTIMATORS LEARNED FROM THE FULL DATA SET WITH TRUE PARAMETERS
EMPCOV EMPIRICALCOVARIANCEFITX
952 CHAPTER 5 EXAMPLES
```

SCIKITLEARN USER GUIDE RELEASE 0213

DISPLAY RESULTS  
FIG PLTFigure  
PLTSUBPLOTSADJUSTHSPACE1 WSPACE4 TOP95 BOTTOM05  
SHOW DATA SET  
SUBFIG1 PLTSUBPLOT3 1 1  
INLIERPLOT SUBFIG1SCATTERX 0 X 1  
COLORBLACK LABELINLIERS  
OUTLIERPLOT SUBFIG1SCATTERX 0NOUTLIERS X 1NOUTLIERS  
COLORRED LABELOUTLIERS  
SUBFIG1SETXLIMSUBFIG1GETXLIM0 11  
SUBFIG1SETTITLEMAHALANOBIS DISTANCES OF A CONTAMINATED DATA SET  
SHOW CONTOURS OF THE DISTANCE FUNCTIONS  
XX YY NPMESHGRIDNPLINSPACEPLTXLIM0 PLTXLIM1 100  
NPLINSPACEPLTYLIM0 PLTYLIM1 100  
ZZ NPCXXRAVEL YYRAVEL  
MAHALEMPCOV EMPCOVMAHALANOBISZZ  
MAHALEMPCOV MAHALEMPCOVRESHAPEXXSHAPE  
EMPCOVCONTOUR SUBFIG1CONTOURXX YY NPSQRTMAHALEMPCOV  
CMAPPLTCMPUBUR  
LINESTYLES DASHED  
MAHALROBUSTCOV ROBUSTCOVMAHALANOBISZZ  
MAHALROBUSTCOV MAHALROBUSTCOVRESHAPEXXSHAPE  
ROBUSTCONTOUR SUBFIG1CONTOURXX YY NPSQRTMAHALROBUSTCOV  
CMAPPLTCMYLORBRR LINESTYLES DOTTED  
SUBFIG1LEGENDEMPCOVCONTOURCOLLECTIONS1 ROBUSTCONTOURCOLLECTIONS1  
INLIERPLOT OUTLIERPLOT  
MLE DIST ROBUST DIST INLIERS OUTLIERS  
LOCUPPER RIGHT BORDERAXESPAD0  
PLXTICKS  
PLTYTICKS  
PLOT THE SCORES FOR EACH POINT  
EMPMahal EMPCOVMAHALANOBISX NPMEANX 0 033  
SUBFIG2 PLTSUBPLOT2 2 3  
SUBFIG2BOXPLOTMPMAHALNOUTLIERS EMPMAHALNOUTLIERS WIDTHS25  
SUBFIG2PLOTNPFULLNSAMPLES NOUTLIERS 126  
EMPMahalNOUTLIERS K MARKEREDGEWIDTH1  
SUBFIG2PLOTNPFULLNOUTLIERS 226  
EMPMahalNOUTLIERS K MARKEREDGEWIDTH1  
SUBFIG2AXESSETXTICKLABELSINLIERS OUTLIERS SIZE15  
SUBFIG2SETYLABELRSQRT3RMMahal DIST SIZE16  
SUBFIG2SETTITLE1 FROM NONROBUST ESTIMATES NMAXIMUM LIKELIHOOD  
PLTYTICKS  
ROBUSTMahal ROBUSTCOVMAHALANOBISX ROBUSTCOVLOCATION 033  
SUBFIG3 PLTSUBPLOT2 2 4  
SUBFIG3BOXPLOTROBUSTMAHALNOUTLIERS ROBUSTMAHALNOUTLIERS  
WIDTHS25  
SUBFIG3PLOTNPFULLNSAMPLES NOUTLIERS 126  
ROBUSTMAHALNOUTLIERS K MARKEREDGEWIDTH1  
SUBFIG3PLOTNPFULLNOUTLIERS 226  
ROBUSTMAHALNOUTLIERS K MARKEREDGEWIDTH1  
58 COVARIANCE ESTIMATION 953

SCIKITLEARN USER GUIDE RELEASE 0213

SUBFIG3AXESSETXTICKLABELSINLIERS OUTLIERS SIZE15

SUBFIG3SETYLABELRSQRT3RMMahal DIST SIZE16

SUBFIG3SETTITLE2 FROM ROBUST ESTIMATES NMINIMUM COVARIANCE DETERMINANT

PLTYTICKS

PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0298 SECONDS

NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE

585 ROBUST VS EMPIRICAL COVARIANCE ESTIMATE

THE USUAL COVARIANCE MAXIMUM LIKELIHOOD ESTIMATE IS VERY SENSITIVE TO THE PRESENCE OF OUTLIERS IN THE DATA SET IN SUCH A CASE IT WOULD BE BETTER TO USE A ROBUST ESTIMATOR OF COVARIANCE TO GUARANTEE THAT THE ESTIMATION IS RESISTANT TO “ERRONEOUS” OBSERVATIONS IN THE DATA SET12

MINIMUM COVARIANCE DETERMINANT ESTIMATOR

THE MINIMUM COVARIANCE DETERMINANT ESTIMATOR IS A ROBUST HIGHBREAKDOWN POINT IE IT CAN BE USED TO ESTIMATE THE COVARIANCE MATRIX OF HIGHLY CONTAMINATED DATASETS UP TO  $\sqrt{\text{SAMPLES}-\text{FEATURES}-1}$

2OUTLIERS ESTIMATOR OF COVARIANCE THE IDEA IS TO FIND  $\sqrt{\text{SAMPLES}-\text{FEATURES}}$ 1

2OBSERVATIONS WHOSE EMPIRICAL COVARIANCE HAS THE SMALLEST DETERMINANT YIELDING A “PURE” SUBSET OF OBSERVATIONS FROM WHICH TO COMPUTE STANDARDS ESTIMATES OF LOCATION AND COVARIANCE AFTER A CORRECTION STEP AIMING AT COMPENSATING THE FACT THAT THE ESTIMATES WERE LEARNED FROM ONLY A PORTION OF THE INITIAL DATA WE END UP WITH ROBUST ESTIMATES OF THE DATA SET LOCATION AND COVARIANCE

THE MINIMUM COVARIANCE DETERMINANT ESTIMATOR MCD HAS BEEN INTRODUCED BY PJROUSSEUW IN3

EVALUATION

IN THIS EXAMPLE WE COMPARE THE ESTIMATION ERRORS THAT ARE MADE WHEN USING VARIOUS TYPES OF LOCATION AND COVARIANCE ESTIMATES ON CONTAMINATED GAUSSIAN DISTRIBUTED DATA SETS

- THE MEAN AND THE EMPIRICAL COVARIANCE OF THE FULL DATASET WHICH BREAK DOWN AS SOON AS THERE ARE OUTLIERS IN THE DATA SET
- THE ROBUST MCD THAT HAS A LOW ERROR PROVIDED  $\sqrt{\text{SAMPLES}-5-\text{FEATURES}}$
- THE MEAN AND THE EMPIRICAL COVARIANCE OF THE OBSERVATIONS THAT ARE KNOWN TO BE GOOD ONES THIS CAN BE CONSIDERED AS A “PERFECT” MCD ESTIMATION SO ONE CAN TRUST OUR IMPLEMENTATION BY COMPARING TO THIS CASE

1JOHANNA HARDIN DAVID M ROCKE THE DISTRIBUTION OF ROBUST DISTANCES JOURNAL OF COMPUTATIONAL AND GRAPHICAL STATISTICS DECEMBER 144 928946

2ZOU BIR A KOIVUNEN V CHAKHCHOUKH Y AND MUMA M 2012 ROBUST ESTIMATION IN SIGNAL PROCESSING A TUTORIALSTYLE TREATMENT OF FU

MENTAL CONCEPTS IEEE SIGNAL PROCESSING MAGAZINE 294 6180

3P J ROUSSEUW LEAST MEDIAN OF SQUARES REGRESSION JOURNAL OF AMERICAN STATISTICAL ASS 79871 1984

954 CHAPTER 5 EXAMPLES



```
SCIKITLEARN USER GUIDE RELEASE 0213
REFERENCES
PRINTDOC
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
IMPORT MATPLOTLIBFONTMANAGER
FROM SKLEARNCOVARIANCE IMPORT EMPIRICALCOVARIANCE MINCOVDET
EXAMPLE SETTINGS
NSAMPLES 80
NFEATURES 5
REPEAT 10
RANGENOUTLIERS NPCONCATENATE
NPLINSPACE0 NSAMPLES 8 5
NPLINSPACENSAMPLES 8 NSAMPLES 2 511ASTYPENPINT
DEFINITION OF ARRAYS TO STORE RESULTS
ERRLOCMCD NPZEROSRANGENOUTLIERSSIZE REPEAT
ERRCOVMCD NPZEROSRANGENOUTLIERSSIZE REPEAT
ERRLOCEMPFULL NPZEROSRANGENOUTLIERSSIZE REPEAT
ERRCOVEMPFULL NPZEROSRANGENOUTLIERSSIZE REPEAT
ERRLOCEMPPURE NPZEROSRANGENOUTLIERSSIZE REPEAT
ERRCOVEMPPURE NPZEROSRANGENOUTLIERSSIZE REPEAT
58 COVARIANCE ESTIMATION 955
```

SCIKITLEARN USER GUIDE RELEASE 0213  
 COMPUTATION  
 FORI NOUTLIERS INENUMERATERANGENOUTLIERS  
 FORJINRANGEREPEAT  
 RNG NPRANDOMRANDOMSTATEI J  
 GENERATE DATA  
 X RNGRANDNNSAMPLES NFEATURES  
 ADD SOME OUTLIERS  
 OUTLIERSINDEX RNGPERMUTATIONNNSAMPLESNOUTLIERS  
 OUTLIERSOFFSET 10  
 NPRANDOMRANDINT2 SIZENOUTLIERS NFEATURES 05  
 XOUTLIERSINDEX OUTLIERSOFFSET  
 INLIERSMASK NPONESNSAMPLESASTYPEBOOL  
 INLIERSMASKOUTLIERSINDEX FALSE  
 FIT A MINIMUM COVARIANCE DETERMINANT MCD ROBUST ESTIMATOR TO DATA  
 MCD MINCOVDETFITX  
 COMPARE RAW ROBUST ESTIMATES WITH THE TRUE LOCATION AND COVARIANCE  
 ERRLOCMCDI J NPSUMMCDLOCATION 2  
 ERRCOVMCDI J MCDERRORNORMNPEYENFEATURES  
 COMPARE ESTIMATORS LEARNED FROM THE FULL DATA SET WITH TRUE  
 PARAMETERS  
 ERRLOCEMPFULLI J NPSUMXMEAN0 2  
 ERRCOVEMPFULLI J EMPIRICALCOVARIANCEFITXERRORNORM  
 NPEYENFEATURES  
 COMPARE WITH AN EMPIRICAL COVARIANCE LEARNED FROM A PURE DATA SET  
 IE PERFECT MCD  
 PUREX XINLIERSMASK  
 PURELOCATION PUREXMEAN0  
 PUREEMPCOV EMPIRICALCOVARIANCEFITPUREX  
 ERRLOCEMPPUREI J NPSUMPURELOCATION 2  
 ERRCOVEMPPUREI J PUREEMPCOVERERRORNORMNPEYENFEATURES  
 DISPLAY RESULTS  
 FONTPROP MATPLOTLIBFONTMANAGERFONTPROPERTIESSIZE11  
 PLTSUBPLOT2 1 1  
 LW 2  
 PLTERRORBARRANGENOUTLIERS ERRLOCMCDMEAN1  
 YERRERRLOCMCDSTD1 NPSQRTREPEAT  
 LABELROBUST LOCATION LWLW COLORM  
 PLTERRORBARRANGENOUTLIERS ERRLOCEMPFULLMEAN1  
 YERRERRLOCEMPFULLSTD1 NPSQRTREPEAT  
 LABELFULL DATA SET MEAN LWLW COLORGREEN  
 PLTERRORBARRANGENOUTLIERS ERRLOCEMPPUREMEAN1  
 YERRERRLOCEMPPURESTD1 NPSQRTREPEAT  
 LABELPURE DATA SET MEAN LWLW COLORBLACK  
 PLTTITLEINFLUENCE OF OUTLIERS ON THE LOCATION ESTIMATION  
 PLTYLABELRERROR MU HATMU22  
 PLTLEGENDLOCUPPER LEFT PROPFONTPROP  
 PLTSUBPLOT2 1 2  
 XSIZE RANGENOUTLIERSSIZE  
 PLTERRORBARRANGENOUTLIERS ERRCOVMCDMEAN1  
 YERRERRCOVMCDSTD1  
 956 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
LABELROBUST COVARIANCE MCD COLORM  
PLTERRORBARRANGENOUTLIERSXSIZE 5 1  
ERRCOVEMPFULLMEAN1XSIZE 5 1  
YERRERRCOVEMPFULLSTD1XSIZE 5 1  
LABELFULL DATA SET EMPIRICAL COVARIANCE COLORGREEN  
PLTPLOTTRANGENOUTLIERSXSIZE 5XSIZE 2 1  
ERRCOVEMPFULLMEAN1XSIZE 5XSIZE 2 1  
COLORGREEN LS  
PLTERRORBARRANGENOUTLIERS ERRCOVEMPPUREMEAN1  
YERRERRCOVEMPPURESTD1  
LABELPURE DATA SET EMPIRICAL COVARIANCE COLORBLACK  
PLTTITLEINFLUENCE OF OUTLIERS ON THE COVARIANCE ESTIMATION  
PLTXLABELAMOUNT OF CONTAMINATION  
PLTYLABELRMSE  
PLTLEGENDLOCUPPER CENTER PROPFONTPROP  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 2985 SECONDS  
59 CROSS DECOMPOSITION  
EXAMPLES CONCERNING THE SKLEARN CROSS DECOMPOSITION MODULE  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
591 COMPARE CROSS DECOMPOSITION METHODS  
SIMPLE USAGE OF VARIOUS CROSS DECOMPOSITION ALGORITHMS PLSCANONICAL PLSREGRESSION WITH MULTIVARIATE RE  
SPONSE AKA PLS2 PLSREGRESSION WITH UNIVARIATE RESPONSE AKA PLS1 CCA  
GIVEN 2 MULTIVARIATE COVARYING TWODIMENSIONAL DATASETS X AND Y PLS EXTRACTS THE 'DIRECTIONS OF COVARIANCE' IE  
THE COMPONENTS OF EACH DATASETS THAT EXPLAIN THE MOST SHARED VARIANCE BETWEEN BOTH DATASETS THIS IS APPARENT ON THE  
SCATTERPLOT MATRIX DISPLAY COMPONENTS 1 IN DATASET X AND DATASET Y ARE MAXIMALLY CORRELATED POINTS LIE AROUND THE  
FIRST DIAGONAL THIS IS ALSO TRUE FOR COMPONENTS 2 IN BOTH DATASET HOWEVER THE CORRELATION ACROSS DATASETS FOR DIFFERENT  
COMPONENTS IS WEAK THE POINT CLOUD IS VERY SPHERICAL  
59 CROSS DECOMPOSITION 957

SCIKITLEARN USER GUIDE RELEASE 0213

OUT

CORRX

1 051 007 005

051 1 011 001

007 011 1 049

005 001 049 1

CORRY

1 048 005 003

048 1 004 012

005 004 1 051

003 012 051 1

TRUE B SUCH THAT Y XB ERR

1 1 1

2 2 2

0 0 0

0 0 0

0 0 0

0 0 0

0 0 0

0 0 0

0 0 0

0 0 0

ESTIMATED B

1 1 1

2 2 2

0 0 0

0 0 0

0 0 0

958 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
0 0 0
0 0 01
0 0 0
0 0 01
0 0 0
ESTIMATED BETAS
1
21
0
0
0
0
0
0
0
0
0
PRINTDOC
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARNCROSSDECOMPOSITION IMPORT PLSCANONICAL PLSREGRESSION CCA

DATASET BASED LATENT VARIABLES MODEL
N 500
2 LATENTS VARS
L1 NPRANDOMNORMALSIZE4
L2 NPRANDOMNORMALSIZE4
LATENTS NPARRAYL1 L1 L2 L2T
X LATENTS NPRANDOMNORMALSIZE4 NRESHAPEN 4
Y LATENTS NPRANDOMNORMALSIZE4 NRESHAPEN 4
XTRAIN XN 2
YTRAIN YN 2
XTEST XN 2
YTEST YN 2
PRINTCORRX
PRINTNPROUNDNPCORRCOEFT 2
PRINTCORRY
PRINTNPROUNDNPCORRCOEFT 2

CANONICAL SYMMETRIC PLS
TRANSFORM DATA

PLSCA PLSCANONICALNCOMPONENTS2
PLSCAFITXTRAIN YTRAIN
XTRAINR YTRAINR PLSCATTRANSFORMXTRAIN YTRAIN
59 CROSS DECOMPOSITION 959
```

SCIKITLEARN USER GUIDE RELEASE 0213

XTESTR YTESTR PLSCATTRANSFORMXTEST YTEST  
SCATTER PLOT OF SCORES

1 ON DIAGONAL PLOT X VS Y SCORES ON EACH COMPONENTS

PLTFIGUREFIGSIZE12 8

PLTSUBPLOT221

PLTSCATTERXTRAINR 0 YTRAINR 0 LABELTRAIN  
MARKERO CB S25

PLTSCATTERXTESTR 0 YTESTR 0 LABELTEST  
MARKERO CR S25

PLTXLABELX SCORES

PLTYLABELY SCORES

PLTTITLECOMP 1 X VS Y TEST CORR 2F

NPCORRCOEFXTESTR 0 YTESTR 00 1

PLXTTICKS

PLTYTICKS

PLTLEGENDLOCBEST

PLTSUBPLOT224

PLTSCATTERXTRAINR 1 YTRAINR 1 LABELTRAIN  
MARKERO CB S25

PLTSCATTERXTESTR 1 YTESTR 1 LABELTEST  
MARKERO CR S25

PLTXLABELX SCORES

PLTYLABELY SCORES

PLTTITLECOMP 2 X VS Y TEST CORR 2F

NPCORRCOEFXTESTR 1 YTESTR 10 1

PLXTTICKS

PLTYTICKS

PLTLEGENDLOCBEST

2 OFF DIAGONAL PLOT COMPONENTS 1 VS 2 FOR X AND Y

PLTSUBPLOT222

PLTSCATTERXTRAINR 0 XTRAINR 1 LABELTRAIN  
MARKER CB S50

PLTSCATTERXTESTR 0 XTESTR 1 LABELTEST  
MARKER CR S50

PLTXLABELX COMP 1

PLTYLABELX COMP 2

PLTTITLEX COMP 1 VS X COMP 2 TEST CORR 2F

NPCORRCOEFXTESTR 0 XTESTR 10 1

PLTLEGENDLOCBEST

PLXTTICKS

PLTYTICKS

PLTSUBPLOT223

PLTSCATTERYTRAINR 0 YTRAINR 1 LABELTRAIN  
MARKER CB S50

PLTSCATTERYTESTR 0 YTESTR 1 LABELTEST  
MARKER CR S50

PLTXLABELY COMP 1

PLTYLABELY COMP 2

PLTTITLEY COMP 1 VS Y COMP 2 TEST CORR 2F

NPCORRCOEFYTESTR 0 YTESTR 10 1

PLTLEGENDLOCBEST

PLXTTICKS

PLTYTICKS

960 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
PLTSHOW

PLS REGRESSION WITH MULTIVARIATE RESPONSE AKA PLS2

```
N 1000
Q 3
P 10
X NPRANDOMNORMALSIZEN PRESHAPEN P
B NPARRAY1 2 0 P 2 QT
EACH YJ 1 X1 2X2 NOIZE
Y NPDOTX B NPRANDOMNORMALSIZEN QRESHAPEN Q 5
PLS2 PLSREGRESSIONNCOMPONENTS3
PLS2FITX Y
PRINTTRUE B SUCH THAT Y XB ERR
PRINTB
COMPARE PLS2COEF WITH B
PRINTESTIMATED B
PRINTNPROUNDPLS2COEF 1
PLS2PREDICTX
```

PLS REGRESSION WITH UNIVARIATE RESPONSE AKA PLS1

```
N 1000
P 10
X NPRANDOMNORMALSIZEN PRESHAPEN P
Y X 0 2 X 1 NPRANDOMNORMALSIZEN 1 5
PLS1 PLSREGRESSIONNCOMPONENTS3
PLS1FITX Y
```

NOTE THAT THE NUMBER OF COMPONENTS EXCEEDS 1 THE DIMENSION OF Y  
PRINTESTIMATED BETAS  
PRINTNPROUNDPLS1COEF 1

CCA PLS MODE B WITH SYMMETRIC DEFLATION

```
CCA CCANCOMPONENTS2
CCAFITXTRAIN YTRAIN
XTRAINR YTRAINR CCATransformXTRAIN YTRAIN
XTESTR YTESTR CCATransformXTEST YTEST
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0099 SECONDS
510 DATASET EXAMPLES
EXAMPLES CONCERNING THE SKLEARNDATASETS MODULE
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE
510 DATASET EXAMPLES 961
```

SCIKITLEARN USER GUIDE RELEASE 0213

5101 THE DIGIT DATASET

THIS DATASET IS MADE UP OF 1797 8X8 IMAGES EACH IMAGE LIKE THE ONE SHOWN BELOW IS OF A HANDWRITTEN DIGIT IN ORDER TO UTILIZE AN 8X8 FIGURE LIKE THIS WE'D HAVE TO FIRST TRANSFORM IT INTO A FEATURE VECTOR WITH LENGTH 64

SEE [HERE](#) FOR MORE INFORMATION ABOUT THIS DATASET

PRINTDOC

```
CODE SOURCE GAËL VAROQUAUX
MODIFIED FOR DOCUMENTATION BY JAQUES GROBLER
LICENSE BSD 3 CLAUSE
FROM SKLEARN IMPORT DATASETS
IMPORT MATPLOTLIBPYPLOT AS PLT
LOAD THE DIGITS DATASET
DIGITS = DATASETSLOADDIGITS
DISPLAY THE FIRST DIGIT
PLTFigure(1,figsize=(3,3))
PLT.imshow(DIGITS[0],cmap=PLT.cm.gray,rasterized=True)
PLT.show()
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0113 SECONDS
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE
```

5102 THE IRIS DATASET

THIS DATA SETS CONSISTS OF 3 DIFFERENT TYPES OF IRISES' SETOSA VERSICOLOUR AND VIRGINICA PETAL AND SEPAL LENGTH STORED IN A 150X4 NUMPYNDARRAY

THE ROWS BEING THE SAMPLES AND THE COLUMNS BEING SEPAL LENGTH SEPAL WIDTH PETAL LENGTH AND PETAL WIDTH

962 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213  
THE BELOW PLOT USES THE FIRST TWO FEATURES SEE [HERE](#) FOR MORE INFORMATION ON THIS DATASET

- 

510 DATASET EXAMPLES 963

SCIKITLEARN USER GUIDE RELEASE 0213

•

```
PRINTDOC
CODE SOURCE GAËL VAROQUAUX
MODIFIED FOR DOCUMENTATION BY JAQUES GROBLER
LICENSE BSD 3 CLAUSE
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM MPLTOOLKITSMPLPLOT3D IMPORT AXES3D
FROM SKLEARN IMPORT DATASETS
FROM SKLEARNDECOMPOSITION IMPORT PCA
IMPORT SOME DATA TO PLAY WITH
IRIS = DATASETSLOADIRIS
X = IRISDATA[2] WE ONLY TAKE THE FIRST TWO FEATURES
Y = IRISTARGET
XMIN XMAX X 0MIN 5 X 0MAX 5
YMIN YMAX X 1MIN 5 X 1MAX 5
PLTFigure2 figsize(8, 6)
PLTCLF
PLOT THE TRAINING POINTS
PLTSCATTER(X[0], X[1], C=Y, CMAP=PLTCMSET1,
EDGE=COLORMAP)
PLTXLABELSEPAL LENGTH
964 CHAPTER 5 EXAMPLES
```

```
SCIKITLEARN USER GUIDE RELEASE 0213
PLTYLABELSEPAL WIDTH
PLTXLIMXMIN XMAX
PLTYLIMYMIN YMAX
PLXTTICKS
PLTYTICKS
    TO GETTER A BETTER UNDERSTANDING OF INTERACTION OF THE DIMENSIONS
    PLOT THE FIRST THREE PCA DIMENSIONS
FIG PLTFigure1 FIGSIZE8 6
AX AXES3DFIG ELEV150 AZIM110
XREDUCED PCANCOMPONENTS3FITTRANSFORMIRISDATA
AXSCATTERXREDUCED 0 XREDUCED 1 XREDUCED 2 CY
CMAPPLTCMSET1 EDGECOLORK S40
AXSETTITLEFIRST THREE PCA DIRECTIONS
AXSETXLABEL1ST EIGENVECTOR
AXWXAXISSETTICKLABELS
AXSETYLABEL2ND EIGENVECTOR
AXWYAXISSETTICKLABELS
AXSETZLABEL3RD EIGENVECTOR
AXWZAXISSETTICKLABELS
PLTSHOW
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0059 SECONDS
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE
5103 PLOT RANDOMLY GENERATED CLASSIFICATION DATASET
PLOT SEVERAL RANDOMLY GENERATED 2D CLASSIFICATION DATASETS THIS EXAMPLE ILLUSTRATES THE DATASETS
MAKECLASSIFICATION DATASETSMAKEBLOBS ANDDATASETSMAKEGAUSSIANQUANTILES FUNC
TIONS
FORMAKECLASSIFICATION THREE BINARY AND TWO MULTICLASS CLASSIFICATION DATASETS ARE GENERATED WITH DIFFERENT
NUMBERS OF INFORMATIVE FEATURES AND CLUSTERS PER CLASS
510 DATASET EXAMPLES 965
```

SCIKITLEARN USER GUIDE RELEASE 0213  
PRINTDOC  
IMPORT MATPLOTLIBPYPLOT AS PLT  
FROM SKLEARNDATASETS IMPORT MAKECLASSIFICATION  
FROM SKLEARNDATASETS IMPORT MAKEBLOBS  
FROM SKLEARNDATASETS IMPORT MAKEGAUSSIANQUANTILES  
PLTFIGUREFIGSIZE8 8  
PLTSUBPLOTSADJUSTBOTTOM05 TOP9 LEFT05 RIGHT95  
PLTSUBPLOT321  
PLTTITLEONE INFORMATIVE FEATURE ONE CLUSTER PER CLASS FONTSIZESMALL  
X1 Y1 MAKECLASSIFICATIONNFEATURES2 NREDUNDANT0 NINFORMATIVE1  
NCLUSTERSPERCLASS1  
966 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
PLTSCATTERX1 0 X1 1 MARKERO CY1  
S25 EDGECOLORK  
PLTSUBPLOT322  
PLTTITLETWO INFORMATIVE FEATURES ONE CLUSTER PER CLASS FONTSIZESMALL  
X1 Y1 MAKECLASSIFICATIONNFEATURES2 NREDUNDANT0 NINFORMATIVE2  
NCLUSTERSPERCLASS1  
PLTSCATTERX1 0 X1 1 MARKERO CY1  
S25 EDGECOLORK  
PLTSUBPLOT323  
PLTTITLETWO INFORMATIVE FEATURES TWO CLUSTERS PER CLASS  
FONTSIZESMALL  
X2 Y2 MAKECLASSIFICATIONNFEATURES2 NREDUNDANT0 NINFORMATIVE2  
PLTSCATTERX2 0 X2 1 MARKERO CY2  
S25 EDGECOLORK  
PLTSUBPLOT324  
PLTTITLEMULTICLASS TWO INFORMATIVE FEATURES ONE CLUSTER  
FONTSIZESMALL  
X1 Y1 MAKECLASSIFICATIONNFEATURES2 NREDUNDANT0 NINFORMATIVE2  
NCLUSTERSPERCLASS1 NCLASSES3  
PLTSCATTERX1 0 X1 1 MARKERO CY1  
S25 EDGECOLORK  
PLTSUBPLOT325  
PLTTITLETHREE BLOBS FONTSIZESMALL  
X1 Y1 MAKEBLOBSNFEATURES2 CENTERS3  
PLTSCATTERX1 0 X1 1 MARKERO CY1  
S25 EDGECOLORK  
PLTSUBPLOT326  
PLTTITLEGAUSSIAN DIVIDED INTO THREE QUANTILES FONTSIZESMALL  
X1 Y1 MAKEGAUSSIANQUANTILESNFEATURES2 NCLASSES3  
PLTSCATTERX1 0 X1 1 MARKERO CY1  
S25 EDGECOLORK  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0097 SECONDS  
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE  
5104 PLOT RANDOMLY GENERATED MULTILABEL DATASET  
THIS ILLUSTRATES THE DATASETSMAKEMULTILABELCLASSIFICATION DATASET GENERATOR EACH SAMPLE CONSISTS  
OF COUNTS OF TWO FEATURES UP TO 50 IN TOTAL WHICH ARE DIFFERENTLY DISTRIBUTED IN EACH OF TWO CLASSES  
POINTS ARE LABELED AS FOLLOWS WHERE Y MEANS THE CLASS IS PRESENT  
510 DATASET EXAMPLES 967

SCIKITLEARN USER GUIDE RELEASE 0213

123COLOR  
Y N N RED  
N Y N BLUE  
N N Y YELLOW  
Y Y N PURPLE  
Y N Y ORANGE  
Y Y N GREEN  
Y Y Y BROWN

A STAR MARKS THE EXPECTED SAMPLE FOR EACH CLASS ITS SIZE REFLECTS THE PROBABILITY OF SELECTING THAT CLASS LABEL  
THE LEFT AND RIGHT EXAMPLES HIGHLIGHT THE NLABELS PARAMETER MORE OF THE SAMPLES IN THE RIGHT PLOT HAVE 2 OR 3  
LABELS

NOTE THAT THIS TWODIMENSIONAL EXAMPLE IS VERY DEGENERATE GENERALLY THE NUMBER OF FEATURES WOULD BE MUCH GREATER  
THAN THE “DOCUMENT LENGTH” WHILE HERE WE HAVE MUCH LARGER DOCUMENTS THAN VOCABULARY SIMILARLY WITH NCLASSES  
NFEATURES IT IS MUCH LESS LIKELY THAT A FEATURE DISTINGUISHES A PARTICULAR CLASS  
OUT

THE DATA WAS GENERATED FROM RANDOMSTATE1013  
CLASS PC PW0C PW1C  
RED 064 097 003  
BLUE 006 060 040  
YELLOW 030 009 091  
IMPORT NUMPY AS NP  
IMPORT MATPLOTLIBPYPLOT AS PLT  
FROM SKLEARNDATASETS IMPORT MAKEMULTILABELCLASSIFICATION ASMAKEMLCLF  
968 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

PRINTDOC  
COLORS NPARRAY  
FF3333 RED  
0198E1 BLUE  
BF5FFF PURPLE  
FCD116 YELLOW  
FF7216 ORANGE  
4DBD33 GREEN  
87421F BROWN

USE SAME RANDOM SEED FOR MULTIPLE CALLS TO MAKEMULTILABELCLASSIFICATION TO  
ENSURE SAME DISTRIBUTIONS  
RANDOMSEED NPRANDOMRANDINT2 10  
DEFPLOT2DAX NLABELS1 NCLASSES3 LENGTH50  
X Y PC PWC MAKEMLCLFNSAMPLES150 NFEATURES2  
NCLASSESNCLASSES NLABELSNLABELS  
LENGTHLENGTH ALLOWUNLABELEDFALSE  
RETURNDISTRIBUTIONSTRUE  
RANDOMSTATERANDOMSEED  
AXSCATTERX 0 X 1 COLORCOLORSTAKE1 2 4  
SUMAXIS1  
MARKER  
AXSCATTERPWC0 LENGTH PWC1 LENGTH  
MARKER LINEWIDTH5 EDGECOLORBLACK  
S20 1500 PC2  
COLORCOLORSTAKE1 2 4  
AXSETXLABELFEATURE 0 COUNT  
RETURNPC PWC  
AX1 AX2 PLTSUBPLOTS1 2 SHAREXROW SHAREYROW FIGSIZE8 4  
PLTSUBPLOTSADJUSTBOTTOM15  
PC PWC PLOT2DAX1 NLABELS1  
AX1SETTITLENLABELS1 LENGTH50  
AX1SETYLABELFEATURE 1 COUNT  
PLOT2DAX2 NLABELS3  
AX2SETTITLENLABELS3 LENGTH50  
AX2SETXLIMLEFT0 AUTOTRUE  
AX2SETYLIMBOTTOM0 AUTOTRUE  
PLTSHOW  
PRINTTHE DATA WAS GENERATED FROM RANDOMSTATE D RANDOMSEED  
PRINTCLASS PC PW0C PW1C SEP T  
FORK P PW INZIPRED BLUE YELLOW PC PWCT  
PRINTST02FT02FT02F K P PW0 PW1  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0085 SECONDS  
510 DATASET EXAMPLES 969

SCIKITLEARN USER GUIDE RELEASE 0213

511 DECOMPOSITION

EXAMPLES CONCERNING THE SKLEARNDECOMPOSITION MODULE

NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)

5111 BETADIVERGENCE LOSS FUNCTIONS

A PLOT THAT COMPARES THE VARIOUS BETADIVERGENCE LOSS FUNCTIONS SUPPORTED BY THE MULTIPLICATIVEUPDATE ‘MU’ SOLVER

INSKLEARNDECOMPOSITIONNMF

IMPORT NUMPY AS NP

IMPORT MATPLOTLIBPYPLOT AS PLT

FROM SKLEARNDECOMPOSITIONNMF IMPORT BETADIVERGENCE

PRINTDOC

X NPLinspace0001 4 1000

Y NPZEROSXSHAPE

COLORS MBGYR

FORJ BETA INENUMERATE0 05 1 15 2

970 CHAPTER 5 EXAMPLES



```
SCIKITLEARN USER GUIDE RELEASE 0213
FORI XIINENUMERATEX
YI BETADIVERGENCE1 XI 1 BETA
NAME BETA 11F BETA
PLTPLOTX Y LABELNAME COLORCOLORSJ
PLTXLABELX
PLTTITLEBETADIVERGENCE1 X
PLTLEGENDLOC0
PLTAXISO 4 0 3
PLTSHOW
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0225 SECONDS
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE
5112 PCA EXAMPLE WITH IRIS DATASET
PRINCIPAL COMPONENT ANALYSIS APPLIED TO THE IRIS DATASET
SEE HERE FOR MORE INFORMATION ON THIS DATASET
PRINTDOC
CODE SOURCE GAËL VAROQUAUX
LICENSE BSD 3 CLAUSE
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM MPLTOOLKITSMPLLOT3D IMPORT AXES3D
FROM SKLEARN IMPORT DECOMPOSITION
FROM SKLEARN IMPORT DATASETS
511 DECOMPOSITION 971
```

SCIKITLEARN USER GUIDE RELEASE 0213  
NPRANDOMSEED5  
CENTERS 1 1 1 1 1 1  
IRIS DATASETSLOADIRIS  
X IRISDATA  
Y IRISTARGET  
FIG PLTFigure1 FIGSize4 3  
PLTCLF  
AX AXES3DFIG RECT0 0 95 1 ELEV48 AZIM134  
PLTCLA  
PCA DECOMPOSITIONPCANCOMPONENTS3  
PCAFITX  
X PCATransformX  
FORNAME LABEL INSETOSA 0 VERSICOLOUR 1 VIRGINICA 2  
AXTEXT3DXY LABEL 0MEAN  
XY LABEL 1MEAN 15  
XY LABEL 2MEAN NAME  
HORIZONTALALIGNMENTCENTER  
BBOXDICTALPHA5 EDGEColorW FACEColorW  
REORDER THE LABELS TO HAVE COLORS MATCHING THE CLUSTER RESULTS  
Y NPCHOOSEY 1 2 0ASTYPENPFloat  
AXSCATTERX 0 X 1 X 2 CY CMAPPLTCMNIPYSPECTRAL  
EDGEColorK  
AXWXAXISSETTICKLABELS  
AXWYAXISSETTICKLABELS  
AXWZAXISSETTICKLABELS  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0186 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
5113 INCREMENTAL PCA  
INCREMENTAL PRINCIPAL COMPONENT ANALYSIS IPCA IS TYPICALLY USED AS A REPLACEMENT FOR PRINCIPAL COMPONENT ANALYSIS  
PCA WHEN THE DATASET TO BE DECOMPOSED IS TOO LARGE TO FIT IN MEMORY IPCA BUILDS A LOWRANK APPROXIMATION FOR THE  
INPUT DATA USING AN AMOUNT OF MEMORY WHICH IS INDEPENDENT OF THE NUMBER OF INPUT DATA SAMPLES IT IS STILL DEPENDENT  
ON THE INPUT DATA FEATURES BUT CHANGING THE BATCH SIZE ALLOWS FOR CONTROL OF MEMORY USAGE  
THIS EXAMPLE SERVES AS A VISUAL CHECK THAT IPCA IS ABLE TO FIND A SIMILAR PROJECTION OF THE DATA TO PCA TO A SIGN FLIP  
WHILE ONLY PROCESSING A FEW SAMPLES AT A TIME THIS CAN BE CONSIDERED A “TOY EXAMPLE” AS IPCA IS INTENDED FOR LARGE  
DATASETS WHICH DO NOT FIT IN MAIN MEMORY REQUIRING INCREMENTAL APPROACHES  
972 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

- 511 DECOMPOSITION 973

SCIKITLEARN USER GUIDE RELEASE 0213

•

```
PRINTDOC
AUTHORS KYLE KASTNER
LICENSE BSD 3 CLAUSE
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARNDATASETS IMPORT LOADIRIS
FROM SKLEARNDECOMPOSITION IMPORT PCA INCREMENTALPCA
IRIS LOADIRIS
X IRISDATA
Y IRISTARGET
NCOMPONENTS 2
IPCA INCREMENTALPCANCOMPONENTSNCOMPONENTS BATCHSIZE10
974 CHAPTER 5 EXAMPLES
```

SCIKITLEARN USER GUIDE RELEASE 0213

XIPCA IPCAFITTRANSFORMX

PCA PCANCOMPONENTSNCOMPONENTS

XPCA PCAFITTRANSFORMX

COLORS NAVY TURQUOISE DARKORANGE

FORXTRANSFORMED TITLE INXIPCA INCREMENTAL PCA XPCA PCA

PLTFIGUREFIGSIZE8 8

FORCOLOR I TARGETNAME INZIPCOLORS 0 1 2 IRISTARGETNAMES

PLTSCATTERXTRANSFORMEDY I 0 XTRANSFORMEDY I 1

COLORCOLOR LW2 LABELTARGETNAME

IFINCREMENTAL INTITLE

ERR NPABSNPABSXPCA NPABSXIPCAMEAN

PLTTITLETITLE OF IRIS DATASET NMEAN ABSOLUTE UNSIGNED ERROR

6F ERR

ELSE

PLTTITLETITLE OF IRIS DATASET

PLTLEGENDLOCBEST SHADOWFALSE SCATTERPOINTS1

PLTAXIS4 4 15 15

PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0143 SECONDS

NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)

5114 COMPARISON OF LDA AND PCA 2D PROJECTION OF IRIS DATASET

THE IRIS DATASET REPRESENTS 3 KIND OF IRIS FLOWERS SETOSA VERSICOLOUR AND VIRGINICA WITH 4 ATTRIBUTES SEPAL LENGTH

SEPAL WIDTH PETAL LENGTH AND PETAL WIDTH

PRINCIPAL COMPONENT ANALYSIS PCA APPLIED TO THIS DATA IDENTIFIES THE COMBINATION OF ATTRIBUTES PRINCIPAL COMPO

NENTS OR DIRECTIONS IN THE FEATURE SPACE THAT ACCOUNT FOR THE MOST VARIANCE IN THE DATA HERE WE PLOT THE DIFFERENT

SAMPLES ON THE 2 FIRST PRINCIPAL COMPONENTS

LINEAR DISCRIMINANT ANALYSIS LDA TRIES TO IDENTIFY ATTRIBUTES THAT ACCOUNT FOR THE MOST VARIANCE BETWEEN CLASSES IN

PARTICULAR LDA IN CONTRAST TO PCA IS A SUPERVISED METHOD USING KNOWN CLASS LABELS

511 DECOMPOSITION 975



SCIKITLEARN USER GUIDE RELEASE 0213

```
•
OUT
EXPLAINED VARIANCE RATIO FIRST TWO COMPONENTS 092461872 005306648
PRINTDOC
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARN IMPORT DATASETS
FROM SKLEARNDECOMPOSITION IMPORT PCA
FROM SKLEARNDISCRIMINANTANALYSIS IMPORT LINEARDISCRIMINANTANALYSIS
IRIS DATASETSLOADIRIS
X IRISDATA
Y IRISTARGET
TARGETNAMES IRISTARGETNAMES
PCA PCANCOMPONENTS2
XR PCAFITXTRANSFORMX
LDA LINEARDISCRIMINANTANALYSISNCOMPONENTS2
511 DECOMPOSITION 977
```

SCIKITLEARN USER GUIDE RELEASE 0213

XR2 LDAFITX YTRANSFORMX

PERCENTAGE OF VARIANCE EXPLAINED FOR EACH COMPONENTS

PRINTEXPLAINED VARIANCE RATIO FIRST TWO COMPONENTS S

STRPCAEXPLAINEDVARIANCERATIO

PLTFigure

COLORS NAVY TURQUOISE DARKORANGE

LW 2

FORCOLOR I TARGETNAME INZIPCOLORS 0 1 2 TARGETNAMES

PLTSCATTERXRY I 0 XRY I 1 COLORCOLOR ALPHA8 LWLW

LABELTARGETNAME

PLTLEGENDLOCBEST SHADOWFALSE SCATTERPOINTS1

PLTTITLEPCA OF IRIS DATASET

PLTFigure

FORCOLOR I TARGETNAME INZIPCOLORS 0 1 2 TARGETNAMES

PLTSCATTERXR2Y I 0 XR2Y I 1 ALPHA8 COLORCOLOR

LABELTARGETNAME

PLTLEGENDLOCBEST SHADOWFALSE SCATTERPOINTS1

PLTTITLELDA OF IRIS DATASET

PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0057 SECONDS

NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE

5115 BLIND SOURCE SEPARATION USING FASTICA

AN EXAMPLE OF ESTIMATING SOURCES FROM NOISY DATA

INDEPENDENT COMPONENT ANALYSIS ICA IS USED TO ESTIMATE SOURCES GIVEN NOISY MEASUREMENTS IMAGINE 3 INSTRUMENTS PLAYING SIMULTANEOUSLY AND 3 MICROPHONES RECORDING THE MIXED SIGNALS ICA IS USED TO RECOVER THE SOURCES IE WHAT IS PLAYED BY EACH INSTRUMENT IMPORTANTLY PCA FAILS AT RECOVERING OUR INSTRUMENTS SINCE THE RELATED SIGNALS REFLECT NONGAUSSIAN PROCESSES

978 CHAPTER 5 EXAMPLES



```
SCIKITLEARN USER GUIDE RELEASE 0213
PRINTDOC
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SCIPY IMPORT SIGNAL
FROM SKLEARNDECOMPOSITION IMPORT FASTICA PCA
```

```
GENERATE SAMPLE DATA
NPRANDOMSEED0
NSAMPLES 2000
TIME NPLinspace0 8 NSAMPLES
S1 NPSIN2 TIME SIGNAL 1 SINUSOIDAL SIGNAL
S2 NPSIGNNPSIN3 TIME SIGNAL 2 SQUARE SIGNAL
S3 SIGNALSAWTOOTH2 NPPITIME SIGNAL 3 SAW TOOTH SIGNAL
S NPCS1 S2 S3
S 02 NPRANDOMNORMALSIZESSHAPE ADD NOISE
S SSTDAXIS0 STANDARDIZE DATA
MIX DATA
A NPARRAY1 1 1 05 2 10 15 10 20 MIXING MATRIX
X NPDOTS AT GENERATE OBSERVATIONS
511 DECOMPOSITION 979
```

SCIKITLEARN USER GUIDE RELEASE 0213  
COMPUTE ICA  
ICA FASTICANCOMPONENTS3  
S ICAFITTRANSFORMX RECONSTRUCT SIGNALS  
A ICAMIXING GET ESTIMATED MIXING MATRIX  
WE CAN PROVE THAT THE ICA MODEL APPLIES BY REVERTING THE UNMIXING  
ASSERTNPALLCLOSEX NPDOTS AT ICAMEAN  
FOR COMPARISON COMPUTE PCA  
PCA PCANCOMPONENTS3  
H PCAFITTRANSFORMX RECONSTRUCT SIGNALS BASED ON ORTHOGONAL COMPONENTS  
  
PLOT RESULTS  
PLTFigure  
MODELS X S S H  
NAMES OBSERVATIONS MIXED SIGNAL  
TRUE SOURCES  
ICA RECOVERED SIGNALS  
PCA RECOVERED SIGNALS  
COLORS RED STEELBLUE ORANGE  
FORII MODEL NAME INENUMERATEZIPMODELS NAMES 1  
PLTSUBPLOT4 1 II  
PLTTITLENAME  
FORSIG COLOR INZIPMODELT COLORS  
PLTPLOTSIG COLORCOLOR  
PLTSUBPLOTSADJUST009 004 094 094 026 046  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0089 SECONDS  
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE  
5116 PRINCIPAL COMPONENTS ANALYSIS PCA  
THESE FIGURES AID IN ILLUSTRATING HOW A POINT CLOUD CAN BE VERY FLAT IN ONE DIRECTION-WHICH IS WHERE PCA COMES IN TO  
CHOOSE A DIRECTION THAT IS NOT FLAT  
980 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

- 
- 

PRINTDOC  
AUTHORS GAELE VAROQUAUX  
JAQUES GROBLER  
KEVIN HUGHES  
LICENSE BSD 3 CLAUSE  
FROM SKLEARNDECOMPOSITION IMPORT PCA  
FROM MPLTOOLKITSMPLPLOT3D IMPORT AXES3D  
IMPORT NUMPY AS NP  
IMPORT MATPLOTLIBPYPLOT AS PLT  
FROM SCIPY IMPORT STATS

CREATE THE DATA  
511 DECOMPOSITION 981

SCIKITLEARN USER GUIDE RELEASE 0213  
E NPEXP1  
NPRANDOMSEED4  
DEFPDFX  
RETURN05STATSNORMSCALE025 EPDFX  
STATSNORMSCALE4 EPDFX  
Y NPRANDOMNORMALSCALE05 SIZE30000  
X NPRANDOMNORMALSCALE05 SIZE30000  
Z NPRANDOMNORMALSCALE01 SIZELENX  
DENSITY PDFX PDFY  
PDFZ PDF5 Z  
DENSITY PDFZ  
A X Y  
B 2Y  
C A B Z  
NORM NPSQRTAVAR BVAR  
A NORM  
B NORM  
  
PLOT THE FIGURES  
DEFPLOTFIGSFIGNUM ELEV AZIM  
FIG PLTFIGUREFIGNUM FIGSIZE4 3  
PLTCLF  
AX AXES3DFIG RECT0 0 95 1 ELEVELEV AZIMAZIM  
AXSCATTERA10 B10 C10 CDENSITY10 MARKER ALPHA4  
Y NPCA B C  
USING SCIPYS SVD THIS WOULD BE  
PCAScore V SCIPYLINALGSVDY FULLMATRICESFALSE  
PCA PCANCOMPONENTS3  
PCAFITY  
PCAScore PCAEXPLAINEDVARIANCERATIO  
V PCACOMPONENTS  
XPCAAXIS YPCAAXIS ZPCAAXIS 3 VT  
XPCAPLANE NPRXPCAAXIS2 XPCAAXIS11  
YPCAPLANE NPRYPCAAXIS2 YPCAAXIS11  
ZPCAPLANE NPRZPCAAXIS2 ZPCAAXIS11  
XPCAPLANESHAPE 2 2  
YPCAPLANESHAPE 2 2  
ZPCAPLANESHAPE 2 2  
AXPLOTSURFACEXPCAPLANE YPCAPLANE ZPCAPLANE  
AXWXAXISSETTICKLABELS  
AXWYAXISSETTICKLABELS  
AXWZAXISSETTICKLABELS  
ELEV 40  
982 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

AZIM 80  
PLOTFIGS1 ELEV AZIM  
ELEV 30  
AZIM 20  
PLOTFIGS2 ELEV AZIM  
PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0114 SECONDS

NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)

5117 FASTICA ON 2D POINT CLOUDS

THIS EXAMPLE ILLUSTRATES VISUALLY IN THE FEATURE SPACE A COMPARISON BY RESULTS USING TWO DIFFERENT COMPONENT ANALYSIS TECHNIQUES

INDEPENDENT COMPONENT ANALYSIS ICA VS PRINCIPAL COMPONENT ANALYSIS PCA

REPRESENTING ICA IN THE FEATURE SPACE GIVES THE VIEW OF 'GEOMETRIC ICA' ICA IS AN ALGORITHM THAT FINDS DIRECTIONS IN THE FEATURE SPACE CORRESPONDING TO PROJECTIONS WITH HIGH NONGAUSSIANITY THESE DIRECTIONS NEED NOT BE ORTHOGONAL IN THE ORIGINAL FEATURE SPACE BUT THEY ARE ORTHOGONAL IN THE WHITENED FEATURE SPACE IN WHICH ALL DIRECTIONS CORRESPOND TO THE SAME VARIANCE

PCA ON THE OTHER HAND FINDS ORTHOGONAL DIRECTIONS IN THE RAW FEATURE SPACE THAT CORRESPOND TO DIRECTIONS ACCOUNTING FOR MAXIMUM VARIANCE

HERE WE SIMULATE INDEPENDENT SOURCES USING A HIGHLY NONGAUSSIAN PROCESS 2 STUDENT T WITH A LOW NUMBER OF DEGREES

OF FREEDOM TOP LEFT FIGURE WE MIX THEM TO CREATE OBSERVATIONS TOP RIGHT FIGURE IN THIS RAW OBSERVATION SPACE

DIRECTIONS IDENTIFIED BY PCA ARE REPRESENTED BY ORANGE VECTORS WE REPRESENT THE SIGNAL IN THE PCA SPACE AFTER

WHITENING BY THE VARIANCE CORRESPONDING TO THE PCA VECTORS LOWER LEFT RUNNING ICA CORRESPONDS TO FINDING A

ROTATION IN THIS SPACE TO IDENTIFY THE DIRECTIONS OF LARGEST NONGAUSSIANITY LOWER RIGHT

511 DECOMPOSITION 983

```
SCIKITLEARN USER GUIDE RELEASE 0213
PRINTDOC
AUTHORS ALEXANDRE GRAMFORT GAELE VAROQUAUX
LICENSE BSD 3 CLAUSE
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARNDECOMPOSITION IMPORT PCA FASTICA

GENERATE SAMPLE DATA
RNG NPRANDOMRANDOMSTATE42
S RNGSTANDARDT15 SIZE20000 2
S 0 2
MIX DATA
A NPARRAY1 1 0 2 MIXING MATRIX
X NPDOTS AT GENERATE OBSERVATIONS
PCA PCA
SPCA PCAFITXTRANSFORMX
ICA FASTICARANDOMSTATERNG
SICA ICAFITXTRANSFORMX ESTIMATE THE SOURCES
984 CHAPTER 5 EXAMPLES
```

SCIKITLEARN USER GUIDE RELEASE 0213  
SICA SICASTDAXISO

PLOT RESULTS  
DEFPLOTSAMPLESS AXISLISTNONE  
PLTSCATTERS 0 S 1 S2 MARKERO ZORDER10  
COLORSTEELBLUE ALPHA05  
IFAXISLIST IS NOTNONE  
COLORS ORANGE RED  
FORCOLOR AXIS INZIPCOLORS AXISLIST  
AXIS AXISSTD  
XAXIS YAXIS AXIS  
TRICK TO GET LEGEND TO WORK  
PLTPLOT01 XAXIS 01 YAXIS LINEWIDTH2 COLORCOLOR  
PLTQUIVER0 0 XAXIS YAXIS ZORDER11 WIDTH001 SCALE6  
COLORCOLOR  
PLTHLINES0 3 3  
PLTVLINES0 3 3  
PLTXLIM3 3  
PLTYLIM3 3  
PLTXLABELX  
PLTYLABELY  
PLTFigure  
PLTSUBPLOT2 2 1  
PLOTSAMPLESS SSTD  
PLTTITLETRUE INDEPENDENT SOURCES  
AXISLIST PCACOMPONENTST ICAMIXING  
PLTSUBPLOT2 2 2  
PLOTSAMPLESX NPSTD X AXISLISTAXISLIST  
LEGEND PLTLEGENDPCA ICA LOCUPPER RIGHT  
LEGENDSETZORDER100  
PLTTITLEOBSERVATIONS  
PLTSUBPLOT2 2 3  
PLOTSAMPLESSPCA NPSTDSPCA AXIS0  
PLTTITLEPCA RECOVERED SIGNALS  
PLTSUBPLOT2 2 4  
PLOTSAMPLESSICA NPSTD SICA  
PLTTITLEICA RECOVERED SIGNALS  
PLTSUBPLOTSADJUST009 004 094 094 026 036  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0357 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
511 DECOMPOSITION 985

```
SCIKITLEARN USER GUIDE RELEASE 0213
5118 KERNEL PCA
THIS EXAMPLE SHOWS THAT KERNEL PCA IS ABLE TO FIND A PROJECTION OF THE DATA THAT MAKES DATA LINEARLY SEPARABLE
PRINTDOC
AUTHORS MATHIEU BLONDEL
ANDREAS MUELLER
LICENSE BSD 3 CLAUSE
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARNDECOMPOSITION IMPORT PCA KERNELPCA
FROM SKLEARNDATASETS IMPORT MAKECIRCLES
NPRANDOMSEED0
X Y MAKECIRCLESNSAMPLES400 FACTOR3 NOISE05
KPCA KERNELPCAKERNELRBF FITINVERSETRANSFORMTRUE GAMMA10
XKPCA KPCAFITTRANSFORMX
XBACK KPCAINVERSETRANSFORMXKPCA
PCA PCA
XPCA PCAFITTRANSFORMX
986 CHAPTER 5 EXAMPLES
```



SCIKITLEARN USER GUIDE RELEASE 0213

PLOT RESULTS

PLTFigure

PLTSUBPLOT2 2 1 ASPECTEQUAL

PLTTITLEORIGINAL SPACE

REDS Y 0

BLUES Y 1

PLTSCATTERXREDS 0 XREDS 1 CRED

S20 EDGECOLORK

PLTSCATTERXBLUES 0 XBLUES 1 CBLUE

S20 EDGECOLORK

PLTXLABELX1

PLTYLABELX2

X1 X2 NPMESHGRIDNPLINSPACE15 15 50 NPLINSPACE15 15 50

XGRID NPARRAYNPRAVELX1 NPRAVELX2T

PROJECTION ON THE FIRST PRINCIPAL COMPONENT IN THE PHI SPACE

ZGRID KPCATTRANSFORMXGRID 0RESHAPEX1SHAPE

PLTCONTOURX1 X2 ZGRID COLORSGREY LINEWIDTHS1 ORIGINLOWER

PLTSUBPLOT2 2 2 ASPECTEQUAL

PLTSCATTERXPCAREDS 0 XPCAREDS 1 CRED

S20 EDGECOLORK

PLTSCATTERXPCABLUES 0 XPCABLUES 1 CBLUE

S20 EDGECOLORK

PLTTITLEPROJECTION BY PCA

PLTXLABEL1ST PRINCIPAL COMPONENT

PLTYLABEL2ND COMPONENT

PLTSUBPLOT2 2 3 ASPECTEQUAL

PLTSCATTERXKPCAREDS 0 XKPCAREDS 1 CRED

S20 EDGECOLORK

PLTSCATTERXKPCABLUES 0 XKPCABLUES 1 CBLUE

S20 EDGECOLORK

PLTTITLEPROJECTION BY KPCA

PLTXLABELR1ST PRINCIPAL COMPONENT IN SPACE INDUCED BY PHI

PLTYLABEL2ND COMPONENT

PLTSUBPLOT2 2 4 ASPECTEQUAL

PLTSCATTERXBACKREDS 0 XBACKREDS 1 CRED

S20 EDGECOLORK

PLTSCATTERXBACKBLUES 0 XBACKBLUES 1 CBLUE

S20 EDGECOLORK

PLTTITLEORIGINAL SPACE AFTER INVERSE TRANSFORM

PLTXLABELX1

PLTYLABELX2

PLTTIGHTLAYOUT

PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0460 SECONDS

NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)

511 DECOMPOSITION 987

SCIKITLEARN USER GUIDE RELEASE 0213

5119 MODEL SELECTION WITH PROBABILISTIC PCA AND FACTOR ANALYSIS FA

PROBABILISTIC PCA AND FACTOR ANALYSIS ARE PROBABILISTIC MODELS THE CONSEQUENCE IS THAT THE LIKELIHOOD OF NEW DATA CAN BE USED FOR MODEL SELECTION AND COVARIANCE ESTIMATION HERE WE COMPARE PCA AND FA WITH CROSSVALIDATION ON LOW RANK DATA CORRUPTED WITH HOMOSCEDASTIC NOISE NOISE VARIANCE IS THE SAME FOR EACH FEATURE OR HETEROSCEDASTIC NOISE NOISE VARIANCE IS THE DIFFERENT FOR EACH FEATURE IN A SECOND STEP WE COMPARE THE MODEL LIKELIHOOD TO THE LIKELIHOODS OBTAINED FROM SHRINKAGE COVARIANCE ESTIMATORS

ONE CAN OBSERVE THAT WITH HOMOSCEDASTIC NOISE BOTH FA AND PCA SUCCEED IN RECOVERING THE SIZE OF THE LOW RANK SUBSPACE THE LIKELIHOOD WITH PCA IS HIGHER THAN FA IN THIS CASE HOWEVER PCA FAILS AND OVERESTIMATES THE RANK WHEN HETEROSCEDASTIC NOISE IS PRESENT UNDER APPROPRIATE CIRCUMSTANCES THE LOW RANK MODELS ARE MORE LIKELY THAN SHRINKAGE MODELS

THE AUTOMATIC ESTIMATION FROM AUTOMATIC CHOICE OF DIMENSIONALITY FOR PCA NIPS 2000 598604 BY THOMAS P MINKA IS ALSO COMPARED

•

988 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

- OUT
- BEST NCOMPONENTS BY PCA CV 10
- BEST NCOMPONENTS BY FACTORANALYSIS CV 10
- BEST NCOMPONENTS BY PCA MLE 10
- BEST NCOMPONENTS BY PCA CV 35
- BEST NCOMPONENTS BY FACTORANALYSIS CV 10
- BEST NCOMPONENTS BY PCA MLE 38
- AUTHORS ALEXANDRE GRAMFORT
- DENIS A ENGEMANN
- LICENSE BSD 3 CLAUSE
- IMPORT NUMPY AS NP
- IMPORT MATPLOTLIBPYPLOT AS PLT
- FROM SCIPY IMPORT LINALG
- FROM SKLEARNDECOMPOSITION IMPORT PCA FACTORANALYSIS
- FROM SKLEARNCOVARIANCE IMPORT SHRUNKCOVARIANCE LEDOITWOLF
- FROM SKLEARNMODELSELECTION IMPORT CROSSVALSCORE
- FROM SKLEARNMODELSELECTION IMPORT GRIDSEARCHCV
- 511 DECOMPOSITION 989

SCIKITLEARN USER GUIDE RELEASE 0213  
PRINTDOC

CREATE THE DATA  
NSAMPLES NFEATURES RANK 1000 50 10  
SIGMA 1  
RNG NPRANDOMRANDOMSTATE42  
U LINALGSVDRNGRANDNNFEATURES NFEATURES  
X NPDOTRNGRANDNNSAMPLES RANK U RANKT  
ADDING HOMOSCEDASTIC NOISE  
XHOMO X SIGMA RNGRANDNNSAMPLES NFEATURES  
ADDING HETEROSCEDASTIC NOISE  
SIGMAS SIGMA RNGRANDNFEATURES SIGMA 2  
XHETERO X RNGRANDNNSAMPLES NFEATURES SIGMAS

FIT THE MODELS  
NCOMPONENTS NPARANGE0 NFEATURES 5 OPTIONS FOR NCOMPONENTS  
DEFCOMPUTESCORES  
PCA PCASVDSOLVERFULL  
FA FACTORANALYSIS  
PCASCORES FASCORES  
FORNINNCOMPONENTS  
PCANCOMPONENTS N  
FANCOMPONENTS N  
PCASCORESAPPENDNPMEANCROSSVALSCOREPCA X CV5  
FASCORESAPPENDNPMEANCROSSVALSCOREFA X CV5  
RETURNPCASCORES FASCORES  
DEFshrunkCOVSCOREX  
SHRINKAGES NPLOGSPACE2 0 30  
CV GRIDSEARCHCVshrunkCOVARIANCE SHRINKAGE SHRINKAGES CV5  
RETURNNPMEANCROSSVALSCORECVFITXBESTESTIMATOR X CV5  
DEFLWSCOREX  
RETURNNPMEANCROSSVALSCORELEDOITWOLF X CV5  
FORX TITLE INXHOMO HOMOSCEDASTIC NOISE  
XHETERO HETEROSCEDASTIC NOISE  
PCASCORES FASCORES COMPUTESCORES  
NCOMPONENTSPCA NCOMPONENTSNPARGMAXPCASCORES  
NCOMPONENTSFa NCOMPONENTSNPARGMAXFASCORES  
PCA PCASVDSOLVERFULL NCOMPONENTSMLE  
PCAFITX  
NCOMPONENTSPCAMLE PCANCOMPONENTS  
990 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
PRINTBEST NCOMPONENTS BY PCA CV D NCOMPONENTSPCA  
PRINTBEST NCOMPONENTS BY FACTORANALYSIS CV D NCOMPONENTSFA  
PRINTBEST NCOMPONENTS BY PCA MLE D NCOMPONENTSPCAMLE  
PLTFigure  
PLTPLOTNCOMPONENTS PCAScores B LABELPCA Scores  
PLTPLOTNCOMPONENTS FAScores R LABELFA Scores  
PLTAXVLINERANK COLORG LABELTRUTH D RANK LINESTYLE  
PLTAXVLINENCOMPONENTSPCA COLORB  
LABELPCA CV D NCOMPONENTSPCA LINESTYLE  
PLTAXVLINENCOMPONENTSFA COLORR  
LABELFACTORANALYSIS CV D NCOMPONENTSFA  
LINESTYLE  
PLTAXVLINENCOMPONENTSPCAMLE COLORK  
LABELPCA MLE D NCOMPONENTSPCAMLE LINESTYLE  
COMPARE WITH OTHER COVARIANCE ESTIMATORS  
PLTAXHLINESHRUNKCOVSCOREX COLORVIOLET  
LABELSHRUNK COVARIANCE MLE LINESTYLE  
PLTAXHLINELWSCOREX COLORORANGE  
LABELLEDOITWOLF MLE NCOMPONENTSPCAMLE LINESTYLE  
PLTXLABELNB OF COMPONENTS  
PLTYLABELCV Scores  
PLTLEGENDLOCLOWER RIGHT  
PLTTITLETITLE  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 17528 SECONDS  
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE  
51110 SPARSE CODING WITH A PRECOMPUTED DICTIONARY  
TRANSFORM A SIGNAL AS A SPARSE COMBINATION OF RICKER WAVELETS THIS EXAMPLE VISUALLY COMPARES DIFFERENT SPARSE CODING  
METHODS USING THE SKLEARNDECOMPOSITIONSPARSECODER ESTIMATOR THE RICKER ALSO KNOWN AS MEXICAN  
HAT OR THE SECOND DERIVATIVE OF A GAUSSIAN IS NOT A PARTICULARLY GOOD KERNEL TO REPRESENT PIECEWISE CONSTANT SIGNALS  
LIKE THIS ONE IT CAN THEREFORE BE SEEN HOW MUCH ADDING DIFFERENT WIDTHS OF ATOMS MATTERS AND IT THEREFORE MOTIVATES  
LEARNING THE DICTIONARY TO BEST FIT YOUR TYPE OF SIGNALS  
THE RICHER DICTIONARY ON THE RIGHT IS NOT LARGER IN SIZE HEAVIER SUBSAMPLING IS PERFORMED IN ORDER TO STAY ON THE SAME  
ORDER OF MAGNITUDE  
511 DECOMPOSITION 991

SCIKITLEARN USER GUIDE RELEASE 0213  
PRINTDOC  
FROM DISTUTILSVERSION IMPORT LOOSEVERSION  
IMPORT NUMPY AS NP  
IMPORT MATPLOTLIBPYPLOT AS PLT  
FROM SKLEARNDECOMPOSITION IMPORT SPARSECODER  
DEFRICKERFUNCTIONRESOLUTION CENTER WIDTH  
DISCRETE SUBSAMPLED RICKER MEXICAN HAT WAVELET  
X NPLinspace0 RESOLUTION 1 RESOLUTION  
X 2 NPSQRT3 WIDTHNPPI1 4  
1 X CENTER 2 WIDTH 2  
NPEXPX CENTER 2 2 WIDTH2  
RETURNX  
DEFRICKERMATRIXWIDTH RESOLUTION NCOMPONENTS  
DICTIONARY OF RICKER MEXICAN HAT WAVELETS  
CENTERS NPLinspace0 RESOLUTION 1 NCOMPONENTS  
D NPEMPTYNCOMPONENTS RESOLUTION  
FORI CENTER INENUMERATECENTERS  
DI RICKERFUNCTIONRESOLUTION CENTER WIDTH  
D NPSQRTNPSUMD 2 AXIS1 NPNEWAXIS  
RETURNND  
RESOLUTION 1024  
SUBSAMPLING 3 SUBSAMPLING FACTOR  
WIDTH 100  
NCOMPONENTS RESOLUTION SUBSAMPLING  
COMPUTE A WAVELET DICTIONARY  
DFIXED RICKERMATRIXWIDTHWIDTH RESOLUTIONRESOLUTION  
NCOMPONENTSNCOMPONENTS  
DMULTI NPRTUPLERICKERMATRIXWIDTHW RESOLUTIONRESOLUTION  
992 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
NCOMPONENTSNCOMPONENTS 5  
FORWIN10 50 100 500 1000  
GENERATE A SIGNAL  
Y NPLinspace0 RESOLUTION 1 RESOLUTION  
FIRSTQUARTER Y RESOLUTION 4  
YFIRSTQUARTER 3  
YNPLOGICALNOTFIRSTQUARTER 1  
LIST THE DIFFERENT SPARSE CODING METHODS IN THE FOLLOWING FORMAT  
TITLE TRANSFORMALGORITHM TRANSFORMALPHA  
TRANSFORMNNOZEROCOEFSCOLOR  
ESTIMATORS OMP OMP NONE 15 NAVY  
LASSO LASSOLARS 2 NONE TURQUOISE  
LW 2  
AVOID FUTUREWARNING ABOUT DEFAULT VALUE CHANGE WHEN NUMPY 114  
LSTSQRCOND NONE IFLOOSEVERSIONNPVERSION 114 ELSE1  
PLTFIGUREFIGSIZE13 6  
FORSUBPLOT D TITLE INENUMERATEZIPDFIXED DMULTI  
FIXED WIDTH MULTIPLE WIDTHS  
PLTSUBPLOT1 2 SUBPLOT 1  
PLTTITLESPARSE CODING AGAINST SDICTIONARY TITLE  
PLTPLOTY LWLW LINESTYLE LABELORIGINAL SIGNAL  
DO A WAVELET APPROXIMATION  
FORTITLE ALGO ALPHA NNONZERO COLOR INESTIMATORS  
CODER SPARSECODERDICTIONARYD TRANSFORMNNONZEROCOEFSSNONZERO  
TRANSFORMALPHAALPHA TRANSFORMALGORITHMALGO  
X CODERTRANSFORMYRESHAPE1 1  
DENSITY LENNPFLATNONZEROX  
X NPRAVELNPDOTX D  
SQUAREDERROR NPSUMY X 2  
PLTPLOTX COLORCOLOR LWLW  
LABELSSNONZERO COEFS N2FERROR  
TITLE DENSITY SQUAREDERROR  
SOFT THRESHOLDING DEBIASING  
CODER SPARSECODERDICTIONARYD TRANSFORMALGORITHMTHRESHOLD  
TRANSFORMALPHA20  
X CODERTRANSFORMYRESHAPE1 1  
IDX NPWHEREX 0  
X0 IDX NPLINALGLSTSQDIDX T Y RCONDLSTSQRCOND  
X NPRAVELNPDOTX D  
SQUAREDERROR NPSUMY X 2  
PLTPLOTX COLORDARKORANGE LWLW  
LABELTHRESHOLDING W DEBIASING NDNONZERO COEFS 2FERROR  
LENIDX SQUAREDERROR  
PLTAXISTIGHT  
PLTLEGENDSHADOWFALSE LOCBEST  
PLTSUBPLOTSADJUST04 07 97 90 09 2  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0238 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
511 DECOMPOSITION 993

SCIKITLEARN USER GUIDE RELEASE 0213

51111 IMAGE DENOISING USING DICTIONARY LEARNING

AN EXAMPLE COMPARING THE EFFECT OF RECONSTRUCTING NOISY FRAGMENTS OF A RACCOON FACE IMAGE USING FIRSTLY ONLINE DICTIONARY LEARNING AND VARIOUS TRANSFORM METHODS

THE DICTIONARY IS FITTED ON THE DISTORTED LEFT HALF OF THE IMAGE AND SUBSEQUENTLY USED TO RECONSTRUCT THE RIGHT HALF NOTE THAT EVEN BETTER PERFORMANCE COULD BE ACHIEVED BY FITTING TO AN UNDISTORTED IE NOISELESS IMAGE BUT HERE WE START FROM THE ASSUMPTION THAT IT IS NOT AVAILABLE

A COMMON PRACTICE FOR EVALUATING THE RESULTS OF IMAGE DENOISING IS BY LOOKING AT THE DIFFERENCE BETWEEN THE RECONSTRUCTION AND THE ORIGINAL IMAGE IF THE RECONSTRUCTION IS PERFECT THIS WILL LOOK LIKE GAUSSIAN NOISE

IT CAN BE SEEN FROM THE PLOTS THAT THE RESULTS OF ORTHOGONAL MATCHING PURSUIT OMP WITH TWO NONZERO COEFFICIENTS IS A BIT LESS BIASED THAN WHEN KEEPING ONLY ONE THE EDGES LOOK LESS PROMINENT IT IS IN ADDITION CLOSER FROM THE GROUND TRUTH IN FROBENIUS NORM

THE RESULT OF LEAST ANGLE REGRESSION IS MUCH MORE STRONGLY BIASED THE DIFFERENCE IS REMINISCENT OF THE LOCAL INTENSITY VALUE OF THE ORIGINAL IMAGE

THRESHOLDING IS CLEARLY NOT USEFUL FOR DENOISING BUT IT IS HERE TO SHOW THAT IT CAN PRODUCE A SUGGESTIVE OUTPUT WITH VERY HIGH SPEED AND THUS BE USEFUL FOR OTHER TASKS SUCH AS OBJECT CLASSIFICATION WHERE PERFORMANCE IS NOT NECESSARILY RELATED TO VISUALISATION

•

994 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213

- 
- 

511 DECOMPOSITION 995

- 
-

SCIKITLEARN USER GUIDE RELEASE 0213

- OUT  
DISTORTING IMAGE  
EXTRACTING REFERENCE PATCHES  
DONE IN 001S  
LEARNING THE DICTIONARY  
DONE IN 722S  
EXTRACTING NOISY PATCHES  
DONE IN 000S  
ORTHOGONAL MATCHING PURSUIT  
1 ATOM  
DONE IN 119S  
ORTHOGONAL MATCHING PURSUIT  
2 ATOMS  
DONE IN 235S  
LEASTANGLE REGRESSION  
5 ATOMS  
DONE IN 1750S  
THRESHOLDING  
ALPHA01  
DONE IN 025S  
PRINTDOC  
FROM TIME IMPORT TIME  
IMPORT MATPLOTLIBPYPLOT AS PLT  
IMPORT NUMPY AS NP  
IMPORT SCIPY AS SP  
FROM SKLEARNDECOMPOSITION IMPORT MINIBATCHDICTIONARYLEARNING  
511 DECOMPOSITION 997

```
SCIKITLEARN USER GUIDE RELEASE 0213
FROM SKLEARNFEATUREEXTRACTIONIMAGE IMPORT EXTRACTPATCHES2D
FROM SKLEARNFEATUREEXTRACTIONIMAGE IMPORT RECONSTRUCTFROMPATCHES2D
TRY SCIPY 016 HAVE FACE IN MISC
FROM SCIPYMISC IMPORT FACE
FACE FACEGRAYTRUE
EXCEPTIMPORTERROR
FACE SPFACEGRAYTRUE
CONVERT FROM UINT8 REPRESENTATION WITH VALUES BETWEEN 0 AND 255 TO
A FLOATING POINT REPRESENTATION WITH VALUES BETWEEN 0 AND 1
FACE FACE 255
DOWNSAMPLE FOR HIGHER SPEED
FACE FACE4 4 FACE14 4 FACE4 14 FACE14 14
FACE 40
HEIGHT WIDTH FACESHAPE
DISTORT THE RIGHT HALF OF THE IMAGE
PRINTDISTORTING IMAGE
DISTORTED FACECOPY
DISTORTED WIDTH 2 0075 NPRANDOMRANDNHEIGHT WIDTH 2
EXTRACT ALL REFERENCE PATCHES FROM THE LEFT HALF OF THE IMAGE
PRINTEXTRACTING REFERENCE PATCHES
TO TIME
PATCHSIZE 7 7
DATA EXTRACTPATCHES2DDISTORTED WIDTH 2 PATCHSIZE
DATA DATARESHAPEDATASHAPE0 1
DATA NPMEANDATA AXIS0
DATA NPSTDDATA AXIS0
PRINTDONE IN 2FS TIME TO

LEARN THE DICTIONARY FROM REFERENCE PATCHES
PRINTLEARNING THE DICTIONARY
TO TIME
DICO MINIBATCHDICTIONARYLEARNINGNCOMPONENTS100 ALPHA1 NITER500
V DICOFITDATACOMPONENTS
DT TIME TO
PRINTDONE IN 2FS DT
PLTFIGUREFIGSIZE42 4
FORI COMP INENUMERATEV100
PLTSUBPLOT10 10 I 1
PLTIMSHOWCOMPRESHAPEPATCHSIZE CMAPPLTCMGRAYR
INTERPOLATIONNEAREST
PLTXTICKS
PLTYTICKS
PLTSUPTITLEDICTIONARY LEARNED FROM FACE PATCHES N
TRAIN TIME 1FS ONDPATCHES DT LENDATA
FONTSIZE16
PLTSUBPLOTSADJUST008 002 092 085 008 023
```

SCIKITLEARN USER GUIDE RELEASE 0213  
  DISPLAY THE DISTORTED IMAGE  
  DEFSHOWWITHDIFFIMAGE REFERENCE TITLE  
  HELPER FUNCTION TO DISPLAY DENOISING  
  PLTFIGUREFIGSIZE5 33  
  PLTSUBPLOT1 2 1  
  PLTTITLEIMAGE  
  PLTIMSHOWIMAGE VMIN0 VMAX1 CMAPPLTCMGRAY  
  INTERPOLATIONNEAREST  
  PLXTICKS  
  PLTYTICKS  
  PLTSUBPLOT1 2 2  
  DIFFERENCE IMAGE REFERENCE  
  PLTTITLEDIFFERENCE NORM 2F NPSQRTNPSUMDIFFERENCE 2  
  PLTIMSHOWDIFFERENCE VMIN05 VMAX05 CMAPPLTCMPUOR  
  INTERPOLATIONNEAREST  
  PLXTICKS  
  PLTYTICKS  
  PLTSUPTITLETITLE SIZE16  
  PLTSUBPLOTSADJUST002 002 098 079 002 02  
  SHOWWITHDIFFDISTORTED FACE DISTORTED IMAGE

  EXTRACT NOISY PATCHES AND RECONSTRUCT THEM USING THE DICTIONARY  
  PRINTEXTRACTING NOISY PATCHES  
  T0 TIME  
  DATA EXTRACTPATCHES2DDISTORTED WIDTH 2 PATCHSIZE  
  DATA DATARESHAPEDATASHAPE0 1  
  INTERCEPT NPMEANDATA AXIS0  
  DATA INTERCEPT  
  PRINTDONE IN 2FS TIME T0  
  TRANSFORMALGORITHMS  
  ORTHOGONAL MATCHING PURSUIT N1 ATOM OMP  
  TRANSFORMMNNONZEROCOEF5 1  
  ORTHOGONAL MATCHING PURSUIT N2 ATOMS OMP  
  TRANSFORMMNNONZEROCOEF5 2  
  LEASTANGLE REGRESSION N5 ATOMS LARS  
  TRANSFORMMNNONZEROCOEF5 5  
  THRESHOLDING NALPHA01 THRESHOLD TRANSFORMALPHA 1  
  RECONSTRUCTIONS  
  FORTITLE TRANSFORMALGORITHM KWARGS INTRANSFORMALGORITHMS  
  PRINTTITLE  
  RECONSTRUCTIONSTITLE FACECOPY  
  T0 TIME  
  DICOSETPARAMSTRANSFORMALGORITHMTRANSFORMALGORITHM KWARGS  
  CODE DICOTRANSFORMDATA  
  PATCHES NPDOTCODE V  
  PATCHES INTERCEPT  
  PATCHES PATCHESRESHAPELENDATA PATCHSIZE  
  IFTRANSFORMALGORITHM THRESHOLD  
  PATCHES PATCHESMIN  
  PATCHES PATCHESMAX  
  511 DECOMPOSITION 999

SCIKITLEARN USER GUIDE RELEASE 0213  
RECONSTRUCTIONSTITLE WIDTH 2 RECONSTRUCTFROMPATCHES2D  
PATCHES HEIGHT WIDTH 2  
DT TIME TO  
PRINTDONE IN 2FS DT  
SHOWWITHDIFFRECONSTRUCTIONSTITLE FACE  
TITLE TIME 1FS DT  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 30244 SECONDS  
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE  
51112 FACES DATASET DECOMPOSITIONS  
THIS EXAMPLE APPLIES TO OLIVETTIFACES DIFFERENT UNSUPERVISED MATRIX DECOMPOSITION DIMENSION REDUCTION METHODS  
FROM THE MODULE SKLEARNDECOMPOSITION SEE THE DOCUMENTATION CHAPTER DECOMPOSING SIGNALS IN COMPONENTS  
MATRIX FACTORIZATION PROBLEMS  
•  
1000 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

- 511 DECOMPOSITION 1001





SCIKITLEARN USER GUIDE RELEASE 0213

- 511 DECOMPOSITION 1003



SCIKITLEARN USER GUIDE RELEASE 0213

- 511 DECOMPOSITION 1005

SCIKITLEARN USER GUIDE RELEASE 0213

- 
- 

1006 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

- 511 DECOMPOSITION 1007



SCIKITLEARN USER GUIDE RELEASE 0213

- 511 DECOMPOSITION 1009







SCIKITLEARN USER GUIDE RELEASE 0213

•

OUT

DATASET CONSISTS OF 400 FACES

EXTRACTING THE TOP 6 EIGENFACES PCA USING RANDOMIZED SVD

DONE IN 0018S

EXTRACTING THE TOP 6 NONNEGATIVE COMPONENTS NMF

DONE IN 0109S

EXTRACTING THE TOP 6 INDEPENDENT COMPONENTS FASTICA

DONE IN 0295S

EXTRACTING THE TOP 6 SPARSE COMP MINIBATCHSPARSEPCA

DONE IN 1129S

EXTRACTING THE TOP 6 MINIBATCHDICTIONARYLEARNING

DONE IN 0979S

EXTRACTING THE TOP 6 CLUSTER CENTERS MINIBATCHKMEANS

DONE IN 0221S

EXTRACTING THE TOP 6 FACTOR ANALYSIS COMPONENTS FA

DONE IN 0206S

EXTRACTING THE TOP 6 DICTIONARY LEARNING

DONE IN 1099S

EXTRACTING THE TOP 6 DICTIONARY LEARNING POSITIVE DICTIONARY

DONE IN 1213S

EXTRACTING THE TOP 6 DICTIONARY LEARNING POSITIVE CODE

DONE IN 0688S

EXTRACTING THE TOP 6 DICTIONARY LEARNING POSITIVE DICTIONARY CODE

DONE IN 0470S

1012 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
PRINTDOC
AUTHORS VLAD NICULAE ALEXANDRE GRAMFORT
LICENSE BSD 3 CLAUSE
IMPORT LOGGING
FROM TIME IMPORT TIME
FROM NUMPYRANDOM IMPORT RANDOMSTATE
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARNDATASETS IMPORT FETCHOLIVETTIFACES
FROM SKLEARNCLUSTER IMPORT MINIBATCHKMEANS
FROM SKLEARN IMPORT DECOMPOSITION
    DISPLAY PROGRESS LOGS ON STDOUT
LOGGINGBASICCONFIGLEVELLOGGINGINFO
FORMAT ASCTIMES LEVELNAMES MESSAGES
NROW NCOL 2 3
NCOMPONENTS NROW NCOL
IMAGESHAPE 64 64
RNG RANDOMSTATE0

LOAD FACES DATA
DATASET FETCHOLIVETTIFACESSHUFFLETRUE RANDOMSTATERNG
FACES DATASETDATA
NSAMPLES NFEATURES FACESSHAPE
GLOBAL CENTERING
FACESCENTERED FACES FACESMEANAXISO
LOCAL CENTERING
FACESCENTERED FACESCENTEREDMEANAXIS1RESHAPENSAMPLES 1
PRINTDATASET CONSISTS OF DFACES NSAMPLES
DEFPLOTGALLERYTITLE IMAGES NCOLNCOL NROWNROW CMAPPLTCMGRAY
PLTFIGUREFIGSIZE2 NCOL 226 NROW
PLTSUPTITLETITLE SIZE16
FORI COMP INENUMERATEIMAGES
PLTSUBPLOTNROW NCOL I 1
VMAX MAXCOMPMAX COMPMIN
PLTIMSHOWCOMPRESHAPEIMAGESHAPE CMAPPCMAP
INTERPOLATIONNEAREST
VMINVMAX VMAXVMAX
PLXTICKS
PLTYTICKS
PLTSUBPLOTSADJUST001 005 099 093 004 0

LIST OF THE DIFFERENT ESTIMATORS WHETHER TO CENTER AND TRANSPOSE THE
PROBLEM AND WHETHER THE TRANSFORMER USES THE CLUSTERING API
ESTIMATORS
EIGENFACES PCA USING RANDOMIZED SVD
DECOMPOSITIONPCANCOMPONENTSNCOMPONENTS SVDSOLVERRANDOMIZED
511 DECOMPOSITION 1013
```

SCIKITLEARN USER GUIDE RELEASE 0213  
WHITENTRUE  
TRUE  
NONNEGATIVE COMPONENTS NMF  
DECOMPOSITIONNMFNCOMPONENTSNSCOMPONENTS INITNNDSDVDA TOL5E3  
FALSE  
INDEPENDENT COMPONENTS FASTICA  
DECOMPOSITIONFASTICANCOMPONENTSNSCOMPONENTS WHITENTRUE  
TRUE  
SPARSE COMP MINIBATCHSPARSEPCA  
DECOMPOSITIONMINIBATCHSPARSEPCANCOMPONENTSNSCOMPONENTS ALPHA08  
NITER100 BATCHSIZE3  
RANDOMSTATERNG  
NORMALIZECOMPONENTSTRUE  
TRUE  
MINIBATCHDICTIONARYLEARNING  
DECOMPOSITIONMINIBATCHDICTIONARYLEARNINGNCOMPONENTS15 ALPHA01  
NITER50 BATCHSIZE3  
RANDOMSTATERNG  
TRUE  
CLUSTER CENTERS MINIBATCHKMEANS  
MINIBATCHKMEANSNCLUSTERSNCOMPONENTS TOL1E3 BATCHSIZE20  
MAXITER50 RANDOMSTATERNG  
TRUE  
FACTOR ANALYSIS COMPONENTS FA  
DECOMPOSITIONFACTORANALYSISNCOMPONENTSNSCOMPONENTS MAXITER20  
TRUE

PLOT A SAMPLE OF THE INPUT DATA  
PLOTGALLERYFIRST CENTERED OLIVETTI FACES FACESCENTEREDNCOMPONENTS

DO THE ESTIMATION AND PLOT IT  
FORNAME ESTIMATOR CENTER INESTIMATORS  
PRINTEXTRACTING THE TOP D S NCOMPONENTS NAME  
TO TIME  
DATA FACES  
IFCENTER  
DATA FACESCENTERED  
ESTIMATORFITDATA  
TRAINTIME TIME TO  
PRINTDONE IN 03FS TRRAINTIME  
IFHASATTRESTIMATOR CLUSTERCENTERS  
COMPONENTS ESTIMATORCLUSTERCENTERS  
ELSE  
COMPONENTS ESTIMATORCOMPONENTS  
PLOT AN IMAGE REPRESENTING THE PIXELWISE VARIANCE PROVIDED BY THE  
1014 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
ESTIMATOR EG ITS NOISEVARIANCE ATTRIBUTE THE EIGENFACES ESTIMATOR  
VIA THE PCA DECOMPOSITION ALSO PROVIDES A SCALAR NOISEVARIANCE  
THE MEAN OF PIXELWISE VARIANCE THAT CANNOT BE DISPLAYED AS AN IMAGE  
SO WE SKIP IT  
IFHASATTRESTIMATOR NOISEVARIANCE AND  
ESTIMATORNOISEVARIANCENDIM 0 SKIP THE EIGENFACES CASE  
PLOTGALLERYPIXELWISE VARIANCE  
ESTIMATORNOISEVARIANCERESHAPE1 1 NCOL1  
NROW1  
PLOTGALLERY S TRAIN TIME 1FS NAME TRAJTIME  
COMPONENTSNCOMPONENTS  
PLTSHOW

VARIOUS POSITIVITY CONSTRAINTS APPLIED TO DICTIONARY LEARNING  
ESTIMATORS  
DICTIONARY LEARNING  
DECOMPOSITIONMINIBATCHDICTIONARYLEARNINGNCOMPONENTS15 ALPHA01  
NITER50 BATCHSIZE3  
RANDOMSTATERNG  
TRUE  
DICTIONARY LEARNING POSITIVE DICTIONARY  
DECOMPOSITIONMINIBATCHDICTIONARYLEARNINGNCOMPONENTS15 ALPHA01  
NITER50 BATCHSIZE3  
RANDOMSTATERNG  
POSITIVEDICTTRUE  
TRUE  
DICTIONARY LEARNING POSITIVE CODE  
DECOMPOSITIONMINIBATCHDICTIONARYLEARNINGNCOMPONENTS15 ALPHA01  
NITER50 BATCHSIZE3  
RANDOMSTATERNG  
POSITIVECODETRUE  
TRUE  
DICTIONARY LEARNING POSITIVE DICTIONARY CODE  
DECOMPOSITIONMINIBATCHDICTIONARYLEARNINGNCOMPONENTS15 ALPHA01  
NITER50 BATCHSIZE3  
RANDOMSTATERNG  
POSITIVEDICTTRUE  
POSITIVECODETRUE  
TRUE

PLOT A SAMPLE OF THE INPUT DATA  
PLOTGALLERYFIRST CENTERED OLIVETTI FACES FACESCENTEREDNCOMPONENTS  
CMAPPLTCMRDBU

DO THE ESTIMATION AND PLOT IT  
FORNAME ESTIMATOR CENTER INESTIMATORS  
PRINTEXTRACTING THE TOP D S NCOMPONENTS NAME  
TO TIME  
DATA FACES  
IFCENTER  
511 DECOMPOSITION 1015

SCIKITLEARN USER GUIDE RELEASE 0213  
DATA FACESCENTERED  
ESTIMATORFITDATA  
TRAINTIME TIME TO  
PRINTDONE IN 03FS TRRAINTIME  
COMPONENTS ESTIMATORCOMPONENTS  
PLOTGALLERYNAME COMPONENTSNCOMPONENTS CMAPPLTCMRDBU  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 8614 SECONDS  
512 ENSEMBLE METHODS  
EXAMPLES CONCERNING THE SKLEARNENSEMBLE MODULE  
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE  
5121 DECISION TREE REGRESSION WITH ADABOOST  
A DECISION TREE IS BOOSTED USING THE ADABOOSTR21ALGORITHM ON A 1D SINUSOIDAL DATASET WITH A SMALL AMOUNT OF  
GAUSSIAN NOISE 299 BOOSTS 300 DECISION TREES IS COMPARED WITH A SINGLE DECISION TREE REGRESSOR AS THE NUMBER OF  
BOOSTS IS INCREASED THE REGRESSOR CAN FIT MORE DETAIL  
1  
8 DRUCKER “IMPROVING REGRESSORS USING BOOSTING TECHNIQUES” 1997  
1016 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
PRINTDOC  
AUTHOR NOEL DAWE NOELDAWEGMAILCOM

LICENSE BSD 3 CLAUSE  
IMPORTING NECESSARY LIBRARIES  
IMPORT NUMPY AS NP  
IMPORT MATPLOTLIBPYPLOT AS PLT  
FROM SKLEARNTREE IMPORT DECISIONTREEREgressor  
FROM SKLEARNENSEMBLE IMPORT ADABOOSTREGRESSOR  
CREATE THE DATASET  
RNG NPRANDOMRANDOMSTATE1  
X NPLinspace0 6 100 NPNEWAXIS  
Y NPSINXRAVEL NPSIN6 XRAVEL RNGNORMAL0 01 XSHAPE0  
FIT REGRESSION MODEL  
REGR1 DECISIONTREEREgressorMAXDEPTH4  
REGR2 ADABOOSTREGRESSORDECISIONTREEREgressorMAXDEPTH4  
NESTIMATORS300 RANDOMSTATERNG  
REGR1FITX Y  
REGR2FITX Y  
512 ENSEMBLE METHODS 1017

SCIKITLEARN USER GUIDE RELEASE 0213

PREDICT

Y1 REGR1PREDICTX

Y2 REGR2PREDICTX

PLOT THE RESULTS

PLTFigure

PLTSCATTERX Y CK LABELTRAINING SAMPLES

PLTPLOTX Y1 CG LABELNESTIMATORS1 LINEWIDTH2

PLTPLOTX Y2 CR LABELNESTIMATORS300 LINEWIDTH2

PLTXLABELDATA

PLTYLABELTARGET

PLTTITLEBOOSTED DECISION TREE REGRESSION

PLTLEGEND

PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0226 SECONDS

NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE

5122 PIXEL IMPORTANCES WITH A PARALLEL FOREST OF TREES

THIS EXAMPLE SHOWS THE USE OF FORESTS OF TREES TO EVALUATE THE IMPORTANCE OF THE PIXELS IN AN IMAGE CLASSIFICATION TASK

FACES THE HOTTER THE PIXEL THE MORE IMPORTANT

THE CODE BELOW ALSO ILLUSTRATES HOW THE CONSTRUCTION AND THE COMPUTATION OF THE PREDICTIONS CAN BE PARALLELIZED WITHIN MULTIPLE JOBS

1018 CHAPTER 5 EXAMPLES



```
SCIKITLEARN USER GUIDE RELEASE 0213
OUT
FITTING EXTRATREESCLASSIFIER ON FACES DATA WITH 1 CORES
DONE IN 1028S
PRINTDOC
FROM TIME IMPORT TIME
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARNDATASETS IMPORT FETCHOLIVETTIFACES
FROM SKLEARNENSEMBLE IMPORT EXTRATREESCLASSIFIER
NUMBER OF CORES TO USE TO PERFORM PARALLEL FITTING OF THE FOREST MODEL
NJOBS 1
LOAD THE FACES DATASET
DATA FETCHOLIVETTIFACES
X DATAIMAGESRESHAPELENDATAIMAGES 1
Y DATATARGET
512 ENSEMBLE METHODS 1019
```

SCIKITLEARN USER GUIDE RELEASE 0213  
MASK Y 5 LIMIT TO 5 CLASSES  
X XMASK  
Y YMASK  
BUILD A FOREST AND COMPUTE THE PIXEL IMPORTANCES  
PRINTFITTING EXTRATREESCLASSIFIER ON FACES DATA WITH DCORES NJOBS  
TO TIME  
FOREST EXTRATREESCLASSIFIERNESTIMATORS1000  
MAXFEATURES128  
NJOBSNJOBS  
RANDOMSTATE0  
FORESTFITX Y  
PRINTDONE IN 03FS TIME TO  
IMPORTANCES FORESTFEATUREIMPORTANCES  
IMPORTANCES IMPORTANCESRESHAPEDATAIMAGES0SHAPE  
PLOT PIXEL IMPORTANCES  
PLTMATSHOWIMPORTANCES CMAPPLTCMHOT  
PLTTITLEPIXEL IMPORTANCES WITH FORESTS OF TREES  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 1147 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
5123 PLOT INDIVIDUAL AND VOTING REGRESSION PREDICTIONS  
PLOT INDIVIDUAL AND AVERAGED REGRESSION PREDICTIONS FOR BOSTON DATASET  
FIRST THREE EXEMPLARY REGRESSORS ARE INITIALIZED GRADIENTBOOSTINGREGRESSOR  
RANDOMFORESTREGRESSOR ANDLINEARREGRESSION AND USED TO INITIALIZE A VOTINGREGRESSOR  
THE RED STARRED DOTS ARE THE AVERAGED PREDICTIONS  
1020 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
PRINTDOC
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARN IMPORT DATASETS
FROM SKLEARNENSEMBLE IMPORT GRADIENTBOOSTINGREGRESSOR
FROM SKLEARNENSEMBLE IMPORT RANDOMFORESTREGRESSOR
FROM SKLEARNLINEARMODEL IMPORT LINEARREGRESSION
FROM SKLEARNENSEMBLE IMPORT VOTINGREGRESSOR
LOADING SOME EXAMPLE DATA
BOSTON DATASETSLOADBOSTON
X BOSTONDATA
Y BOSTONTARGET
TRAINING CLASSIFIERS
REG1 GRADIENTBOOSTINGREGRESSORRANDOMSTATE1 NESTIMATORS10
REG2 RANDOMFORESTREGRESSORRANDOMSTATE1 NESTIMATORS10
REG3 LINEARREGRESSION
EREG VOTINGREGRESSORGB REG1 RF REG2 LR REG3
REG1FITX Y
REG2FITX Y
REG3FITX Y
EREGFITX Y
XT X20
512 ENSEMBLE METHODS 1021
```

SCIKITLEARN USER GUIDE RELEASE 0213

PLTFigure

PLTPLOTREG1PREDICTXT GD LABELGRADIENTBOOSTINGREGRESSOR

PLTPLOTREG2PREDICTXT B LABELRANDOMFORESTREGRESSOR

PLTPLOTREG3PREDICTXT YS LABELLINEARREGRESSION

PLTPLOTREGPREDICTXT R LABELVOTINGREGRESSOR

PLTTICKPARAMSAXISX WHICHBOTH BOTTOMFALSE TOPFALSE

LABELBOTTOMFALSE

PLTYLABELPREDICTED

PLTXLABELTRAINING SAMPLES

PLTLEGENDLOCBEST

PLTTITLECOMPARISON OF INDIVIDUAL PREDICTIONS WITH AVERAGED

PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0117 SECONDS

NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)

5124 FEATURE IMPORTANCES WITH FORESTS OF TREES

THIS EXAMPLES SHOWS THE USE OF FORESTS OF TREES TO EVALUATE THE IMPORTANCE OF FEATURES ON AN ARTIFICIAL CLASSIFICATION TASK THE RED BARS ARE THE FEATURE IMPORTANCES OF THE FOREST ALONG WITH THEIR INTERTREES VARIABILITY AS EXPECTED THE PLOT SUGGESTS THAT 3 FEATURES ARE INFORMATIVE WHILE THE REMAINING ARE NOT

1022 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
OUT  
FEATURE RANKING  
1 FEATURE 1 0295902  
2 FEATURE 2 0208351  
3 FEATURE 0 0177632  
4 FEATURE 3 0047121  
5 FEATURE 6 0046303  
6 FEATURE 8 0046013  
7 FEATURE 7 0045575  
8 FEATURE 4 0044614  
9 FEATURE 9 0044577  
10 FEATURE 5 0043912  
PRINTDOC  
IMPORT NUMPY AS NP  
IMPORT MATPLOTLIBPYPLOT AS PLT  
FROM SKLEARNDATASETS IMPORT MAKECLASSIFICATION  
FROM SKLEARNENSEMBLE IMPORT EXTRATREESCLASSIFIER  
512 ENSEMBLE METHODS 1023

SCIKITLEARN USER GUIDE RELEASE 0213

BUILD A CLASSIFICATION TASK USING 3 INFORMATIVE FEATURES

X Y MAKECLASSIFICATIONNSAMPLES1000

NFEATURES10

NINFORMATIVE3

NREDUNDANT0

NREPEATED0

NCLASSES2

RANDOMSTATE0

SHUFFLEFALSE

BUILD A FOREST AND COMPUTE THE FEATURE IMPORTANCES

FOREST EXTRATREESCLASSIFIERNESTIMATORS250

RANDOMSTATE0

FORESTFITX Y

IMPORTANCES FORESTFEATUREIMPORTANCES

STD NPSTDTREEFEATUREIMPORTANCES FORTREEINFORESTESTIMATORS

AXISO

INDICES NPARGSORTIMPORTANCES1

PRINT THE FEATURE RANKING

PRINTFEATURE RANKING

FORFINRANGEXSHAPE1

PRINTD FEATURE DF F 1 INDICESF IMPORTANCESINDICESF

PLOT THE FEATURE IMPORTANCES OF THE FOREST

PLTFigure

PLTTITLEFEATURE IMPORTANCES

PLTBARRANGEXSHAPE1 IMPORTANCESINDICES

COLORR YERRSTDINDICES ALIGNCENTER

PLXTICKSRANGEXSHAPE1 INDICES

PLTXLIM1 XSHAPE1

PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0343 SECONDS

NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE

5125 ISOLATIONFOREST EXAMPLE

AN EXAMPLE USING SKLEARNENSEMBLEISOLATIONFOREST FOR ANOMALY DETECTION

THE ISOLATIONFOREST 'ISOLATES' OBSERVATIONS BY RANDOMLY SELECTING A FEATURE AND THEN RANDOMLY SELECTING A SPLIT VALUE BETWEEN THE MAXIMUM AND MINIMUM VALUES OF THE SELECTED FEATURE

SINCE RECURSIVE PARTITIONING CAN BE REPRESENTED BY A TREE STRUCTURE THE NUMBER OF SPLITTINGS REQUIRED TO ISOLATE A SAMPLE IS EQUIVALENT TO THE PATH LENGTH FROM THE ROOT NODE TO THE TERMINATING NODE

THIS PATH LENGTH AVERAGED OVER A FOREST OF SUCH RANDOM TREES IS A MEASURE OF NORMALITY AND OUR DECISION FUNCTION

RANDOM PARTITIONING PRODUCES NOTICEABLE SHORTER PATHS FOR ANOMALIES HENCE WHEN A FOREST OF RANDOM TREES COLLECTIVELY PRODUCE SHORTER PATH LENGTHS FOR PARTICULAR SAMPLES THEY ARE HIGHLY LIKELY TO BE ANOMALIES

1024 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
PRINTDOC
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARNENSEMBLE IMPORT ISOLATIONFOREST
RNG NPRANDOMRANDOMSTATE42
GENERATE TRAIN DATA
X 03 RNGRANDN100 2
XTRAIN NPRX 2 X 2
GENERATE SOME REGULAR NOVEL OBSERVATIONS
X 03 RNGRANDN20 2
XTEST NPRX 2 X 2
GENERATE SOME ABNORMAL NOVEL OBSERVATIONS
XOUTLIERS RNGUNIFORMLOW4 HIGH4 SIZE20 2
FIT THE MODEL
CLF ISOLATIONFORESTBEHAVIOURNEW MAXSAMPLES100
RANDOMSTATERNG CONTAMINATIONAUTO
CLFFITXTRAIN
YPREDTRAIN CLFPREDICTXTRAIN
YPREDTEST CLFPREDICTXTEST
YPREDOUTLIERS CLFPREDICTXOUTLIERS
PLOT THE LINE THE SAMPLES AND THE NEAREST VECTORS TO THE PLANE
512 ENSEMBLE METHODS 1025
```

SCIKITLEARN USER GUIDE RELEASE 0213  
XX YY NPMESHGRIDNPLINSPACE5 5 50 NPLINSPACE5 5 50  
Z CLFDECISIONFUNCTIONNPCXXRAVEL YYRAVEL  
Z ZRESHAPEXXSHAPE  
PLTTITLEISOLATIONFOREST  
PLTCONTOURFXX YY Z CMAPPLTCMBLUESR  
B1 PLTSCATTERXTRAIN 0 XTRAIN 1 CWHITE  
S20 EDGECOLORK  
B2 PLTSCATTERXTEST 0 XTEST 1 CGREEN  
S20 EDGECOLORK  
C PLTSCATTERXOUTLIERS 0 XOUTLIERS 1 CRED  
S20 EDGECOLORK  
PLTAXISTIGHT  
PLTXLIM5 5  
PLTYLIM5 5  
PLTLEGENDB1 B2 C  
TRAINING OBSERVATIONS  
NEW REGULAR OBSERVATIONS NEW ABNORMAL OBSERVATIONS  
LOCUPPER LEFT  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0237 SECONDS  
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE  
5126 PLOT THE DECISION BOUNDARIES OF A VOTINGCLASSIFIER  
PLOT THE DECISION BOUNDARIES OF A VOTINGCLASSIFIER FOR TWO FEATURES OF THE IRIS DATASET  
PLOT THE CLASS PROBABILITIES OF THE FIRST SAMPLE IN A TOY DATASET PREDICTED BY THREE DIFFERENT CLASSIFIERS AND AVERAGED BY  
THEVOTINGCLASSIFIER  
FIRST THREE EXEMPLARY CLASSIFIERS ARE INITIALIZED DECISIONTREECLASSIFIER KNEIGHBORSCLASSIFIER AND  
SVC AND USED TO INITIALIZE A SOFTVOTING VOTINGCLASSIFIER WITH WEIGHTS 2 1 2 WHICH MEANS THAT THE  
PREDICTED PROBABILITIES OF THE DECISIONTREECLASSIFIER ANDSVC COUNT 5 TIMES AS MUCH AS THE WEIGHTS OF THE  
KNEIGHBORSCLASSIFIER CLASSIFIER WHEN THE AVERAGED PROBABILITY IS CALCULATED  
1026 CHAPTER 5 EXAMPLES



```
SCIKITLEARN USER GUIDE RELEASE 0213
PRINTDOC
FROM ITERTOOLS IMPORT PRODUCT
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARN IMPORT DATASETS
FROM SKLEARNTREE IMPORT DECISIONTREECLASSIFIER
FROM SKLEARNNEIGHBORS IMPORT KNEIGHBORSCLASSIFIER
FROM SKLEARN SVM IMPORT SVC
FROM SKLEARNENSEMBLE IMPORT VOTINGCLASSIFIER
LOADING SOME EXAMPLE DATA
IRIS DATASETSLOADIRIS
X IRISDATA 0 2
Y IRISTARGET
TRAINING CLASSIFIERS
CLF1 DECISIONTREECLASSIFIERMAXDEPTH4
CLF2 KNEIGHBORSCLASSIFIERNNEIGHBORS7
CLF3 SVCGAMMA1 KERNELRBF PROBABILITYTRUE
ECLF VOTINGCLASSIFIERESTIMATORS DT CLF1 KNN CLF2
512 ENSEMBLE METHODS 1027
```

SCIKITLEARN USER GUIDE RELEASE 0213

SVC CLF3

VOTINGSOFT WEIGHTS2 1 2

CLF1FITX Y

CLF2FITX Y

CLF3FITX Y

ECLFFITX Y

PLOTTING DECISION REGIONS

XMIN XMAX X 0MIN 1 X 0MAX 1

YMIN YMAX X 1MIN 1 X 1MAX 1

XX YY NPMESHGRIDNPARANGEXMIN XMAX 01

NPARANGEYMIN YMAX 01

F AXARR PLTSUBPLOTS2 2 SHAREXCOL SHAREYROW FIGSIZE10 8

FORIDX CLF TT INZIPPRODUCT0 1 0 1

CLF1 CLF2 CLF3 ECLF

DECISION TREE DEPTH4 KNN K7

KERNEL SVM SOFT VOTING

Z CLFPREDICTNPCXXRAVEL YYRAVEL

Z ZRESHAPEXXSHAPE

AXARRIDX0 IDX1CONTOURFXX YY Z ALPHA04

AXARRIDX0 IDX1SCATTERX 0 X 1 CY

S20 EDGECOLORK

AXARRIDX0 IDX1SETTITLET

PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0202 SECONDS

NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)

5127 COMPARING RANDOM FORESTS AND THE MULTIOUTPUT META ESTIMATOR

AN EXAMPLE TO COMPARE MULTIOUTPUT REGRESSION WITH RANDOM FOREST AND THE MULTIOUTPUTMULTIOUTPUTREGRESSOR META ESTIMATOR

THIS EXAMPLE ILLUSTRATES THE USE OF THE MULTIOUTPUTMULTIOUTPUTREGRESSOR METAESTIMATOR TO PERFORM MULTIOUTPUT REGRESSION A RANDOM FOREST REGRESSOR IS USED WHICH SUPPORTS MULTIOUTPUT REGRESSION NATIVELY SO THE RESULTS CAN BE COMPARED

THE RANDOM FOREST REGRESSOR WILL ONLY EVER PREDICT VALUES WITHIN THE RANGE OF OBSERVATIONS OR CLOSER TO ZERO FOR EACH OF THE TARGETS AS A RESULT THE PREDICTIONS ARE BIASED TOWARDS THE CENTRE OF THE CIRCLE

USING A SINGLE UNDERLYING FEATURE THE MODEL LEARNS BOTH THE X AND Y COORDINATE AS OUTPUT

1028 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
PRINTDOC  
AUTHOR TIM HEAD BETATIMGMAILCOM

LICENSE BSD 3 CLAUSE  
IMPORT NUMPY AS NP  
IMPORT MATPLOTLIBPYPLOT AS PLT  
FROM SKLEARNENSEMBLE IMPORT RANDOMFORESTREGRESSOR  
FROM SKLEARNMODELSELECTION IMPORT TRAINTESTSPLIT  
FROM SKLEARNMULTIOUTPUT IMPORT MULTIOUTPUTREGRESSOR  
CREATE A RANDOM DATASET  
RNG NPRANDOMRANDOMSTATE1  
X NPSORT200 RNGRAND600 1 100 AXIS0  
Y NPARRAYNPPI NPSINXRAVEL NPPI NPCOSXRAVELT  
Y 05 RNGRAND YSHAPE  
XTRAIN XTEST YTRAIN YTEST TRAINTESTSPLIT  
X Y TRAINSIZE400 TESTSIZE200 RANDOMSTATE4  
MAXDEPTH 30  
REGRMULTIRF MULTIOUTPUTREGRESSORRANDOMFORESTREGRESSORNESTIMATORS100  
MAXDEPTHMAXDEPTH  
RANDOMSTATE0  
512 ENSEMBLE METHODS 1029

SCIKITLEARN USER GUIDE RELEASE 0213  
REGRMULTIRFFITXTRAIN YTRAIN  
REGRRF RANDOMFORESTREGRESSORNESTIMATORS100 MAXDEPTHMAXDEPTH  
RANDOMSTATE2  
REGRRFFITXTRAIN YTRAIN  
PREDICT ON NEW DATA  
YMULTIRF REGRMULTIRFPREDICTXTEST  
YRF REGRRFPREDICTXTEST  
PLOT THE RESULTS  
PLTFigure  
S 50  
A 04  
PLTSCATTERYTEST 0 YTEST 1 EDGEColorK  
CNAVY SS MARKERS ALPHAa LABELDATA  
PLTSCATTERYMULTIRF 0 YMULTIRF 1 EDGEColorK  
CCORNflowerBLUE SS ALPHAa  
LABELMULTI RF SCORE 2F REGRMULTIRFSCOREXTEST YTEST  
PLTSCATTERYRF 0 YRF 1 EDGEColorK  
CC SS MARKER ALPHAa  
LABELRF SCORE 2F REGRRFSCOREXTEST YTEST  
PLTXLIM6 6  
PLTYLIM6 6  
PLTXLABELTARGET 1  
PLTYLABELTARGET 2  
PLTTITLECOMPARING RANDOM FORESTS AND THE MULTIOUTPUT META ESTIMATOR  
PLTLEGEND  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0299 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
5128 PREDICTION INTERVALS FOR GRADIENT BOOSTING REGRESSION  
THIS EXAMPLE SHOWS HOW QUANTILE REGRESSION CAN BE USED TO CREATE PREDICTION INTERVALS  
1030 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARNENSEMBLE IMPORT GRADIENTBOOSTINGREGRESSOR
NPRANDOMSEED1
DEFFX
THE FUNCTION TO PREDICT
RETURNXNPSINX

FIRST THE NOISELESS CASE
X NPATLEAST2DNPRANDOMUNIFORM0 100 SIZE100T
X XASTYPENPFLOAT32
OBSERVATIONS
Y FXRAVEL
DY 15 10 NPRANDOMRANDOMMYSHAPE
NOISE NPRANDOMNORMAL0 DY
Y NOISE
Y YASTYPENPFLOAT32
MESH THE INPUT SPACE FOR EVALUATIONS OF THE REAL FUNCTION THE PREDICTION AND
512 ENSEMBLE METHODS 1031
```

SCIKITLEARN USER GUIDE RELEASE 0213  
ITS MSE  
XX NPATLEAST2DNPLINSPACE0 10 1000T  
XX XXASTYPENPFLOAT32  
ALPHA 095  
CLF GRADIENTBOOSTINGREGRESSORLOSSQUANTILE ALPHAALPHA  
NESTIMATORS250 MAXDEPTH3  
LEARNINGRATE1 MINSAMPLESLEAF9  
MINSAMPLESSPLIT9  
CLFFITX Y  
MAKE THE PREDICTION ON THE MESHED XAXIS  
YUPPER CLFPREDICTXX  
CLFSETPARAMSALPHA10 ALPHA  
CLFFITX Y  
MAKE THE PREDICTION ON THE MESHED XAXIS  
YLOWER CLFPREDICTXX  
CLFSETPARAMSLOSSLS  
CLFFITX Y  
MAKE THE PREDICTION ON THE MESHED XAXIS  
YPRED CLFPREDICTXX  
PLOT THE FUNCTION THE PREDICTION AND THE 90 CONFIDENCE INTERVAL BASED ON  
THE MSE  
FIG PLTFigure  
PLTPLOTXX FXX G LABELRFX XSINX  
PLTPLOTX Y B MARKERSIZE10 LABELUOBSERVATIONS  
PLTPLOTXX YPRED R LABELUPREDICTION  
PLTPLOTXX YUPPER K  
PLTPLOTXX YLOWER K  
PLTFILLNPCONCATENATEXX XX1  
NPCONCATENATEYUPPER YLOWER1  
ALPHA5 FCB ECNONE LABEL90 PREDICTION INTERVAL  
PLTXLABELX  
PLTYLABELFX  
PLTYLIM10 20  
PLTLEGENDLOCUPPER LEFT  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0275 SECONDS  
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE  
5129 GRADIENT BOOSTING REGULARIZATION  
ILLUSTRATION OF THE EFFECT OF DIFFERENT REGULARIZATION STRATEGIES FOR GRADIENT BOOSTING THE EXAMPLE IS TAKEN FROM HASTIE  
ET AL 20091  
1T HASTIE R TIBSHIRANI AND J FRIEDMAN “ELEMENTS OF STATISTICAL LEARNING ED 2” SPRINGER 2009  
1032 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
THE LOSS FUNCTION USED IS BINOMIAL DEVIANCE REGULARIZATION VIA SHRINKAGE LEARNINGRATE 10 IMPROVES  
PERFORMANCE CONSIDERABLY IN COMBINATION WITH SHRINKAGE STOCHASTIC GRADIENT BOOSTING SUBSAMPLE 10 CAN  
PRODUCE MORE ACCURATE MODELS BY REDUCING THE VARIANCE VIA BAGGING SUBSAMPLING WITHOUT SHRINKAGE USUALLY DOES  
POORLY ANOTHER STRATEGY TO REDUCE THE VARIANCE IS BY SUBSAMPLING THE FEATURES ANALOGOUS TO THE RANDOM SPLITS IN  
RANDOM FORESTS VIA THE MAXFEATURES PARAMETER  
PRINTDOC  
AUTHOR PETER PRETTENHOFER PETERPRETTENHOFERGMAILCOM

LICENSE BSD 3 CLAUSE  
IMPORT NUMPY AS NP  
IMPORT MATPLOTLIBPYPLOT AS PLT  
FROM SKLEARN IMPORT ENSEMBLE  
FROM SKLEARN IMPORT DATASETS  
X Y DATASETSMAKEHASTIE102NSAMPLES12000 RANDOMSTATE1  
X XASTYPENPFLOAT32  
MAP LABELS FROM 1 1 TO 0 1  
LABELS Y NPUNIQUEY RETURNINVERSETRUE  
XTRAIN XTEST X2000 X2000  
512 ENSEMBLE METHODS 1033

SCIKITLEARN USER GUIDE RELEASE 0213  
YTRAIN YTEST Y2000 Y2000  
ORIGINALPARAMS NESTIMATORS 1000 MAXLEAFNODES 4 MAXDEPTH NONE  
↩→RANDOMSTATE 2  
MINSAMPLESSPLIT 5  
PLTFigure  
FORLABEL COLOR SETTING INNO SHRINKAGE ORANGE  
LEARNINGRATE 10 SUBSAMPLE 10  
LEARNINGRATE01 TURQUOISE  
LEARNINGRATE 01 SUBSAMPLE 10  
SUBSAMPLE05 BLUE  
LEARNINGRATE 10 SUBSAMPLE 05  
LEARNINGRATE01 SUBSAMPLE05 GRAY  
LEARNINGRATE 01 SUBSAMPLE 05  
LEARNINGRATE01 MAXFEATURES2 MAGENTA  
LEARNINGRATE 01 MAXFEATURES 2  
PARAMS DICTORIGINALPARAMS  
PARAMSUPDATESSETTING  
CLF ENSEMBLEGRADIENTBOOSTINGCLASSIFIER PARAMS  
CLFFITXTRAIN YTRAIN  
COMPUTE TEST SET DEVIANCE  
TESTDEVIANCE NPZEROPARAMSNESTIMATORS DTYPENPFLOAT64  
FORI YPRED INENUMERATECLFSTAGEDDECISIONFUNCTIONXTEST  
CLFLOSS ASSUMES THAT YTESTI IN 0 1  
TESTDEVIANCEI CLFLOSSYTEST YPRED  
PLTPLOTNPARANGETESTDEVIANCESHAPE0 15 TESTDEVIANCES5  
COLORCOLOR LABELLABEL  
PLTLEGENDLOCUPPER LEFT  
PLTXLABELBOOSTING ITERATIONS  
PLTYLABELTEST SET DEVIANCE  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 10180 SECONDS  
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE  
51210 PLOT CLASS PROBABILITIES CALCULATED BY THE VOTINGCLASSIFIER  
PLOT THE CLASS PROBABILITIES OF THE FIRST SAMPLE IN A TOY DATASET PREDICTED BY THREE DIFFERENT CLASSIFIERS AND AVERAGED BY  
THEVOTINGCLASSIFIER  
FIRST THREE EXAMPLARY CLASSIFIERS ARE INITIALIZED LOGISTICREGRESSION GAUSSIANNB AND  
RANDOMFORESTCLASSIFIER AND USED TO INITIALIZE A SOFTVOTING VOTINGCLASSIFIER WITH WEIGHTS 1 1  
5 WHICH MEANS THAT THE PREDICTED PROBABILITIES OF THE RANDOMFORESTCLASSIFIER COUNT 5 TIMES AS MUCH AS THE  
WEIGHTS OF THE OTHER CLASSIFIERS WHEN THE AVERAGED PROBABILITY IS CALCULATED  
TO VISUALIZE THE PROBABILITY WEIGHTING WE FIT EACH CLASSIFIER ON THE TRAINING SET AND PLOT THE PREDICTED CLASS PROBABILITIES  
1034 CHAPTER 5 EXAMPLES



```
SCIKITLEARN USER GUIDE RELEASE 0213
FOR THE FIRST SAMPLE IN THIS EXAMPLE DATASET
PRINTDOC
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARNLINEARMODEL IMPORT LOGISTICREGRESSION
FROM SKLEARNNAIVEBAYES IMPORT GAUSSIANNB
FROM SKLEARNENSEMBLE IMPORT RANDOMFORESTCLASSIFIER
FROM SKLEARNENSEMBLE IMPORT VOTINGCLASSIFIER
CLF1 LOGISTICREGRESSIONSOLVERLBFGS MAXITER1000 RANDOMSTATE123
CLF2 RANDOMFORESTCLASSIFIERNESTIMATORS100 RANDOMSTATE123
CLF3 GAUSSIANNB
X NPARRAY10 10 12 14 34 22 11 12
Y NPARRAY1 1 2 2
ECLF VOTINGCLASSIFIERESTIMATORSCLR CLF1 RF CLF2 GNB CLF3
VOTINGSOFT
WEIGHTS1 1 5
PREDICT CLASS PROBABILITIES FOR ALL CLASSIFIERS
PROBAS CFITX YPREDICTPROBAX FORCINCLF1 CLF2 CLF3 ECLF
GET CLASS PROBABILITIES FOR THE FIRST SAMPLE IN THE DATASET
512 ENSEMBLE METHODS 1035
```

SCIKITLEARN USER GUIDE RELEASE 0213  
CLASS11 PRO 0 FORPRINPROBAS  
CLASS21 PRO 1 FORPRINPROBAS  
PLOTING  
N 4 NUMBER OF GROUPS  
IND NPARANGEN GROUP POSITIONS  
WIDTH 035 BAR WIDTH  
FIG AX PLTSUBPLOTS  
BARS FOR CLASSIFIER 13  
P1 AXBARIND NPHSTACKCLASS111 0 WIDTH  
COLORGREEN EDGECOLORK  
P2 AXBARIND WIDTH NPHSTACKCLASS211 0 WIDTH  
COLORLIGHTGREEN EDGECOLORK  
BARS FOR VOTINGCLASSIFIER  
P3 AXBARIND 0 0 0 CLASS111 WIDTH  
COLORBLUE EDGECOLORK  
P4 AXBARIND WIDTH 0 0 0 CLASS211 WIDTH  
COLORSTEELBLUE EDGECOLORK  
PLOT ANNOTATIONS  
PLTAXVLINE28 COLORK LINESTYLEDASHED  
AXSETXTICKSIND WIDTH  
AXSETXTICKLABELSLOGISTICREGRESSION NWEIGHT 1  
GAUSSIANNB NWEIGHT 1  
RANDOMFORESTCLASSIFIER NWEIGHT 5  
VOTINGCLASSIFIER NAVERAGE PROBABILITIES  
ROTATION40  
HARIGHT  
PLTYLIM0 1  
PLTTITLECLASS PROBABILITIES FOR SAMPLE 1 BY DIFFERENT CLASSIFIERS  
PLTLEGENDP10 P20 CLASS 1 CLASS 2 LOCUPPER LEFT  
PLTTIGHTLAYOUT  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0245 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
51211 GRADIENT BOOSTING REGRESSION  
DEMONSTRATE GRADIENT BOOSTING ON THE BOSTON HOUSING DATASET  
THIS EXAMPLE FITS A GRADIENT BOOSTING MODEL WITH LEAST SQUARES LOSS AND 500 REGRESSION TREES OF DEPTH 4  
1036 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
OUT  
MSE 65493  
PRINTDOC  
AUTHOR PETER PRETTENHOFER PETERPRETTENHOFERGMAILCOM

LICENSE BSD 3 CLAUSE  
IMPORT NUMPY AS NP  
IMPORT MATPLOTLIBPYPLOT AS PLT  
FROM SKLEARN IMPORT ENSEMBLE  
FROM SKLEARN IMPORT DATASETS  
FROM SKLEARNUTILS IMPORT SHUFFLE  
FROM SKLEARNMETRICS IMPORT MEANSQUAREDERROR  
  
LOAD DATA  
BOSTON DATASETSLOADBOSTON  
X Y SHUFFLEBOSTONDATA BOSTONTARGET RANDOMSTATE13  
X XASTYPENPFLOAT32  
OFFSET INTXSHAPE0 09  
XTRAIN YTRAIN XOFFSET YOFFSET  
XTEST YTEST XOFFSET YOFFSET

FIT REGRESSION MODEL  
PARAMS NESTIMATORS 500 MAXDEPTH 4 MINSAMPLESSPLIT 2  
LEARNINGRATE 001 LOSS LS  
512 ENSEMBLE METHODS 1037

SCIKITLEARN USER GUIDE RELEASE 0213  
CLF ENSEMBLEGRADIENTBOOSTINGREGRESSOR PARAMS  
CLFFITXTRAIN YTRAIN  
MSE MEANSQUAREDERRORYTEST CLFPREDICTXTEST  
PRINTMSE4F MSE

PLOT TRAINING DEVIANCE  
COMPUTE TEST SET DEVIANCE  
TESTSCORE NPZEROPARAMSNESTIMATORS DTYENPFLOAT64  
FORI YPRED INENUMERATECLFSTAGEDPREDICTXTEST  
TESTSCOREI CLFLOSSYTEST YPRED  
PLTFIGUREFIGSIZE12 6  
PLTSUBPLOT1 2 1  
PLTTITLEDEVIANCE  
PLTPLOTNPARANGEPARAMSNESTIMATORS 1 CLFTRAINSORE B  
LABELTRAINING SET DEVIANCE  
PLTPLOTNPARANGEPARAMSNESTIMATORS 1 TESTSCORE R  
LABELTEST SET DEVIANCE  
PLTLEGENDLOCUPPER RIGHT  
PLTXLABELBOOSTING ITERATIONS  
PLTYLABELDEVIANCE

PLOT FEATURE IMPORTANCE  
FEATUREIMPORTANCE CLFFEATUREIMPORTANCES  
MAKE IMPORTANCES RELATIVE TO MAX IMPORTANCE  
FEATUREIMPORTANCE 1000 FEATUREIMPORTANCE FEATUREIMPORTANCEMAX  
SORTEDIDX NPARGSORTFEATUREIMPORTANCE  
POS NPARANGESORTEDIDXSHAPE0 5  
PLTSUBPLOT1 2 2  
PLTBARHPOS FEATUREIMPORTANCESORTEDIDX ALIGNCENTER  
PLTYTICKSPOS BOSTONFEATURENAMESSORTEDIDX  
PLTXLABELRELATIVE IMPORTANCE  
PLTTITLEVARIABLE IMPORTANCE  
PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0416 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
51212 TWOCLASS ADABOOST

THIS EXAMPLE FITS AN ADABOOSTED DECISION STUMP ON A NONLINEARLY SEPARABLE CLASSIFICATION DATASET COMPOSED OF TWO  
“GAUSSIAN QUANTILES” CLUSTERS SEE SKLEARNDATASETSMAKEGAUSSIANQUANTILES AND PLOTS THE DECISION  
BOUNDARY AND DECISION SCORES THE DISTRIBUTIONS OF DECISION SCORES ARE SHOWN SEPARATELY FOR SAMPLES OF CLASS A AND B  
THE PREDICTED CLASS LABEL FOR EACH SAMPLE IS DETERMINED BY THE SIGN OF THE DECISION SCORE SAMPLES WITH DECISION SCORES  
GREATER THAN ZERO ARE CLASSIFIED AS B AND ARE OTHERWISE CLASSIFIED AS A THE MAGNITUDE OF A DECISION SCORE DETERMINES  
THE DEGREE OF LIKENESS WITH THE PREDICTED CLASS LABEL ADDITIONALLY A NEW DATASET COULD BE CONSTRUCTED CONTAINING A  
DESIRED PURITY OF CLASS B FOR EXAMPLE BY ONLY SELECTING SAMPLES WITH A DECISION SCORE ABOVE SOME VALUE  
1038 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
PRINTDOC  
AUTHOR NOEL DAWE NOELDAWEGMAILCOM

LICENSE BSD 3 CLAUSE  
IMPORT NUMPY AS NP  
IMPORT MATPLOTLIBPYPLOT AS PLT  
FROM SKLEARNENSEMBLE IMPORT ADABOOSTCLASSIFIER  
FROM SKLEARNTREE IMPORT DECISIONTREECLASSIFIER  
FROM SKLEARNDATASETS IMPORT MAKEGAUSSIANQUANTILES  
CONSTRUCT DATASET  
X1 Y1 MAKEGAUSSIANQUANTILESCOV2  
NSAMPLES200 NFEATURES2  
NCLASSES2 RANDOMSTATE1  
X2 Y2 MAKEGAUSSIANQUANTILESMEAN3 3 COV15  
NSAMPLES300 NFEATURES2  
NCLASSES2 RANDOMSTATE1  
X NPCONCATENATEX1 X2  
Y NPCONCATENATEY1 Y2 1  
CREATE AND FIT AN ADABOOSTED DECISION TREE  
BDT ADABOOSTCLASSIFIERDECISIONTREECLASSIFIERMAXDEPTH1  
ALGORITHMSAMME  
NESTIMATORS200  
BDTFITX Y  
PLOTCOLORS BR  
PLOTSTEP 002  
CLASSNAMES AB  
PLTFIGUREFIGSIZE10 5  
512 ENSEMBLE METHODS 1039

SCIKITLEARN USER GUIDE RELEASE 0213  
PLOT THE DECISION BOUNDARIES  
PLTSUBPLOT121  
XMIN XMAX X 0MIN 1 X 0MAX 1  
YMIN YMAX X 1MIN 1 X 1MAX 1  
XX YY NPMESHGRIDNPARANGEXMIN XMAX PLOTSTEP  
NPARANGEYMIN YMAX PLOTSTEP  
Z BDTPREDICTNPCXXRAVEL YYRAVEL  
Z ZRESHAPEXXSHAPE  
CS PLTCONTOURFXX YY Z CMAPPLTCMPAIED  
PLTAXISTIGHT  
PLOT THE TRAINING POINTS  
FORI N C INZIPRANGE2 CLASSNAMES PLOTCOLORS  
IDX NPWHEREY I  
PLTSCATTERXIDX 0 XIDX 1  
CC CMAPPLTCMPAIED  
S20 EDGECOLORK  
LABELCLASS S N  
PLTXLIMXMIN XMAX  
PLTYLIMYMIN YMAX  
PLTLEGENDLOCUPPER RIGHT  
PLTXLABELX  
PLTYLABELY  
PLTTITLEDECISION BOUNDARY  
PLOT THE TWOCLASS DECISION SCORES  
TWOCLASSOUTPUT BDTDECISIONFUNCTIONX  
PLOT RANGE TWOCLASSOUTPUTMIN TWOCLASSOUTPUTMAX  
PLTSUBPLOT122  
FORI N C INZIPRANGE2 CLASSNAMES PLOTCOLORS  
PLTHISTTWOCLASSOUTPUTY I  
BINS10  
RANGEPLOT RANGE  
FACECOLORC  
LABELCLASS S N  
ALPHA5  
EDGECOLORK  
X1 X2 Y1 Y2 PLTAXIS  
PLTAXISX1 X2 Y1 Y2 12  
PLTLEGENDLOCUPPER RIGHT  
PLTYLABELSAMPLES  
PLTXLABELSCORE  
PLTTITLEDECISION SCORES  
PLTTIGHTLAYOUT  
PLTSUBPLOTSADJUSTWSPACE035  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 2255 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
1040 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

51213 OOB ERRORS FOR RANDOM FORESTS

THE RANDOMFORESTCLASSIFIER IS TRAINED USING BOOTSTRAP AGGREGATION WHERE EACH NEW TREE IS FIT FROM A BOOTSTRAP SAMPLE OF THE TRAINING OBSERVATIONS. THE OUTFBAG OOB ERROR IS THE AVERAGE ERROR FOR EACH CALCULATED USING PREDICTIONS FROM THE TREES THAT DO NOT CONTAIN IN THEIR RESPECTIVE BOOTSTRAP SAMPLE. THIS ALLOWS THE RANDOMFORESTCLASSIFIER TO BE FIT AND VALIDATED WHILST BEING TRAINED.

THE EXAMPLE BELOW DEMONSTRATES HOW THE OOB ERROR CAN BE MEASURED AT THE ADDITION OF EACH NEW TREE DURING TRAINING. THE RESULTING PLOT ALLOWS A PRACTITIONER TO APPROXIMATE A SUITABLE VALUE OF NESTIMATORS AT WHICH THE ERROR STABILIZES.

```
import matplotlib.pyplot as plt
from collections import OrderedDict
from sklearn.datasets import make_classification
from sklearn.ensemble import RandomForestClassifier
author = 'Kian Ho Hui, Kian Hogg, Gilles Louppe, Andreas Mueller'
gilles = 'Gilles Louppe, Gilles Louppe'
andreas = 'Andreas Mueller, Andreas Mueller'

license = 'BSD 3 Clause'
printdoc = True
hastie = 'Hastie, Tibshirani, Friedman: Elements of Statistical Learning, 2nd ed, Springer, 2009'
ensemble = 'Ensemble Methods, 1041'
```

SCIKITLEARN USER GUIDE RELEASE 0213

RANDOMSTATE 123

GENERATE A BINARY CLASSIFICATION DATASET

X Y MAKECLASSIFICATIONNSAMPLES500 NFEATURES25

NCLUSTERSPERCLASS1 NINFORMATIVE15

RANDOMSTATERANDOMSTATE

NOTE SETTING THE WARMSTART CONSTRUCTION PARAMETER TO TRUE DISABLES

SUPPORT FOR PARALLELIZED ENSEMBLES BUT IS NECESSARY FOR TRACKING THE OOB

ERROR TRAJECTORY DURING TRAINING

ENSEMBLECLFS

RANDOMFORESTCLASSIFIER MAXFEATURESSQRT

RANDOMFORESTCLASSIFIERNESTIMATORS100

WARMSTARTTRUE OOBSCORETRUE

MAXFEATURESSQRT

RANDOMSTATERANDOMSTATE

RANDOMFORESTCLASSIFIER MAXFEATURESLOG2

RANDOMFORESTCLASSIFIERNESTIMATORS100

WARMSTARTTRUE MAXFEATURESLOG2

OOBSCORETRUE

RANDOMSTATERANDOMSTATE

RANDOMFORESTCLASSIFIER MAXFEATURESNONE

RANDOMFORESTCLASSIFIERNESTIMATORS100

WARMSTARTTRUE MAXFEATURESNONE

OOBSCORETRUE

RANDOMSTATERANDOMSTATE

MAP A CLASSIFIER NAME TO A LIST OF NESTIMATORS ERROR RATE PAIRS

ERRORRATE ORDEREDDICTLABEL FORLABEL INENSEMBLECLFS

RANGE OF NESTIMATORS VALUES TO EXPLORE

MINESTIMATORS 15

MAXESTIMATORS 175

FORLABEL CLF INENSEMBLECLFS

FORIINRANGEMINESTIMATORS MAXESTIMATORS 1

CLFSETPARAMSNESTIMATORSI

CLFFITX Y

RECORD THE OOB ERROR FOR EACH NESTIMATORSI SETTING

OOBERROR 1 CLFOOBSCORE

ERRORRATELABELAPPENDI OOBERROR

GENERATE THE OOB ERROR RATE VS NESTIMATORS PLOT

FORLABEL CLFERR INERRORRATEITEMS

XS YS ZIP CLFERR

PLTPLOTXS YS LABELLABEL

PLTXLIMMINESTIMATORS MAXESTIMATORS

PLTXLABELNESTIMATORS

PLTYLABELOOB ERROR RATE

PLTLEGENDLOCUPPER RIGHT

PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 5517 SECONDS

1042 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213

NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)

51214 HASHING FEATURE TRANSFORMATION USING TOTALLY RANDOM TREES

RANDOMTREEEMBEDDING PROVIDES A WAY TO MAP DATA TO A VERY HIGHDIMENSIONAL SPARSE REPRESENTATION WHICH MIGHT BE BENEFICIAL FOR CLASSIFICATION THE MAPPING IS COMPLETELY UNSUPERVISED AND VERY EFFICIENT

THIS EXAMPLE VISUALIZES THE PARTITIONS GIVEN BY SEVERAL TREES AND SHOWS HOW THE TRANSFORMATION CAN ALSO BE USED FOR NONLINEAR DIMENSIONALITY REDUCTION OR NONLINEAR CLASSIFICATION

POINTS THAT ARE NEIGHBORING OFTEN SHARE THE SAME LEAF OF A TREE AND THEREFORE SHARE LARGE PARTS OF THEIR HASHED REPRESENTATION THIS ALLOWS TO SEPARATE TWO CONCENTRIC CIRCLES SIMPLY BASED ON THE PRINCIPAL COMPONENTS OF THE TRANSFORMED DATA WITH TRUNCATED SVD

IN HIGHDIMENSIONAL SPACES LINEAR CLASSIFIERS OFTEN ACHIEVE EXCELLENT ACCURACY FOR SPARSE BINARY DATA BERNOULLINB IS PARTICULARLY WELLSUITED THE BOTTOM ROW COMPARES THE DECISION BOUNDARY OBTAINED BY BERNOULLINB IN THE TRANSFORMED SPACE WITH AN EXTRATREESCLASSIFIER FORESTS LEARNED ON THE ORIGINAL DATA

512 ENSEMBLE METHODS 1043

```
SCIKITLEARN USER GUIDE RELEASE 0213
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARNDATASETS IMPORT MAKECIRCLES
FROM SKLEARNENSEMBLE IMPORT RANDOMTREEEMBEDDING EXTRATREESCLASSIFIER
FROM SKLEARNDECOMPOSITION IMPORT TRUNCATEDSVD
FROM SKLEARNNAIVEBAYES IMPORT BERNOULLINB
MAKE A SYNTHETIC DATASET
X Y MAKECIRCLESFACTOR05 RANDOMSTATE0 NOISE005
USE RANDOMTREEEMBEDDING TO TRANSFORM DATA
HASHER RANDOMTREEEMBEDDINGNESTIMATORS10 RANDOMSTATE0 MAXDEPTH3
XTRANSFORMED HASHERFITTRANSFORMX
VISUALIZE RESULT AFTER DIMENSIONALITY REDUCTION USING TRUNCATED SVD
SVD TRUNCATEDSVDNCOMPONENTS2
XREDUCED SVDFITTRANSFORMXTRANSFORMED
LEARN A NAIVE BAYES CLASSIFIER ON THE TRANSFORMED DATA
NB BERNOULLINB
NBFITXTRANSFORMED Y
LEARN AN EXTRATREESCLASSIFIER FOR COMPARISON
TREES EXTRATREESCLASSIFIERMAXDEPTH3 NESTIMATORS10 RANDOMSTATE0
TREESFITX Y
SCATTER PLOT OF ORIGINAL AND REDUCED DATA
FIG PLTFIGUREFIGSIZE9 8
AX PLTSUBPLOT221
AXSCATTERX 0 X 1 CY S50 EDGECOLORK
AXSETTITLEORIGINAL DATA 2D
AXSETXTICKS
AXSETYTICKS
AX PLTSUBPLOT222
AXSCATTERXREDUCED 0 XREDUCED 1 CY S50 EDGECOLORK
AXSETTITLETRUNCATED SVD REDUCTION 2D OF TRANSFORMED DATA DD
XTRANSFORMEDSHAPE1
AXSETXTICKS
AXSETYTICKS
PLOT THE DECISION IN ORIGINAL SPACE FOR THAT WE WILL ASSIGN A COLOR
TO EACH POINT IN THE MESH XMIN XMAXXYMIN YMAX
H 01
XMIN XMAX X 0MIN 5 X 0MAX 5
YMIN YMAX X 1MIN 5 X 1MAX 5
XX YY NPMESHGRIDNPARANGEXMIN XMAX H NPARANGEYMIN YMAX H
TRANSFORM GRID USING RANDOMTREEEMBEDDING
TRANSFORMEDGRID HASHERTRANSFORMNPCXXRAVEL YYRAVEL
YGRIDPRED NBPREDICTPROBATTRANSFORMEDGRID 1
AX PLTSUBPLOT223
AXSETTITLENAIVE BAYES ON TRANSFORMED DATA
1044 CHAPTER 5 EXAMPLES
```

SCIKITLEARN USER GUIDE RELEASE 0213  
AXPCOLORMESHXX YY YGRIDPREDRESHAPEXXSHAPE  
AXSCATTERX 0 X 1 CY S50 EDGECOLORK  
AXSETYLIM14 14  
AXSETXLIM14 14  
AXSETXTICKS  
AXSETYTICKS  
TRANSFORM GRID USING EXTRATREESCLASSIFIER  
YGRIDPRED TREESPREDICTPROBANPCXXRAVEL YYRAVEL 1  
AX PLTSUBPLOT224  
AXSETTITLEEXTRATREES PREDICTIONS  
AXPCOLORMESHXX YY YGRIDPREDRESHAPEXXSHAPE  
AXSCATTERX 0 X 1 CY S50 EDGECOLORK  
AXSETYLIM14 14  
AXSETXLIM14 14  
AXSETXTICKS  
AXSETYTICKS  
PLTTIGHTLAYOUT  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0250 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
51215 MULTICLASS ADABOOSTED DECISION TREES  
THIS EXAMPLE REPRODUCES FIGURE 1 OF ZHU ET AL1AND SHOWS HOW BOOSTING CAN IMPROVE PREDICTION ACCURACY ON A MULTI  
CLASS PROBLEM THE CLASSIFICATION DATASET IS CONSTRUCTED BY TAKING A TENDIMENSIONAL STANDARD NORMAL DISTRIBUTION AND  
DEFINING THREE CLASSES SEPARATED BY NESTED CONCENTRIC TENDIMENSIONAL SPHERES SUCH THAT ROUGHLY EQUAL NUMBERS OF  
SAMPLES ARE IN EACH CLASS QUANTILES OF THE 2DISTRIBUTION  
THE PERFORMANCE OF THE SAMME AND SAMMER1ALGORITHMS ARE COMPARED SAMMER USES THE PROBABILITY ESTIMATES  
TO UPDATE THE ADDITIVE MODEL WHILE SAMME USES THE CLASSIFICATIONS ONLY AS THE EXAMPLE ILLUSTRATES THE SAMMER  
ALGORITHM TYPICALLY CONVERGES FASTER THAN SAMME ACHIEVING A LOWER TEST ERROR WITH FEWER BOOSTING ITERATIONS THE  
ERROR OF EACH ALGORITHM ON THE TEST SET AFTER EACH BOOSTING ITERATION IS SHOWN ON THE LEFT THE CLASSIFICATION ERROR ON THE  
TEST SET OF EACH TREE IS SHOWN IN THE MIDDLE AND THE BOOST WEIGHT OF EACH TREE IS SHOWN ON THE RIGHT ALL TREES HAVE A  
WEIGHT OF ONE IN THE SAMMER ALGORITHM AND THEREFORE ARE NOT SHOWN  
1  
10 ZHU H ZOU S ROSSET T HASTIE “MULTICLASS ADABOOST” 2009  
512 ENSEMBLE METHODS 1045

SCIKITLEARN USER GUIDE RELEASE 0213  
PRINTDOC  
AUTHOR NOEL DAWE NOELDAWEGMAILCOM

LICENSE BSD 3 CLAUSE  
IMPORT MATPLOTLIBPYPLOT AS PLT  
FROM SKLEARNDATASETS IMPORT MAKEGAUSSIANQUANTILES  
FROM SKLEARNENSEMBLE IMPORT ADABOOSTCLASSIFIER  
FROM SKLEARNMETRICS IMPORT ACCURACYScore  
FROM SKLEARNTREE IMPORT DECISIONTREECLASSIFIER  
X Y MAKEGAUSSIANQUANTILES NSAMPLES13000 NFEATURES10  
NCLASSES3 RANDOMSTATE1  
NSPLIT 3000  
XTRAIN XTEST XNSPLIT XNSPLIT  
YTRAIN YTEST YNSPLIT YNSPLIT  
BDTREAL ADABOOSTCLASSIFIER  
DECISIONTREECLASSIFIERMAXDEPTH2  
NESTIMATORS600  
LEARNINGRATE1  
BDTDISCRETE ADABOOSTCLASSIFIER  
DECISIONTREECLASSIFIERMAXDEPTH2  
NESTIMATORS600  
LEARNINGRATE15  
ALGORITHMSAMME  
BDTREALFITXTRAIN YTRAIN  
BDTDISCRETEFITXTRAIN YTRAIN  
REALTESTERRORS  
DISCRETETESTERRORS  
FORREALTESTPREDICT DISCRETETRAINPREDICT INZIP  
BDTREALSTAGEDPREDICTXTEST BDTDISCRETETAGEDPREDICTXTEST  
REALTESTERRORSAPPEND  
1 ACCURACYScoreREALTESTPREDICT YTEST  
1046 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
DISCRETETESTERRORSAPPEND  
1 ACCURACYSCOREDISCRETETRAINPREDICT YTEST  
NTREESDISCRETE LENBDTDISCRETE  
NTREESREAL LENBDTREAL  
BOOSTING MIGHT TERMINATE EARLY BUT THE FOLLOWING ARRAYS ARE ALWAYS  
NESTIMATORS LONG WE CROP THEM TO THE ACTUAL NUMBER OF TREES HERE  
DISCRETEESTIMATORERRORS BDTDISCRETEESTIMATORERRORSNTREESDISCRETE  
REALESTIMATORERRORS BDTREALESTIMATORERRORSNTREESREAL  
DISCRETEESTIMATORWEIGHTS BDTDISCRETEESTIMATORWEIGHTSNTREESDISCRETE  
PLTFIGUREFIGSIZE15 5  
PLTSUBPLOT131  
PLTPLOTRANGE1 NTREESDISCRETE 1  
DISCRETETESTERRORS CBLACK LABELSAMME  
PLTPLOTRANGE1 NTREESREAL 1  
REALTESTERRORS CBLACK  
LINESTYLEDASHED LABELSAMMER  
PLTLEGEND  
PLTYLIM018 062  
PLTYLABELTEST ERROR  
PLTXLABELNUMBER OF TREES  
PLTSUBPLOT132  
PLTPLOTRANGE1 NTREESDISCRETE 1 DISCRETEESTIMATORERRORS  
B LABELSAMME ALPHA5  
PLTPLOTRANGE1 NTREESREAL 1 REALESTIMATORERRORS  
R LABELSAMMER ALPHA5  
PLTLEGEND  
PLTYLABELERROR  
PLTXLABELNUMBER OF TREES  
PLTYLIM2  
MAXREALESTIMATORERRORSMAX  
DISCRETEESTIMATORERRORSMAX 12  
PLTXLIM20 LENBDTDISCRETE 20  
PLTSUBPLOT133  
PLTPLOTRANGE1 NTREESDISCRETE 1 DISCRETEESTIMATORWEIGHTS  
B LABELSAMME  
PLTLEGEND  
PLTYLABELWEIGHT  
PLTXLABELNUMBER OF TREES  
PLTYLIM0 DISCRETEESTIMATORWEIGHTSMAX 12  
PLTXLIM20 NTREESDISCRETE 20  
PREVENT OVERLAPPING YAXIS LABELS  
PLTSUBPLOTSADJUSTWSPACE025  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 11401 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
512 ENSEMBLE METHODS 1047

SCIKITLEARN USER GUIDE RELEASE 0213

51216 DISCRETE VERSUS REAL ADABOOST

THIS EXAMPLE IS BASED ON FIGURE 102 FROM HASTIE ET AL 2009<sup>1</sup>AND ILLUSTRATES THE DIFFERENCE IN PERFORMANCE BETWEEN THE DISCRETE SAMME<sup>2</sup>BOOSTING ALGORITHM AND REAL SAMMER BOOSTING ALGORITHM BOTH ALGORITHMS ARE EVALUATED ON A BINARY CLASSIFICATION TASK WHERE THE TARGET Y IS A NONLINEAR FUNCTION OF 10 INPUT FEATURES

DISCRETE SAMME ADABOOST ADAPTS BASED ON ERRORS IN PREDICTED CLASS LABELS WHEREAS REAL SAMMER USES THE PREDICTED CLASS PROBABILITIES

PRINTDOC

AUTHOR PETER PRETTENHOFER PETERPRETTENHOFER@GMAIL.COM

NOEL DAWE NOELDAWEG@GMAIL.COM

LICENSE BSD 3 CLAUSE

IMPORT NUMPY AS NP

IMPORT MATPLOTLIBPYPLOT AS PLT

FROM SKLEARN IMPORT DATASETS

FROM SKLEARNTREE IMPORT DECISIONTREECLASSIFIER

FROM SKLEARNMETRICS IMPORT ZEROONELOSS

<sup>1</sup>T HASTIE R TIBSHIRANI AND J FRIEDMAN “ELEMENTS OF STATISTICAL LEARNING ED 2” SPRINGER 2009

<sup>2</sup>

<sup>10</sup> ZHU H ZOU S ROSSET T HASTIE “MULTICLASS ADABOOST” 2009

1048 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
 FROM SKLEARNENSEMBLE IMPORT ADABOOSTCLASSIFIER  
 NESTIMATORS 400  
 A LEARNING RATE OF 1 MAY NOT BE OPTIMAL FOR BOTH SAMME AND SAMMER  
 LEARNINGRATE 1  
 X Y DATASETSMAKEHASTIE102NSAMPLES12000 RANDOMSTATE1  
 XTEST YTEST X2000 Y2000  
 XTRAIN YTRAIN X2000 Y2000  
 DTSTUMP DECISIONTREECLASSIFIERMAXDEPTH1 MINSAMPLESLEAF1  
 DTSTUMPFITXTRAIN YTRAIN  
 DTSTUMPERR 10 DTSTUMPSCOREXTEST YTEST  
 DT DECISIONTREECLASSIFIERMAXDEPTH9 MINSAMPLESLEAF1  
 DTFITXTRAIN YTRAIN  
 DTERR 10 DTSCOREXTEST YTEST  
 ADADISCRETE ADABOOSTCLASSIFIER  
 BASEESTIMATORDTSTUMP  
 LEARNINGRATELEARNINGRATE  
 NESTIMATORSNESTIMATORS  
 ALGORITHMSAMME  
 ADADISCRETEFITXTRAIN YTRAIN  
 ADAREAL ADABOOSTCLASSIFIER  
 BASEESTIMATORDTSTUMP  
 LEARNINGRATELEARNINGRATE  
 NESTIMATORSNESTIMATORS  
 ALGORITHMSAMMER  
 ADAREALFITXTRAIN YTRAIN  
 FIG PLTFigure  
 AX FIGADDSUBPLOT111  
 AXPLOT1 NESTIMATORS DTSTUMPERR 2 K  
 LABELDECISION STUMP ERROR  
 AXPLOT1 NESTIMATORS DTERR 2 K  
 LABELDECISION TREE ERROR  
 ADADISCRETEERR NPZEROSNESTIMATORS  
 FORI YPRED INENUMERATEADADISCRETESTAGEDPREDICTXTEST  
 ADADISCRETEERRI ZEROONELOSSYPRED YTEST  
 ADADISCRETEERRTRAIN NPZEROSNESTIMATORS  
 FORI YPRED INENUMERATEADADISCRETESTAGEDPREDICTXTRAIN  
 ADADISCRETEERRTRAINI ZEROONELOSSYPRED YTRAIN  
 ADAREALERR NPZEROSNESTIMATORS  
 FORI YPRED INENUMERATEADAREALSTAGEDPREDICTXTEST  
 ADAREALERRI ZEROONELOSSYPRED YTEST  
 ADAREALERRTRAIN NPZEROSNESTIMATORS  
 FORI YPRED INENUMERATEADAREALSTAGEDPREDICTXTRAIN  
 ADAREALERRTRAINI ZEROONELOSSYPRED YTRAIN  
 512 ENSEMBLE METHODS 1049

SCIKITLEARN USER GUIDE RELEASE 0213  
AXPLOTNPARANGENESTIMATORS 1 ADADISCRETEERR  
LABELDISCRETE ADABOOST TEST ERROR  
COLORRED  
AXPLOTNPARANGENESTIMATORS 1 ADADISCRETEERRTRAIN  
LABELDISCRETE ADABOOST TRAIN ERROR  
COLORBLUE  
AXPLOTNPARANGENESTIMATORS 1 ADAREALERR  
LABELREAL ADABOOST TEST ERROR  
COLORORANGE  
AXPLOTNPARANGENESTIMATORS 1 ADAREALERRTRAIN  
LABELREAL ADABOOST TRAIN ERROR  
COLORGREEN  
AXSETYLIM00 05  
AXSETXLABELNESTIMATORS  
AXSETYLABELERROR RATE  
LEG AXLEGENDLOCUPPER RIGHT FANCYBOXTRUE  
LEGGETFRAMESETALPHA07  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 4579 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
51217 EARLY STOPPING OF GRADIENT BOOSTING  
GRADIENT BOOSTING IS AN ENSEMBLING TECHNIQUE WHERE SEVERAL WEAK LEARNERS REGRESSION TREES ARE COMBINED TO YIELD A  
POWERFUL SINGLE MODEL IN AN ITERATIVE FASHION  
EARLY STOPPING SUPPORT IN GRADIENT BOOSTING ENABLES US TO FIND THE LEAST NUMBER OF ITERATIONS WHICH IS SUFFICIENT TO  
BUILD A MODEL THAT GENERALIZES WELL TO UNSEEN DATA  
THE CONCEPT OF EARLY STOPPING IS SIMPLE WE SPECIFY A VALIDATIONFRACTION WHICH DENOTES THE FRACTION OF THE  
WHOLE DATASET THAT WILL BE KEPT ASIDE FROM TRAINING TO ASSESS THE VALIDATION LOSS OF THE MODEL THE GRADIENT BOOSTING  
MODEL IS TRAINED USING THE TRAINING SET AND EVALUATED USING THE VALIDATION SET WHEN EACH ADDITIONAL STAGE OF REGRESSION  
TREE IS ADDED THE VALIDATION SET IS USED TO SCORE THE MODEL THIS IS CONTINUED UNTIL THE SCORES OF THE MODEL IN THE LAST  
NITERNOCHANGE STAGES DO NOT IMPROVE BY ATLEAST TOL AFTER THAT THE MODEL IS CONSIDERED TO HAVE CONVERGED  
AND FURTHER ADDITION OF STAGES IS “STOPPED EARLY”  
THE NUMBER OF STAGES OF THE FINAL MODEL IS AVAILABLE AT THE ATTRIBUTE NESTIMATORS  
THIS EXAMPLE ILLUSTRATES HOW THE EARLY STOPPING CAN USED IN THE SKLEARNENSEMBLE  
GRADIENTBOOSTINGCLASSIFIER MODEL TO ACHIEVE ALMOST THE SAME ACCURACY AS COMPARED TO A MODEL  
BUILT WITHOUT EARLY STOPPING USING MANY FEWER ESTIMATORS THIS CAN SIGNIFICANTLY REDUCE TRAINING TIME MEMORY USAGE  
AND PREDICTION LATENCY  
AUTHORS VIGHNESH BIRODKAR VIGHNESHBIRODKARNYUEDU  
RAGHAV RV RVRAGHAV93GMAILCOM  
LICENSE BSD 3 CLAUSE  
IMPORT TIME  
IMPORT NUMPY AS NP  
1050 CHAPTER 5 EXAMPLES



```
SCIKITLEARN USER GUIDE RELEASE 0213
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARN IMPORT ENSEMBLE
FROM SKLEARN IMPORT DATASETS
FROM SKLEARNMODELSELECTION IMPORT TRAINTESTSPLIT
PRINTDOC
DATALIST DATASETSLOADIRIS DATASETSLOADDIGITS
DATALIST DDATA DTARGET FORDINDATALIST
DATALIST DATASETSMAKEHASTIE102
NAMES IRIS DATA DIGITS DATA HASTIE DATA
NGB
SCOREGB
TIMEGB
NGBES
SCOREGBES
TIMEGBES
NESTIMATORS 500
FORX YINDATALIST
XTRAIN XTEST YTRAIN YTEST TRAINTESTSPLITX Y TESTSIZE02
RANDOMSTATE0
WE SPECIFY THAT IF THE SCORES DONT IMPROVE BY ATLEAST 001 FOR THE LAST
10 STAGES STOP FITTING ADDITIONAL STAGES
GBES ENSEMBLEGRADIENTBOOSTINGCLASSIFIERNESTIMATORSNESTIMATORS
VALIDATIONFRACTION02
NITERNOCHANGE5 TOL001
RANDOMSTATE0
GB ENSEMBLEGRADIENTBOOSTINGCLASSIFIERNESTIMATORSNESTIMATORS
RANDOMSTATE0
START TIMETIME
GBFITXTRAIN YTRAIN
TIMEGBAPPENDTIMETIME START
START TIMETIME
GBESFITXTRAIN YTRAIN
TIMEGBESAPPENDTIMETIME START
SCOREGBAPPENDGBSCOREXTEST YTEST
SCOREGBESAPPENDGBESSCOREXTEST YTEST
NGBAPPENDGBNESTIMATORS
NGBESAPPENDGBESNESTIMATORS
BARWIDTH 02
N LENDATALIST
INDEX NPARANGE0 N BARWIDTH BARWIDTH 25
INDEX INDEX0N
512 ENSEMBLE METHODS 1051
```

SCIKITLEARN USER GUIDE RELEASE 0213  
COMPARE SCORES WITH AND WITHOUT EARLY STOPPING  
PLTFIGUREFIGSIZE9 5  
BAR1 PLTBARINDEX SCOREGB BARWIDTH LABELWITHOUT EARLY STOPPING  
COLORCRIMSON  
BAR2 PLTBARINDEX BARWIDTH SCOREGBES BARWIDTH  
LABELWITH EARLY STOPPING COLORCORAL  
PLXTICKSINDEX BARWIDTH NAMES  
PLTYTICKSNPARANGE0 13 01  
DEFAUTOLABELRECTS NESTIMATORS

ATTACH A TEXT LABEL ABOVE EACH BAR DISPLAYING NESTIMATORS OF EACH MODEL

FORI RECT INENUMERATERECTS  
PLTTEXTRECTGETX RECTGETWIDTH 2  
105RECTGETHEIGHT NEST D NESTIMATORSI  
HACENTER VABOTTOM  
AUTOLABELBAR1 NGB  
AUTOLABELBAR2 NGBES  
PLTYLIM0 13  
PLTLEGENDLOCBEST  
PLTGRIDTRUE  
PLTXLABELDATASETS  
PLTYLABELTEST SCORE  
PLTSHOW  
1052 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
COMPARE FIT TIMES WITH AND WITHOUT EARLY STOPPING  
PLTFIGUREFIGSIZE9 5  
BAR1 PLTBARINDEX TIMEGB BARWIDTH LABELWITHOUT EARLY STOPPING  
COLORCRIMSON  
BAR2 PLTBARINDEX BARWIDTH TIMEGBES BARWIDTH  
LABELWITH EARLY STOPPING COLORCORAL  
MAXY NPAMAXNPMAXIMUMTIMEGB TIMEGBES  
PLTXTICKSINDEX BARWIDTH NAMES  
PLTYTICKSNPLINSPACE0 13 MAXY 13  
AUTOLABELBAR1 NGB  
AUTOLABELBAR2 NGBES  
PLTYLIM0 13 MAXY  
PLTLEGENDLOCBEST  
PLTGRIDTRUE  
PLTXLABELDATASETS  
PLTYLABELFIT TIME  
PLTSHOW  
512 ENSEMBLE METHODS 1053

SCIKITLEARN USER GUIDE RELEASE 0213  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 15622 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
51218 FEATURE TRANSFORMATIONS WITH ENSEMBLES OF TREES  
TRANSFORM YOUR FEATURES INTO A HIGHER DIMENSIONAL SPARSE SPACE THEN TRAIN A LINEAR MODEL ON THESE FEATURES  
FIRST FIT AN ENSEMBLE OF TREES TOTALLY RANDOM TREES A RANDOM FOREST OR GRADIENT BOOSTED TREES ON THE TRAINING SET THEN  
EACH LEAF OF EACH TREE IN THE ENSEMBLE IS ASSIGNED A FIXED ARBITRARY FEATURE INDEX IN A NEW FEATURE SPACE THESE LEAF  
INDICES ARE THEN ENCODED IN A ONEHOT FASHION  
EACH SAMPLE GOES THROUGH THE DECISIONS OF EACH TREE OF THE ENSEMBLE AND ENDS UP IN ONE LEAF PER TREE THE SAMPLE IS  
ENCODED BY SETTING FEATURE VALUES FOR THESE LEAVES TO 1 AND THE OTHER FEATURE VALUES TO 0  
THE RESULTING TRANSFORMER HAS THEN LEARNED A SUPERVISED SPARSE HIGHDIMENSIONAL CATEGORICAL EMBEDDING OF THE DATA  
1054 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213

• AUTHOR TIM HEAD BETATIMGMAILCOM

LICENSE BSD 3 CLAUSE  
IMPORT NUMPY AS NP  
NPRANDOMSEED10  
IMPORT MATPLOTLIBPYPLOT AS PLT  
FROM SKLEARNDATASETS IMPORT MAKECLASSIFICATION  
FROM SKLEARNLINEARMODEL IMPORT LOGISTICREGRESSION  
FROM SKLEARNENSEMBLE IMPORT RANDOMTREESEMBEDDING RANDOMFORESTCLASSIFIER  
GRADIENTBOOSTINGCLASSIFIER  
FROM SKLEARNPREPROCESSING IMPORT ONEHOTENCODER  
FROM SKLEARNMODELSELECTION IMPORT TRAINTESTSPLIT  
FROM SKLEARNMETRICS IMPORT ROCCURVE  
FROM SKLEARNPIPELINE IMPORT MAKEPIPELINE  
NESTIMATOR 10  
X Y MAKECLASSIFICATIONNSAMPLES80000  
XTRAIN XTEST YTRAIN YTEST TRAINTESTSPLITX Y TESTSIZE05  
IT IS IMPORTANT TO TRAIN THE ENSEMBLE OF TREES ON A DIFFERENT SUBSET  
OF THE TRAINING DATA THAN THE LINEAR REGRESSION MODEL TO AVOID  
OVERFITTING IN PARTICULAR IF THE TOTAL NUMBER OF LEAVES IS  
SIMILAR TO THE NUMBER OF TRAINING SAMPLES  
XTRAIN XTRAINLR YTRAIN YTRAINLR TRAINTESTSPLIT  
1056 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

XTRAIN YTRAIN TESTSIZE05

UNSUPERVISED TRANSFORMATION BASED ON TOTALLY RANDOM TREES

RT RANDOMTREEEMBEDDINGMAXDEPTH3 NESTIMATORSNESTIMATOR

RANDOMSTATE0

RTLM LOGISTICREGRESSIONSOLVERLBFGS MAXITER1000

PIPELINE MAKEPIPELINERT RTLM

PIPELINEFITXTRAIN YTRAIN

YPREDRT PIPELINEPREDICTPROBAXTEST 1

FPRRTLM TPRRTLM ROCCURVEYTEST YPREDRT

SUPERVISED TRANSFORMATION BASED ON RANDOM FORESTS

RF RANDOMFORESTCLASSIFIERMAXDEPTH3 NESTIMATORSNESTIMATOR

RFENC ONEHOTENCODERCATEGORIESAUTO

RFLM LOGISTICREGRESSIONSOLVERLBFGS MAXITER1000

RFFITXTRAIN YTRAIN

RFENCFITRFAPPLYXTRAIN

RFLMFITRFENCTransformRFAPPLYXTRAINLR YTRAINLR

YPREDRFLM RFLMPREDICTPROBARFENCTransformRFAPPLYXTEST 1

FPRRFLM TPRRFLM ROCCURVEYTEST YPREDRFLM

SUPERVISED TRANSFORMATION BASED ON GRADIENT BOOSTED TREES

GRD GRADIENTBOOSTINGCLASSIFIERNESTIMATORSNESTIMATOR

GRDENC ONEHOTENCODERCATEGORIESAUTO

GRDLM LOGISTICREGRESSIONSOLVERLBFGS MAXITER1000

GRDFITXTRAIN YTRAIN

GRDENCFITGRDAPPLYXTRAIN 0

GRDLMFITGRDENCTransformGRDAPPLYXTRAINLR 0 YTRAINLR

YPREDGRDLM GRDLMpredictPROBA

GRDENCTransformGRDAPPLYXTEST 0 1

FPRGRDLM TPRGRDLM ROCCURVEYTEST YPREDGRDLM

THE GRADIENT BOOSTED MODEL BY ITSELF

YPREDGRD GRDPREDICTPROBAXTEST 1

FPRGRD TPRGRD ROCCURVEYTEST YPREDGRD

THE RANDOM FOREST MODEL BY ITSELF

YPREDRF RFPREDICTPROBAXTEST 1

FPRRF TPRRF ROCCURVEYTEST YPREDRF

PLTFigure1

PLTPLOT0 1 0 1 K

PLTPLOTfprRTLM TPRRTLM LABELRT LR

PLTPLOTfprRF TPRRF LABELRF

PLTPLOTfprRFLM TPRRFLM LABELRF LR

PLTPLOTfprGRD TPRGRD LABELGBT

PLTPLOTfprGRDLM TPRGRDLM LABELGBT LR

PLTXLABELFALSE POSITIVE RATE

PLTYLABELTRUE POSITIVE RATE

PLTTITLEROC CURVE

PLTLEGENDLOCBEST

PLTSHOW

PLTFigure2

PLTXLIM0 02

512 ENSEMBLE METHODS 1057

SCIKITLEARN USER GUIDE RELEASE 0213

PLTYLIM08 1

PLTPLOT0 1 0 1 K

PLTPLOTFPRRTLM TPRRTLM LABELRT LR

PLTPLOTFPRRF TPRRF LABELRF

PLTPLOTFPRRFLM TPRRFLM LABELRF LR

PLTPLOTFPRGRD TPRGRD LABELGBT

PLTPLOTFPRGRDLM TPRGRDLM LABELGBT LR

PLTXLABELFALSE POSITIVE RATE

PLTYLABELTRUE POSITIVE RATE

PLTTITLEROC CURVE ZOOMED IN AT TOP LEFT

PLTLEGENDLOCBEST

PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 2261 SECONDS

NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)

51219 GRADIENT BOOSTING OUTFBAG ESTIMATES

OUTFBAG OOB ESTIMATES CAN BE A USEFUL HEURISTIC TO ESTIMATE THE “OPTIMAL” NUMBER OF BOOSTING ITERATIONS OOB ESTIMATES ARE ALMOST IDENTICAL TO CROSSVALIDATION ESTIMATES BUT THEY CAN BE COMPUTED ONTHEFLY WITHOUT THE NEED FOR REPEATED MODEL FITTING OOB ESTIMATES ARE ONLY AVAILABLE FOR STOCHASTIC GRADIENT BOOSTING IE SUBSAMPLE 1

0 THE ESTIMATES ARE DERIVED FROM THE IMPROVEMENT IN LOSS BASED ON THE EXAMPLES NOT INCLUDED IN THE BOOTSTRAP SAMPLE THE SOCALLED OUTFBAG EXAMPLES THE OOB ESTIMATOR IS A PESSIMISTIC ESTIMATOR OF THE TRUE TEST LOSS BUT REMAINS A FAIRLY GOOD APPROXIMATION FOR A SMALL NUMBER OF TREES

THE FIGURE SHOWS THE CUMULATIVE SUM OF THE NEGATIVE OOB IMPROVEMENTS AS A FUNCTION OF THE BOOSTING ITERATION AS YOU CAN SEE IT TRACKS THE TEST LOSS FOR THE FIRST HUNDRED ITERATIONS BUT THEN DIVERGES IN A PESSIMISTIC WAY THE FIGURE ALSO SHOWS THE PERFORMANCE OF 3FOLD CROSS VALIDATION WHICH USUALLY GIVES A BETTER ESTIMATE OF THE TEST LOSS BUT IS COMPUTATIONALLY MORE DEMANDING

1058 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213  
OUT  
ACCURACY 06840  
PRINTDOC  
AUTHOR PETER PRETTENHOFER PETERPRETTENHOFERGMAILCOM

LICENSE BSD 3 CLAUSE  
IMPORT NUMPY AS NP  
IMPORT MATPLOTLIBPYPLOT AS PLT  
FROM SKLEARN IMPORT ENSEMBLE  
FROM SKLEARNMODELSELECTION IMPORT KFOLD  
FROM SKLEARNMODELSELECTION IMPORT TRAINTESTSPLIT  
FROM SCIPYSPECIAL IMPORT EXPIT  
GENERATE DATA ADAPTED FROM G RIDGEWAYS GBM EXAMPLE  
NSAMPLES 1000  
512 ENSEMBLE METHODS 1059

SCIKITLEARN USER GUIDE RELEASE 0213  
 RANDOMSTATE NPRANDOMRANDOMSTATE13  
 X1 RANDOMSTATEUNIFORMSIZENSAMPLES  
 X2 RANDOMSTATEUNIFORMSIZENSAMPLES  
 X3 RANDOMSTATERANDINT0 4 SIZENSAMPLES  
 P EXPITNPSIN3 X1 4 X2 X3  
 Y RANDOMSTATEBINOMIAL1 P SIZENSAMPLES  
 X NPCX1 X2 X3  
 X XASTYPENPFLOAT32  
 XTRAIN XTEST YTRAIN YTEST TRAINTESTSPLITX Y TESTSIZE05  
 RANDOMSTATE9  
 FIT CLASSIFIER WITH OUTFBAG ESTIMATES  
 PARAMS NESTIMATORS 1200 MAXDEPTH 3 SUBSAMPLE 05  
 LEARNINGRATE 001 MINSAMPLESLEAF 1 RANDOMSTATE 3  
 CLF ENSEMBLEGRADIENTBOOSTINGCLASSIFIER PARAMS  
 CLFFITXTRAIN YTRAIN  
 ACC CLFSCOREXTEST YTEST  
 PRINTACCURACY 4FFORMATAACC  
 NESTIMATORS PARAMSNESTIMATORS  
 X NPARANGENESTIMATORS 1  
 DEFHELDOUTSCORECLF XTEST YTEST  
 COMPUTE DEVIANCE SCORES ON XTEST AND YTEST  
 SCORE NPZEROSNESTIMATORS DTYPENPFLOAT64  
 FORI YPRED INENUMERATECLFSTAGEDDECISIONFUNCTIONXTEST  
 SCOREI CLFLOSSYTEST YPRED  
 RETURNSCORE  
 DEFCVESTIMATENSPLITSNONE  
 CV KFOLDNSPLITSNSPLITS  
 CVCLF ENSEMBLEGRADIENTBOOSTINGCLASSIFIER PARAMS  
 VALSCORES NPZEROSNESTIMATORS DTYPENPFLOAT64  
 FORTRAIN TEST INCVSPLITXTRAIN YTRAIN  
 CVCLFFITXTRAINTRAIN YTRAINTRAIN  
 VALSCORES HELDOUTSCORECVCLF XTRAINTEST YTRAINTEST  
 VALSCORES NSPLITS  
 RETURNVALSCORES  
 ESTIMATE BEST NESTIMATOR USING CROSSVALIDATION  
 CVSCORE CVESTIMATE3  
 COMPUTE BEST NESTIMATOR FOR TEST DATA  
 TESTSCORE HELDOUTSCORECLF XTEST YTEST  
 NEGATIVE CUMULATIVE SUM OF OOB IMPROVEMENTS  
 CUMSUM NPCUMSUMCLFOOBIMPROVEMENT  
 MIN LOSS ACCORDING TO OOB  
 OOBBESTITER XNPARGMINCUMSUM  
 1060 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
MIN LOSS ACCORDING TO TEST NORMALIZE SUCH THAT FIRST LOSS IS 0  
TESTSCORE TESTSCORE0  
TESTBESTITER XNPARGMINTESTSCORE  
MIN LOSS ACCORDING TO CV NORMALIZE SUCH THAT FIRST LOSS IS 0  
CVSCORE CVSCORE0  
CVBESTITER XNPARGMINCVSCORE  
COLOR BREW FOR THE THREE CURVES  
OOBCOLOR LISTMAP LAMBDAX X 2560 190 174 212  
TESTCOLOR LISTMAP LAMBDAX X 2560 127 201 127  
CVCOLOR LISTMAP LAMBDAX X 2560 253 192 134  
PLOT CURVES AND VERTICAL LINES FOR BEST ITERATIONS  
PLTPLOTX CUMSUM LABELOOB LOSS COLOROOBCOLOR  
PLTPLOTX TESTSCORE LABELTEST LOSS COLORTESTCOLOR  
PLTPLOTX CVSCORE LABELCV LOSS COLORCVCOLOR  
PLTAXVLINEXOOBBESTITER COLOROOBCOLOR  
PLTAXVLINEXTESTBESTITER COLORTESTCOLOR  
PLTAXVLINEXCVBESTITER COLORCVCOLOR  
ADD THREE VERTICAL LINES TO XTICKS  
XTICKS PLTXTICKS  
XTICKSPOS NPARRAYXTICKS0TOLIST  
OOBBESTITER CVBESTITER TESTBESTITER  
XTICKSLABEL NPARRAYLISTMAP LAMBDAT INTT XTICKS0  
OOB CV TEST  
IND NPARGSORTXTICKSPOS  
XTICKSPOS XTICKSPOSIND  
XTICKSLABEL XTICKSLABELIND  
PLTXTICKSXTICKSPOS XTICKSLABEL  
PLTLEGENDLOCUPPER RIGHT  
PLTYLABELNORMALIZED LOSS  
PLTXLABELNUMBER OF ITERATIONS  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 2759 SECONDS  
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE  
51220 SINGLE ESTIMATOR VERSUS BAGGING BIASVARIANCE DECOMPOSITION  
THIS EXAMPLE ILLUSTRATES AND COMPARES THE BIASVARIANCE DECOMPOSITION OF THE EXPECTED MEAN SQUARED ERROR OF A SINGLE ESTIMATOR AGAINST A BAGGING ENSEMBLE  
IN REGRESSION THE EXPECTED MEAN SQUARED ERROR OF AN ESTIMATOR CAN BE DECOMPOSED IN TERMS OF BIAS VARIANCE AND NOISE ON AVERAGE OVER DATASETS OF THE REGRESSION PROBLEM THE BIAS TERM MEASURES THE AVERAGE AMOUNT BY WHICH THE PREDICTIONS OF THE ESTIMATOR DIFFER FROM THE PREDICTIONS OF THE BEST POSSIBLE ESTIMATOR FOR THE PROBLEM IE THE BAYES MODEL THE VARIANCE TERM MEASURES THE VARIABILITY OF THE PREDICTIONS OF THE ESTIMATOR WHEN FIT OVER DIFFERENT INSTANCES LS OF THE PROBLEM FINALLY THE NOISE MEASURES THE IRREDUCIBLE PART OF THE ERROR WHICH IS DUE THE VARIABILITY IN THE DATA THE UPPER LEFT FIGURE ILLUSTRATES THE PREDICTIONS IN DARK RED OF A SINGLE DECISION TREE TRAINED OVER A RANDOM DATASET LS THE BLUE DOTS OF A TOY 1D REGRESSION PROBLEM IT ALSO ILLUSTRATES THE PREDICTIONS IN LIGHT RED OF OTHER SINGLE DECISION  
512 ENSEMBLE METHODS 1061

SCIKITLEARN USER GUIDE RELEASE 0213

TREES TRAINED OVER OTHER AND DIFFERENT RANDOMLY DRAWN INSTANCES LS OF THE PROBLEM INTUITIVELY THE VARIANCE TERM HERE CORRESPONDS TO THE WIDTH OF THE BEAM OF PREDICTIONS IN LIGHT RED OF THE INDIVIDUAL ESTIMATORS THE LARGER THE VARIANCE THE MORE SENSITIVE ARE THE PREDICTIONS FOR XTO SMALL CHANGES IN THE TRAINING SET THE BIAS TERM CORRESPONDS TO THE DIFFERENCE BETWEEN THE AVERAGE PREDICTION OF THE ESTIMATOR IN CYAN AND THE BEST POSSIBLE MODEL IN DARK BLUE ON THIS PROBLEM WE CAN THUS OBSERVE THAT THE BIAS IS QUITE LOW BOTH THE CYAN AND THE BLUE CURVES ARE CLOSE TO EACH OTHER WHILE THE VARIANCE IS LARGE THE RED BEAM IS RATHER WIDE

THE LOWER LEFT FIGURE PLOTS THE POINTWISE DECOMPOSITION OF THE EXPECTED MEAN SQUARED ERROR OF A SINGLE DECISION TREE IT CONFIRMS THAT THE BIAS TERM IN BLUE IS LOW WHILE THE VARIANCE IS LARGE IN GREEN IT ALSO ILLUSTRATES THE NOISE PART OF THE ERROR WHICH AS EXPECTED APPEARS TO BE CONSTANT AND AROUND 001

THE RIGHT FIGURES CORRESPOND TO THE SAME PLOTS BUT USING INSTEAD A BAGGING ENSEMBLE OF DECISION TREES IN BOTH FIGURES WE CAN OBSERVE THAT THE BIAS TERM IS LARGER THAN IN THE PREVIOUS CASE IN THE UPPER RIGHT FIGURE THE DIFFERENCE BETWEEN THE AVERAGE PREDICTION IN CYAN AND THE BEST POSSIBLE MODEL IS LARGER EG NOTICE THE OFFSET AROUND X2 IN THE LOWER RIGHT FIGURE THE BIAS CURVE IS ALSO SLIGHTLY HIGHER THAN IN THE LOWER LEFT FIGURE IN TERMS OF VARIANCE HOWEVER THE BEAM OF PREDICTIONS IS NARROWER WHICH SUGGESTS THAT THE VARIANCE IS LOWER INDEED AS THE LOWER RIGHT FIGURE CONFIRMS THE VARIANCE TERM IN GREEN IS LOWER THAN FOR SINGLE DECISION TREES OVERALL THE BIAS VARIANCE DECOMPOSITION IS THEREFORE NO LONGER THE SAME THE TRADEOFF IS BETTER FOR BAGGING AVERAGING SEVERAL DECISION TREES FIT ON BOOTSTRAP COPIES OF THE DATASET SLIGHTLY INCREASES THE BIAS TERM BUT ALLOWS FOR A LARGER REDUCTION OF THE VARIANCE WHICH RESULTS IN A LOWER OVERALL MEAN SQUARED ERROR COMPARE THE RED CURVES INT THE LOWER FIGURES THE SCRIPT OUTPUT ALSO CONFIRMS THIS INTUITION THE TOTAL ERROR OF THE BAGGING ENSEMBLE IS LOWER THAN THE TOTAL ERROR OF A SINGLE DECISION TREE AND THIS DIFFERENCE INDEED MAINLY STEMS FROM A REDUCED VARIANCE

FOR FURTHER DETAILS ON BIASVARIANCE DECOMPOSITION SEE SECTION 73 OF1

1T HASTIE R TIBSHIRANI AND J FRIEDMAN “ELEMENTS OF STATISTICAL LEARNING” SPRINGER 2009

1062 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
REFERENCES  
OUT  
TREE 00255 ERROR 00003 BIAS2 00152 VAR 00098 NOISE  
BAGGINGTREE 00196 ERROR 00004 BIAS2 00092 VAR 00098 NOISE  
PRINTDOC  
AUTHOR GILLES LOUPPE GLOUPPEGMAILCOM  
LICENSE BSD 3 CLAUSE  
IMPORT NUMPY AS NP  
IMPORT MATPLOTLIBPYPLOT AS PLT  
FROM SKLEARNENSEMBLE IMPORT BAGGINGREGRESSOR  
FROM SKLEARNTREE IMPORT DECISIONTREEREGRESSOR  
512 ENSEMBLE METHODS 1063

```
SCIKITLEARN USER GUIDE RELEASE 0213
SETTINGS
NREPEAT 50 NUMBER OF ITERATIONS FOR COMPUTING EXPECTATIONS
NTRAIN 50 SIZE OF THE TRAINING SET
NTEST 1000 SIZE OF THE TEST SET
NOISE 01 STANDARD DEVIATION OF THE NOISE
NPRANDOMSEED0
CHANGE THIS FOR EXPLORING THE BIASVARIANCE DECOMPOSITION OF OTHER
ESTIMATORS THIS SHOULD WORK WELL FOR ESTIMATORS WITH HIGH VARIANCE EG
DECISION TREES OR KNN BUT POORLY FOR ESTIMATORS WITH LOW VARIANCE EG
LINEAR MODELS
ESTIMATORS TREE DECISIONTREEREgressor
BAGGINGTREE BAGGINGREGRESSORDECISIONTREEREgressor
NESTIMATORS LENESTIMATORS
GENERATE DATA
DEFFX
X XRAVEL
RETURNNPXPX 2 15 NXPX 2 2
DEFGENERATENSAMPLES NOISE NREPEAT1
X NPRANDOMRANDNSAMPLES 10 5
X NPSORTX
IFNREPEAT 1
Y FX NPRANDOMNORMAL00 NOISE NSAMPLES
ELSE
Y NPZEROSNSAMPLES NREPEAT
FORIINRANGENREPEAT
Y I FX NPRANDOMNORMAL00 NOISE NSAMPLES
X XRESHAPENSAMPLES 1
RETURNX Y
XTRAIN
YTRAIN
FORIINRANGENREPEAT
X Y GENERATENSAMPLESNTRAIN NOISENOISE
XTRAINAPPENDX
YTRAINAPPENDY
XTEST YTEST GENERATENSAMPLESNTEST NOISENOISE NREPEATNREPEAT
PLTFIGUREFIGSIZE10 8
LOOP OVER ESTIMATORS TO COMPARE
FORN NAME ESTIMATOR INENUMERATEESTIMATORS
COMPUTE PREDICTIONS
YPREDICT NPZEROSNTEST NREPEAT
1064 CHAPTER 5 EXAMPLES
```

SCIKITLEARN USER GUIDE RELEASE 0213  
FORIINRANGENREPEAT  
ESTIMATORFITXTRAINI YTRAINI  
YPREDICT I ESTIMATORPREDICTXTEST  
BIAS2 VARIANCE NOISE DECOMPOSITION OF THE MEAN SQUARED ERROR  
YERROR NPZEROSNTEST  
FORIINRANGENREPEAT  
FORJINRANGENREPEAT  
YERROR YTEST J YPREDICT I 2  
YERROR NREPEAT NREPEAT  
YNOISE NPVARYTEST AXIS1  
YBIAS FXTEST NPMEANYPREDICT AXIS1 2  
YVAR NPVARYPREDICT AXIS1  
PRINT0 14F ERROR 24F BIAS2  
34F VAR 44F NOISEFORMATNAME  
NPMEANYERROR  
NPMEANYBIAS  
NPMEANYVAR  
NPMEANYNOISE  
PLOT FIGURES  
PLTSUBPLOT2 NESTIMATORS N 1  
PLTPLOTXTEST FXTEST B LABELFX  
PLTPLOTXTRAIN0 YTRAIN0 B LABELLS Y FXNOISE  
FORIINRANGENREPEAT  
IFI 0  
PLTPLOTXTEST YPREDICT I R LABELRYX  
ELSE  
PLTPLOTXTEST YPREDICT I R ALPHA005  
PLTPLOTXTEST NPMEANYPREDICT AXIS1 C  
LABELRMATHBBELS YX  
PLTXLIM5 5  
PLTTITLENAME  
IFN NESTIMATORS 1  
PLTLEGENDLOC11 5  
PLTSUBPLOT2 NESTIMATORS NESTIMATORS N 1  
PLTPLOTXTEST YERROR R LABELERRORX  
PLTPLOTXTEST YBIAS B LABELBIAS2X  
PLTPLOTXTEST YVAR G LABELVARIANCEX  
PLTPLOTXTEST YNOISE C LABELNOISEX  
PLTXLIM5 5  
PLTYLIM0 01  
IFN NESTIMATORS 1  
PLTLEGENDLOC11 5  
PLTSUBPLOTSADJUSTRIGHT75  
512 ENSEMBLE METHODS 1065

SCIKITLEARN USER GUIDE RELEASE 0213

PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0515 SECONDS

NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)

51221 PLOT THE DECISION SURFACES OF ENSEMBLES OF TREES ON THE IRIS DATASET

PLOT THE DECISION SURFACES OF FORESTS OF RANDOMIZED TREES TRAINED ON PAIRS OF FEATURES OF THE IRIS DATASET

THIS PLOT COMPARES THE DECISION SURFACES LEARNED BY A DECISION TREE CLASSIFIER FIRST COLUMN BY A RANDOM FOREST CLASSIFIER SECOND COLUMN BY AN EXTRA TREES CLASSIFIER THIRD COLUMN AND BY AN ADABOOST CLASSIFIER FOURTH COLUMN

IN THE FIRST ROW THE CLASSIFIERS ARE BUILT USING THE SEPAL WIDTH AND THE SEPAL LENGTH FEATURES ONLY ON THE SECOND ROW USING THE PETAL LENGTH AND SEPAL LENGTH ONLY AND ON THE THIRD ROW USING THE PETAL WIDTH AND THE PETAL LENGTH ONLY

IN DESCENDING ORDER OF QUALITY WHEN TRAINED OUTSIDE OF THIS EXAMPLE ON ALL 4 FEATURES USING 30 ESTIMATORS AND SCORED USING 10 FOLD CROSS VALIDATION WE SEE

EXTRATREESCLASSIFIER 095 SCORE

RANDOMFORESTCLASSIFIER 094 SCORE

ADABOOSTDECISIONTREEMAXDEPTH3 094 SCORE

DECISIONTREEMAXDEPTH NONE 094 SCORE

INCREASING MAXDEPTH FOR ADABOOST LOWERS THE STANDARD DEVIATION OF THE SCORES BUT THE AVERAGE SCORE DOES NOT IMPROVE

SEE THE CONSOLE'S OUTPUT FOR FURTHER DETAILS ABOUT EACH MODEL

IN THIS EXAMPLE YOU MIGHT TRY TO

1 VARY THE MAXDEPTH FOR THE DECISIONTREECLASSIFIER ANDADABOOSTCLASSIFIER

PERHAPS TRY MAXDEPTH3 FOR THEDECISIONTREECLASSIFIER ORMAXDEPTHNONE FOR ADABOOSTCLASSIFIER

2 VARYNESTIMATORS

IT IS WORTH NOTING THAT RANDOMFORESTS AND EXTRATREES CAN BE FITTED IN PARALLEL ON MANY CORES AS EACH TREE IS BUILT INDEPENDENTLY OF THE OTHERS ADABOOST'S SAMPLES ARE BUILT SEQUENTIALLY AND SO DO NOT USE MULTIPLE CORES

1066 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213  
OUT  
DECISIONTREE WITH FEATURES 0 1 HAS A SCORE OF 092666666666666666  
RANDOMFOREST WITH 30 ESTIMATORS WITH FEATURES 0 1 HAS A SCORE OF 092666666666666666  
EXTRATREES WITH 30 ESTIMATORS WITH FEATURES 0 1 HAS A SCORE OF 092666666666666666  
ADABOOST WITH 30 ESTIMATORS WITH FEATURES 0 1 HAS A SCORE OF 084  
DECISIONTREE WITH FEATURES 0 2 HAS A SCORE OF 099333333333333333  
RANDOMFOREST WITH 30 ESTIMATORS WITH FEATURES 0 2 HAS A SCORE OF 099333333333333333  
EXTRATREES WITH 30 ESTIMATORS WITH FEATURES 0 2 HAS A SCORE OF 099333333333333333  
ADABOOST WITH 30 ESTIMATORS WITH FEATURES 0 2 HAS A SCORE OF 099333333333333333  
DECISIONTREE WITH FEATURES 2 3 HAS A SCORE OF 099333333333333333  
RANDOMFOREST WITH 30 ESTIMATORS WITH FEATURES 2 3 HAS A SCORE OF 099333333333333333  
EXTRATREES WITH 30 ESTIMATORS WITH FEATURES 2 3 HAS A SCORE OF 099333333333333333  
ADABOOST WITH 30 ESTIMATORS WITH FEATURES 2 3 HAS A SCORE OF 099333333333333333  
PRINTDOC  
IMPORT NUMPY AS NP  
IMPORT MATPLOTLIBPYPLOT AS PLT  
FROM MATPLOTLIBCOLORS IMPORT LISTEDCOLORMAP  
512 ENSEMBLE METHODS 1067

SCIKITLEARN USER GUIDE RELEASE 0213  
FROM SKLEARNDATASETS IMPORT LOADIRIS  
FROM SKLEARNENSEMBLE IMPORT RANDOMFORESTCLASSIFIER EXTRATREESCLASSIFIER  
ADABOOSTCLASSIFIER  
FROM SKLEARNTREE IMPORT DECISIONTREECLASSIFIER  
PARAMETERS  
NCLASSES 3  
NESTIMATORS 30  
CMAP PLTCMRDYLBU  
PLOTSTEP 002 FINE STEP WIDTH FOR DECISION SURFACE CONTOURS  
PLOTSTEPCOARSER 05 STEP WIDTHS FOR COARSE CLASSIFIER GUESSES  
RANDOMSEED 13 FIX THE SEED ON EACH ITERATION  
LOAD DATA  
IRIS LOADIRIS  
PLOTIDX 1  
MODELS DECISIONTREECLASSIFIERMAXDEPTHNONE  
RANDOMFORESTCLASSIFIERNESTIMATORSNESTIMATORS  
EXTRATREESCLASSIFIERNESTIMATORSNESTIMATORS  
ADABOOSTCLASSIFIERDECISIONTREECLASSIFIERMAXDEPTH3  
NESTIMATORSNESTIMATORS  
FORPAIRIN0 1 0 2 2 3  
FORMODELINMODELS  
WE ONLY TAKE THE TWO CORRESPONDING FEATURES  
X IRISDATA PAIR  
Y IRISTARGET  
SHUFFLE  
IDX NPARANGEXSHAPE0  
NPRANDOMSEEDRANDOMSEED  
NPRANDOMSHUFFLEIDX  
X XIDX  
Y YIDX  
STANDARDIZE  
MEAN XMEANAXIS0  
STD XSTDAXIS0  
X X MEAN STD  
TRAIN  
MODELFITX Y  
SCORES MODELSCOREX Y  
CREATE A TITLE FOR EACH COLUMN AND THE CONSOLE BY USING STR AND  
SLICING AWAY USELESS PARTS OF THE STRING  
MODELTITLE STRTYPEMODELSPLIT  
12LENCLASSIFIER  
MODELDETAILS MODELTITLE  
IFHASATTRMODEL ESTIMATORS  
MODELDETAILS WITH ESTIMATORSFORMAT  
LENMODELESTIMATORS  
PRINTMODELDETAILS WITH FEATURES PAIR  
HAS A SCORE OF SCORES  
1068 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
PLTSUBPLOT3 4 PLOTIDX  
IFPLOTIDX LENMODELS  
ADD A TITLE AT THE TOP OF EACH COLUMN  
PLTTITLEMODELTITLE FONTSIZE9  
NOW PLOT THE DECISION BOUNDARY USING A FINE MESH AS INPUT TO A  
FILLED CONTOUR PLOT  
XMIN XMAX X 0MIN 1 X 0MAX 1  
YMIN YMAX X 1MIN 1 X 1MAX 1  
XX YY NPMESHGRIDNPARANGEXMIN XMAX PLOTSTEP  
NPARANGEYMIN YMAX PLOTSTEP  
PLOT EITHER A SINGLE DECISIONTREECLASSIFIER OR ALPHA BLEND THE  
DECISION SURFACES OF THE ENSEMBLE OF CLASSIFIERS  
IFISINSTANCEMODEL DECISIONTREECLASSIFIER  
Z MODELPREDICTNPCXXRAVEL YYRAVEL  
Z ZRESHAPEXXSHAPE  
CS PLTCONTOURFXX YY Z CMAPCMAP  
ELSE  
CHOOSE ALPHA BLEND LEVEL WITH RESPECT TO THE NUMBER  
OF ESTIMATORS  
THAT ARE IN USE NOTING THAT ADABOOST CAN USE FEWER ESTIMATORS  
THAN ITS MAXIMUM IF IT ACHIEVES A GOOD ENOUGH FIT EARLY ON  
ESTIMATORALPHA 10 LENMODELESTIMATORS  
FORTREEINMODELESTIMATORS  
Z TREEPREDICTNPCXXRAVEL YYRAVEL  
Z ZRESHAPEXXSHAPE  
CS PLTCONTOURFXX YY Z ALPHAESTIMATORALPHA CMAPCMAP  
BUILD A COARSER GRID TO PLOT A SET OF ENSEMBLE CLASSIFICATIONS  
TO SHOW HOW THESE ARE DIFFERENT TO WHAT WE SEE IN THE DECISION  
SURFACES THESE POINTS ARE REGULARLY SPACE AND DO NOT HAVE A  
BLACK OUTLINE  
XXCOARSER YYCOARSER NPMESHGRID  
NPARANGEXMIN XMAX PLOTSTEPCOARSER  
NPARANGEYMIN YMAX PLOTSTEPCOARSER  
ZPOINTSCOARSER MODELPREDICTNPCXXCOARSERRAVEL  
YYCOARSERRAVEL  
RESHAPEXXCOASERSHAPE  
CSPOINTS PLTSCATTERXXCOARSER YYCOARSER S15  
CZPOINTSCOARSER CMAPCMAP  
EDGECOLORSNONE  
PLOT THE TRAINING POINTS THESE ARE CLUSTERED TOGETHER AND HAVE A  
BLACK OUTLINE  
PLTSCATTERX 0 X 1 CY  
CMAPLISTEDCOLORMAPR Y B  
EDGECOLORK S20  
PLOTIDX 1 MOVE ON TO THE NEXT PLOT IN SEQUENCE  
PLTSUPTITLECLASSIFIERS ON FEATURE SUBSETS OF THE IRIS DATASET FONTSIZE12  
PLTAXISTIGHT  
PLTTIGHTLAYOUTHPAD02 WPAD02 PAD25  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 6177 SECONDS  
512 ENSEMBLE METHODS 1069

SCIKITLEARN USER GUIDE RELEASE 0213  
513 TUTORIAL EXERCISES  
EXERCISES FOR THE TUTORIALS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
5131 DIGITS CLASSIFICATION EXERCISE  
A TUTORIAL EXERCISE REGARDING THE USE OF CLASSIFICATION TECHNIQUES ON THE DIGITS DATASET  
THIS EXERCISE IS USED IN THE CLASSIFICATION PART OF THE SUPERVISED LEARNING PREDICTING AN OUTPUT VARIABLE FROM HIGH  
DIMENSIONAL OBSERVATIONS SECTION OF THE A TUTORIAL ON STATISTICALLEARNING FOR SCIENTIFIC DATA PROCESSING  
OUT  
KNN SCORE 0961111  
LOGISTICREGRESSION SCORE 0933333  
PRINTDOC  
FROM SKLEARN IMPORT DATASETS NEIGHBORS LINEARMODEL  
DIGITS DATASETSLOADDIGITS  
XDIGITS DIGITSDATA DIGITSDATAMAX  
YDIGITS DIGITSTARGET  
NSAMPLES LENXDIGITS  
XTRAIN XDIGITSINT9 NSAMPLES  
YTRAIN YDIGITSINT9 NSAMPLES  
XTEST XDIGITSINT9 NSAMPLES  
YTEST YDIGITSINT9 NSAMPLES  
KNN NEIGHBORSKNEIGHBORSCLASSIFIER  
LOGISTIC LINEARMODELLOGISTICREGRESSIONSOLVERLBFSGS MAXITER1000  
MULTICLASSMULTINOMIAL  
PRINTKNN SCORE F KNNFITXTRAIN YTRAINSCOREXTEST YTEST  
PRINTLOGISTICREGRESSION SCORE F  
LOGISTICFITXTRAIN YTRAINSCOREXTEST YTEST  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0432 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
1070 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
5132 CROSSVALIDATION ON DIGITS DATASET EXERCISE
A TUTORIAL EXERCISE USING CROSSVALIDATION WITH AN SVM ON THE DIGITS DATASET
THIS EXERCISE IS USED IN THE CROSSVALIDATION GENERATORS PART OF THE MODEL SELECTION CHOOSING ESTIMATORS AND THEIR
PARAMETERS SECTION OF THE A TUTORIAL ON STATISTICALLEARNING FOR SCIENTIFIC DATA PROCESSING
PRINTDOC
IMPORT NUMPY AS NP
FROM SKLEARNMODELSELECTION IMPORT CROSSVALSCORE
FROM SKLEARN IMPORT DATASETS SVM
DIGITS DATASETSLOADDIGITS
X DIGITSDATA
Y DIGITSTARGET
SVC SVMKVCKERNELLINER
CS NPLOGSPACE10 0 10
SCORES LIST
SCORESSTD LIST
FORCINCS
SVCC C
THISSCORES CROSSVALSCORESVC X Y CV5 NJOBS1
513 TUTORIAL EXERCISES 1071
```

```
SCIKITLEARN USER GUIDE RELEASE 0213
SCORESAPPENDNPMEANTHISSCORES
SCORESSTDAPPENDNPSTDTHISSCORES
DO THE PLOTTING
IMPORT MATPLOTLIBPYPLOT AS PLT
PLTFigure
PLTSEMILOGXCS SCORES
PLTSEMILOGXCS NPARRAYSCORES NPARRAYSCORESSTD B
PLTSEMILOGXCS NPARRAYSCORES NPARRAYSCORESSTD B
LOCS LABELS PLTYTICKS
PLTYTICKSLOCS LISTMAP LAMBDA X G X LOCS
PLTYLABELCV SCORE
PLTXLABELPARAMETER C
PLTYLIM0 11
PLTSHOW
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 8826 SECONDS
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE
5133 SVM EXERCISE
A TUTORIAL EXERCISE FOR USING DIFFERENT SVM KERNELS
THIS EXERCISE IS USED IN THE USING KERNELS PART OF THE SUPERVISED LEARNING PREDICTING AN OUTPUT VARIABLE FROM HIGH
DIMENSIONAL OBSERVATIONS SECTION OF THE A TUTORIAL ON STATISTICAL LEARNING FOR SCIENTIFIC DATA PROCESSING
1072 CHAPTER 5 EXAMPLES
```

SCIKITLEARN USER GUIDE RELEASE 0213

- 513 TUTORIAL EXERCISES 1073





SCIKITLEARN USER GUIDE RELEASE 0213

•

```
PRINTDOC
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARN IMPORT DATASETS SVM
IRIS DATASETSLOADIRIS
X IRISDATA
Y IRISTARGET
X XY 0 2
Y YY 0
NSAMPLE LENX
NPRANDOMSEED0
ORDER NPRANDOMPERMUTATIONNSAMPLE
X XORDER
Y YORDERASTYPENPFLOAT
XTRAIN XINT9 NSAMPLE
YTRAIN YINT9 NSAMPLE
XTEST XINT9 NSAMPLE
YTEST YINT9 NSAMPLE
FIT THE MODEL
513 TUTORIAL EXERCISES 1075
```

SCIKITLEARN USER GUIDE RELEASE 0213  
FORKERNELINLINEAR RBF POLY  
CLF SVMSVCKERNELKERNEL GAMMA10  
CLFFITXTRAIN YTRAIN  
PLTFigure  
PLTCLF  
PLTSCATTERX 0 X 1 CY ZORDER10 CMAPPLTCMPAIED  
EDGECOLORK S20  
CIRCLE OUT THE TEST DATA  
PLTSCATTERXTEST 0 XTEST 1 S80 FACECOLORSNONE  
ZORDER10 EDGECOLORK  
PLTAXISTIGHT  
XMIN X OMIN  
XMAX X OMAX  
YMIN X 1MIN  
YMAX X 1MAX  
XX YY NPMGRIDXMINXMAX200J YMINYMAX200J  
Z CLFDECISIONFUNCTIONNPCXXRAVEL YYRAVEL  
PUT THE RESULT INTO A COLOR PLOT  
Z ZRESHAPEXXSHAPE  
PLTPCOLORMESHXX YY Z 0 CMAPPLTCMPAIED  
PLTCONTOURXX YY Z COLORSK K K  
LINESTYLES LEVELS5 0 5  
PLTTITLEKERNEL  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 5320 SECONDS  
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE  
5134 CROSSVALIDATION ON DIABETES DATASET EXERCISE  
A TUTORIAL EXERCISE WHICH USES CROSSVALIDATION WITH LINEAR MODELS  
THIS EXERCISE IS USED IN THE CROSSVALIDATED ESTIMATORS PART OF THE MODEL SELECTION CHOOSING ESTIMATORS AND THEIR  
PARAMETERS SECTION OF THE A TUTORIAL ON STATISTICALLEARNING FOR SCIENTIFIC DATA PROCESSING  
1076 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
OUT  
ANSWER TO THE BONUS QUESTION HOW MUCH CAN YOU TRUST THE SELECTION OF ALPHA  
ALPHA PARAMETERS MAXIMISING THE GENERALIZATION SCORE ON DIFFERENT  
SUBSETS OF THE DATA  
FOLD 0 ALPHA 005968 SCORE 054209  
FOLD 1 ALPHA 004520 SCORE 015523  
FOLD 2 ALPHA 007880 SCORE 045193  
ANSWER NOT VERY MUCH SINCE WE OBTAINED DIFFERENT ALPHAS FOR DIFFERENT  
SUBSETS OF THE DATA AND MOREOVER THE SCORES FOR THESE ALPHAS DIFFER  
QUITE SUBSTANTIALLY  
PRINTDOC  
IMPORT NUMPY AS NP  
IMPORT MATPLOTLIBPYPLOT AS PLT  
FROM SKLEARN IMPORT DATASETS  
FROM SKLEARNLINEARMODEL IMPORT LASSOCV  
513 TUTORIAL EXERCISES 1077

```

SCIKITLEARN USER GUIDE RELEASE 0213
FROM SKLEARNLINEARMODEL IMPORT LASSO
FROM SKLEARNMODELSELECTION IMPORT KFOLD
FROM SKLEARNMODELSELECTION IMPORT GRIDSEARCHCV
DIABETES  DATASETSLOADDIABETES
X  DIABETESDATA150
Y  DIABETESTARGET150
LASSO  LASSORANDOMSTATE0 MAXITER10000
ALPHAS  NPLOGSPACE4 05 30
TUNEDPARAMETERS  ALPHA ALPHAS
NFOLDS  5
CLF  GRIDSEARCHCVLASSO TUNEDPARAMETERS CVNFOLDS REFITFALSE
CLFFITX Y
SCORES  CLFCVRESULTSMEANTESTSCORE
SCORESSTD  CLFCVRESULTSSTDTESTSCORE
PLTFIGURESETSIZEINCHES8 6
PLTSEMILOGXALPHAS SCORES
  PLOT ERROR LINES SHOWING  STD ERRORS OF THE SCORES
STDERROR  SCORESSTD  NPSQRTNFOLDS
PLTSEMILOGXALPHAS SCORES  STDERROR B
PLTSEMILOGXALPHAS SCORES  STDERROR B
  ALPHA02 CONTROLS THE TRANSLUCENCY OF THE FILL COLOR
PLTFILLBETWEENALPHAS SCORES  STDERROR SCORES  STDERROR ALPHA02
PLTYLABELCV SCORE  STD ERROR
PLTXLABELALPHA
PLTAXHLINENPMAXScores LINESTYLE COLOR5
PLTXLIMALPHAS0 ALPHAS1

```

```

BONUS HOW MUCH CAN YOU TRUST THE SELECTION OF ALPHA
TO ANSWER THIS QUESTION WE USE THE LASSOCV OBJECT THAT SETS ITS ALPHA
PARAMETER AUTOMATICALLY FROM THE DATA BY INTERNAL CROSSVALIDATION IE IT
PERFORMS CROSSVALIDATION ON THE TRAINING DATA IT RECEIVES
WE USE EXTERNAL CROSSVALIDATION TO SEE HOW MUCH THE AUTOMATICALLY OBTAINED
ALPHAS DIFFER ACROSS DIFFERENT CROSSVALIDATION FOLDS
LASSOCV  LASSOCVALPHASALPHAS CV5 RANDOMSTATE0 MAXITER10000
KFOLD  KFOLD3
PRINTANSWER TO THE BONUS QUESTION
HOW MUCH CAN YOU TRUST THE SELECTION OF ALPHA
PRINT
PRINTALPHA PARAMETERS MAXIMISING THE GENERALIZATION SCORE ON DIFFERENT
PRINTSUBSETS OF THE DATA
FORK TRAIN TEST INENUMERATEKFOLDSPLITX Y
LASSOCVFITXTRAIN YTRAIN
PRINTFOLD 0 ALPHA 15F SCORE 25F
FORMATK LASSOCVALPHA LASSOCVSCOREXTEST YTEST
PRINT
PRINTANSWER NOT VERY MUCH SINCE WE OBTAINED DIFFERENT ALPHAS FOR DIFFERENT
PRINTSUBSETS OF THE DATA AND MOREOVER THE SCORES FOR THESE ALPHAS DIFFER
1078 CHAPTER 5 EXAMPLES

```

SCIKITLEARN USER GUIDE RELEASE 0213  
PRINTQUITE SUBSTANTIALY  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0294 SECONDS  
514 FEATURE SELECTION  
EXAMPLES CONCERNING THE SKLEARNFEATURESELECTION MODULE  
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE  
5141 RECURSIVE FEATURE ELIMINATION  
A RECURSIVE FEATURE ELIMINATION EXAMPLE SHOWING THE RELEVANCE OF PIXELS IN A DIGIT CLASSIFICATION TASK  
NOTE SEE ALSO RECURSIVE FEATURE ELIMINATION WITH CROSSVALIDATION  
514 FEATURE SELECTION 1079

SCIKITLEARN USER GUIDE RELEASE 0213

PRINTDOC

```
FROM SKLEARN SVM IMPORT SVC
FROM SKLEARN DATASETS IMPORT LOADDIGITS
FROM SKLEARN FEATURE SELECTION IMPORT RFE
IMPORT MATPLOTLIB PYPLOT AS PLT
LOAD THE DIGITS DATASET
DIGITS LOADDIGITS
X DIGITS IMAGES RESHAPE LE NDIGITS IMAGES 1
Y DIGIT TARGET
CREATE THE RFE OBJECT AND RANK EACH PIXEL
SVC SVCKERNEL LINEAR C1
RFE RFE ESTIMATOR SVC NFEATURES TO SELECT 1 STEP 1
RFE FIT X Y
RANKING RFE RANKING RESHAPED DIGITS IMAGES 0 SHAPE
PLOT PIXEL RANKING
PLT MAT SHOW RANKING CM APPLT CM BLUES
PLT COLOR BAR
PLT TITLE RANKING OF PIXELS WITH RFE
PLT SHOW
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 3436 SECONDS
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE
5142 COMPARISON OF FTEST AND MUTUAL INFORMATION
THIS EXAMPLE ILLUSTRATES THE DIFFERENCES BETWEEN UNIVARIATE FTEST STATISTICS AND MUTUAL INFORMATION
WE CONSIDER 3 FEATURES X1 X2 X3 DISTRIBUTED UNIFORMLY OVER 0 1 THE TARGET DEPENDS ON THEM AS FOLLOWS
Y X1 SIN 6 PI X2 01 NO 1 THAT IS THE THIRD FEATURES IS COMPLETELY IRRELEVANT
THE CODE BELOW PLOTS THE DEPENDENCY OF Y AGAINST INDIVIDUAL XI AND NORMALIZED VALUES OF UNIVARIATE FTESTS STATISTICS
AND MUTUAL INFORMATION
AS FTEST CAPTURES ONLY LINEAR DEPENDENCY IT RATES X1 AS THE MOST DISCRIMINATIVE FEATURE ON THE OTHER HAND MUTUAL
INFORMATION CAN CAPTURE ANY KIND OF DEPENDENCY BETWEEN VARIABLES AND IT RATES X2 AS THE MOST DISCRIMINATIVE FEATURE
WHICH PROBABLY AGREES BETTER WITH OUR INTUITIVE PERCEPTION FOR THIS EXAMPLE BOTH METHODS CORRECTLY MARKS X3 AS
IRRELEVANT
1080 CHAPTER 5 EXAMPLES
```

```
SCIKITLEARN USER GUIDE RELEASE 0213
PRINTDOC
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARNFEATURESELECTION IMPORT FREGRESSION MUTUALINFOREGRESSION
NPRANDOMSEED0
X NPRANDOMRAND1000 3
Y X 0 NPSIN6 NPPIX 1 01 NPRANDOMRANDN1000
FTEST FREGRESSIONX Y
FTEST NPMAXFTEST
MI MUTUALINFOREGRESSIONX Y
MI NPMAXMI
PLTFIGUREFIGSIZE15 5
FORIINRANGE3
PLTSUBPLOT1 3 I 1
PLTSCATTERX I Y EDGECOLORBLACK S20
PLTXLABELXFORMATI 1 FONTSIZE14
IFI 0
PLTYLABELY FONTSIZE14
PLTTITLEFTEST2F MI2FFORMATFTESTI MII
FONTSIZE16
PLTSHOW
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0062 SECONDS
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE
5143 PIPELINE ANOVA SVM
SIMPLE USAGE OF PIPELINE THAT RUNS SUCCESSIVELY A UNIVARIATE FEATURE SELECTION WITH ANOVA AND THEN A SVM OF THE
SELECTED FEATURES
USING A SUBPIPELINE THE FITTED COEFFICIENTS CAN BE MAPPED BACK INTO THE ORIGINAL FEATURE SPACE
OUT
514 FEATURE SELECTION 1081
```

SCIKITLEARN USER GUIDE RELEASE 0213  
 PRECISION RECALL F1SCORE SUPPORT  
 0 075 050 060 6  
 1 067 100 080 6  
 2 067 080 073 5  
 3 100 075 086 8  
 ACCURACY 076 25  
 MACRO AVG 077 076 075 25  
 WEIGHTED AVG 079 076 076 25  
 023912131 0 0 0 03236911 0  
 0 0 0 0 0 0  
 010836648 0 0 0 0 0  
 0 0  
 043878747 0 0 0 051415652 0  
 0 0 0 0 0 0  
 004845652 0 0 0 0 0  
 0 0  
 065382998 0 0 0 057962856 0  
 0 0 0 0 0 0  
 004736524 0 0 0 0 0  
 0 0  
 054403412 0 0 0 058478491 0  
 0 0 0 0 0 0  
 011344659 0 0 0 0 0  
 0 0  
 FROM SKLEARN IMPORT SVM  
 FROM SKLEARNDATASETS IMPORT SAMPLESGENERATOR  
 FROM SKLEARNFEATURESELECTION IMPORT SELECTKBEST FREGRESSION  
 FROM SKLEARNPIPELINE IMPORT MAKEPIPELINE  
 FROM SKLEARNMODELSELECTION IMPORT TRAINTESTSPLIT  
 FROM SKLEARNMETRICS IMPORT CLASSIFICATIONREPORT  
 PRINTDOC  
 IMPORT SOME DATA TO PLAY WITH  
 X Y SAMPLESGENERATORMAKECLASSIFICATION  
 NFEATURES20 NINFORMATIVE3 NREDUNDANT0 NCLASSES4  
 NCLUSTERSPERCLASS2  
 XTRAIN XTEST YTRAIN YTEST TRAINTESTSPLITX Y RANDOMSTATE42  
 ANOVA SVMC  
 1 ANOVA FILTER TAKE 3 BEST RANKED FEATURES  
 ANOVAFILTER SELECTKBESTFREGRESSION K3  
 2 SVM  
 CLF SVMLINEARSVC  
 ANOVASVM MAKEPIPELINEANOVAFILTER CLF  
 ANOVASVMFITXTRAIN YTRAIN  
 YPRED ANOVASVMPREDICTXTEST  
 1082 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213  
PRINTCLASSIFICATIONREPORTYTEST YPRED  
COEF ANOVASVM1INVERSETRANSFORMANOVASVMLINEARSVCCOEF  
PRINTCOEF  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0008 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
5144 RECURSIVE FEATURE ELIMINATION WITH CROSSVALIDATION  
A RECURSIVE FEATURE ELIMINATION EXAMPLE WITH AUTOMATIC TUNING OF THE NUMBER OF FEATURES SELECTED WITH CROSSVALIDATION  
OUT  
OPTIMAL NUMBER OF FEATURES 3  
514 FEATURE SELECTION 1083

SCIKITLEARN USER GUIDE RELEASE 0213

PRINTDOC

IMPORT MATPLOTLIBPYPLOT AS PLT

FROM SKLEARN SVM IMPORT SVC

FROM SKLEARNMODELSELECTION IMPORT STRATIFIEDKFOLD

FROM SKLEARNFEATURESELECTION IMPORT RFECV

FROM SKLEARNDATASETS IMPORT MAKECLASSIFICATION

BUILD A CLASSIFICATION TASK USING 3 INFORMATIVE FEATURES

X Y MAKECLASSIFICATIONNNSAMPLES1000 NFEATURES25 NINFORMATIVE3

NREDUNDANT2 NREPEATED0 NCLASSES8

NCLUSTERSPERCLASS1 RANDOMSTATE0

CREATE THE RFE OBJECT AND COMPUTE A CROSSVALIDATED SCORE

SVC SVCKERNELLINEAR

THE ACCURACY SCORING IS PROPORTIONAL TO THE NUMBER OF CORRECT CLASSIFICATIONS

RFECV RFECVESTIMATORSVC STEP1 CVSTRATIFIEDKFOLD2

SCORINGACCURACY

RFECVFITX Y

PRINTOPTIMAL NUMBER OF FEATURES D RFECVNFEATURES

PLOT NUMBER OF FEATURES VS CROSSVALIDATION SCORES

PLTFigure

PLTXLABELNUMBER OF FEATURES SELECTED

PLTYLABELCROSS VALIDATION SCORE NB OF CORRECT CLASSIFICATIONS

PLTPLOT RANGE1 LENRFECVGRIDScores 1 RFECVGRIDScores

PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 1806 SECONDS

NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)

5145 FEATURE SELECTION USING SELECTFROMMODEL AND LASSOCV

USE SELECTFROMMODEL METATransformer ALONG WITH LASSO TO SELECT THE BEST COUPLE OF FEATURES FROM THE BOSTON DATASET

1084 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
AUTHOR MANOJ KUMAR MKS542NYUEDU
LICENSE BSD 3 CLAUSE
PRINTDOC
IMPORT MATPLOTLIBPYPLOT AS PLT
IMPORT NUMPY AS NP
FROM SKLEARNDATASETS IMPORT LOADBOSTON
FROM SKLEARNFEATURESELECTION IMPORT SELECTFROMMODEL
FROM SKLEARNLINEARMODEL IMPORT LASSOCV
LOAD THE BOSTON DATASET
BOSTON LOADBOSTON
X Y BOSTONDATA BOSTONTARGET
WE USE THE BASE ESTIMATOR LASSOCV SINCE THE L1 NORM PROMOTES SPARSITY OF FEATURES
CLF LASSOCVCV5
SET A MINIMUM THRESHOLD OF 025
SFM SELECTFROMMODELCLF THRESHOLD025
SFMFITX Y
NFEATURES SFMTRANSFORMXSHAPE1
RESET THE THRESHOLD TILL THE NUMBER OF FEATURES EQUALS TWO
NOTE THAT THE ATTRIBUTE CAN BE SET DIRECTLY INSTEAD OF REPEATEDLY
514 FEATURE SELECTION 1085
```

SCIKITLEARN USER GUIDE RELEASE 0213  
FITTING THE METATRANSFORMER  
WHILENFEATURES 2  
SFMTHRESHOLD 01  
XTRANSFORM SFMTRANSFORMX  
NFEATURES XTRANSFORMSHAPE1  
PLOT THE SELECTED TWO FEATURES FROM X  
PLTTITLE  
FEATURES SELECTED FROM BOSTON USING SELECTFROMMODEL WITH  
THRESHOLD 03F SFMTHRESHOLD  
FEATURE1 XTRANSFORM 0  
FEATURE2 XTRANSFORM 1  
PLTPLOTFEATURE1 FEATURE2 R  
PLTXLABELFEATURE NUMBER 1  
PLTYLABELFEATURE NUMBER 2  
PLTYLIMNPMINFEATURE2 NPMAXFEATURE2  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0056 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
5146 TEST WITH PERMUTATIONS THE SIGNIFICANCE OF A CLASSIFICATION SCORE  
IN ORDER TO TEST IF A CLASSIFICATION SCORE IS SIGNIFICATIVE A TECHNIQUE IN REPEATING THE CLASSIFICATION PROCEDURE AFTER RAN  
DOMIZING PERMUTING THE LABELS THE PVALUE IS THEN GIVEN BY THE PERCENTAGE OF RUNS FOR WHICH THE SCORE OBTAINED IS  
GREATER THAN THE CLASSIFICATION SCORE OBTAINED IN THE FIRST PLACE  
1086 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
OUT  
CLASSIFICATION SCORE 051333333333333333 PVALUE 0009900990099009901  
AUTHOR ALEXANDRE GRAMFORT ALEXANDREGRAMFORTINRIA.FR  
LICENSE BSD 3 CLAUSE  
PRINTDOC  
IMPORT NUMPY AS NP  
IMPORT MATPLOTLIBPY.PLOT AS PLT  
FROM SKLEARN.SVM IMPORT SVC  
FROM SKLEARN.MODEL.SELECTION IMPORT STRATIFIEDKFOLD  
FROM SKLEARN.MODEL.SELECTION IMPORT PERMUTATIONTESTSCORE  
FROM SKLEARN IMPORT DATASETS  
  
LOADING A DATASET  
IRIS DATASETS.LOAD.IRIS  
514 FEATURE SELECTION 1087

SCIKITLEARN USER GUIDE RELEASE 0213

X IRISDATA  
Y IRISTARGET  
NCLASSES NPUNIQUEYSIZE  
SOME NOISY DATA NOT CORRELATED  
RANDOM NPRANDOMRANDOMSTATESEED0  
E RANDOMNORMALSIZELENX 2200  
ADD NOISY DATA TO THE INFORMATIVE FEATURES FOR MAKE THE TASK HARDER  
X NPCX E  
SVM SVCKERNELLINEAR  
CV STRATIFIEDKFOLD2  
SCORE PERMUTATIONScores PVALUE PERMUTATIONTESTSCORE  
SVM X Y SCORINGACCURACY CVCV NPERMUTATIONS100 NJOBS1  
PRINTCLASSIFICATION SCORE SPVALUE S SCORE PVALUE

VIEW HISTOGRAM OF PERMUTATION SCORES  
PLTHISTPERMUTATIONScores 20 LABELPERMUTATION SCORES  
EDGEcolorBLACK  
YLIM PLTYLIM  
BUG Vlines linestyle FAILS ON OLDER VERSIONS OF MATPLOTLIB  
PLTVLINESScore YLIM0 YLIM1 LINESTYLE  
COLORG LINEWIDTH3 LABELCLASSIFICATION SCORE  
PVALUE S PVALUE  
PLTVLINES10 NCLASSES YLIM0 YLIM1 LINESTYLE  
COLORK LINEWIDTH3 LABELLUCK  
PLTPLOT2 SCORE YLIM G LINEWIDTH3  
LABELCLASSIFICATION SCORE  
PVALUE S PVALUE  
PLTPLOT2 1 NCLASSES YLIM K LINEWIDTH3 LABELLUCK  
PLTYLIMYLIM  
PLTLEGEND  
PLTXLABELSCORE  
PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 7883 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
5147 UNIVARIATE FEATURE SELECTION  
AN EXAMPLE SHOWING UNIVARIATE FEATURE SELECTION  
NOISY NON INFORMATIVE FEATURES ARE ADDED TO THE IRIS DATA AND UNIVARIATE FEATURE SELECTION IS APPLIED FOR EACH FEATURE  
WE PLOT THE PVALUES FOR THE UNIVARIATE FEATURE SELECTION AND THE CORRESPONDING WEIGHTS OF AN SVM WE CAN SEE THAT  
UNIVARIATE FEATURE SELECTION SELECTS THE INFORMATIVE FEATURES AND THAT THESE HAVE LARGER SVM WEIGHTS  
IN THE TOTAL SET OF FEATURES ONLY THE 4 FIRST ONES ARE SIGNIFICANT WE CAN SEE THAT THEY HAVE THE HIGHEST SCORE WITH  
UNIVARIATE FEATURE SELECTION THE SVM ASSIGNS A LARGE WEIGHT TO ONE OF THESE FEATURES BUT ALSO SELECTS MANY OF THE  
NONINFORMATIVE FEATURES APPLYING UNIVARIATE SELECTION BEFORE THE SVM INCREASES THE SVM WEIGHT ATTRIBUTED  
1088 CHAPTER 5 EXAMPLES

```

SCIKITLEARN USER GUIDE RELEASE 0213
TO THE SIGNIFICANT FEATURES AND WILL THUS IMPROVE CLASSIFICATION
PRINTDOC
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARN IMPORT DATASETS SVM
FROM SKLEARNFEATURESELECTION IMPORT SELECTPERCENTILE FCLASSIF

IMPORT SOME DATA TO PLAY WITH
THE IRIS DATASET
IRIS DATASETSLOADIRIS
SOME NOISY DATA NOT CORRELATED
E NPRANDOMUNIFORM0 01 SIZELENIRISDATA 20
ADD THE NOISY DATA TO THE INFORMATIVE FEATURES
X NPHSTACKIRISDATA E
Y IRISTARGET
PLTFigure1
PLTCLF
514 FEATURE SELECTION 1089

```

SCIKITLEARN USER GUIDE RELEASE 0213  
XINDICES NPARANGEXSHAPE1

UNIVARIATE FEATURE SELECTION WITH FTEST FOR FEATURE SCORING  
WE USE THE DEFAULT SELECTION FUNCTION THE 10 MOST SIGNIFICANT FEATURES  
SELECTOR SELECTPERCENTILEFCCLASSIF PERCENTILE10  
SELECTORFITX Y  
SCORES NPLOG10SELECTORPVALUES  
SCORES SCORESMAX  
PLTBARXINDICES 45 SCORES WIDTH2  
LABELRUNIVARIATE SCORE LOGPVALUE COLORDARKORANGE  
EDGECOLORBLACK

COMPARE TO THE WEIGHTS OF AN SVM  
CLF SVMVCKERNELLINEAR  
CLFFITX Y  
SVMWEIGHTS CLFCOEF 2SUMAXIS0  
SVMWEIGHTS SVMWEIGHTSMAX  
PLTBARXINDICES 25 SVMWEIGHTS WIDTH2 LABELSVM WEIGHT  
COLORNAVY EDGECOLORBLACK  
CLFSELECTED SVMVCKERNELLINEAR  
CLFSELECTEDFITSELECTORTTRANSFORMX Y  
SVMWEIGHTSSELECTED CLFSELECTEDCOEF 2SUMAXIS0  
SVMWEIGHTSSELECTED SVMWEIGHTSSELECTEDMAX  
PLTBARXINDICESSELECTORGETSUPPORT 05 SVMWEIGHTSSELECTED  
WIDTH2 LABELSVM WEIGHTS AFTER SELECTION COLORC  
EDGECOLORBLACK  
PLTTITLECOMPARING FEATURE SELECTION  
PLTXLABELFEATURE NUMBER  
PLTYTICKS  
PLTAXISTIGHT  
PLTLEGENDLOCUPPER RIGHT  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0045 SECONDS  
515 GAUSSIAN PROCESS FOR MACHINE LEARNING  
EXAMPLES CONCERNING THE SKLEARNGAUSSIANPROCESS MODULE  
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE  
1090 CHAPTER 5 EXAMPLES



```
SCIKITLEARN USER GUIDE RELEASE 0213
5151 ILLUSTRATION OF GAUSSIAN PROCESS CLASSIFICATION GPC ON THE XOR DATASET
THIS EXAMPLE ILLUSTRATES GPC ON XOR DATA COMPARED ARE A STATIONARY ISOTROPIC KERNEL RBF AND A NONSTATIONARY
KERNEL DOTPRODUCT ON THIS PARTICULAR DATASET THE DOTPRODUCT KERNEL OBTAINS CONSIDERABLY BETTER RESULTS BECAUSE THE
CLASSBOUNDARIES ARE LINEAR AND COINCIDE WITH THE COORDINATE AXES IN GENERAL STATIONARY KERNELS OFTEN OBTAIN BETTER
RESULTS
PRINTDOC
AUTHORS JAN HENDRIK METZEN JHMINFORMATIKUNIBREMENDE

LICENSE BSD 3 CLAUSE
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARNGAUSSIANPROCESS IMPORT GAUSSIANPROCESSCLASSIFIER
FROM SKLEARNGAUSSIANPROCESSKERNELS IMPORT RBF DOTPRODUCT
XX YY NPMESHGRIDNPLINSPACE3 3 50
NPLINSPACE3 3 50
RNG NPRANDOMRANDOMSTATE0
X RNGRANDN200 2
Y NPLOGICALXORX 0 0 X 1 0
FIT THE MODEL
PLTFIGUREFIGSIZE10 5
KERNELS 10 RBFLLENGTHSCALE10 10 DOTPRODUCTSIGMA010 2
FORI KERNEL INENUMERATEKERNELS
CLF GAUSSIANPROCESSCLASSIFIERKERNELKERNEL WARMSTARTTRUEFITX Y
PLOT THE DECISION FUNCTION FOR EACH DATAPOINT ON THE GRID
Z CLFPREDICTPROBANPVSTACKXXRAVEL YYRAVELT 1
Z ZRESHAPEXXSHAPE
515 GAUSSIAN PROCESS FOR MACHINE LEARNING 1091
```

SCIKITLEARN USER GUIDE RELEASE 0213  
PLTSUBPLOT1 2 I 1  
IMAGE PLTIMSHOWZ INTERPOLATIONNEAREST  
EXTENTXXMIN XXMAX YYMIN YYMAX  
ASPECTAUTO ORIGINLOWER CMAPPLTCMPUORR  
CONTOURS PLTCONTOURXX YY Z LEVELS05 LINEWIDTHS2  
COLORSK  
PLTSCATTERX 0 X 1 S30 CY CMAPPLTCMPAIED  
EDGECOLORS0 0 0  
PLTXTICKS  
PLTYTICKS  
PLTAXIS3 3 3 3  
PLTCOLORBARIMAGE  
PLTTITLE SNLOGMARGINALLIKELIHOOD 3F  
CLFKERNEL CLFLOGMARGINALLIKELIHOODCLFKERNELTHETA  
FONTSIZE12  
PLTTIGHTLAYOUT  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0686 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
5152 GAUSSIAN PROCESS CLASSIFICATION GPC ON IRIS DATASET  
THIS EXAMPLE ILLUSTRATES THE PREDICTED PROBABILITY OF GPC FOR AN ISOTROPIC AND ANISOTROPIC RBF KERNEL ON A TWO  
DIMENSIONAL VERSION FOR THE IRISDATASET THE ANISOTROPIC RBF KERNEL OBTAINS SLIGHTLY HIGHER LOGMARGINALLIKELIHOOD  
BY ASSIGNING DIFFERENT LENGTHSCALES TO THE TWO FEATURE DIMENSIONS  
PRINTDOC  
IMPORT NUMPY AS NP  
1092 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARN IMPORT DATASETS
FROM SKLEARNGAUSSIANPROCESS IMPORT GAUSSIANPROCESSCLASSIFIER
FROM SKLEARNGAUSSIANPROCESSKERNELS IMPORT RBF
IMPORT SOME DATA TO PLAY WITH
IRIS DATASETSLOADIRIS
X IRISDATA 2 WE ONLY TAKE THE FIRST TWO FEATURES
Y NPARRAYIRISTARGET DTYPEINT
H 02 STEP SIZE IN THE MESH
KERNEL 10 RBF10
GPCRBFISOTROPIC GAUSSIANPROCESSCLASSIFIERKERNELKERNELFITX Y
KERNEL 10 RBF10 10
GPCRBFANISOTROPIC GAUSSIANPROCESSCLASSIFIERKERNELKERNELFITX Y
CREATE A MESH TO PLOT IN
XMIN XMAX X 0MIN 1 X 0MAX 1
YMIN YMAX X 1MIN 1 X 1MAX 1
XX YY NPMESHGRIDNPARANGEXMIN XMAX H
NPARANGEYMIN YMAX H
TITLES ISOTROPIC RBF ANISOTROPIC RBF
PLTFIGUREFIGSIZE10 5
FORI CLFINENUMERATEGPCRBFISOTROPIC GPCRBFANISOTROPIC
PLOT THE PREDICTED PROBABILITIES FOR THAT WE WILL ASSIGN A COLOR TO
EACH POINT IN THE MESH XMIN MMAXXYMIN YMAX
PLTSUBPLOT1 2 I 1
Z CLFPREDICTPROBANPCXXRAVEL YYRAVEL
PUT THE RESULT INTO A COLOR PLOT
Z ZRESHAPEXXSHAPE0 XXSHAPE1 3
PLTIMSHOWZ EXTENTXMIN XMAX YMIN YMAX ORIGINLOWER
PLOT ALSO THE TRAINING POINTS
PLTSCATTERX 0 X 1 CNPARRAYR G BY
EDGECOLORS0 0 0
PLTXLABELSEPAL LENGTH
PLTYLABELSEPAL WIDTH
PLTXLIMXXMIN XXMAX
PLTYLIMYYMIN YYMAX
PLTXTICKS
PLTYTICKS
PLTTITLE S LML3F
TITLES I CLFLOGMARGINALLIKELIHOODCLFKERNELTHETA
PLTTIGHTLAYOUT
PLTSHOW
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 4265 SECONDS
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE
515 GAUSSIAN PROCESS FOR MACHINE LEARNING 1093
```

SCIKITLEARN USER GUIDE RELEASE 0213

5153 COMPARISON OF KERNEL RIDGE AND GAUSSIAN PROCESS REGRESSION

BOTH KERNEL RIDGE REGRESSION KRR AND GAUSSIAN PROCESS REGRESSION GPR LEARN A TARGET FUNCTION BY EMPLOYING INTERNALLY THE “KERNEL TRICK” KRR LEARNS A LINEAR FUNCTION IN THE SPACE INDUCED BY THE RESPECTIVE KERNEL WHICH CORRESPONDS TO A NONLINEAR FUNCTION IN THE ORIGINAL SPACE THE LINEAR FUNCTION IN THE KERNEL SPACE IS CHOSEN BASED ON THE MEANSQUARED ERROR LOSS WITH RIDGE REGULARIZATION GPR USES THE KERNEL TO DEFINE THE COVARIANCE OF A PRIOR DISTRIBUTION OVER THE TARGET FUNCTIONS AND USES THE OBSERVED TRAINING DATA TO DEFINE A LIKELIHOOD FUNCTION BASED ON BAYES THEOREM A GAUSSIAN POSTERIOR DISTRIBUTION OVER TARGET FUNCTIONS IS DEFINED WHOSE MEAN IS USED FOR PREDICTION A MAJOR DIFFERENCE IS THAT GPR CAN CHOOSE THE KERNEL’S HYPERPARAMETERS BASED ON GRADIENTASCENT ON THE MARGINAL LIKELIHOOD FUNCTION WHILE KRR NEEDS TO PERFORM A GRID SEARCH ON A CROSSVALIDATED LOSS FUNCTION MEANSQUARED ERROR LOSS A FURTHER DIFFERENCE IS THAT GPR LEARNS A GENERATIVE PROBABILISTIC MODEL OF THE TARGET FUNCTION AND CAN THUS PROVIDE MEANINGFUL CONFIDENCE INTERVALS AND POSTERIOR SAMPLES ALONG WITH THE PREDICTIONS WHILE KRR ONLY PROVIDES PREDICTIONS

THIS EXAMPLE ILLUSTRATES BOTH METHODS ON AN ARTIFICIAL DATASET WHICH CONSISTS OF A SINUSOIDAL TARGET FUNCTION AND STRONG NOISE THE FIGURE COMPARES THE LEARNED MODEL OF KRR AND GPR BASED ON A EXPSINESQUARED KERNEL WHICH IS SUITED FOR LEARNING PERIODIC FUNCTIONS THE KERNEL’S HYPERPARAMETERS CONTROL THE SMOOTHNESS L AND PERIODICITY OF THE KERNEL P MOREOVER THE NOISE LEVEL OF THE DATA IS LEARNED EXPLICITLY BY GPR BY AN ADDITIONAL WHITEKERNEL COMPONENT IN THE KERNEL AND BY THE REGULARIZATION PARAMETER ALPHA OF KRR

THE FIGURE SHOWS THAT BOTH METHODS LEARN REASONABLE MODELS OF THE TARGET FUNCTION GPR CORRECTLY IDENTIFIES THE PERIODICITY OF THE FUNCTION TO BE ROUGHLY 2PI 628 WHILE KRR CHOOSES THE DOUBLED PERIODICITY 4PI BESIDES THAT GPR PROVIDES REASONABLE CONFIDENCE BOUNDS ON THE PREDICTION WHICH ARE NOT AVAILABLE FOR KRR A MAJOR DIFFERENCE BETWEEN THE TWO METHODS IS THE TIME REQUIRED FOR FITTING AND PREDICTING WHILE FITTING KRR IS FAST IN PRINCIPLE THE GRIDSEARCH FOR HYPERPARAMETER OPTIMIZATION SCALES EXPONENTIALLY WITH THE NUMBER OF HYPERPARAMETERS “CURSE OF DIMENSIONALITY” THE GRADIENTBASED OPTIMIZATION OF THE PARAMETERS IN GPR DOES NOT SUFFER FROM THIS EXPONENTIAL SCALING AND IS THUS CONSIDERABLE FASTER ON THIS EXAMPLE WITH 3DIMENSIONAL HYPERPARAMETER SPACE THE TIME FOR PREDICTING IS SIMILAR HOWEVER GENERATING THE VARIANCE OF THE PREDICTIVE DISTRIBUTION OF GPR TAKES CONSIDERABLE LONGER THAN JUST PREDICTING THE MEAN

OUT

TIME FOR KRR FITTING 3180

TIME FOR GPR FITTING 0096

TIME FOR KRR PREDICTION 0009

1094 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
TIME FOR GPR PREDICTION 0010
TIME FOR GPR PREDICTION WITH STANDARDDEVIATION 0014
PRINTDOC
AUTHORS JAN HENDRIK METZEN JHMINFORMATIKUNIBREMENDE
LICENSE BSD 3 CLAUSE
IMPORT TIME
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARNKERNELRIDGE IMPORT KERNELRIDGE
FROM SKLEARNMODELSELECTION IMPORT GRIDSEARCHCV
FROM SKLEARNGAUSSIANPROCESS IMPORT GAUSSIANPROCESSREGRESSOR
FROM SKLEARNGAUSSIANPROCESSKERNELS IMPORT WHITEKERNEL EXPSINESQUARED
RNG NPRANDOMRANDOMSTATE0
GENERATE SAMPLE DATA
X 15RNGRAND100 1
Y NPSINXRABEL
Y 305 RNGRANDXSHAPE0 ADD NOISE
FIT KERNELRIDGE WITH PARAMETER SELECTION BASED ON 5FOLD CROSS VALIDATION
PARAMGRID ALPHA 1E0 1E1 1E2 1E3
KERNEL EXPSINESQUARED L P
FORLINNPLOGSPACE2 2 10
FORPINNPLOGSPACE0 2 10
KR GRIDSEARCHCVKERNELRIDGE CV5 PARAMGRIDPARAMGRID
STIME TIMETIME
KRFITX Y
PRINTTIME FOR KRR FITTING 3F TIMETIME STIME
GPKERNEL EXPSINESQUARED10 50 PERIODICITYBOUNDS1E2 1E1
WHITEKERNEL1E1
GPR GAUSSIANPROCESSREGRESSORKERNELGPKERNEL
STIME TIMETIME
GPRFITX Y
PRINTTIME FOR GPR FITTING 3F TIMETIME STIME
PREDICT USING KERNEL RIDGE
XPLOT NPLINSPACE0 20 10000 NONE
STIME TIMETIME
YKR KRPREDICTXPLOT
PRINTTIME FOR KRR PREDICTION 3F TIMETIME STIME
PREDICT USING GAUSSIAN PROCESS REGRESSOR
STIME TIMETIME
YGPR GPRPREDICTXPLOT RETURNSTDFALSE
515 GAUSSIAN PROCESS FOR MACHINE LEARNING 1095
```

SCIKITLEARN USER GUIDE RELEASE 0213  
PRINTTIME FOR GPR PREDICTION 3F TIMETIME STIME  
STIME TIMETIME  
YGPR YSTD GPRPREDICTXPLOT RETURNSTDTRUE  
PRINTTIME FOR GPR PREDICTION WITH STANDARDDEVIATION 3F  
TIMETIME STIME  
PLOT RESULTS  
PLTFIGUREFIGSIZE10 5  
LW 2  
PLTSCATTERX Y CK LABELDATA  
PLTPLOTXPLOT NPSINXPLOT COLORNAVY LWLW LABELTRUE  
PLTPLOTXPLOT YKR COLORTURQUOISE LWLW  
LABELKRR S KRBESTPARAMS  
PLTPLOTXPLOT YGPR COLORDARKORANGE LWLW  
LABELGPR S GPRKERNEL  
PLTFILLBETWEENXPLOT 0 YGPR YSTD YGPR YSTD COLORDARKORANGE  
ALPHA02  
PLTXLABELDATA  
PLTYLABELTARGET  
PLTXLIM0 20  
PLTYLIM4 4  
PLTTITLEGPR VERSUS KERNEL RIDGE  
PLTLEGENDLOCBEST SCATTERPOINTS1 PROPSIZE 8  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 3377 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
5154 ILLUSTRATION OF PRIOR AND POSTERIOR GAUSSIAN PROCESS FOR DIFFERENT KERNELS  
THIS EXAMPLE ILLUSTRATES THE PRIOR AND POSTERIOR OF A GPR WITH DIFFERENT KERNELS MEAN STANDARD DEVIATION AND 10  
SAMPLES ARE SHOWN FOR BOTH PRIOR AND POSTERIOR  
1096 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

- 515 GAUSSIAN PROCESS FOR MACHINE LEARNING 1097





SCIKITLEARN USER GUIDE RELEASE 0213

- 515 GAUSSIAN PROCESS FOR MACHINE LEARNING 1099



SCIKITLEARN USER GUIDE RELEASE 0213

•

PRINTDOC

AUTHORS JAN HENDRIK METZEN JHMINFORMATIKUNIBREMENDE

LICENSE BSD 3 CLAUSE

IMPORT NUMPY AS NP

FROM MATPLOTLIB IMPORT PYPLOTASPLT

FROM SKLEARNGAUSSIANPROCESS IMPORT GAUSSIANPROCESSREGRESSOR

FROM SKLEARNGAUSSIANPROCESSKERNELS IMPORT RBF MATERN RATIONALQUADRATIC

EXPSINESQUARED DOTPRODUCT

CONSTANTKERNEL

KERNELS 10 RBFLENGTHSCALE10 LENGTHSCALEBOUNDS1E1 100

515 GAUSSIAN PROCESS FOR MACHINE LEARNING 1101

SCIKITLEARN USER GUIDE RELEASE 0213  
10RATIONALQUADRATICLENGTHSCALE10 ALPHA01  
10EXPSINESQUAREDLENGTHSCALE10 PERIODICITY30  
LENGTHSCALEBOUNDS01 100  
PERIODICITYBOUNDS10 100  
CONSTANTKERNEL01 001 100  
DOTPRODUCTSIGMA010 SIGMA0BOUNDS01 100 2  
10MATERNLENGTHSCALE10 LENGTHSCALEBOUNDS1E1 100  
NU15  
FORKERNELINKERNELS  
SPECIFY GAUSSIAN PROCESS  
GP GAUSSIANPROCESSREGRESSORKERNELKERNEL  
PLOT PRIOR  
PLTFIGUREFIGSIZE8 8  
PLTSUBPLOT2 1 1  
X NPLinspace0 5 100  
YMEAN YSTD GPPREDICTX NPNEWAXIS RETURNSTDTRUE  
PLTPLOTX YMEAN K LW3 ZORDER9  
PLTFILLBETWEENX YMEAN YSTD YMEAN YSTD  
ALPHA02 COLORK  
YSAMPLES GPSAMPLEYX NPNEWAXIS 10  
PLTPLOTX YSAMPLES LW1  
PLTXLIM0 5  
PLTYLIM3 3  
PLTTITLEPRIOR KERNEL S KERNEL FONTSIZE12  
GENERATE DATA AND FIT GP  
RNG NPRANDOMRANDOMSTATE4  
X RNGUNIFORM0 5 10 NPNEWAXIS  
Y NPSINX 0 25 2  
GPFITX Y  
PLOT POSTERIOR  
PLTSUBPLOT2 1 2  
X NPLinspace0 5 100  
YMEAN YSTD GPPREDICTX NPNEWAXIS RETURNSTDTRUE  
PLTPLOTX YMEAN K LW3 ZORDER9  
PLTFILLBETWEENX YMEAN YSTD YMEAN YSTD  
ALPHA02 COLORK  
YSAMPLES GPSAMPLEYX NPNEWAXIS 10  
PLTPLOTX YSAMPLES LW1  
PLTSCATTERX 0 Y CR S50 ZORDER10 EDGECOLORS0 0 0  
PLTXLIM0 5  
PLTYLIM3 3  
PLTTITLEPOSTERIOR KERNEL SNLOGLIKELIHOOD 3F  
GPKERNEL GPLOGMARGINALLIKELIHOODGPKERNELTHETA  
FONTSIZE12  
PLTTIGHTLAYOUT  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 1458 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
1102 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
5155 ISOPROBABILITY LINES FOR GAUSSIAN PROCESSES CLASSIFICATION GPC  
A TWODIMENSIONAL CLASSIFICATION EXAMPLE SHOWING ISOPROBABILITY LINES FOR THE PREDICTED PROBABILITIES  
OUT  
LEARNED KERNEL 00256 2DOTPRODUCTSIGMA0572 2  
PRINTDOC  
AUTHOR VINCENT DUBOURG VINCENTDUBOURGGMAILCOM  
ADAPTED TO GAUSSIANPROCESSCLASSIFIER  
JAN HENDRIK METZEN JHMINFORMATIKUNIBREMENDE  
LICENSE BSD 3 CLAUSE  
IMPORT NUMPY AS NP  
FROM MATPLOTLIB IMPORT PYPLASPLT  
FROM MATPLOTLIB IMPORT CM  
FROM SKLEARNPROCESS IMPORT GAUSSIANPROCESSCLASSIFIER  
515 GAUSSIAN PROCESS FOR MACHINE LEARNING 1103

SCIKITLEARN USER GUIDE RELEASE 0213  
FROM SKLEARNGAUSSIANPROCESSKERNELS IMPORT DOTPRODUCT CONSTANTKERNEL ASC  
A FEW CONSTANTS  
LIM 8  
DEFGX  
THE FUNCTION TO PREDICT CLASSIFICATION WILL THEN CONSIST IN PREDICTING  
WHETHER GX 0 OR NOT  
RETURN5 X 1 5 X 0 2  
DESIGN OF EXPERIMENTS  
X NPARRAY461611719 600099547  
410469096 532782448  
000000000 050000000  
617289014 46984743  
13109306 693271427  
503823144 310584743  
287600388 674310541  
521301203 426386883  
OBSERVATIONS  
Y NPARRAYGX 0 DTYPEINT  
INSTANTIATE AND FIT GAUSSIAN PROCESS MODEL  
KERNEL C01 1E5 NPINF DOTPRODUCTSIGMA001 2  
GP GAUSSIANPROCESSCLASSIFIERKERNELKERNEL  
GPFITX Y  
PRINTLEARNED KERNEL S GPKERNEL  
EVALUATE REAL FUNCTION AND THE PREDICTED PROBABILITY  
RES 50  
X1 X2 NPMESHGRIDNPLINSPACE LIM LIM RES  
NPLINSPACE LIM LIM RES  
XX NPVSTACKX1RESHAPEX1SIZE X2RESHAPEX2SIZET  
YTRUE GXX  
YPROB GPPREDICTPROBAXX 1  
YTRUE YTRUERESHAPERES RES  
YPROB YPROBRESHAPERES RES  
PLOT THE PROBABILISTIC CLASSIFICATION ISOVALUES  
FIG PLTFigure1  
AX FIGGCA  
AXAXESSETASPECTEQUAL  
PLXTICKS  
PLTYTICKS  
AXSETXTICKLABELS  
AXSETYTICKLABELS  
PLTXLABELX1  
PLTYLABELX2  
CAX PLTIMSHOWYPROB CMAPCMGRAYR ALPHA08  
EXTENTLIM LIM LIM LIM  
NORM PLTMATPLOTLIBCOLORSNORMALIZEVMIN0 VMAX09  
CB PLTCOLORBARCAX TICKS0 02 04 06 08 1 NORMNORM  
CBSETLABELRRM MATHBBPLEFTWIDEHATGMATHBFX LEQ 0RIGHT  
PLTCLIM0 1  
1104 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
PLTPLOTXY 0 0 XY 0 1 R MARKERSIZE12
PLTPLOTXY 0 0 XY 0 1 B MARKERSIZE12
PLTCONTOURX1 X2 YTRUE 0 COLORSK LINESTYLES DASHDOT
CS PLTCONTOURX1 X2 YPROB 0666 COLORSB
LINESTYLESSOLID
PLTCLABELCS FONTSIZE11
CS PLTCONTOURX1 X2 YPROB 05 COLORSK
LINESTYLESDASHED
PLTCLABELCS FONTSIZE11
CS PLTCONTOURX1 X2 YPROB 0334 COLORSR
LINESTYLESSOLID
PLTCLABELCS FONTSIZE11
PLTSHOW
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0098 SECONDS
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE
5156 PROBABILISTIC PREDICTIONS WITH GAUSSIAN PROCESS CLASSIFICATION GPC
THIS EXAMPLE ILLUSTRATES THE PREDICTED PROBABILITY OF GPC FOR AN RBF KERNEL WITH DIFFERENT CHOICES OF THE HYPERPARAMETERS THE FIRST FIGURE SHOWS THE PREDICTED PROBABILITY OF GPC WITH ARBITRARILY CHOSEN HYPERPARAMETERS AND WITH THE HYPERPARAMETERS CORRESPONDING TO THE MAXIMUM LOGMARGINALLIKELIHOOD LML WHILE THE HYPERPARAMETERS CHOSEN BY OPTIMIZING LML HAVE A CONSIDERABLE LARGER LML THEY PERFORM SLIGHTLY WORSE ACCORDING TO THE LOGLOSS ON TEST DATA THE FIGURE SHOWS THAT THIS IS BECAUSE THEY EXHIBIT A STEEP CHANGE OF THE CLASS PROBABILITIES AT THE CLASS BOUNDARIES WHICH IS GOOD BUT HAVE PREDICTED PROBABILITIES CLOSE TO 05 FAR AWAY FROM THE CLASS BOUNDARIES WHICH IS BAD THIS UNDESIRABLE EFFECT IS CAUSED BY THE LAPLACE APPROXIMATION USED INTERNALLY BY GPC
THE SECOND FIGURE SHOWS THE LOGMARGINALLIKELIHOOD FOR DIFFERENT CHOICES OF THE KERNEL'S HYPERPARAMETERS HIGHLIGHTING THE TWO CHOICES OF THE HYPERPARAMETERS USED IN THE FIRST FIGURE BY BLACK DOTS
515 GAUSSIAN PROCESS FOR MACHINE LEARNING 1105
```





SCIKITLEARN USER GUIDE RELEASE 0213

- 

OUT

LOG MARGINAL LIKELIHOOD INITIAL 17598

LOG MARGINAL LIKELIHOOD OPTIMIZED 3875

ACCURACY 1000 INITIAL 1000 OPTIMIZED

LOGLOSS 0214 INITIAL 0319 OPTIMIZED

PRINTDOC

AUTHORS JAN HENDRIK METZEN JHMINFORMATIKUNIBREMENDE

LICENSE BSD 3 CLAUSE

IMPORT NUMPY AS NP

FROM MATPLOTLIB IMPORT PYPLASPLT

FROM SKLEARNMETRICSClassification IMPORT ACCURACYScore LOGLOSS

FROM SKLEARNGAUSSIANPROCESS IMPORT GAUSSIANPROCESSClassifier

FROM SKLEARNGAUSSIANPROCESSKernels IMPORT RBF

515 GAUSSIAN PROCESS FOR MACHINE LEARNING 1107

SCIKITLEARN USER GUIDE RELEASE 0213

GENERATE DATA

TRAINSIZ 50

RNG NPRANDOMRANDOMSTATE0

X RNGUNIFORM0 5 100 NPNEWAXIS

Y NPARRAYX 0 25 DTYPEINT

SPECIFY GAUSSIAN PROCESSES WITH FIXED AND OPTIMIZED HYPERPARAMETERS

GPFIX GAUSSIANPROCESSCLASSIFIERKERNEL10 RBLENGTHSCALE10

OPTIMIZERNONE

GPFIXFITXTRAINSIZ YTRAINSIZ

GPOPT GAUSSIANPROCESSCLASSIFIERKERNEL10 RBLENGTHSCALE10

GPOPTFITXTRAINSIZ YTRAINSIZ

PRINTLOG MARGINAL LIKELIHOOD INITIAL 3F

GPFIXLOGMARGINALLIKELIHOODGPFIXKERNELTHETA

PRINTLOG MARGINAL LIKELIHOOD OPTIMIZED 3F

GPOPTLOGMARGINALLIKELIHOODGPOPTKERNELTHETA

PRINTACCURACY 3FINITIAL 3FOPTIMIZED

ACCURACYSCOREYTRAINSIZ GPFIXPREDICTXTRAINSIZ

ACCURACYSCOREYTRAINSIZ GPOPTPREDICTXTRAINSIZ

PRINTLOGLOSS 3FINITIAL 3FOPTIMIZED

LOGLOSSYTRAINSIZ GPFIXPREDICTPROBAXTRAINSIZ 1

LOGLOSSYTRAINSIZ GPOPTPREDICTPROBAXTRAINSIZ 1

PLOT POSTERIOR

PLTFigure

PLTSCATTERXTRAINSIZ 0 YTRAINSIZ CK LABELTRAIN DATA

EDGECOLORS0 0 0

PLTSCATTERXTRAINSIZ 0 YTRAINSIZ CG LABELTEST DATA

EDGECOLORS0 0 0

X NPLINSPACE0 5 100

PLTPLOTX GPFIXPREDICTPROBAX NPNEWAXIS 1 R

LABELINITIAL KERNEL S GPFIXKERNEL

PLTPLOTX GPOPTPREDICTPROBAX NPNEWAXIS 1 B

LABELOPTIMIZED KERNEL S GPOPTKERNEL

PLTXLABELFEATURE

PLTYLABELCLASS 1 PROBABILITY

PLTXLIM0 5

PLTYLIM025 15

PLTLEGENDLOCBEST

PLOT LML LANDSCAPE

PLTFigure

THETA0 NPLOGSPACE0 8 30

THETA1 NPLOGSPACE1 1 29

THETA0 THETA1 NPMESHGRIDTHETA0 THETA1

LML GPOPTLOGMARGINALLIKELIHOODNPLOGTHETA0I J THETA1I J

FORIINRANGETHETA0SHAPE0 FORJINRANGETHETA0SHAPE1

LML NPARRAYLMLT

PLTPLOTNPXPGPFIXKERNELTHETA0 NPXPGPFIXKERNELTHETA1

KO ZORDER10

PLTPLOTNPXPGPPOPTKERNELTHETA0 NPXPGPPOPTKERNELTHETA1

KO ZORDER10

PLTPCOLORTHETA0 THETA1 LML

PLTXSCALELOG

1108 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

PLTYSCALELOG

PLTCOLORBAR

PLTXLABELMAGNITUDE

PLTYLABELLENGTHSCALE

PLTTITLELOGMARGINALLIKELIHOOD

PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 2514 SECONDS

NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)

5157 GAUSSIAN PROCESS REGRESSION GPR WITH NOISELEVEL ESTIMATION

THIS EXAMPLE ILLUSTRATES THAT GPR WITH A SUMKERNEL INCLUDING A WHITEKERNEL CAN ESTIMATE THE NOISE LEVEL OF DATA AN ILLUSTRATION OF THE LOGMARGINALLIKELIHOOD LML LANDSCAPE SHOWS THAT THERE EXIST TWO LOCAL MAXIMA OF LML THE FIRST CORRESPONDS TO A MODEL WITH A HIGH NOISE LEVEL AND A LARGE LENGTH SCALE WHICH EXPLAINS ALL VARIATIONS IN THE DATA BY NOISE THE SECOND ONE HAS A SMALLER NOISE LEVEL AND SHORTER LENGTH SCALE WHICH EXPLAINS MOST OF THE VARIATION BY THE NOISEFREE FUNCTIONAL RELATIONSHIP THE SECOND MODEL HAS A HIGHER LIKELIHOOD HOWEVER DEPENDING ON THE INITIAL VALUE FOR THE HYPERPARAMETERS THE GRADIENTBASED OPTIMIZATION MIGHT ALSO CONVERGE TO THE HIGHNOISE SOLUTION IT IS THUS IMPORTANT TO REPEAT THE OPTIMIZATION SEVERAL TIMES FOR DIFFERENT INITIALIZATIONS

•

515 GAUSSIAN PROCESS FOR MACHINE LEARNING 1109



SCIKITLEARN USER GUIDE RELEASE 0213

•  
PRINTDOC  
AUTHORS JAN HENDRIK METZEN JHMINFORMATIKUNIBREMENDE

```
LICENSE BSD 3 CLAUSE
IMPORT NUMPY AS NP
FROM MATPLOTLIB IMPORT PYPLLOTASPLT
FROM MATPLOTLIBCOLORS IMPORT LOGNORM
FROM SKLEARNGAUSSIANPROCESS IMPORT GAUSSIANPROCESSREGRESSOR
FROM SKLEARNGAUSSIANPROCESSKERNELS IMPORT RBF WHITEKERNEL
RNG NPRANDOMRANDOMSTATE0
X RNGUNIFORM0 5 20 NPNEWAXIS
Y 05 NPSIN3 X 0 RNGNORMAL0 05 XSHAPE0
FIRST RUN
PLTFigure
KERNEL 10 RBFLengthScale1000 LengthScaleBounds1E2 1E3
WHITEKernelNoiseLevel1 NoiseLevelBounds1E10 1E1
GP GAUSSIANPROCESSREGRESSORKernelKernel
ALPHA00FITX Y
X NPLinspace0 5 100
YMEAN YCOV GPPREDICTX NPNEWAXIS RETURNCOVTRUE
515 GAUSSIAN PROCESS FOR MACHINE LEARNING 1111
```

SCIKITLEARN USER GUIDE RELEASE 0213  
PLTPLOTX YMEAN K LW3 ZORDER9  
PLTFILLBETWEENX YMEAN NPSQRTNPDIAGYCOV  
YMEAN NPSQRTNPDIAGYCOV  
ALPHA05 COLORK  
PLTPLOTX 05 NPSIN3 X R LW3 ZORDER9  
PLTSCATTERX 0 Y CR S50 ZORDER10 EDGECOLORS0 0 0  
PLTTITLEINITIAL SNOPTIMUM SNLOGMARGINALLIKELIHOOD S  
KERNEL GPKERNEL  
GPLOGMARGINALLIKELIHOODGPKERNELTHETA  
PLTTIGHTLAYOUT  
SECOND RUN  
PLTFigure  
KERNEL 10 RBFLENGTHSCALE10 LENGTHSCALEBOUNDS1E2 1E3  
WHITEKERNELNOISELEVEL1E5 NOISELEVELBOUNDS1E10 1E1  
GP GAUSSIANPROCESSREGRESSORKERNELKERNEL  
ALPHA00FITX Y  
X NPLINSPACE0 5 100  
YMEAN YCOV GPPREDICTX NPNEWAXIS RETURNCOVTRUE  
PLTPLOTX YMEAN K LW3 ZORDER9  
PLTFILLBETWEENX YMEAN NPSQRTNPDIAGYCOV  
YMEAN NPSQRTNPDIAGYCOV  
ALPHA05 COLORK  
PLTPLOTX 05 NPSIN3 X R LW3 ZORDER9  
PLTSCATTERX 0 Y CR S50 ZORDER10 EDGECOLORS0 0 0  
PLTTITLEINITIAL SNOPTIMUM SNLOGMARGINALLIKELIHOOD S  
KERNEL GPKERNEL  
GPLOGMARGINALLIKELIHOODGPKERNELTHETA  
PLTTIGHTLAYOUT  
PLOT LML LANDSCAPE  
PLTFigure  
THETA0 NPLOGSPACE2 3 49  
THETA1 NPLOGSPACE2 0 50  
THETA0 THETA1 NPMESHGRIDTHETA0 THETA1  
LML GPLOGMARGINALLIKELIHOODNPLOG036 THETA0I J THETA1I J  
FORIINRANGETHETA0SHAPE0 FORJINRANGETHETA0SHAPE1  
LML NPARRAYLMLT  
VMIN VMAX LMLMIN LMLMAX  
VMAX 50  
LEVEL NPAROUNDNPLOGSPACENPLOG10VMIN NPLOG10VMAX 50 DECIMALS1  
PLTCONTOURTHETA0 THETA1 LML  
LEVELSLEVEL NORMLOGNORMVMINVMIN VMAXVMAX  
PLTCOLORBAR  
PLTXSCALELOG  
PLTYSCALELOG  
PLTXLABELLENGTHSCALE  
PLTYLABELNOISELEVEL  
PLTTITLELOGMARGINALLIKELIHOOD  
PLTTIGHTLAYOUT  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 2874 SECONDS  
1112 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)

5158 GAUSSIAN PROCESSES REGRESSION BASIC INTRODUCTORY EXAMPLE

A SIMPLE ONEDIMENSIONAL REGRESSION EXAMPLE COMPUTED IN TWO DIFFERENT WAYS

1 A NOISEFREE CASE

2 A NOISY CASE WITH KNOWN NOISELEVEL PER DATAPOINT

IN BOTH CASES THE KERNEL’S PARAMETERS ARE ESTIMATED USING THE MAXIMUM LIKELIHOOD PRINCIPLE

THE FIGURES ILLUSTRATE THE INTERPOLATING PROPERTY OF THE GAUSSIAN PROCESS MODEL AS WELL AS ITS PROBABILISTIC NATURE IN THE FORM OF A POINTWISE 95 CONFIDENCE INTERVAL

NOTE THAT THE PARAMETER ALPHA IS APPLIED AS A TIKHONOV REGULARIZATION OF THE ASSUMED COVARIANCE BETWEEN THE TRAINING POINTS

- 

515 GAUSSIAN PROCESS FOR MACHINE LEARNING 1113

SCIKITLEARN USER GUIDE RELEASE 0213

```
•
PRINTDOC
AUTHOR VINCENT DUBOURG VINCENTDUBOURGGMAILCOM
JAKE VANDERPLAS VANDERPLASASTROWASHINGTONEDU
JAN HENDRIK METZEN JHMINFORMATIKUNIBREMENDES
LICENSE BSD 3 CLAUSE
IMPORT NUMPY AS NP
FROM MATPLOTLIB IMPORT PYPLASPLT
FROM SKLEARNGAUSSIANPROCESS IMPORT GAUSSIANPROCESSREGRESSOR
FROM SKLEARNGAUSSIANPROCESSKERNELS IMPORT RBF CONSTANTKERNEL ASC
NPRANDOMSEED1
DEFFX
THE FUNCTION TO PREDICT
RETURNXNPSINX

FIRST THE NOISELESS CASE
X NPATLEAST2D1 3 5 6 7 8T
OBSERVATIONS
Y FXRAVEL
1114 CHAPTER 5 EXAMPLES
```



```
SCIKITLEARN USER GUIDE RELEASE 0213
MESH THE INPUT SPACE FOR EVALUATIONS OF THE REAL FUNCTION THE PREDICTION AND
ITS MSE
X NPATLEAST2DNPLINSPACE0 10 1000T
INSTANTIATE A GAUSSIAN PROCESS MODEL
KERNEL C10 1E3 1E3 RBF10 1E2 1E2
GP GAUSSIANPROCESSREGRESSORKERNELKERNEL NRESTARTSOPTIMIZER9
FIT TO DATA USING MAXIMUM LIKELIHOOD ESTIMATION OF THE PARAMETERS
GPFITX Y
MAKE THE PREDICTION ON THE MESHED XAXIS ASK FOR MSE AS WELL
YPRED SIGMA GPPREDICTX RETURNSTDTRUE
PLOT THE FUNCTION THE PREDICTION AND THE 95 CONFIDENCE INTERVAL BASED ON
THE MSE
PLTFigure
PLTPLOTX FX R LABELRFX XSINX
PLTPLOTX Y R MARKERSIZE10 LABELOBSERVATIONS
PLTPLOTX YPRED B LABELPREDICTION
PLTFILLNPConcatenateX X1
NPConcatenateYPRED 19600 SIGMA
YPRED 19600 SIGMA1
ALPHA5 FCB ECNONE LABEL95 CONFIDENCE INTERVAL
PLTXLABELX
PLTYLABELFX
PLTYLIM10 20
PLTLEGENDLOCUPPER LEFT

NOW THE NOISY CASE
X NPLINSPACE01 99 20
X NPATLEAST2DXT
OBSERVATIONS AND NOISE
Y FXRAVEL
DY 05 10 NPRANDOMRANDOMMYSHAPE
NOISE NPRANDOMNORMAL0 DY
Y NOISE
INSTANTIATE A GAUSSIAN PROCESS MODEL
GP GAUSSIANPROCESSREGRESSORKERNELKERNEL ALPHADY 2
NRESTARTSOPTIMIZER10
FIT TO DATA USING MAXIMUM LIKELIHOOD ESTIMATION OF THE PARAMETERS
GPFITX Y
MAKE THE PREDICTION ON THE MESHED XAXIS ASK FOR MSE AS WELL
YPRED SIGMA GPPREDICTX RETURNSTDTRUE
PLOT THE FUNCTION THE PREDICTION AND THE 95 CONFIDENCE INTERVAL BASED ON
THE MSE
PLTFigure
PLTPLOTX FX R LABELRFX XSINX
PLTERRORBARXRAVEL Y DY FMTR MARKERSIZE10 LABELOBSERVATIONS
PLTPLOTX YPRED B LABELPREDICTION
PLTFILLNPConcatenateX X1
NPConcatenateYPRED 19600 SIGMA
515 GAUSSIAN PROCESS FOR MACHINE LEARNING 1115
```

SCIKITLEARN USER GUIDE RELEASE 0213  
YPRED 19600 SIGMA1  
ALPHA5 FCB ECNONE LABEL95 CONFIDENCE INTERVAL  
PLTXLABELX  
PLTYLABELFX  
PLTYLIM10 20  
PLTLEGENDLOCUPPER LEFT  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0284 SECONDS  
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE  
5159 GAUSSIAN PROCESS REGRESSION GPR ON MAUNA LOA CO2 DATA  
THIS EXAMPLE IS BASED ON SECTION 543 OF “GAUSSIAN PROCESSES FOR MACHINE LEARNING” RW2006 IT ILLUSTRATES AN  
EXAMPLE OF COMPLEX KERNEL ENGINEERING AND HYPERPARAMETER OPTIMIZATION USING GRADIENT ASCENT ON THE LOGMARGINAL  
LIKELIHOOD THE DATA CONSISTS OF THE MONTHLY AVERAGE ATMOSPHERIC CO2 CONCENTRATIONS IN PARTS PER MILLION BY VOLUME  
PPMV COLLECTED AT THE MAUNA LOA OBSERVATORY IN HAWAII BETWEEN 1958 AND 2001 THE OBJECTIVE IS TO MODEL THE  
CO2 CONCENTRATION AS A FUNCTION OF THE TIME T  
THE KERNEL IS COMPOSED OF SEVERAL TERMS THAT ARE RESPONSIBLE FOR EXPLAINING DIFFERENT PROPERTIES OF THE SIGNAL  
• A LONG TERM SMOOTH RISING TREND IS TO BE EXPLAINED BY AN RBF KERNEL THE RBF KERNEL WITH A LARGE LENGTHSCALE  
ENFORCES THIS COMPONENT TO BE SMOOTH IT IS NOT ENFORCED THAT THE TREND IS RISING WHICH LEAVES THIS CHOICE TO THE  
GP THE SPECIFIC LENGTHSCALE AND THE AMPLITUDE ARE FREE HYPERPARAMETERS  
• A SEASONAL COMPONENT WHICH IS TO BE EXPLAINED BY THE PERIODIC EXPSINESQUARED KERNEL WITH A FIXED PERIODICITY  
OF 1 YEAR THE LENGTHSCALE OF THIS PERIODIC COMPONENT CONTROLLING ITS SMOOTHNESS IS A FREE PARAMETER IN ORDER  
TO ALLOW DECAYING AWAY FROM EXACT PERIODICITY THE PRODUCT WITH AN RBF KERNEL IS TAKEN THE LENGTHSCALE OF THIS  
RBF COMPONENT CONTROLS THE DECAY TIME AND IS A FURTHER FREE PARAMETER  
• SMALLER MEDIUM TERM IRREGULARITIES ARE TO BE EXPLAINED BY A RATIONALQUADRATIC KERNEL COMPONENT WHOSE LENGTH  
SCALE AND ALPHA PARAMETER WHICH DETERMINES THE DIFFUSENESS OF THE LENGTHSCALES ARE TO BE DETERMINED AC  
CORDING TO RW2006 THESE IRREGULARITIES CAN BETTER BE EXPLAINED BY A RATIONALQUADRATIC THAN AN RBF KERNEL  
COMPONENT PROBABLY BECAUSE IT CAN ACCOMMODATE SEVERAL LENGTHSCALES  
• A “NOISE” TERM CONSISTING OF AN RBF KERNEL CONTRIBUTION WHICH SHALL EXPLAIN THE CORRELATED NOISE COMPONENTS  
SUCH AS LOCAL WEATHER PHENOMENA AND A WHITEKERNEL CONTRIBUTION FOR THE WHITE NOISE THE RELATIVE AMPLITUDES  
AND THE RBF’S LENGTH SCALE ARE FURTHER FREE PARAMETERS  
MAXIMIZING THE LOGMARGINALLIKELIHOOD AFTER SUBTRACTING THE TARGET’S MEAN YIELDS THE FOLLOWING KERNEL WITH AN LML  
OF 83214  
3442RBFLENGTHSCALE418  
3272RBFLENGTHSCALE180 EXPSINESQUAREDLENGTHSCALE144  
PERIODICITY1  
04462RATIONALQUADRATICALPHA177 LENGTHSCALE0957  
01972RBFLENGTHSCALE0138 WHITEKERNELNOISELEVEL00336  
THUS MOST OF THE TARGET SIGNAL 344PPM IS EXPLAINED BY A LONGTERM RISING TREND LENGTHSCALE 418 YEARS THE  
PERIODIC COMPONENT HAS AN AMPLITUDE OF 327PPM A DECAY TIME OF 180 YEARS AND A LENGTHSCALE OF 144 THE LONG  
DECAY TIME INDICATES THAT WE HAVE A LOCALLY VERY CLOSE TO PERIODIC SEASONAL COMPONENT THE CORRELATED NOISE HAS AN  
AMPLITUDE OF 0197PPM WITH A LENGTH SCALE OF 0138 YEARS AND A WHITENOISE CONTRIBUTION OF 0197PPM THUS THE  
1116 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
OVERALL NOISE LEVEL IS VERY SMALL INDICATING THAT THE DATA CAN BE VERY WELL EXPLAINED BY THE MODEL THE FIGURE SHOWS  
ALSO THAT THE MODEL MAKES VERY CONFIDENT PREDICTIONS UNTIL AROUND 2015  
OUT  
GPML KERNEL 66 2RBFLENGTHSCALE67 24 2RBFLENGTHSCALE90  
↪EXPSINESQUAREDLENGTHSCALE13 PERIODICITY1 066 2  
↪RATIONALQUADRATICAPHA078 LENGTHSCALE12 018 2RBFLENGTHSCALE0134  
↪ WHITEKERNELNOISELEVEL00361  
LOGMARGINALLIKELIHOOD 117023  
LEARNED KERNEL 448 2RBFLENGTHSCALE516 264 2RBFLENGTHSCALE915  
↪EXPSINESQUAREDLENGTHSCALE148 PERIODICITY1 0536 2  
↪RATIONALQUADRATICAPHA289 LENGTHSCALE0968 0188 2RBFLENGTHSCALE0  
↪122 WHITEKERNELNOISELEVEL00367  
LOGMARGINALLIKELIHOOD 115050  
AUTHORS JAN HENDRIK METZEN JHMINFORMATIKUNIBREMENDE  
  
LICENSE BSD 3 CLAUSE  
515 GAUSSIAN PROCESS FOR MACHINE LEARNING 1117

```
SCIKITLEARN USER GUIDE RELEASE 0213
IMPORT NUMPY AS NP
FROM MATPLOTLIB IMPORT PYPLLOTASPLT
FROM SKLEARNDATASETS IMPORT FETCHOPENML
FROM SKLEARNGAUSSIANPROCESS IMPORT GAUSSIANPROCESSREGRESSOR
FROM SKLEARNGAUSSIANPROCESSKERNELS
IMPORTRBF WHITEKERNEL RATIONALQUADRATIC EXPSINESQUARED
PRINTDOC
DEFLOADMAUNALOAATMOSPHERICCO2
MLDATA FETCHOPENMLDATAID41187
MONTHS
PPMVSUMS
COUNTS
Y MLDATADATA 0
M MLDATADATA 1
MONTHFLOAT Y M 1 12
PPMVS MLDATATARGET
FORMONTH PPMV INZIPMONTHFLOAT PPMVS
IF NOTMONTHSORMONTH MONTHS1
MONTHSAPPENDMONTH
PPMVSUMSAPPENDPPMV
COUNTSAPPEND1
ELSE
  AGGREGATE MONTHLY SUM TO PRODUCE AVERAGE
PPMVSUMS1 PPMV
COUNTS1 1
MONTHS NPASARRAYMONTHSRESHAPE1 1
AVGPPMVS NPASARRAYPPMVSUMS COUNTS
RETURNMONTHS AVGPPMVS
X Y LOADMAUNALOAATMOSPHERICCO2
  KERNEL WITH PARAMETERS GIVEN IN GPML BOOK
K1 660 2RBFLENGTHSCALE670 LONG TERM SMOOTH RISING TREND
K2 24 2RBFLENGTHSCALE900
EXPSINESQUAREDLENGTHSCALE13 PERIODICITY10 SEASONAL COMPONENT
  MEDIUM TERM IRREGULARITY
K3 066 2
RATIONALQUADRATICLENGTHSCALE12 ALPHA078
K4 018 2RBFLENGTHSCALE0134
  WHITEKERNELNOISELEVEL019 2 NOISE TERMS
KERNELGPML K1 K2 K3 K4
GP GAUSSIANPROCESSREGRESSORKERNELKERNELGPML ALPHA0
OPTIMIZERNONE NORMALIZEYTRUE
GPFITX Y
PRINTGPML KERNEL S GPKERNEL
PRINTLOGMARGINALLIKELIHOOD 3F
  GPLOGMARGINALLIKELIHOODGPKERNELTHETA
1118 CHAPTER 5 EXAMPLES
```

SCIKITLEARN USER GUIDE RELEASE 0213

KERNEL WITH OPTIMIZED PARAMETERS

K1 500 2RBFLENGTHSCALE500 LONG TERM SMOOTH RISING TREND

K2 20 2RBFLENGTHSCALE1000

EXPSINESQUAREDLENGTHSCALE10 PERIODICITY10

PERIODICITYBOUNDSFIXED SEASONAL COMPONENT

MEDIUM TERM IRREGULARITIES

K3 05 2RATIONALQUADRATICLENGTHSCALE10 ALPHA10

K4 01 2RBFLENGTHSCALE01

WHITEKERNELNOISELEVEL01 2

NOISELEVELBOUNDS1E3 NPINF NOISE TERMS

KERNEL K1 K2 K3 K4

GP GAUSSIANPROCESSREGRESSORKERNELKERNEL ALPHA0

NORMALIZEYTRUE

GPFITX Y

PRINTNLEARNED KERNEL S GPKERNEL

PRINTLOGMARGINALLIKELIHOOD 3F

GPLOGMARGINALLIKELIHOODGPKERNELTHETA

X NPLINSPACEXMIN XMAX 30 1000 NPNEWAXIS

YPRED YSTD GPPREDICTX RETURNSTDTRUE

ILLUSTRATION

PLTSCATTERX Y CK

PLTPLOTX YPRED

PLTFILLBETWEENX 0 YPRED YSTD YPRED YSTD

ALPHA05 COLORK

PLTXLIMXMIN XMAX

PLTXLABELYEAR

PLTYLABELRCO2 IN PPM

PLTTITLERATMOSPHERIC CO2 CONCENTRATION AT MAUNA LOA

PLTTIGHTLAYOUT

PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 11550 SECONDS

516 MISSING VALUE IMPUTATION

EXAMPLES CONCERNING THE SKLEARNIMPUTE MODULE

NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)

5161 IMPUTING MISSING VALUES WITH VARIANTS OF ITERATIVEIMPUTER

THESKLEARNIMPUTEITERATIVEIMPUTER CLASS IS VERY FLEXIBLE IT CAN BE USED WITH A VARIETY OF ESTIMATORS

TO DO ROUNDROBIN REGRESSION TREATING EVERY VARIABLE AS AN OUTPUT IN TURN

IN THIS EXAMPLE WE COMPARE SOME ESTIMATORS FOR THE PURPOSE OF MISSING FEATURE IMPUTATION WITH SKLEARNIMPUTE

ITERATIVEIMPUTER

•BAYESIANRIDGE REGULARIZED LINEAR REGRESSION

516 MISSING VALUE IMPUTATION 1119

SCIKITLEARN USER GUIDE RELEASE 0213

- DECISIONTREEREgressor NONLINEAR REGRESSION
- EXTRATREESRegressor SIMILAR TO MISSFOREST IN R
- KNEIGHBORSRegressor COMPARABLE TO OTHER KNN IMPUTATION APPROACHES

OF PARTICULAR INTEREST IS THE ABILITY OF SKLEARNIMPUTEITERATIVEIMPUTER TO MIMIC THE BEHAVIOR OF MISS FOREST A POPULAR IMPUTATION PACKAGE FOR R IN THIS EXAMPLE WE HAVE CHOSEN TO USE SKLEARNENSEMBLE EXTRATREESRegressor INSTEAD OF SKLEARNENSEMBLERANDOMFORESTRegressor AS IN MISSFOREST DUE TO ITS INCREASED SPEED

NOTE THAT SKLEARNNEIGHBORSKNEIGHBORSRegressor IS DIFFERENT FROM KNN IMPUTATION WHICH LEARNS FROM SAMPLES WITH MISSING VALUES BY USING A DISTANCE METRIC THAT ACCOUNTS FOR MISSING VALUES RATHER THAN IMPUTING THEM

THE GOAL IS TO COMPARE DIFFERENT ESTIMATORS TO SEE WHICH ONE IS BEST FOR THE SKLEARNIMPUTE ITERATIVEIMPUTER WHEN USING A SKLEARNLINEARMODEL BAYESIANRIDGE ESTIMATOR ON THE CALIFORNIA HOUSING DATASET WITH A SINGLE VALUE RANDOMLY REMOVED FROM EACH ROW

FOR THIS PARTICULAR PATTERN OF MISSING VALUES WE SEE THAT SKLEARNENSEMBLEEXTRATREESRegressor AND SKLEARNLINEARMODEL BAYESIANRIDGE GIVE THE BEST RESULTS

PRINTDOC

IMPORT NUMPY AS NP

IMPORT MATPLOTLIBPYPLT AS PLT

IMPORT PANDAS AS PD

TO USE THIS EXPERIMENTAL FEATURE WE NEED TO EXPLICITLY ASK FOR IT

FROM SKLEARNEXPERIMENTAL IMPORT ENABLEITERATIVEIMPUTER NOQA

FROM SKLEARNDATASETS IMPORT FETCHCALIFORNIAHOUSING

FROM SKLEARNIMPUTE IMPORT SIMPLEIMPUTER

FROM SKLEARNIMPUTE IMPORT ITERATIVEIMPUTER

FROM SKLEARNLINEARMODEL IMPORT BAYESIANRIDGE

FROM SKLEARNTREE IMPORT DECISIONTREEREgressor

FROM SKLEARNENSEMBLE IMPORT EXTRATREESRegressor

FROM SKLEARNNEIGHBORS IMPORT KNEIGHBORSRegressor

FROM SKLEARNPIPELINE IMPORT MAKEPIPELINE

FROM SKLEARNMODELSELECTION IMPORT CROSSVALSCORE

1120 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
NSPLITS 5
RNG NPRANDOMRANDOMSTATE0
XFULL YFULL FETCHCALIFORNIAHOUSINGRETURNXYTRUE
2K SAMPLES IS ENOUGH FOR THE PURPOSE OF THE EXAMPLE
REMOVE THE FOLLOWING TWO LINES FOR A SLOWER RUN WITH DIFFERENT ERROR BARS
XFULL XFULL10
YFULL YFULL10
NSAMPLES NFEATURES XFULLSHAPE
ESTIMATE THE SCORE ON THE ENTIRE DATASET WITH NO MISSING VALUES
BRESTIMATOR BAYESIANRIDGE
SCOREFULLDATA PDDATAFRAME
CROSSVALSCORE
BRESTIMATOR XFULL YFULL SCORINGNEGMEANSQUAREDERROR
CVNSPLITS

COLUMNSFULL DATA

ADD A SINGLE MISSING VALUE TO EACH ROW
XMISSING XFULLCOPY
YMISSING YFULL
MISSINGSAMPLES NPARANGENSAMPLES
MISSINGFEATURES RNGCHOICENFEATURES NSAMPLES REPLACETRUE
XMISSINGMISSINGSAMPLES MISSINGFEATURES NPNAN
ESTIMATE THE SCORE AFTER IMPUTATION MEAN AND MEDIAN STRATEGIES
SCORESIMPLEIMPUTER PDDATAFRAME
FORSTRATEGY INMEAN MEDIAN
ESTIMATOR MAKEPIPELINE
SIMPLEIMPUTERMISSINGVALUESNPNAN STRATEGYSTRATEGY
BRESTIMATOR

SCORESIMPLEIMPUTERSTRATEGY CROSSVALSCORE
ESTIMATOR XMISSING YMISSING SCORINGNEGMEANSQUAREDERROR
CVNSPLITS

ESTIMATE THE SCORE AFTER ITERATIVE IMPUTATION OF THE MISSING VALUES
WITH DIFFERENT ESTIMATORS
ESTIMATORS
BAYESIANRIDGE
DECISIONTREEREGRESSORMAXFEATURESSORT RANDOMSTATE0
EXTRATREESREGRESSORNESTIMATORS10 RANDOMSTATE0
KNEIGHBORSREGRESSORNNEIGHBORS15

SCOREITERATIVEIMPUTER PDDATAFRAME
FORIMPUTEESTIMATOR INESTIMATORS
ESTIMATOR MAKEPIPELINE
ITERATIVEIMPUTERRANDOMSTATE0 ESTIMATORIMPUTEESTIMATOR
BRESTIMATOR

SCOREITERATIVEIMPUTERIMPUTEESTIMATORCLASSNAME
CROSSVALSCORE
ESTIMATOR XMISSING YMISSING SCORINGNEGMEANSQUAREDERROR
CVNSPLITS
516 MISSING VALUE IMPUTATION 1121
```

SCORES PDCONCAT  
SCOREFULLDATA SCORESIMPLEIMPUTER SCOREITERATIVEIMPUTER  
KEYSORIGINAL SIMPLEIMPUTER ITERATIVEIMPUTER AXIS1

PLOT BOSTON RESULTS  
FIG AX PLTSUBPLOTSFIGSIZE13 6  
MEANS SCORESMEAN  
ERRORS SCORESSTD  
MEANSPLOTBARHXERRERRORS AXAX  
AXSETTITLECALIFORNIA HOUSING REGRESSION WITH DIFFERENT IMPUTATION METHODS  
AXSETXLABELMSE SMALLER IS BETTER  
AXSETYTICKSNPARANGEMEANSSSHAPE0  
AXSETYTICKLABELS W JOINLABEL FORLABELINMEANSINDEXGETVALUES  
PLTTIGHTLAYOUTPAD1  
PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 19017 SECONDS  
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE  
5162 IMPUTING MISSING VALUES BEFORE BUILDING AN ESTIMATOR  
MISSING VALUES CAN BE REPLACED BY THE MEAN THE MEDIAN OR THE MOST FREQUENT VALUE USING THE BASIC SKLEARN  
IMPUTESIMPLEIMPUTER THE MEDIAN IS A MORE ROBUST ESTIMATOR FOR DATA WITH HIGH MAGNITUDE VARIABLES WHICH  
COULD DOMINATE RESULTS OTHERWISE KNOWN AS A 'LONG TAIL'  
ANOTHER OPTION IS THE SKLEARNIMPUTEITERATIVEIMPUTER THIS USES ROUNDROBIN LINEAR REGRESSION TREATING  
EVERY VARIABLE AS AN OUTPUT IN TURN THE VERSION IMPLEMENTED ASSUMES GAUSSIAN OUTPUT VARIABLES IF YOUR FEATURES ARE  
OBVIOUSLY NONNORMAL CONSIDER TRANSFORMING THEM TO LOOK MORE NORMAL SO AS TO POTENTIALLY IMPROVE PERFORMANCE  
IN ADDITION OF USING AN IMPUTING METHOD WE CAN ALSO KEEP AN INDICATION OF THE MISSING INFORMATION USING SKLEARN  
IMPUTEMISSINGINDICATOR WHICH MIGHT CARRY SOME INFORMATION

1122 CHAPTER 5 EXAMPLES



```
SCIKITLEARN USER GUIDE RELEASE 0213
PRINTDOC
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
  TO USE THE EXPERIMENTAL ITERATIVEIMPUTER WE NEED TO EXPLICITLY ASK FOR IT
FROM SKLEARNEXPERIMENTAL IMPORT ENABLEITERATIVEIMPUTER  NOQA
FROM SKLEARNDATASETS IMPORT LOADDIABETES
FROM SKLEARNDATASETS IMPORT LOADBOSTON
FROM SKLEARNENSEMBLE IMPORT RANDOMFORESTREGRESSOR
FROM SKLEARNPIPELINE IMPORT MAKEPIPELINE MAKEUNION
FROM SKLEARNIMPUTE IMPORT SIMPLEIMPUTER ITERATIVEIMPUTER MISSINGINDICATOR
FROM SKLEARNMODELSELECTION IMPORT CROSSVALSCORE
RNG  NPRANDOMRANDOMSTATE0
NSPLITS  5
REGRESSOR  RANDOMFORESTREGRESSORRANDOMSTATE0 NESTIMATORS100
DEFGETSCORESFORIMPUTERIMPUTER XMISSING YMISSING
ESTIMATOR  MAKEPIPELINE
MAKEUNIONIMPUTER MISSINGINDICATORMISSINGVALUES0
REGRESSOR
IMPUTESCORES  CROSSVALSCOREESTIMATOR XMISSING YMISSING
SCORINGNEGMEANSQUAREDERROR
CVNSPLITS
RETURNIMPUTESCORES
DEFGETRESULTS DATASET
XFULL YFULL  DATASETDATA DATASETTARGET
NSAMPLES  XFULLSHAPE0
NFEATURES  XFULLSHAPE1
  ESTIMATE THE SCORE ON THE ENTIRE DATASET WITH NO MISSING VALUES
516 MISSING VALUE IMPUTATION 1123
```

SCIKITLEARN USER GUIDE RELEASE 0213  
FULLSCORES CROSSVALSCOREREREgressor XFULL YFULL  
SCORINGNEGMEANSQUAREDERROR  
CVNSPLITS  
ADD MISSING VALUES IN 75 OF THE LINES  
MISSINGRATE 075  
NMISSINGSAMPLES INTNPFLOORNSAMPLES MISSINGRATE  
MISSINGSAMPLES NPHSTACKNPZEROSNSAMPLES NMISSINGSAMPLES  
DTYPENPBOOL  
NPONESNMISSINGSAMPLES  
DTYPENPBOOL  
RNGSHUFFLEMISSINGSAMPLES  
MISSINGFEATURES RNGRANDINTO NFEATURES NMISSINGSAMPLES  
XMISSING XFULLCOPY  
XMISSINGNPWHEREMISSINGSAMPLES0 MISSINGFEATURES 0  
YMISSING YFULLCOPY  
ESTIMATE THE SCORE AFTER REPLACING MISSING VALUES BY 0  
IMPUTER SIMPLEIMPUTERMISsingVALUES0  
STRATEGYCONSTANT  
FILLVALUE0  
ZEROIMPUTEScores GETScoresFORIMPUTERIMPUTER XMISSING YMISSING  
ESTIMATE THE SCORE AFTER IMPUTATION MEAN STRATEGY OF THE MISSING VALUES  
IMPUTER SIMPLEIMPUTERMISsingVALUES0 STRATEGYMEAN  
MEANIMPUTEScores GETScoresFORIMPUTERIMPUTER XMISSING YMISSING  
ESTIMATE THE SCORE AFTER ITERATIVE IMPUTATION OF THE MISSING VALUES  
IMPUTER ITERATIVEIMPUTERMISsingVALUES0  
RANDOMSTATE0  
NNEARESTFEATURES5  
ITERATIVEIMPUTEScores GETScoresFORIMPUTERIMPUTER  
XMISSING  
YMISSING  
RETURNFULLSCORESMEAN FULLSCORESSTD  
ZEROIMPUTEScoresMEAN ZEROIMPUTEScoresSTD  
MEANIMPUTEScoresMEAN MEANIMPUTEScoresSTD  
ITERATIVEIMPUTEScoresMEAN ITERATIVEIMPUTEScoresSTD  
RESULTSdiabetes NPARRAYGETRESULTSLOADdiabetes  
MSESDiabetes RESULTSdiabetes 0 1  
STDSDiabetes RESULTSdiabetes 1  
RESULTSboston NPARRAYGETRESULTSLOADboston  
MSESBoston RESULTSboston 0 1  
STDSBoston RESULTSboston 1  
NBARS LENMSESDiabetes  
XVAL NPARANGENBARS  
XLABELS FULL DATA  
ZERO IMPUTATION  
MEAN IMPUTATION  
MULTIVARIATE IMPUTATION  
COLORS R G B ORANGE  
1124 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

PLOT DIABETES RESULTS

PLTFIGUREFIGSIZE12 6

AX1 PLTSUBPLOT121

FORJINXVAL

AX1BARHJ MSESIDIABETESJ XERRSTDSDIABETESJ

COLORCOLORSJ ALPHA06 ALIGNCENTER

AX1SETTITLEIMPUTATION TECHNIQUES WITH DIABETES DATA

AX1SETXLIMLEFTNPMINMSESIDIABETES 09

RIGHTNPMAXMSESIDIABETES 11

AX1SETYTICKSXVAL

AX1SETXLABELMSE

AX1INVERTYAXIS

AX1SETYTICKLABELSXLABELS

PLOT BOSTON RESULTS

AX2 PLTSUBPLOT122

FORJINXVAL

AX2BARHJ MSESIBOSTONJ XERRSTDIBOSTONJ

COLORCOLORSJ ALPHA06 ALIGNCENTER

AX2SETTITLEIMPUTATION TECHNIQUES WITH BOSTON DATA

AX2SETYTICKSXVAL

AX2SETXLABELMSE

AX2INVERTYAXIS

AX2SETYTICKLABELS NBARS

PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 11569 SECONDS

517 INSPECTION

EXAMPLES RELATED TO THE SKLEARNINSPECTION MODULE

NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE

5171 PARTIAL DEPENDENCE PLOTS

PARTIAL DEPENDENCE PLOTS SHOW THE DEPENDENCE BETWEEN THE TARGET FUNCTION2AND A SET OF ‘TARGET’ FEATURES MARGINALIZING OVER THE VALUES OF ALL OTHER FEATURES THE COMPLEMENT FEATURES DUE TO THE LIMITS OF HUMAN PERCEPTION THE SIZE OF THE TARGET FEATURE SET MUST BE SMALL USUALLY ONE OR TWO THUS THE TARGET FEATURES ARE USUALLY CHOSEN AMONG THE MOST IMPORTANT FEATURES

THIS EXAMPLE SHOWS HOW TO OBTAIN PARTIAL DEPENDENCE PLOTS FROM A MLPREGRESSOR AND A GRADIENTBOOSTINGREGRESSOR TRAINED ON THE CALIFORNIA HOUSING DATASET THE EXAMPLE IS TAKEN FROM1

THE PLOTS SHOW FOUR 1WAY AND TWO 1WAY PARTIAL DEPENDENCE PLOTS OMMITTED FOR MLPREGRESSOR DUE TO COMPUTA TION TIME THE TARGET VARIABLES FOR THE ONEWAY PDP ARE MEDIAN INCOME MEDINC AVERAGE OCCUPANTS PER HOUSEHOLD AVGOCUP MEDIAN HOUSE AGE HOUSEAGE AND AVERAGE ROOMS PER HOUSEHOLD AVEROOMS

2FOR CLASSIFICATION YOU CAN THINK OF IT AS THE REGRESSION SCORE BEFORE THE LINK FUNCTION

1T HASTIE R TIBSHIRANI AND J FRIEDMAN “ELEMENTS OF STATISTICAL LEARNING ED 2” SPRINGER 2009

517 INSPECTION 1125

SCIKITLEARN USER GUIDE RELEASE 0213

WE CAN CLEARLY SEE THAT THE MEDIAN HOUSE PRICE SHOWS A LINEAR RELATIONSHIP WITH THE MEDIAN INCOME TOP LEFT AND THAT THE HOUSE PRICE DROPS WHEN THE AVERAGE OCCUPANTS PER HOUSEHOLD INCREASES TOP MIDDLE THE TOP RIGHT PLOT SHOWS THAT THE HOUSE AGE IN A DISTRICT DOES NOT HAVE A STRONG INFLUENCE ON THE MEDIAN HOUSE PRICE SO DOES THE AVERAGE ROOMS PER HOUSEHOLD THE TICK MARKS ON THE XAXIS REPRESENT THE DECILES OF THE FEATURE VALUES IN THE TRAINING DATA WE ALSO OBSERVE THAT MLPREGRESSOR HAS MUCH SMOOTHER PREDICTIONS THAN GRADIENTBOOSTINGREGRESSOR FOR THE PLOTS TO BE COMPARABLE IT IS NECESSARY TO SUBTRACT THE AVERAGE VALUE OF THE TARGET Y THE 'RECURSION' METHOD USED BY DEFAULT FOR GRADIENTBOOSTINGREGRESSOR DOES NOT ACCOUNT FOR THE INITIAL PREDICTOR IN OUR CASE THE AVERAGE TARGET SETTING THE TARGET AVERAGE TO 0 AVOIDS THIS BIAS PARTIAL DEPENDENCE PLOTS WITH TWO TARGET FEATURES ENABLE US TO VISUALIZE INTERACTIONS AMONG THEM THE TOWWAY PARTIAL DEPENDENCE PLOT SHOWS THE DEPENDENCE OF MEDIAN HOUSE PRICE ON JOINT VALUES OF HOUSE AGE AND AVERAGE OCCUPANTS PER HOUSEHOLD WE CAN CLEARLY SEE AN INTERACTION BETWEEN THE TWO FEATURES FOR AN AVERAGE OCCUPANCY GREATER THAN TWO THE HOUSE PRICE IS NEARLY INDEPENDENT OF THE HOUSE AGE WHEREAS FOR VALUES LESS THAN TWO THERE IS A STRONG DEPENDENCE ON AGE ON A THIRD FIGURE WE HAVE PLOTTED THE SAME PARTIAL DEPENDENCE PLOT THIS TIME IN 3 DIMENSIONS

SCIKITLEARN USER GUIDE RELEASE 0213

- 517 INSPECTION 1127

SCIKITLEARN USER GUIDE RELEASE 0213

•

OUT

TRAINING MLPREGRESSOR

COMPUTING PARTIAL DEPENDENCE PLOTS

TRAINING GRADIENTBOOSTINGREGRESSOR

COMPUTING PARTIAL DEPENDENCE PLOTS

CUSTOM 3D PLOT VIA PARTIALDEPENDENCE

PRINTDOC

IMPORT NUMPY AS NP

IMPORT MATPLOTLIBPYPLOT AS PLT

FROM MPLTOOLKITSMPLPLOT3D IMPORT AXES3D

FROM SKLEARNINSPECTION IMPORT PARTIALDEPENDENCE

FROM SKLEARNINSPECTION IMPORT PLOTPARTIALDEPENDENCE

FROM SKLEARNENSEMBLE IMPORT GRADIENTBOOSTINGREGRESSOR

FROM SKLEARNNEURALNETWORK IMPORT MLPREGRESSOR

FROM SKLEARNDATASETS CALIFORNIAHOUSING IMPORT FETCHCALIFORNIAHOUSING

DEFMAIN

1128 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
CALHOUSING FETCHCALIFORNIAHOUSING  
X Y CALHOUSINGDATA CALHOUSINGTARGET  
NAMES CALHOUSINGFEATURENAMES  
CENTER TARGET TO AVOID GRADIENT BOOSTING INIT BIAS GRADIENT BOOSTING  
WITH THE RECURSION METHOD DOES NOT ACCOUNT FOR THE INITIAL ESTIMATOR  
HERE THE AVERAGE TARGET BY DEFAULT  
Y YMEAN  
PRINTTRAINING MLPREGRESSOR  
EST MLPREGRESSORACTIVATIONLOGISTIC  
ESTFITX Y  
PRINTCOMPUTING PARTIAL DEPENDENCE PLOTS  
WE DONT COMPUTE THE 2WAY PDP 5 1 HERE BECAUSE IT IS A LOT SLOWER  
WITH THE BRUTE METHOD  
FEATURES 0 5 1 2  
PLOTPARTIALDEPENDENCEEST X FEATURES FEATURENAMESNAMES  
NJOBS3 GRIDRESOLUTION50  
FIG PLTGCF  
FIGSUPTITLEPARTIAL DEPENDENCE OF HOUSE VALUE ON NONLOCATION FEATURES N  
FOR THE CALIFORNIA HOUSING DATASET WITH MLPREGRESSOR  
PLTSUBPLOTSADJUSTTOP09 TIGHTLAYOUT CAUSES OVERLAP WITH SUPTITLE  
PRINTTRAINING GRADIENTBOOSTINGREGRESSOR  
EST GRADIENTBOOSTINGREGRESSORNESTIMATORS100 MAXDEPTH4  
LEARNINGRATE01 LOSSHUBER  
RANDOMSTATE1  
ESTFITX Y  
PRINTCOMPUTING PARTIAL DEPENDENCE PLOTS  
FEATURES 0 5 1 2 5 1  
PLOTPARTIALDEPENDENCEEST X FEATURES FEATURENAMESNAMES  
NJOBS3 GRIDRESOLUTION50  
FIG PLTGCF  
FIGSUPTITLEPARTIAL DEPENDENCE OF HOUSE VALUE ON NONLOCATION FEATURES N  
FOR THE CALIFORNIA HOUSING DATASET WITH GRADIENT BOOSTING  
PLTSUBPLOTSADJUSTTOP09  
PRINTCUSTOM 3D PLOT VIA PARTIALDEPENDENCE  
FIG PLTFigure  
TARGETFEATURE 1 5  
PDP AXES PARTIALDEPENDENCEEST X TARGETFEATURE  
GRIDRESOLUTION50  
XX YY NPMESHGRIDAXES0 AXES1  
Z PDP0T  
AX AXES3DFIG  
SURF AXPLOTSURFACEXX YY Z RSTRIDE1 CSTRIDE1  
CMAPPLTCMBUPU EDGECOLORK  
AXSETXLABELNAMESTARGETFEATURE0  
AXSETYLABELNAMESTARGETFEATURE1  
AXSETZLABELPARTIAL DEPENDENCE  
PRETTY INIT VIEW  
AXVIEWINITELEV22 AZIM122  
PLTCOLORBARSURF  
PLTSUPTITLEPARTIAL DEPENDENCE OF HOUSE VALUE ON MEDIAN N  
AGE AND AVERAGE OCCUPANCY WITH GRADIENT BOOSTING  
PLTSUBPLOTSADJUSTTOP09  
517 INSPECTION 1129

SCIKITLEARN USER GUIDE RELEASE 0213  
PLTSHOW  
NEEDED ON WINDOWS BECAUSE PLOTPARTIALDEPENDENCE USES MULTIPROCESSING  
IFNAME MAIN  
MAIN  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 24838 SECONDS  
518 GENERALIZED LINEAR MODELS  
EXAMPLES CONCERNING THE SKLEARNLINEARMODEL MODULE  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
5181 LASSO PATH USING LARS  
COMPUTES LASSO PATH ALONG THE REGULARIZATION PARAMETER USING THE LARS ALGORITHM ON THE DIABETES DATASET EACH  
COLOR REPRESENTS A DIFFERENT FEATURE OF THE COEFFICIENT VECTOR AND THIS IS DISPLAYED AS A FUNCTION OF THE REGULARIZATION  
PARAMETER  
1130 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213  
OUT  
COMPUTING REGULARIZATION PATH USING THE LARS

PRINTDOC  
AUTHOR FABIAN PEDREGOSA FABIANPEDREGOSAINRIA  
ALEXANDRE GRAMFORT ALEXANDREGRAMFORTINRIA  
LICENSE BSD 3 CLAUSE

IMPORT NUMPY AS NP  
IMPORT MATPLOTLIBPYPLOT AS PLT  
FROM SKLEARN IMPORT LINEARMODEL  
FROM SKLEARN IMPORT DATASETS  
DIABETES DATASETSLOADDIABETES  
X DIABETESDATA  
Y DIABETESTARGET

518 GENERALIZED LINEAR MODELS 1131

SCIKITLEARN USER GUIDE RELEASE 0213  
PRINTCOMPUTING REGULARIZATION PATH USING THE LARS  
COEFS LINEARMODELLARSPATHX Y METHODLASSO VERBOSETRUE  
XX NPSUMNPABSCOEFS1  
XX XX1  
PLTPLOTXX COEFST  
YMIN YMAX PLTYLIM  
PLTVLINESXX YMIN YMAX LINESTYLEDASHED  
PLTXLABELCOEF MAXCOEF  
PLTYLABELCOEFFICIENTS  
PLTTITLELASSO PATH  
PLTAXISTIGHT  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0025 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
5182 PLOT RIDGE COEFFICIENTS AS A FUNCTION OF THE REGULARIZATION  
SHOWS THE EFFECT OF COLLINEARITY IN THE COEFFICIENTS OF AN ESTIMATOR  
RIDGE REGRESSION IS THE ESTIMATOR USED IN THIS EXAMPLE EACH COLOR REPRESENTS A DIFFERENT FEATURE OF THE COEFFICIENT  
VECTOR AND THIS IS DISPLAYED AS A FUNCTION OF THE REGULARIZATION PARAMETER  
THIS EXAMPLE ALSO SHOWS THE USEFULNESS OF APPLYING RIDGE REGRESSION TO HIGHLY ILLCONDITIONED MATRICES FOR SUCH  
MATRICES A SLIGHT CHANGE IN THE TARGET VARIABLE CAN CAUSE HUGE VARIANCES IN THE CALCULATED WEIGHTS IN SUCH CASES IT IS  
USEFUL TO SET A CERTAIN REGULARIZATION ALPHA TO REDUCE THIS VARIATION NOISE  
WHEN ALPHA IS VERY LARGE THE REGULARIZATION EFFECT DOMINATES THE SQUARED LOSS FUNCTION AND THE COEFFICIENTS TEND TO  
ZERO AT THE END OF THE PATH AS ALPHA TENDS TOWARD ZERO AND THE SOLUTION TENDS TOWARDS THE ORDINARY LEAST SQUARES  
COEFFICIENTS EXHIBIT BIG OSCILLATIONS IN PRACTISE IT IS NECESSARY TO TUNE ALPHA IN SUCH A WAY THAT A BALANCE IS MAINTAINED  
BETWEEN BOTH  
1132 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
AUTHOR FABIAN PEDREGOSA  FABIANPEDREGOSAINRIAFR
LICENSE BSD 3 CLAUSE
PRINTDOC
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARN IMPORT LINEARMODEL
X IS THE 10X10 HILBERT MATRIX
X 1 NPARANGE1 11 NPARANGE0 10 NPNEWAXIS
Y NPONES10

COMPUTE PATHS
NALPHAS 200
ALPHAS NPLOGSPACE10 2 NALPHAS
COEFS
FORAINALPHAS
RIDGE LINEARMODELRIDGEALPHA FITINTERCEPTFALSE
RIDGEFITX Y
COEFSAPPENDRIDGECOEf

518 GENERALIZED LINEAR MODELS 1133
```

SCIKITLEARN USER GUIDE RELEASE 0213

  DISPLAY RESULTS

AX  PLTGCA

AXPLOTALPHAS COEFS

AXSETXSCALELOG

AXSETXLIMAXGETXLIM1  REVERSE AXIS

PLTXLABELALPHA

PLTYLABELWEIGHTS

PLTTITLERIDGE COEFFICIENTS AS A FUNCTION OF THE REGULARIZATION

PLTAXISTIGHT

PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT  0 MINUTES 0161 SECONDS

NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)

5183 SGD MAXIMUM MARGIN SEPARATING HYPERPLANE

PLOT THE MAXIMUM MARGIN SEPARATING HYPERPLANE WITHIN A TWOCCLASS SEPARABLE DATASET USING A LINEAR SUPPORT VECTOR

MACHINES CLASSIFIER TRAINED USING SGD

1134 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
PRINTDOC
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARNLINEARMODEL IMPORT SGDCLASSIFIER
FROM SKLEARNDATASETSSAMPLESGENERATOR IMPORT MAKEBLOBS
    WE CREATE 50 SEPARABLE POINTS
X Y MAKEBLOBSNSAMPLES50 CENTERS2 RANDOMSTATE0 CLUSTERSTD060
FIT THE MODEL
CLF SGDCLASSIFIERLOSSHINGE ALPHA001 MAXITER200
FITINTERCEPTTRUE TOL1E3
CLFFITX Y
    PLOT THE LINE THE POINTS AND THE NEAREST VECTORS TO THE PLANE
XX NPLinspace1 5 10
YY NPLinspace1 5 10
X1 X2 NPMESHGRIDXX YY
Z NPEMPTYX1SHAPE
FOR J VAL INNPNDENUMERATEX1
X1 VAL
X2 X2I J
P CLFDECISIONFUNCTIONX1 X2
ZI J P0
LEVELS 10 00 10
LINESTYLES DASHED SOLID DASHED
COLORS K
PLTCONTOURX1 X2 Z LEVELS COLORSCOLORS LINESTYLESLINESTYLES
PLTSCATTERX 0 X 1 CY CMAPPLTCMPAIED
EDGECOLORBLACK S20
PLTAXISTIGHT
PLTSHOW
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0020 SECONDS
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE
5184 SGD CONVEX LOSS FUNCTIONS
A PLOT THAT COMPARES THE VARIOUS CONVEX LOSS FUNCTIONS SUPPORTED BY SKLEARNLINEARMODEL
SGDCLASSIFIER
518 GENERALIZED LINEAR MODELS 1135
```

```
SCIKITLEARN USER GUIDE RELEASE 0213
PRINTDOC
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
DEFMODIFIEDHUBERLOSSYTRUE YPRED
Z`YPRED YTRUE
LOSS`4`Z
LOSSZ`1`1`ZZ`1`2
LOSSZ`1`0
RETURNLOSS
XMIN XMAX`4`4
XX`NPLINSPACEXMIN XMAX 100
LW`2
PLTPLOTXMIN 0 0 XMAX 1 1 0 0 COLORGOLD LWLW
LABELZEROONE LOSS
PLTPLOTXX NPWHEREXX`1`1`XX 0 COLORTEAL LWLW
LABELHINGE LOSS
PLTPLOTXX NPMINIMUMXX 0 COLORYELLOWGREEN LWLW
LABELPERCEPTRON LOSS
PLTPLOTXX NPLOG21`NPEXPXX COLORCORNFLOWERBLUE LWLW
LABELLOG LOSS
PLTPLOTXX NPWHEREXX`1`1`XX 0 2 COLORORANGE LWLW
1136 CHAPTER 5 EXAMPLES
```

SCIKITLEARN USER GUIDE RELEASE 0213

LABELSQUARED HINGE LOSS

PLTPLOTXX MODIFIEDHUBERLOSSXX 1 COLORDARKORCHID LWLW

LINESTYLE LABELMODIFIED HUBER LOSS

PLTYLIM0 8

PLTLEGENDLOCUPPER RIGHT

PLTXLABELRDECISION FUNCTION FX

PLTYLABELLY1 FX

PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0019 SECONDS

NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE

5185 ORDINARY LEAST SQUARES AND RIDGE REGRESSION VARIANCE

DUE TO THE FEW POINTS IN EACH DIMENSION AND THE STRAIGHT LINE THAT LINEAR REGRESSION USES TO FOLLOW THESE POINTS AS WELL AS IT CAN NOISE ON THE OBSERVATIONS WILL CAUSE GREAT VARIANCE AS SHOWN IN THE FIRST PLOT EVERY LINE’S SLOPE CAN VARY QUITE A BIT FOR EACH PREDICTION DUE TO THE NOISE INDUCED IN THE OBSERVATIONS

RIDGE REGRESSION IS BASICALLY MINIMIZING A PENALISED VERSION OF THE LEASTSQUARED FUNCTION THE PENALISING SHRINKS THE VALUE OF THE REGRESSION COEFFICIENTS DESPITE THE FEW DATA POINTS IN EACH DIMENSION THE SLOPE OF THE PREDICTION IS MUCH MORE STABLE AND THE VARIANCE IN THE LINE ITSELF IS GREATLY REDUCED IN COMPARISON TO THAT OF THE STANDARD LINEAR REGRESSION

- 

518 GENERALIZED LINEAR MODELS 1137

```
•
PRINTDOC
CODE SOURCE GAËL VAROQUAUX
MODIFIED FOR DOCUMENTATION BY JAQUES GROBLER
LICENSE BSD 3 CLAUSE
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARN IMPORT LINEARMODEL
XTRAIN NPC5 1T
YTRAIN 5 1
XTEST NPC0 2T
NPRANDOMSEED0
CLASSIFIERS DICTOLSLINEARMODELLINEARREGRESSION
RIDGE LINEARMODEL RIDGEALPHA1
FORNAME CLF INCLASSIFIERSITEMS
FIG AX PLTSUBPLOTSFIGSIZE4 3
FORINRANGE6
THISX 1 NPRANDOMNORMALSIZE2 1 XTRAIN
CLFFITTHISX YTRAIN
AXPLOTXTEST CLFPREDICTXTEST COLORGRAY
AXSCATTERTHISX YTRAIN S3 CGRAY MARKERO ZORDER10
CLFFITXTRAIN YTRAIN
AXPLOTXTEST CLFPREDICTXTEST LINEWIDTH2 COLORBLUE
AXSCATTERXTRAIN YTRAIN S30 CRED MARKER ZORDER10
AXSETTITLENAME
AXSETXLIM0 2
1138 CHAPTER 5 EXAMPLES
```



SCIKITLEARN USER GUIDE RELEASE 0213

AXSETYLIM0 16

AXSETXLABELX

AXSETYLABELY

FIGTIGHTLAYOUT

PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0130 SECONDS

NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)

5186 PLOT RIDGE COEFFICIENTS AS A FUNCTION OF THE L2 REGULARIZATION

RIDGE REGRESSION IS THE ESTIMATOR USED IN THIS EXAMPLE EACH COLOR IN THE LEFT PLOT REPRESENTS ONE DIFFERENT DIMENSION OF THE COEFFICIENT VECTOR AND THIS IS DISPLAYED AS A FUNCTION OF THE REGULARIZATION PARAMETER THE RIGHT PLOT SHOWS HOW EXACT THE SOLUTION IS THIS EXAMPLE ILLUSTRATES HOW A WELL DEFINED SOLUTION IS FOUND BY RIDGE REGRESSION AND HOW REGULARIZATION AFFECTS THE COEFFICIENTS AND THEIR VALUES THE PLOT ON THE RIGHT SHOWS HOW THE DIFFERENCE OF THE COEFFICIENTS FROM THE ESTIMATOR CHANGES AS A FUNCTION OF REGULARIZATION

IN THIS EXAMPLE THE DEPENDENT VARIABLE  $Y$  IS SET AS A FUNCTION OF THE INPUT FEATURES  $XW$   $C$  THE COEFFICIENT VECTOR  $W$  IS RANDOMLY SAMPLED FROM A NORMAL DISTRIBUTION WHEREAS THE BIAS TERM  $C$  IS SET TO A CONSTANT

AS  $\alpha$  TENDS TOWARD ZERO THE COEFFICIENTS FOUND BY RIDGE REGRESSION STABILIZE TOWARDS THE RANDOMLY SAMPLED VECTOR  $W$  FOR BIG  $\alpha$  STRONG REGULARISATION THE COEFFICIENTS ARE SMALLER EVENTUALLY CONVERGING AT 0 LEADING TO A SIMPLER AND BIASED SOLUTION THESE DEPENDENCIES CAN BE OBSERVED ON THE LEFT PLOT

THE RIGHT PLOT SHOWS THE MEAN SQUARED ERROR BETWEEN THE COEFFICIENTS FOUND BY THE MODEL AND THE CHOSEN VECTOR  $W$  LESS REGULARISED MODELS RETRIEVE THE EXACT COEFFICIENTS ERROR IS EQUAL TO 0 STRONGER REGULARISED MODELS INCREASE THE ERROR

PLEASE NOTE THAT IN THIS EXAMPLE THE DATA IS NONNOISY HENCE IT IS POSSIBLE TO EXTRACT THE EXACT COEFFICIENTS

AUTHOR KORNEL KIELCZEWSKI KORNELKPLUSNETPL

PRINTDOC

IMPORT MATPLOTLIBPYPLOT AS PLT

IMPORT NUMPY AS NP

FROM SKLEARNDATASETS IMPORT MAKEREGRESSION

FROM SKLEARNLINEARMODEL IMPORT RIDGE

518 GENERALIZED LINEAR MODELS 1139

SCIKITLEARN USER GUIDE RELEASE 0213  
FROM SKLEARNMETRICS IMPORT MEANSQUAREDERROR  
CLF RIDGE  
X Y W MAKEREGRESSIONNSAMPLES10 NFEATURES10 COEFTRUE  
RANDOMSTATE1 BIAS35  
COEFS  
ERRORS  
ALPHAS NPLOGSPACE6 6 200  
TRAIN THE MODEL WITH DIFFERENT REGULARISATION STRENGTHS  
FORAINALPHAS  
CLFSETPARAMSALPHAA  
CLFFITX Y  
COEFSAPPENDCLFCOEF  
ERRORSAPPENDMEANSQUAREDERRORCLFCOEF W  
DISPLAY RESULTS  
PLTFIGUREFIGSIZE20 6  
PLTSUBPLOT121  
AX PLTGCA  
AXPLOTALPHAS COEFS  
AXSETXSCALELOG  
PLTXLABELALPHA  
PLTYLABELWEIGHTS  
PLTTITLERIDGE COEFFICIENTS AS A FUNCTION OF THE REGULARIZATION  
PLTAXISTIGHT  
PLTSUBPLOT122  
AX PLTGCA  
AXPLOTALPHAS ERRORS  
AXSETXSCALELOG  
PLTXLABELALPHA  
PLTYLABELERROR  
PLTTITLECOEFFICIENT ERROR AS A FUNCTION OF THE REGULARIZATION  
PLTAXISTIGHT  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0126 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
5187 SGD PENALTIES  
CONTOURS OF WHERE THE PENALTY IS EQUAL TO 1 FOR THE THREE PENALTIES L1 L2 AND ELASTICNET  
ALL OF THE ABOVE ARE SUPPORTED BY SKLEARNLINEARMODELSTOCHASTICGRADIENT  
1140 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
PRINTDOC
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
L1COLOR NAVY
L2COLOR C
ELASTICNETCOLOR DARKORANGE
LINE NPLinspace15 15 1001
XX YY NPMESHGRIDLINE LINE
L2 XX 2 YY2
L1 NPABSXX NPABSY
RHO 05
518 GENERALIZED LINEAR MODELS 1141
```

SCIKITLEARN USER GUIDE RELEASE 0213  
ELASTICNET RHO L1 1 RHO L2  
PLTFIGUREFIGSIZE10 10 DPI100  
AX PLTGCA  
ELASTICNETCONTOUR PLTCONTOURXX YY ELASTICNET LEVELS1  
COLORSELASTICNETCOLOR  
L2CONTOUR PLTCONTOURXX YY L2 LEVELS1 COLORSL2COLOR  
L1CONTOUR PLTCONTOURXX YY L1 LEVELS1 COLORSL1COLOR  
AXSETASPECTEQUAL  
AXSPINESLEFTSETPOSITIONCENTER  
AXSPINESRIGHTSETCOLORNONE  
AXSPINESBOTTOMSETPOSITIONCENTER  
AXSPINESTOPSETCOLORNONE  
PLTCLABELELASTICNETCONTOUR INLINE1 FONTSIZE18  
FMT10 ELASTICNET MANUAL1 1  
PLTCLABELL2CONTOUR INLINE1 FONTSIZE18  
FMT10 L2 MANUAL1 1  
PLTCLABELL1CONTOUR INLINE1 FONTSIZE18  
FMT10 L1 MANUAL1 1  
PLTTIGHTLAYOUT  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0247 SECONDS  
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE  
5188 REGULARIZATION PATH OF L1 LOGISTIC REGRESSION  
TRAIN L1PENALIZED LOGISTIC REGRESSION MODELS ON A BINARY CLASSIFICATION PROBLEM DERIVED FROM THE IRIS DATASET  
THE MODELS ARE ORDERED FROM STRONGEST REGULARIZED TO LEAST REGULARIZED THE 4 COEFFICIENTS OF THE MODELS ARE COLLECTED  
AND PLOTTED AS A “REGULARIZATION PATH” ON THE LEFTHAND SIDE OF THE FIGURE STRONG REGULARIZERS ALL THE COEFFICIENTS ARE  
EXACTLY 0 WHEN REGULARIZATION GETS PROGRESSIVELY LOOSER COEFFICIENTS CAN GET NONZERO VALUES ONE AFTER THE OTHER  
HERE WE CHOOSE THE SAGA SOLVER BECAUSE IT CAN EFFICIENTLY OPTIMIZE FOR THE LOGISTIC REGRESSION LOSS WITH A NON  
SMOOTH SPARSITY INDUCING L1 PENALTY  
ALSO NOTE THAT WE SET A LOW VALUE FOR THE TOLERANCE TO MAKE SURE THAT THE MODEL HAS CONVERGED BEFORE COLLECTING THE  
COEFFICIENTS  
WE ALSO USE WARMSTARTTRUE WHICH MEANS THAT THE COEFFICIENTS OF THE MODELS ARE REUSED TO INITIALIZE THE NEXT MODEL  
FIT TO SPEEDUP THE COMPUTATION OF THE FULLPATH  
1142 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
OUT  
COMPUTING REGULARIZATION PATH  
THIS TOOK 2008S  
PRINTDOC  
AUTHOR ALEXANDRE GRAMFORT ALEXANDREGAMFORTINRIA.FR  
LICENSE BSD 3 CLAUSE  
FROM TIME IMPORT TIME  
IMPORT NUMPY AS NP  
IMPORT MATPLOTLIBPYPLOT AS PLT  
FROM SKLEARN IMPORT LINEARMODEL  
FROM SKLEARN IMPORT DATASETS  
FROM SKLEARN.SVM IMPORT L1MINC  
IRIS DATASETSLOADIRIS  
X IRISDATA  
Y IRISTARGET  
518 GENERALIZED LINEAR MODELS 1143

SCIKITLEARN USER GUIDE RELEASE 0213

X XY 2  
Y YY 2  
X XMAX NORMALIZE X TO SPEEDUP CONVERGENCE

DEMO PATH FUNCTIONS  
CS L1MINCX Y LOSSLOG NPLOGSPACE0 7 16  
PRINTCOMPUTING REGULARIZATION PATH  
START TIME  
CLF LINEARMODELLOGISTICREGRESSIONPENALTYL1 SOLVERSAGA  
TOL1E6 MAXITERINT1E6  
WARMSTARTTRUE  
COEFS  
FORCINCS  
CLFSETPARAMSCC  
CLFFITX Y  
COEFSAPPENDCLFCOEFRAVELCOPY  
PRINTTHIS TOOK 03FS TIME START  
COEFS NPARRAYCOEFS  
PLTPLOTNPLOG10CS COEFS MARKERO  
YMIN YMAX PLTYLIM  
PLTXLABELLOGC  
PLTYLABELCOEFFICIENTS  
PLTTITLELOGISTIC REGRESSION PATH  
PLTAXISTIGHT  
PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 2022 SECONDS

NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE

5189 POLYNOMIAL INTERPOLATION

THIS EXAMPLE DEMONSTRATES HOW TO APPROXIMATE A FUNCTION WITH A POLYNOMIAL OF DEGREE NDEGREE BY USING RIDGE REGRESSION CONCRETELY FROM NSAMPLES 1D POINTS IT SUFFICES TO BUILD THE VANDERMONDE MATRIX WHICH IS NSAMPLES X NDEGREE1 AND HAS THE FOLLOWING FORM

1 X1 X1 2 X1 3 1 X2 X2 2 X2 3

INTUITIVELY THIS MATRIX CAN BE INTERPRETED AS A MATRIX OF PSEUDO FEATURES THE POINTS RAISED TO SOME POWER THE MATRIX IS AKIN TO BUT DIFFERENT FROM THE MATRIX INDUCED BY A POLYNOMIAL KERNEL

THIS EXAMPLE SHOWS THAT YOU CAN DO NONLINEAR REGRESSION WITH A LINEAR MODEL USING A PIPELINE TO ADD NONLINEAR FEATURES KERNEL METHODS EXTEND THIS IDEA AND CAN INDUCE VERY HIGH EVEN INFINITE DIMENSIONAL FEATURE SPACES

SCIKITLEARN USER GUIDE RELEASE 0213  
PRINTDOC  
AUTHOR MATHIEU BLONDEL  
JAKE VANDERPLAS  
LICENSE BSD 3 CLAUSE  
IMPORT NUMPY AS NP  
IMPORT MATPLOTLIBPYPLOT AS PLT  
FROM SKLEARNLINEARMODEL IMPORT RIDGE  
FROM SKLEARNPREPROCESSING IMPORT POLYNOMIALFEATURES  
FROM SKLEARNPIPELINE IMPORT MAKEPIPELINE  
DEFFX  
FUNCTION TO APPROXIMATE BY POLYNOMIAL INTERPOLATION  
RETURNXNPSINX  
GENERATE POINTS USED TO PLOT  
XPLOT NPLinspace0 10 100  
GENERATE POINTS AND KEEP A SUBSET OF THEM  
X NPLinspace0 10 100  
RNG NPRandomRandomState0  
RNGSHUFFLEX  
518 GENERALIZED LINEAR MODELS 1145

SCIKITLEARN USER GUIDE RELEASE 0213  
X NPSORTX20  
Y FX  
CREATE MATRIX VERSIONS OF THESE ARRAYS  
X X NPNEWAXIS  
XPLOT XPLOT NPNEWAXIS  
COLORS TEAL YELLOWGREEN GOLD  
LW 2  
PLTPLOTXPLOT FXPLOT COLORCORNFLOWERBLUE LINEWIDTHLW  
LABELGROUND TRUTH  
PLTSCATTERX Y COLORNAVY S30 MARKERO LABELTRAINING POINTS  
FORCOUNT DEGREE INENUMERATE3 4 5  
MODEL MAKEPIPELINEPOLYNOMIALFEATURESDEGREE RIDGE  
MODELFITX Y  
YPLOT MODEL PREDICTXPLOT  
PLTPLOTXPLOT YPLOT COLORCOLORSCOUNT LINEWIDTHLW  
LABELDEGREE D DEGREE  
PLTLEGENDLOCLOWER LEFT  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0020 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
51810 LOGISTIC FUNCTION  
SHOWN IN THE PLOT IS HOW THE LOGISTIC REGRESSION WOULD IN THIS SYNTHETIC DATASET CLASSIFY VALUES AS EITHER 0 OR 1 IE  
CLASS ONE OR TWO USING THE LOGISTIC CURVE  
1146 CHAPTER 5 EXAMPLES



```
SCIKITLEARN USER GUIDE RELEASE 0213
PRINTDOC
CODE SOURCE GAEI VAROQUAUX
LICENSE BSD 3 CLAUSE
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARN IMPORT LINEARMODEL
FROM SCIPYSPECIAL IMPORT EXPIT
GENERAL A TOY DATASETS ITS JUST A STRAIGHT LINE WITH SOME GAUSSIAN NOISE
XMIN XMAX 5 5
NSAMPLES 100
NPRANDOMSEED0
X NPRANDOMNORMALSIZENSAMPLES
Y X 0ASTYPENPFLOAT
XX 0 4
X 3 NPRANDOMNORMALSIZENSAMPLES
X X NPNEWAXIS
FIT THE CLASSIFIER
CLF LINEARMODELLOGISTICREGRESSIONC1E5 SOLVERLBFGS
CLFFITX Y
AND PLOT THE RESULT
PLTFigure1 FIGSIZE4 3
PLTCLF
PLTSCATTERXRavel Y COLORBLACK ZORDER20
XTEST NPLINSPACE5 10 300
LOSS EXPITXTEST CLFCOEF CLFINTERCEPTRavel
PLTPLOTXTEST LOSS COLORRED LINEWIDTH3
OLS LINEARMODELLINEARREGRESSION
OLSFITX Y
PLTPLOTXTEST OLSOEF XTEST OLSINTERCEPT LINEWIDTH1
PLTAXHLINE5 COLOR5
PLTYLABELY
PLTXLABELX
PLXTICKSRANGE5 10
PLTYTICKS0 05 1
PLTYLIM25 125
PLTXLIM4 10
PLTLEGENDLOGISTIC REGRESSION MODEL LINEAR REGRESSION MODEL
LOCLOWER RIGHT FONTSIZESMALL
PLTTIGHTLAYOUT
PLTSHOW
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0046 SECONDS
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE
518 GENERALIZED LINEAR MODELS 1147
```

```
SCIKITLEARN USER GUIDE RELEASE 0213
51811 SGD WEIGHTED SAMPLES
PLOT DECISION FUNCTION OF A WEIGHTED DATASET WHERE THE SIZE OF POINTS IS PROPORTIONAL TO ITS WEIGHT
PRINTDOC
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARN IMPORT LINEARMODEL
    WE CREATE 20 POINTS
NPRANDOMSEED0
X NPRNPRANDOMRANDN10 2 1 1 NPRANDOMRANDN10 2
Y 1 10 1 10
SAMPLEWEIGHT 100 NPABSNPRANDOMRANDN20
    AND ASSIGN A BIGGER WEIGHT TO THE LAST 10 SAMPLES
SAMPLEWEIGHT10 10
    PLOT THE WEIGHTED DATA POINTS
XX YY NPMESHGRIDNPLINSPACE4 5 500 NPLINSPACE4 5 500
PLTFigure
PLTSCATTERX 0 X 1 CY SSAMPLEWEIGHT ALPHA09
CMAPPLTCMBONE EDGEColorBLACK
    FIT THE UNWEIGHTED MODEL
CLF LINEARMODELSGDCLASSIFIERALPHA001 MAXITER100 TOL1E3
1148 CHAPTER 5 EXAMPLES
```

```
SCIKITLEARN USER GUIDE RELEASE 0213
CLFFITX Y
Z CLFDECISIONFUNCTIONNPCXXRAVEL YYRAVEL
Z ZRESHAPEXXSHAPE
NOWEIGHTS PLTCONTOURXX YY Z LEVELS0 LINESYLESSOLID
FIT THE WEIGHTED MODEL
CLF LINEARMODELSGDCLASSIFIERALPHA001 MAXITER100 TOL1E3
CLFFITX Y SAMPLEWEIGHTSAMPWEIGHT
Z CLFDECISIONFUNCTIONNPCXXRAVEL YYRAVEL
Z ZRESHAPEXXSHAPE
SAMPLESWEIGHTS PLTCONTOURXX YY Z LEVELS0 LINESYLEDASHED
PLTLEGENDNOWEIGHTSCOLLECTIONS0 SAMPLESWEIGHTSCOLLECTIONS0
NO WEIGHTS WITH WEIGHTS LOCLOWER LEFT
PLXTXTICKS
PLTYTICKS
PLTSHOW
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0086 SECONDS
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE
51812 LOGISTIC REGRESSION 3CLASS CLASSIFIER
SHOW BELOW IS A LOGISTICREGRESSION CLASSIFIERS DECISION BOUNDARIES ON THE FIRST TWO DIMENSIONS SEPAL LENGTH AND WIDTH
OF THE IRIS DATASET THE DATAPOINTS ARE COLORED ACCORDING TO THEIR LABELS
PRINTDOC
CODE SOURCE GAËL VAROQUAUX
MODIFIED FOR DOCUMENTATION BY JAQUES GROBLER
LICENSE BSD 3 CLAUSE
518 GENERALIZED LINEAR MODELS 1149
```

SCIKITLEARN USER GUIDE RELEASE 0213

```
import numpy as np
import matplotlib.pyplot as plt
from sklearn.linear_model import LogisticRegression
from sklearn import datasets
import sys

iris = datasets.load_iris()
X = iris.data[:, 2:]  # we only take the first two features
y = iris.target

logreg = LogisticRegression(C=1e5, solver='lbfgs', multi_class='multinomial')
logreg.fit(X, y)

# Plot the decision boundary for that we will assign a color to each
# point in the mesh
xmin, xmax, ymin, ymax = X.min(), X.max(), y.min(), y.max()
h = 0.02  # step size in the mesh
xx, yy = np.meshgrid(np.arange(xmin, xmax, h), np.arange(ymin, ymax, h))
z = logreg.predict(np.c_[xx.ravel(), yy.ravel()])
# Put the result into a color plot
z = z.reshape(xx.shape)
plt.figure(1, figsize=(4, 3))
plt.pcolormesh(xx, yy, z, cmap=plt.cm.Paired)
# Plot also the training points
plt.scatter(X[:, 0], X[:, 1], c=y, edgecolors='k', cmap=plt.cm.Paired)
plt.xlabel('sepal length')
plt.ylabel('sepal width')
plt.xlim(xx.min(), xx.max())
plt.ylim(yy.min(), yy.max())
plt.xticks(())
plt.yticks(())
plt.show()
```

TOTAL RUNNING TIME OF THE SCRIPT: 0 MINUTES 0083 SECONDS

[NOTE: CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)

51813 LINEAR REGRESSION EXAMPLE

THIS EXAMPLE USES THE ONLY THE FIRST FEATURE OF THE DIABETES DATASET IN ORDER TO ILLUSTRATE A TWODIMENSIONAL PLOT OF THIS REGRESSION TECHNIQUE. THE STRAIGHT LINE CAN BE SEEN IN THE PLOT SHOWING HOW LINEAR REGRESSION ATTEMPTS TO DRAW A STRAIGHT LINE THAT WILL BEST MINIMIZE THE RESIDUAL SUM OF SQUARES BETWEEN THE OBSERVED RESPONSES IN THE DATASET AND THE RESPONSES PREDICTED BY THE LINEAR APPROXIMATION. THE COEFFICIENTS, THE RESIDUAL SUM OF SQUARES AND THE VARIANCE SCORE ARE ALSO CALCULATED.

1150 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
OUT  
COEFFICIENTS  
93823786125  
MEAN SQUARED ERROR 254807  
VARIANCE SCORE 047  
PRINTDOC  
CODE SOURCE JAQUES GROBLER  
LICENSE BSD 3 CLAUSE  
IMPORT MATPLOTLIBPYPLOT AS PLT  
IMPORT NUMPY AS NP  
FROM SKLEARN IMPORT DATASETS LINEARMODEL  
FROM SKLEARNMETRICS IMPORT MEANSQUAREDERROR R2SCORE  
LOAD THE DIABETES DATASET  
DIABETES DATASETSLOADDIABETES  
518 GENERALIZED LINEAR MODELS 1151

SCIKITLEARN USER GUIDE RELEASE 0213  
USE ONLY ONE FEATURE  
DIABETESX DIABETESDATA NPNEWAXIS 2  
SPLIT THE DATA INTO TRAININGTESTING SETS  
DIABETESXTRAIN DIABETESX20  
DIABETESXTEST DIABETESX20  
SPLIT THE TARGETS INTO TRAININGTESTING SETS  
DIABETESYTRAIN DIABETESTARGET20  
DIABETESYTEST DIABETESTARGET20  
CREATE LINEAR REGRESSION OBJECT  
REGR LINEARMODELLINEARREGRESSION  
TRAIN THE MODEL USING THE TRAINING SETS  
REGRFITDIABETESXTRAIN DIABETESYTRAIN  
MAKE PREDICTIONS USING THE TESTING SET  
DIABETESYPRED REGRPREDICTDIABETESXTEST  
THE COEFFICIENTS  
PRINTCOEFFICIENTS N REGRCOEF  
THE MEAN SQUARED ERROR  
PRINTMEAN SQUARED ERROR 2F  
MEANSQUAREDERRORDIABETESYTEST DIABETESYPRED  
EXPLAINED VARIANCE SCORE 1 IS PERFECT PREDICTION  
PRINTVARIANCE SCORE 2F R2SCOREDIABETESYTEST DIABETESYPRED  
PLOT OUTPUTS  
PLTSCATTERDIABETESXTEST DIABETESYTEST COLORBLACK  
PLTPLOTDIABETESXTEST DIABETESYPRED COLORBLUE LINEWIDTH3  
PLXTTICKS  
PLTYTICKS  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0175 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
51814 ROBUST LINEAR MODEL ESTIMATION USING RANSAC  
IN THIS EXAMPLE WE SEE HOW TO ROBUSTLY FIT A LINEAR MODEL TO FAULTY DATA USING THE RANSAC ALGORITHM  
1152 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
OUT
ESTIMATED COEFFICIENTS TRUE LINEAR REGRESSION RANSAC
821903908407869 5417236387 8208533159
IMPORT NUMPY AS NP
FROM MATPLOTLIB IMPORT PYPLLOTASPLT
FROM SKLEARN IMPORT LINEARMODEL DATASETS
NSAMPLES 1000
NOUTLIERS 50
X Y COEF DATASETSMAKEREGRESSIONNSAMPLESNSAMPLES NFEATURES1
NINFORMATIVE1 NOISE10
COEFTRUE RANDOMSTATE0
ADD OUTLIER DATA
NPRANDOMSEED0
518 GENERALIZED LINEAR MODELS 1153
```

SCIKITLEARN USER GUIDE RELEASE 0213

XNOUTLIERS 3 05 NPRANDOMNORMALSIZENOUTLIERS 1

YNOUTLIERS 3 10 NPRANDOMNORMALSIZENOUTLIERS

FIT LINE USING ALL DATA

LR LINEARMODELLINEARREGRESSION

LRFITX Y

ROBUSTLY FIT LINEAR MODEL WITH RANSAC ALGORITHM

RANSAC LINEARMODEL RANSACREGRESSOR

RANSACFITX Y

INLIERMASK RANSACINLIERMASK

OUTLIERMASK NPLOGICALNOTINLIERMASK

PREDICT DATA OF ESTIMATED MODELS

LINEX NPARANGEXMIN XMAX NPNEWAXIS

LINEY LRPREDICTLINEX

LINEYRANSAC RANSACPREDICTLINEX

COMPARE ESTIMATED COEFFICIENTS

PRINTESTIMATED COEFFICIENTS TRUE LINEAR REGRESSION RANSAC

PRINTCOEF LR COEF RANSACESTIMATOR COEF

LW 2

PLTSCATTERXINLIERMASK YINLIERMASK COLORYELLOWGREEN MARKER

LABELINLIERS

PLTSCATTERXOUTLIERMASK YOUTLIERMASK COLORGOLD MARKER

LABELOUTLIERS

PLTPLOTLINEX LINEY COLORNAVY LINEWIDTHLW LABELLINEAR REGRESSOR

PLTPLOTLINEX LINEYRANSAC COLORCORNFLOWERBLUE LINEWIDTHLW

LABELRANSAC REGRESSOR

PLTLEGENDLOCLOWER RIGHT

PLTXLABELINPUT

PLTYLABELRESPONSE

PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0028 SECONDS

NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE

51815 SPARSITY EXAMPLE FITTING ONLY FEATURES 1 AND 2

FEATURES 1 AND 2 OF THE DIABETESDATASET ARE FITTED AND PLOTTED BELOW IT ILLUSTRATES THAT ALTHOUGH FEATURE 2 HAS A STRONG COEFFICIENT ON THE FULL MODEL IT DOES NOT GIVE US MUCH REGARDING Y WHEN COMPARED TO JUST FEATURE 1

1154 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213

- 
- 

518 GENERALIZED LINEAR MODELS 1155

SCIKITLEARN USER GUIDE RELEASE 0213

```
•
PRINTDOC
CODE SOURCE GAËL VAROQUAUX
MODIFIED FOR DOCUMENTATION BY JAQUES GROBLER
LICENSE BSD 3 CLAUSE
IMPORT MATPLOTLIBPYPLOT AS PLT
IMPORT NUMPY AS NP
FROM MPLTOOLKITSMPLPLOT3D IMPORT AXES3D
FROM SKLEARN IMPORT DATASETS LINEARMODEL
DIABETES DATASETSLOADDIABETES
INDICES 0 1
XTRAIN DIABETESDATA20 INDICES
XTEST DIABETESDATA20 INDICES
YTRAIN DIABETESTARGET20
YTEST DIABETESTARGET20
OLS LINEARMODELLINEARREGRESSION
OLSFITXTRAIN YTRAIN

PLOT THE FIGURE
DEFPLOTFIGSFIGNUM ELEV AZIM XTRAIN CLF
FIG PLTFIGUREFIGNUM FIGSIZE4 3
PLTCLF
AX AXES3DFIG ELEVELEV AZIMAZIM
AXSCATTERXTRAIN 0 XTRAIN 1 YTRAIN CK MARKER
AXPLOTSURFACENPARRAY1 1 15 15
NPARRAY1 15 1 15
CLFPREDICTNPARRAY1 1 15 15
1 15 1 15T
RESHAPE2 2
1156 CHAPTER 5 EXAMPLES
```

SCIKITLEARN USER GUIDE RELEASE 0213  
ALPHA5  
AXSETXLABELX1  
AXSETYLABELX2  
AXSETZLABELY  
AXWXAXISSETTICKLABELS  
AXWYAXISSETTICKLABELS  
AXWZAXISSETTICKLABELS  
GENERATE THE THREE DIFFERENT FIGURES FROM DIFFERENT VIEWS  
ELEV 435  
AZIM 110  
PLOTFIGS1 ELEV AZIM XTRAIN OLS  
ELEV 5  
AZIM 0  
PLOTFIGS2 ELEV AZIM XTRAIN OLS  
ELEV 5  
AZIM 90  
PLOTFIGS3 ELEV AZIM XTRAIN OLS  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0197 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
51816 LASSO ON DENSE AND SPARSE DATA  
WE SHOW THAT LINEARMODELLASSO PROVIDES THE SAME RESULTS FOR DENSE AND SPARSE DATA AND THAT IN THE CASE OF SPARSE  
DATA THE SPEED IS IMPROVED  
OUT  
DENSE MATRICES  
SPARSE LASSO DONE IN 0186067S  
DENSE LASSO DONE IN 0034536S  
DISTANCE BETWEEN COEFFICIENTS 9043732562018544E14  
SPARSE MATRICES  
MATRIX DENSITY 062630000000000001  
SPARSE LASSO DONE IN 0284597S  
DENSE LASSO DONE IN 0955330S  
DISTANCE BETWEEN COEFFICIENTS 7344760355532163E12  
PRINTDOC  
FROM TIME IMPORT TIME  
FROM SCIPY IMPORT SPARSE  
FROM SCIPY IMPORT LINALG  
518 GENERALIZED LINEAR MODELS 1157

SCIKITLEARN USER GUIDE RELEASE 0213  
FROM SKLEARN DATASET SAMPLES GENERATOR  
IMPORT MAKE REGRESSION  
FROM SKLEARN LINEAR MODEL IMPORT LASSO

THE TWO LASSO IMPLEMENTATIONS ON DENSE DATA  
PRINT DENSE MATRICES  
X Y MAKE REGRESSION N SAMPLES 200 N FEATURES 5000 RANDOM STATE 0  
XSP SPARSE COO MATRIX X  
ALPHA 1  
SPARSE LASSO LASSO ALPHA ALPHA FIT INTERCEPT FALSE MAX ITER 1000  
DENSE LASSO LASSO ALPHA ALPHA FIT INTERCEPT FALSE MAX ITER 1000  
T0 TIME  
SPARSE LASSO FIT XSP Y  
PRINT SPARSE LASSO DONE IN FS TIME T0  
T0 TIME  
DENSE LASSO FIT X Y  
PRINT DENSE LASSO DONE IN FS TIME T0  
PRINT DISTANCE BETWEEN COEFFICIENTS S  
LINALGNORM SPARSE LASSO COEF DENSE LASSO COEF

THE TWO LASSO IMPLEMENTATIONS ON SPARSE DATA  
PRINT SPARSE MATRICES  
XS X COPY  
XSXS 25 00  
XS SPARSE COO MATRIX XS  
XS X STO CSC  
PRINT MATRIX DENSITY S XSNNZ FLOAT XS SIZE 100  
ALPHA 01  
SPARSE LASSO LASSO ALPHA ALPHA FIT INTERCEPT FALSE MAX ITER 10000  
DENSE LASSO LASSO ALPHA ALPHA FIT INTERCEPT FALSE MAX ITER 10000  
T0 TIME  
SPARSE LASSO FIT XS Y  
PRINT SPARSE LASSO DONE IN FS TIME T0  
T0 TIME  
DENSE LASSO FIT X STO ARRAY Y  
PRINT DENSE LASSO DONE IN FS TIME T0  
PRINT DISTANCE BETWEEN COEFFICIENTS S  
LINALGNORM SPARSE LASSO COEF DENSE LASSO COEF  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 1565 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
1158 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

51817 HUBERREGRESSOR VS RIDGE ON DATASET WITH STRONG OUTLIERS

FIT RIDGE AND HUBERREGRESSOR ON A DATASET WITH OUTLIERS

THE EXAMPLE SHOWS THAT THE PREDICTIONS IN RIDGE ARE STRONGLY INFLUENCED BY THE OUTLIERS PRESENT IN THE DATASET THE HUBER REGRESSOR IS LESS INFLUENCED BY THE OUTLIERS SINCE THE MODEL USES THE LINEAR LOSS FOR THESE AS THE PARAMETER EPSILON IS INCREASED FOR THE HUBER REGRESSOR THE DECISION FUNCTION APPROACHES THAT OF THE RIDGE

AUTHORS MANOJ KUMAR MKS542NYUEDU

LICENSE BSD 3 CLAUSE

PRINTDOC

IMPORT NUMPY AS NP

IMPORT MATPLOTLIBPYPLOT AS PLT

FROM SKLEARNDATASETS IMPORT MAKEREGRESSION

FROM SKLEARNLINEARMODEL IMPORT HUBERREGRESSOR RIDGE

GENERATE TOY DATA

RNG NPRANDOMRANDOMSTATE0

X Y MAKEREGRESSIONNNSAMPLES20 NFEATURES1 RANDOMSTATE0 NOISE40

BIAS1000

ADD FOUR STRONG OUTLIERS TO THE DATASET

XOUTLIERS RNGNORMAL0 05 SIZE4 1

518 GENERALIZED LINEAR MODELS 1159

SCIKITLEARN USER GUIDE RELEASE 0213

YOUTLIERS RNGNORMAL0 20 SIZE4

XOUTLIERS2 XMAX XMEAN 4

XOUTLIERS2 XMIN XMEAN 4

YOUTLIERS2 YMIN YMEAN 4

YOUTLIERS2 YMAX YMEAN 4

X NPVSTACKX XOUTLIERS

Y NPCONCATENATEY YOUTLIERS

PLTPLOTX Y B

FIT THE HUBER REGRESSOR OVER A SERIES OF EPSILON VALUES

COLORS R B Y M

X NPLINSPACEXMIN XMAX 7

EPSILONVALUES 135 15 175 19

FORK EPSILON INENUMERATEEPSILONVALUES

HUBER HUBERREGRESSORFITINTERCEPTTRUE ALPHA00 MAXITER100

EPSILONEPSILON

HUBERFITX Y

COEF HUBERCOEF X HUBERINTERCEPT

PLTPLOTX COEF COLORSK LABELHUBER LOSS S EPSILON

FIT A RIDGE REGRESSOR TO COMPARE IT TO HUBER REGRESSOR

RIDGE RIDGEFITINTERCEPTTRUE ALPHA00 RANDOMSTATE0 NORMALIZETRUE

RIDGEFITX Y

COEFRIDGE RIDGECOEf

COEF RIDGECOEf X RIDGEINTERCEPT

PLTPLOTX COEF G LABELRIDGE REGRESSION

PLTTITLECOMPARISON OF HUBERREGRESSOR VS RIDGE

PLTXLABELX

PLTYLABELY

PLTLEGENDLOCO

PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0032 SECONDS

NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE

51818 JOINT FEATURE SELECTION WITH MULTITASK LASSO

THE MULTITASK LASSO ALLOWS TO FIT MULTIPLE REGRESSION PROBLEMS JOINTLY ENFORCING THE SELECTED FEATURES TO BE THE SAME ACROSS TASKS THIS EXAMPLE SIMULATES SEQUENTIAL MEASUREMENTS EACH TASK IS A TIME INSTANT AND THE RELEVANT FEATURES VARY IN AMPLITUDE OVER TIME WHILE BEING THE SAME THE MULTITASK LASSO IMPOSES THAT FEATURES THAT ARE SELECTED AT ONE TIME POINT ARE SELECT FOR ALL TIME POINT THIS MAKES FEATURE SELECTION BY THE LASSO MORE STABLE

1160 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

- 
- 

518 GENERALIZED LINEAR MODELS 1161

```
SCIKITLEARN USER GUIDE RELEASE 0213
PRINTDOC
AUTHOR ALEXANDRE GRAMFORT ALEXANDREGRAMFORTINRIA.FR
LICENSE BSD 3 CLAUSE
IMPORT MATPLOTLIB.PY.PLOT AS PLT
IMPORT NUMPY AS NP
FROM SKLEARN.LINEAR.MODEL IMPORT MULTITASK.LASSO LASSO
RNG = NPY.RANDOM.RANDOM.STATE42
GENERATE SOME 2D COEFFICIENTS WITH SINE WAVES WITH RANDOM FREQUENCY AND PHASE
NSAMPLES, NFEATURES, NTASKS = 100, 30, 40
NRELEVANTFEATURES = 5
COEF = NP.ZEROS(NTASKS, NFEATURES)
TIMES = NP.Linspace(0, 2, NP.I, NTASKS)
FOR K IN RANGE(NRELEVANTFEATURES):
    COEF[K, :] = NP.SIN(1 + RNG.RANDN(1) * TIMES) * 3 * RNG.RANDN(1)
X = RNG.RANDN(NSAMPLES, NFEATURES)
Y = NP.DOT(X, COEF) + RNG.RANDN(NSAMPLES, NTASKS)
COEF.LASSO = NP.ARRAY.LASSO.ALPHA05.FIT(X, Y, COEF)
COEF.MULTITASK.LASSO = MULTITASK.LASSO.ALPHA1.FIT(X, Y, COEF)

PLOT SUPPORT AND TIME SERIES
FIG, PLT.FIGURE.FIGSIZES(8, 5)
PLT.SUBPLOT(1, 2, 1)
PLT.SPY(COEF.LASSO)
PLT.XLABEL('FEATURE')
PLT.YLABEL('TIME OR TASK')
PLT.TEXT(10, 5, 'LASSO')
PLT.SUBPLOT(1, 2, 2)
PLT.SPY(COEF.MULTITASK.LASSO)
PLT.XLABEL('FEATURE')
PLT.YLABEL('TIME OR TASK')
PLT.TEXT(10, 5, 'MULTITASK.LASSO')
FIG.SUPP.TITLE('COEFFICIENT NONZERO LOCATION')
FEATURE_TO_PLOT = 0
PLT.FIGURE
LW = 2
PLT.PLOT(COEF, FEATURE_TO_PLOT, COLOR='SEAGREEN', LINEWIDTH=LW)
LABEL_GROUND_TRUTH
PLT.PLOT(COEF.LASSO, FEATURE_TO_PLOT, COLOR='CORNFLOWERBLUE', LINEWIDTH=LW)
LABEL_LASSO
PLT.PLOT(COEF.MULTITASK.LASSO, FEATURE_TO_PLOT, COLOR='GOLD', LINEWIDTH=LW)
LABEL_MULTITASK.LASSO
PLT.LEGEND(LOC='UPPER CENTER')
PLT.AXIS.TIGHT()
PLT.YLIM(1, 11)
PLT.SHOW()
TOTAL RUNNING TIME OF THE SCRIPT: 0 MINUTES 0061 SECONDS
1162 CHAPTER 5 EXAMPLES
```



SCIKITLEARN USER GUIDE RELEASE 0213  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
51819 COMPARING VARIOUS ONLINE SOLVERS  
AN EXAMPLE SHOWING HOW DIFFERENT ONLINE SOLVERS PERFORM ON THE HANDWRITTEN DIGITS DATASET  
OUT  
TRAINING SGD  
TRAINING ASGD  
TRAINING PERCEPTRON  
TRAINING PASSIVEAGGRESSIVE I  
TRAINING PASSIVEAGGRESSIVE II  
TRAINING SAG  
AUTHOR ROB ZINKOV ROB AT ZINKOV DOT COM  
LICENSE BSD 3 CLAUSE  
518 GENERALIZED LINEAR MODELS 1163

SCIKITLEARN USER GUIDE RELEASE 0213

```
import numpy as np
import matplotlib.pyplot as plt
from sklearn import datasets
from sklearn.model_selection import train_test_split
from sklearn.linear_model import SGDClassifier, Perceptron
from sklearn.linear_model import PassiveAggressiveClassifier
from sklearn.linear_model import LogisticRegression

heldout = 0.95
n_rounds = 20
digits = datasets.load_digits()
x, y, digit_target = digits.data, digits.target, digits.target
classifiers = [
    SGDClassifier(max_iter=100, tol=1e-3),
    SGDClassifier(average=True, max_iter=1000, tol=1e-3),
    Perceptron(tol=1e-3),
    PassiveAggressiveClassifier(loss='hinge', C=10, tol=1e-4),
    PassiveAggressiveClassifier(loss='squared_hinge', C=10, tol=1e-4),
    LogisticRegression(solver='sag', tol=1e-1, C=1e4, max_iter=1000, multi_class='auto')
]

xx = 1
np.random.seed(42)
for i in range(n_rounds):
    x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.1, random_state=42)
    clf = classifiers[xx % len(classifiers)]
    clf.fit(x_train, y_train)
    y_pred = clf.predict(x_test)
    np.mean(y_pred == y_test)
    plt.plot(xx, np.mean(y_pred == y_test), label=clf.__class__.__name__)
    plt.legend(loc='upper right')
    plt.xlabel('Proportion Train')
    plt.ylabel('Test Error Rate')
    plt.show()

total_running_time = 0
note_click_here_to_download_the_full_example_code = 51820
orthogonal_matching_pursuit = 51820
orthogonal_matching_pursuit_for_recovering_a_sparse_signal_from_a_noisy_measurement_encoded_with_a_dictionary = 1164
chapter_5_examples = 1164
```

```
SCIKITLEARN USER GUIDE RELEASE 0213
PRINTDOC
IMPORT MATPLOTLIBPYPLOT AS PLT
IMPORT NUMPY AS NP
FROM SKLEARNLINEARMODEL IMPORT ORTHOGONALMATCHINGPURSUIT
FROM SKLEARNLINEARMODEL IMPORT ORTHOGONALMATCHINGPURSUITCV
FROM SKLEARNDATASETS IMPORT MAKESPARSECODEDSIGNAL
NCOMPONENTS NFEATURES 512 100
NNONZEROCOEFS 17
GENERATE THE DATA
Y XW
X0 NNONZEROCOEFS
518 GENERALIZED LINEAR MODELS 1165
```

SCIKITLEARN USER GUIDE RELEASE 0213  
Y X W MAKESPARSECODEDSIGNALNSAMPLES1  
NCOMPONENTSNCOMPONENTS  
NFEATURESNFEATURES  
NNONZEROCOEFSSNNONZEROCOEFSS  
RANDOMSTATE0  
IDX WNONZERO  
DISTORT THE CLEAN SIGNAL  
YNOISY Y 005 NPRANDOMRANDNLENY  
PLOT THE SPARSE SIGNAL  
PLTFIGUREFIGSIZE7 7  
PLTSUBPLOT4 1 1  
PLTXLIM0 512  
PLTTITLESPARSE SIGNAL  
PLTSTEMIDX WIDX  
PLOT THE NOISEFREE RECONSTRUCTION  
OMP ORTHOGONALMATCHINGPURSUITNNONZEROCOEFSSNNONZEROCOEFSS  
OMPFITX Y  
COEF OMPCOEF  
IDXR COEFNONZERO  
PLTSUBPLOT4 1 2  
PLTXLIM0 512  
PLTTITLERECOVERED SIGNAL FROM NOISEFREE MEASUREMENTS  
PLTSTEMIDXR COEFIDXR  
PLOT THE NOISY RECONSTRUCTION  
OMPFITX YNOISY  
COEF OMPCOEF  
IDXR COEFNONZERO  
PLTSUBPLOT4 1 3  
PLTXLIM0 512  
PLTTITLERECOVERED SIGNAL FROM NOISY MEASUREMENTS  
PLTSTEMIDXR COEFIDXR  
PLOT THE NOISY RECONSTRUCTION WITH NUMBER OF NONZEROS SET BY CV  
OMPCV ORTHOGONALMATCHINGPURSUITCVCV5  
OMPCVFITX YNOISY  
COEF OMPCVCOEF  
IDXR COEFNONZERO  
PLTSUBPLOT4 1 4  
PLTXLIM0 512  
PLTTITLERECOVERED SIGNAL FROM NOISY MEASUREMENTS WITH CV  
PLTSTEMIDXR COEFIDXR  
PLTSUBPLOTSADJUST006 004 094 090 020 038  
PLTSUPTITLESPARSE SIGNAL RECOVERY WITH ORTHOGONAL MATCHING PURSUIT  
FONTSIZE16  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0499 SECONDS  
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE  
1166 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

51821 MNIST CLASSFICATION USING MULTINOMIAL LOGISTIC L1

HERE WE FIT A MULTINOMIAL LOGISTIC REGRESSION WITH L1 PENALTY ON A SUBSET OF THE MNIST DIGITS CLASSIFICATION TASK WE USE THE SAGA ALGORITHM FOR THIS PURPOSE THIS A SOLVER THAT IS FAST WHEN THE NUMBER OF SAMPLES IS SIGNIFICANTLY LARGER THAN THE NUMBER OF FEATURES AND IS ABLE TO FINELY OPTIMIZE NONSMOOTH OBJECTIVE FUNCTIONS WHICH IS THE CASE WITH THE L1PENALTY TEST ACCURACY REACHES 08 WHILE WEIGHT VECTORS REMAINS SPARSE AND THEREFORE MORE EASILY INTERPRETABLE NOTE THAT THIS ACCURACY OF THIS L1PENALIZED LINEAR MODEL IS SIGNIFICANTLY BELOW WHAT CAN BE REACHED BY AN L2PENALIZED LINEAR MODEL OR A NONLINEAR MULTILAYER PERCEPTRON MODEL ON THIS DATASET

OUT

SPARSITY WITH L1 PENALTY 8080

TEST SCORE WITH L1 PENALTY 08351

EXAMPLE RUN IN 28654 S

IMPORT TIME

IMPORT MATPLOTLIBPYPLOT AS PLT

IMPORT NUMPY AS NP

FROM SKLEARNDATASETS IMPORT FETCHOPENML

FROM SKLEARNLINEARMODEL IMPORT LOGISTICREGRESSION

FROM SKLEARNMODELSELECTION IMPORT TRAINTESTSPLIT

FROM SKLEARNPREPROCESSING IMPORT STANDARDSCALER

FROM SKLEARNUTILS IMPORT CHECKRANDOMSTATE

PRINTDOC

AUTHOR ARTHUR MENSCH ARTHURMENSCHM4XORG

LICENSE BSD 3 CLAUSE

518 GENERALIZED LINEAR MODELS 1167

SCIKITLEARN USER GUIDE RELEASE 0213  
TURN DOWN FOR FASTER CONVERGENCE  
T0 TIMETIME  
TRAINSAMPLES 5000  
LOAD DATA FROM HTTPSWWWOPENMLORG554  
X Y FETCHOPENMLMNIST784 VERSION1 RETURNXYTRUE  
RANDOMSTATE CHECKRANDOMSTATE0  
PERMUTATION RANDOMSTATEPERMUTATIONXSHAPE0  
X XPERMUTATION  
Y YPERMUTATION  
X XRESHAPEXSHAPE0 1  
XTRAIN XTEST YTRAIN YTEST TRAINTESTSPLIT  
X Y TRAINSIZETRAINSSAMPLES TESTSIZE10000  
SCALER STANDARDSCALER  
XTRAIN SCALERFITTRANSFORMXTRAIN  
XTEST SCALERTRANSFORMXTEST  
TURN UP TOLERANCE FOR FASTER CONVERGENCE  
CLF LOGISTICREGRESSIONC50 TRAINSSAMPLES  
MULTICLASSMULTINOMIAL  
PENALTYL1 SOLVERSAGA TOL01  
CLFFITXTRAIN YTRAIN  
SPARSITY NPMEANCLFCOEFF 0 100  
SCORE CLFSCOREXTEST YTEST  
PRINTBEST C 4F CLFC  
PRINTSPARSITY WITH L1 PENALTY 2F SPARSITY  
PRINTTEST SCORE WITH L1 PENALTY 4F SCORE  
COEFF CLFCOEFCOPY  
PLTFIGUREFIGSIZE10 5  
SCALE NPABSCOEFFMAX  
FORIINRANGE10  
L1PLOT PLTSUBPLOT2 5 I 1  
L1PLOTIMSHOWCOEFFIRESHAPE28 28 INTERPOLATIONNEAREST  
CMAPPLTCMRDBU VMINSIZE VMAXSCALE  
L1PLOTSETXTICKS  
L1PLOTSETYXTICKS  
L1PLOTSETXLABELCLASS I I  
PLTSUPTITLECLASSIFICATION VECTOR FOR  
RUNTIME TIMETIME T0  
PRINTEXAMPLE RUN IN 3FS RUNTIME  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 28655 SECONDS  
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE  
51822 PLOT MULTICLASS SGD ON THE IRIS DATASET  
PLOT DECISION SURFACE OF MULTICLASS SGD ON IRIS DATASET THE HYPERPLANES CORRESPONDING TO THE THREE ONEVERSUSALL  
OV A CLASSIFIERS ARE REPRESENTED BY THE DASHED LINES  
1168 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
PRINTDOC
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARN IMPORT DATASETS
FROM SKLEARNLINEARMODEL IMPORT SGDCLASSIFIER
IMPORT SOME DATA TO PLAY WITH
IRIS DATASETSLOADIRIS
WE ONLY TAKE THE FIRST TWO FEATURES WE COULD
AVOID THIS UGLY SLICING BY USING A TWODIM DATASET
X IRISDATA 2
Y IRISTARGET
COLORS BRY
SHUFFLE
IDX NPARANGEXSHAPE0
NPRANDOMSEED13
NPRANDOMSHUFFLEIDX
X XIDX
Y YIDX
STANDARDIZE
MEAN XMEANAXIS0
STD XSTDAXIS0
518 GENERALIZED LINEAR MODELS 1169
```

```
SCIKITLEARN USER GUIDE RELEASE 0213
X X MEAN STD
H 02 STEP SIZE IN THE MESH
CLF SGDCLASSIFIERALPHA0001 MAXITER100 TOL1E3FITX Y
  CREATE A MESH TO PLOT IN
XMIN XMAX X 0MIN 1 X 0MAX 1
YMIN YMAX X 1MIN 1 X 1MAX 1
XX YY NPMESHGRIDNPARANGEXMIN XMAX H
NPARANGEYMIN YMAX H
  PLOT THE DECISION BOUNDARY FOR THAT WE WILL ASSIGN A COLOR TO EACH
  POINT IN THE MESH XMIN XMAXXYMIN YMAX
Z CLFPREDICTNPCXXRAVEL YYRAVEL
  PUT THE RESULT INTO A COLOR PLOT
Z ZRESHAPEXXSHAPE
CS PLTCONTOURFXX YY Z CMAPPLTCMPAIED
PLTAXISTIGHT
  PLOT ALSO THE TRAINING POINTS
FORI COLOR INZIPCLFCLASSES COLORS
IDX NPWHEREY I
PLTSCATTERXIDX 0 XIDX 1 CCOLOR LABELIRISTARGETNAMESI
CMAPPLTCMPAIED EDGECOLORBLACK S20
PLTTITLEDECISION SURFACE OF MULTICLASS SGD
PLTAXISTIGHT
  PLOT THE THREE ONEAGAINSTALL CLASSIFIERS
XMIN XMAX PLTXLIM
YMIN YMAX PLTYLIM
COEF CLFCOEF
INTERCEPT CLFINTERCEPT
DEFPLOTHYPERPLANEC COLOR
DEFLINEX0
RETURNX0COEFC 0 INTERCEPTC COEFC 1
PLTPLOTXMIN XMAX LINEXMIN LINEXMAX
LS COLORCOLOR
FORI COLOR INZIPCLFCLASSES COLORS
PLOTHYPERPLANEI COLOR
PLTLEGEND
PLTSHOW
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0050 SECONDS
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE
1170 CHAPTER 5 EXAMPLES
```



```
SCIKITLEARN USER GUIDE RELEASE 0213
51823 LASSO AND ELASTIC NET FOR SPARSE SIGNALS
ESTIMATES LASSO AND ELASTICNET REGRESSION MODELS ON A MANUALLY GENERATED SPARSE SIGNAL CORRUPTED WITH AN ADDITIVE
NOISE ESTIMATED COEFFICIENTS ARE COMPARED WITH THE GROUNDTRUTH
OUT
LASSOALPHA01
R2 ON TEST DATA 0658064
ELASTICNETALPHA01 L1RATIO07
R2 ON TEST DATA 0642515
PRINTDOC
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARNMETRICS IMPORT R2SCORE

GENERATE SOME SPARSE DATA TO PLAY WITH
518 GENERALIZED LINEAR MODELS 1171
```

```
SCIKITLEARN USER GUIDE RELEASE 0213
NPRANDOMSEED42
NSAMPLES NFEATURES 50 100
X NPRANDOMRANDNNSAMPLES NFEATURES
  DECREASING COEF W ALTERNATED SIGNS FOR VISUALIZATION
IDX NPARANGENFEATURES
COEF 1 IDXNPEXPIDX 10
COEF10 0 SPARSIFY COEF
Y NPDOTX COEF
  ADD NOISE
Y 001 NPRANDOMNORMALSIZENSAMPLES
  SPLIT DATA IN TRAIN SET AND TEST SET
NSAMPLES XSHAPE0
XTRAIN YTRAIN XNSAMPLES 2 YNSAMPLES 2
XTEST YTEST XNSAMPLES 2 YNSAMPLES 2

  LASSO
FROM SKLEARNLINEARMODEL IMPORT LASSO
ALPHA 01
LASSO LASSOALPHAALPHA
YPREDLASSO LASSOFITXTRAIN YTRAINPREDICTXTEST
R2SCORELASSO R2SCOREYTEST YPREDLASSO
PRINTLASSO
PRINTR2 ON TEST DATA F R2SCORELASSO

  ELASTICNET
FROM SKLEARNLINEARMODEL IMPORT ELASTICNET
ENET ELASTICNETALPHAALPHA L1RATIO07
YPREDENET ENETFITXTRAIN YTRAINPREDICTXTEST
R2SCOREENET R2SCOREYTEST YPREDENET
PRINTENET
PRINTR2 ON TEST DATA F R2SCOREENET
M S PLTSTEMNPWHEREENETCOEF0 ENETCOEFENETCOEF 0
MARKERFMTX LABELELASTIC NET COEFFICIENTS
PLTSETPM S COLOR2CA02C
M S PLTSTEMNPWHERELASSOCOEF0 LASSOCOEFCLASSOCOEF 0
MARKERFMTX LABELLASSO COEFFICIENTS
PLTSETPM S COLORFF7F0E
PLTSTEMNPWHERECOEF0 COEFCOEF 0 LABELTRUE COEFFICIENTS
MARKERFMTBX
PLTLEGENDLOCBEST
PLTTITLELASSO R2 3F ELASTIC NET R2 3F
  R2SCORELASSO R2SCOREENET
PLTSHOW
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0071 SECONDS
1172 CHAPTER 5 EXAMPLES
```

SCIKITLEARN USER GUIDE RELEASE 0213

NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)

51824 THEILSEN REGRESSION

COMPUTES A THEILSEN REGRESSION ON A SYNTHETIC DATASET

SEETHEILSEN ESTIMATOR GENERALIZEDMEDIANBASED ESTIMATOR FOR MORE INFORMATION ON THE REGRESSOR

COMPARED TO THE OLS ORDINARY LEAST SQUARES ESTIMATOR THE THEILSEN ESTIMATOR IS ROBUST AGAINST OUTLIERS IT HAS A BREAKDOWN POINT OF ABOUT 293 IN CASE OF A SIMPLE LINEAR REGRESSION WHICH MEANS THAT IT CAN TOLERATE ARBITRARY CORRUPTED DATA OUTLIERS OF UP TO 293 IN THE TWODIMENSIONAL CASE

THE ESTIMATION OF THE MODEL IS DONE BY CALCULATING THE SLOPES AND INTERCEPTS OF A SUBPOPULATION OF ALL POSSIBLE COMBINATIONS OF P SUBSAMPLE POINTS IF AN INTERCEPT IS FITTED P MUST BE GREATER THAN OR EQUAL TO NFEATURES 1 THE FINAL SLOPE AND INTERCEPT IS THEN DEFINED AS THE SPATIAL MEDIAN OF THESE SLOPES AND INTERCEPTS

IN CERTAIN CASES THEILSEN PERFORMS BETTER THAN RANSAC WHICH IS ALSO A ROBUST METHOD THIS IS ILLUSTRATED IN THE SECOND EXAMPLE BELOW WHERE OUTLIERS WITH RESPECT TO THE XAXIS PERTURB RANSAC TUNING THE RESIDUALTHRESHOLD

PARAMETER OF RANSAC REMEDIES THIS BUT IN GENERAL A PRIORI KNOWLEDGE ABOUT THE DATA AND THE NATURE OF THE OUTLIERS IS NEEDED DUE TO THE COMPUTATIONAL COMPLEXITY OF THEILSEN IT IS RECOMMENDED TO USE IT ONLY FOR SMALL PROBLEMS IN TERMS OF NUMBER OF SAMPLES AND FEATURES FOR LARGER PROBLEMS THE MAXSUBPOPULATION PARAMETER RESTRICTS THE MAGNITUDE OF ALL POSSIBLE COMBINATIONS OF P SUBSAMPLE POINTS TO A RANDOMLY CHOSEN SUBSET AND THEREFORE ALSO LIMITS THE RUNTIME THEREFORE THEILSEN IS APPLICABLE TO LARGER PROBLEMS WITH THE DRAWBACK OF LOSING SOME OF ITS MATHEMATICAL PROPERTIES SINCE IT THEN WORKS ON A RANDOM SUBSET

- 

518 GENERALIZED LINEAR MODELS 1173

SCIKITLEARN USER GUIDE RELEASE 0213

```
•
AUTHOR FLORIAN WILHELM FLORIANWILHELMGMAILCOM
LICENSE BSD 3 CLAUSE
IMPORT TIME
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARNLINEARMODEL IMPORT LINEARREGRESSION THEILSENREGRESSOR
FROM SKLEARNLINEARMODEL IMPORT RANSACREGRESSOR
PRINTDOC
ESTIMATORS OLS LINEARREGRESSION
THEILSEN THEILSENREGRESSORRANDOMSTATE42
RANSAC RANSACREGRESSORRANDOMSTATE42
COLORS OLS TURQUOISE THEILSEN GOLD RANSAC LIGHTGREEN
LW 2
```

```
OUTLIERS ONLY IN THE Y DIRECTION
NPRANDOMSEED0
NSAMPLES 200
LINEAR MODEL Y 3 X N2 01 2
X NPRANDOMRANDNNSAMPLES
W 3
C 2
NOISE 01 NPRANDOMRANDNNSAMPLES
1174 CHAPTER 5 EXAMPLES
```

SCIKITLEARN USER GUIDE RELEASE 0213  
Y WX C NOISE  
10 OUTLIERS  
Y20 20 X20  
X X NPNEWAXIS  
PLTSCATTERX Y COLORINDIGO MARKERX S40  
LINEX NPARRAY3 3  
FORNAME ESTIMATOR INESTIMATORS  
T0 TIMETIME  
ESTIMATORFITX Y  
ELAPSEDTIME TIMETIME T0  
YPRED ESTIMATORPREDICTLINEXRESHAPE2 1  
PLTPLOTLINEX YPRED COLORCOLORSNAME LINEWIDTHLW  
LABELSFIT TIME 2FS NAME ELAPSEDTIME  
PLTAXISTIGHT  
PLTLEGENDLOCUPPER LEFT  
PLTTITLECORRUPT Y

OUTLIERS IN THE X DIRECTION  
NPRANDOMSEED0  
LINEAR MODEL Y 3 X N2 01 2  
X NPRANDOMRANDNNSAMPLES  
NOISE 01 NPRANDOMRANDNNSAMPLES  
Y 3X 2 NOISE  
10 OUTLIERS  
X20 99  
Y20 22  
X X NPNEWAXIS  
PLTFigure  
PLTSCATTERX Y COLORINDIGO MARKERX S40  
LINEX NPARRAY3 10  
FORNAME ESTIMATOR INESTIMATORS  
T0 TIMETIME  
ESTIMATORFITX Y  
ELAPSEDTIME TIMETIME T0  
YPRED ESTIMATORPREDICTLINEXRESHAPE2 1  
PLTPLOTLINEX YPRED COLORCOLORSNAME LINEWIDTHLW  
LABELSFIT TIME 2FS NAME ELAPSEDTIME  
PLTAXISTIGHT  
PLTLEGENDLOCUPPER LEFT  
PLTTITLECORRUPT X  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 1359 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
518 GENERALIZED LINEAR MODELS 1175

SCIKITLEARN USER GUIDE RELEASE 0213

51825 PLOT MULTINOMIAL AND ONEVSREST LOGISTIC REGRESSION

PLOT DECISION SURFACE OF MULTINOMIAL AND ONEVSREST LOGISTIC REGRESSION THE HYPERPLANES CORRESPONDING TO THE THREE ONEVSREST OVR CLASSIFIERS ARE REPRESENTED BY THE DASHED LINES

- 

1176 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

•

OUT  
TRAINING SCORE 0995 MULTINOMIAL  
TRAINING SCORE 0976 OVR  
PRINTDOC  
AUTHORS TOM DUPRE LA TOUR TOMDUPRELATOURM4XORG  
LICENSE BSD 3 CLAUSE  
IMPORT NUMPY AS NP  
IMPORT MATPLOTLIBPYPLOT AS PLT  
FROM SKLEARNDATASETS IMPORT MAKEBLOBS  
FROM SKLEARNLINEARMODEL IMPORT LOGISTICREGRESSION  
MAKE 3CLASS DATASET FOR CLASSIFICATION  
CENTERS 5 0 0 15 5 1  
X Y MAKEBLOBSNSAMPLES1000 CENTERSCENTERS RANDOMSTATE40  
TRANSFORMATION 04 02 04 12  
X NPDOTX TRANSFORMATION  
FORMULTICLASS INMULTINOMIAL OVR  
CLF LOGISTICREGRESSIONSOLVERSAG MAXITER100 RANDOMSTATE42  
518 GENERALIZED LINEAR MODELS 1177

SCIKITLEARN USER GUIDE RELEASE 0213  
MULTICLASSMULTICLASSFITX Y  
PRINT THE TRAINING SCORES  
PRINTTRAINING SCORE 3FS CLFSCOREX Y MULTICLASS  
CREATE A MESH TO PLOT IN  
H 02 STEP SIZE IN THE MESH  
XMIN XMAX X 0MIN 1 X 0MAX 1  
YMIN YMAX X 1MIN 1 X 1MAX 1  
XX YY NPMESHGRIDNPARANGEXMIN XMAX H  
NPARANGEYMIN YMAX H  
PLOT THE DECISION BOUNDARY FOR THAT WE WILL ASSIGN A COLOR TO EACH  
POINT IN THE MESH XMIN XMAXXYMIN YMAX  
Z CLFPREDICTNPCXXRAVEL YYRAVEL  
PUT THE RESULT INTO A COLOR PLOT  
Z ZRESHAPEXXSHAPE  
PLTFigure  
PLTCONTOURFXX YY Z CMAPPLTCMPAIED  
PLTTITLEDECISION SURFACE OF LOGISTICREGRESSION S MULTICLASS  
PLTAXISTIGHT  
PLOT ALSO THE TRAINING POINTS  
COLORS BRY  
FORI COLOR INZIPCLFCLASSES COLORS  
IDX NPWHEREY I  
PLTSCATTERXIDX 0 XIDX 1 CCOLOR CMAPPLTCMPAIED  
EDGECOLORBLACK S20  
PLOT THE THREE ONEAGAINSTALL CLASSIFIERS  
XMIN XMAX PLTXLIM  
YMIN YMAX PLTYLIM  
COEF CLFCOEF  
INTERCEPT CLFINTERCEPT  
DEFPLOTHYPERPLANEI COLOR  
DEFLINEX0  
RETURNX0COEFC 0 INTERCEPTC COEFC 1  
PLTPLOTXMIN XMAX LINEXMIN LINEXMAX  
LS COLORCOLOR  
FORI COLOR INZIPCLFCLASSES COLORS  
PLOTHYPERPLANEI COLOR  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0262 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
51826 ROBUST LINEAR ESTIMATOR FITTING  
HERE A SINE FUNCTION IS FIT WITH A POLYNOMIAL OF ORDER 3 FOR VALUES CLOSE TO ZERO  
ROBUST FITTING IS DEMOED IN DIFFERENT SITUATIONS  
1178 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213

- NO MEASUREMENT ERRORS ONLY MODELLING ERRORS FITTING A SINE WITH A POLYNOMIAL
- MEASUREMENT ERRORS IN X
- MEASUREMENT ERRORS IN Y

THE MEDIAN ABSOLUTE DEVIATION TO NON CORRUPT NEW DATA IS USED TO JUDGE THE QUALITY OF THE PREDICTION  
WHAT WE CAN SEE THAT

- RANSAC IS GOOD FOR STRONG OUTLIERS IN THE Y DIRECTION
- THEILSEN IS GOOD FOR SMALL OUTLIERS BOTH IN DIRECTION X AND Y BUT HAS A BREAK POINT ABOVE WHICH IT PERFORMS WORSE THAN OLS
- THE SCORES OF HUBERREGRESSOR MAY NOT BE COMPARED DIRECTLY TO BOTH THEILSEN AND RANSAC BECAUSE IT DOES NOT ATTEMPT TO COMPLETELY FILTER THE OUTLIERS BUT LESSEN THEIR EFFECT
- 

518 GENERALIZED LINEAR MODELS 1179

SCIKITLEARN USER GUIDE RELEASE 0213

- 
- 

1180 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

- 
- 

FROM MATPLOTLIB IMPORT PYPLLOTASPLT  
IMPORT NUMPY AS NP  
FROM SKLEARNLINEARMODEL IMPORT  
518 GENERALIZED LINEAR MODELS 1181

SCIKITLEARN USER GUIDE RELEASE 0213  
 LINEARREGRESSION THEILSENREGRESSOR RANSACREGRESSOR HUBERREGRESSOR  
 FROM SKLEARNMETRICS IMPORT MEANSQUAREDERROR  
 FROM SKLEARNPREPROCESSING IMPORT POLYNOMIALFEATURES  
 FROM SKLEARNPIPELINE IMPORT MAKEPIPELINE  
 NPRANDOMSEED42  
 X NPRANDOMNORMALSIZE400  
 Y NPSINX  
 MAKE SURE THAT IT X IS 2D  
 X X NPNEWAXIS  
 XTEST NPRANDOMNORMALSIZE200  
 YTEST NPSINXTEST  
 XTEST XTEST NPNEWAXIS  
 YERRORS YCOPY  
 YERRORS3 3  
 XERRORS XCOPY  
 XERRORS3 3  
 YERRORSLARGE YCOPY  
 YERRORSLARGE3 10  
 XERRORSLARGE XCOPY  
 XERRORSLARGE3 10  
 ESTIMATORS OLS LINEARREGRESSION  
 THEILSEN THEILSENREGRESSORRANDOMSTATE42  
 RANSAC RANSACREGRESSORRANDOMSTATE42  
 HUBERREGRESSOR HUBERREGRESSOR  
 COLORS OLS TURQUOISE THEILSEN GOLD RANSAC LIGHTGREEN  
 ↳HUBERREGRESSOR BLACK  
 LINESTYLE OLS THEILSEN RANSAC HUBERREGRESSOR  
 LW 3  
 XPLOT NPLINSPACEXMIN XMAX  
 FORTITLE THISX THISY IN  
 MODELING ERRORS ONLY X Y  
 CORRUPT X SMALL DEVIANTS XERRORS Y  
 CORRUPT Y SMALL DEVIANTS X YERRORS  
 CORRUPT X LARGE DEVIANTS XERRORSLARGE Y  
 CORRUPT Y LARGE DEVIANTS X YERRORSLARGE  
 PLTFIGUREFIGSIZE5 4  
 PLTPLOTTHISX 0 THISY B  
 FORNAME ESTIMATOR INESTIMATORS  
 MODEL MAKEPIPELINEPOLYNOMIALFEATURES3 ESTIMATOR  
 MODELFITTHISX THISY  
 MSE MEANSQUAREDERRORMODELPREDICTXTEST YTEST  
 YPLOT MODELPREDICTXPLOT NPNEWAXIS  
 PLTPLOTXPLOT YPLOT COLORCOLORSNAME LINESTYLELINESTYLENAME  
 LINEWIDTHLW LABEL S ERROR 3F NAME MSE  
 LEGENDTITLE ERROR OF MEAN NABSOLUTE DEVIATION NTO NONCORRUPT DATA  
 LEGEND PLTLEGENDLOCUPPER RIGHT FRAMEONFALSE TITLELEGENDTITLE  
 PROPDICTSIZEXSMALL  
 1182 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
PLTXLIM4 102  
PLTYLIM2 102  
PLTTITLETITLE  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 3950 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
51827 L1 PENALTY AND SPARSITY IN LOGISTIC REGRESSION  
COMPARISON OF THE SPARSITY PERCENTAGE OF ZERO COEFFICIENTS OF SOLUTIONS WHEN L1 L2 AND ELASTICNET PENALTY ARE USED  
FOR DIFFERENT VALUES OF C WE CAN SEE THAT LARGE VALUES OF C GIVE MORE FREEDOM TO THE MODEL CONVERSELY SMALLER VALUES  
OF C CONSTRAIN THE MODEL MORE IN THE L1 PENALTY CASE THIS LEADS TO SPARSER SOLUTIONS AS EXPECTED THE ELASTICNET  
PENALTY SPARSITY IS BETWEEN THAT OF L1 AND L2  
WE CLASSIFY 8X8 IMAGES OF DIGITS INTO TWO CLASSES 04 AGAINST 59 THE VISUALIZATION SHOWS COEFFICIENTS OF THE MODELS  
FOR VARYING C  
OUT  
C100  
SPARSITY WITH L1 PENALTY 625  
518 GENERALIZED LINEAR MODELS 1183

SCIKITLEARN USER GUIDE RELEASE 0213  
SPARSITY WITH ELASTICNET PENALTY 469  
SPARSITY WITH L2 PENALTY 469  
SCORE WITH L1 PENALTY 090  
SCORE WITH ELASTICNET PENALTY 090  
SCORE WITH L2 PENALTY 090  
C010  
SPARSITY WITH L1 PENALTY 2969  
SPARSITY WITH ELASTICNET PENALTY 1250  
SPARSITY WITH L2 PENALTY 469  
SCORE WITH L1 PENALTY 090  
SCORE WITH ELASTICNET PENALTY 090  
SCORE WITH L2 PENALTY 090  
C001  
SPARSITY WITH L1 PENALTY 8438  
SPARSITY WITH ELASTICNET PENALTY 6875  
SPARSITY WITH L2 PENALTY 469  
SCORE WITH L1 PENALTY 086  
SCORE WITH ELASTICNET PENALTY 088  
SCORE WITH L2 PENALTY 089  
PRINTDOC  
AUTHORS ALEXANDRE GRAMFORT ALEXANDREGRAMFORTINRIA  
MATHIEU BLONDEL MATHIEU MBONDELORG  
ANDREAS MUELLER AMUELLERAISUNIBONNDE  
LICENSE BSD 3 CLAUSE  
IMPORT NUMPY AS NP  
IMPORT MATPLOTLIBPYPLOT AS PLT  
FROM SKLEARNLINEARMODEL IMPORT LOGISTICREGRESSION  
FROM SKLEARN IMPORT DATASETS  
FROM SKLEARNPREPROCESSING IMPORT STANDARDSCALER  
DIGITS DATASETSLOADDIGITS  
X Y DIGITSDATA DIGITSTARGET  
X STANDARDSCALERFITTRANSFORMX  
CLASSIFY SMALL AGAINST LARGE DIGITS  
Y Y 4ASTYPENPINT  
L1RATIO 05 L1 WEIGHT IN THE ELASTICNET REGULARIZATION  
FIG AXES PLTSUBPLOTS3 3  
SET REGULARIZATION PARAMETER  
FORI C AXESROW INENUMERATEZIP1 01 001 AXES  
TURN DOWN TOLERANCE FOR SHORT TRAINING TIME  
CLFL1LR LOGISTICREGRESSIONCC PENALTYL1 TOL001 SOLVERSAGA  
CLFL2LR LOGISTICREGRESSIONCC PENALTYL2 TOL001 SOLVERSAGA  
CLFENLR LOGISTICREGRESSIONCC PENALTYELASTICNET SOLVERSAGA  
L1RATIO1RATIO TOL001  
1184 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
CLFL1LRFITX Y  
CLFL2LRFITX Y  
CLFENLRFITX Y  
COEFL1LR CLFL1LRCOEFRAVEL  
COEFL2LR CLFL2LRCOEFRAVEL  
COEFENLR CLFENLRCOEFRAVEL  
COEFL1LR CONTAINS ZEROS DUE TO THE  
L1 SPARSITY INDUCING NORM  
SPARSITYL1LR NPMEANCOEFL1LR 0 100  
SPARSITYL2LR NPMEANCOEFL2LR 0 100  
SPARSITYENLR NPMEANCOEFENLR 0 100  
PRINTC2F C  
PRINT40 2FFORMATSPARSITY WITH L1 PENALTY SPARSITYL1LR  
PRINT40 2FFORMATSPARSITY WITH ELASTICNET PENALTY  
SPARSITYENLR  
PRINT40 2FFORMATSPARSITY WITH L2 PENALTY SPARSITYL2LR  
PRINT40 2FFORMATSCORE WITH L1 PENALTY  
CLFL1LRSCOREX Y  
PRINT40 2FFORMATSCORE WITH ELASTICNET PENALTY  
CLFENLRSCOREX Y  
PRINT40 2FFORMATSCORE WITH L2 PENALTY  
CLFL2LRSCOREX Y  
IFI 0  
AXESROW0SETTITLEL1 PENALTY  
AXESROW1SETTITLEELASTICNET NL1RATIO S L1RATIO  
AXESROW2SETTITLEL2 PENALTY  
FORAX COEFS INZIPAXESROW COEFL1LR COEFENLR COEFL2LR  
AXIMSHOWNPABSCOEFSRESHAPE8 8 INTERPOLATIONNEAREST  
CMAPBINARY VMAX1 VMIN0  
AXSETXTICKS  
AXSETYTICKS  
AXESROW0SETYLABELC S C  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0659 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
51828 LASSO AND ELASTIC NET  
LASSO AND ELASTIC NET L1 AND L2 PENALISATION IMPLEMENTED USING A COORDINATE DESCENT  
THE COEFFICIENTS CAN BE FORCED TO BE POSITIVE  
518 GENERALIZED LINEAR MODELS 1185







SCIKITLEARN USER GUIDE RELEASE 0213

•

OUT

COMPUTING REGULARIZATION PATH USING THE LASSO  
COMPUTING REGULARIZATION PATH USING THE POSITIVE LASSO  
COMPUTING REGULARIZATION PATH USING THE ELASTIC NET  
COMPUTING REGULARIZATION PATH USING THE POSITIVE ELASTIC NET  
PRINTDOC

AUTHOR ALEXANDRE GRAMFORT ALEXANDREGRAMFORTINRIA.FR

LICENSE BSD 3 CLAUSE

FROM ITERTOOLS IMPORT CYCLE

IMPORT NUMPY AS NP

IMPORT MATPLOTLIBPY.PLOT AS PLT

FROM SKLEARN.LINEAR\_MODEL IMPORT LASSO\_PATH ENET\_PATH

FROM SKLEARN IMPORT DATASETS

DIABETES = DATASETS.LOAD\_DIABETES

X = DIABETES.DATA

Y = DIABETES.TARGET

1188 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

X XSTDAXISO STANDARDIZE DATA EASIER TO SET THE L1RATIO PARAMETER  
COMPUTE PATHS

EPS 5E3 THE SMALLER IT IS THE LONGER IS THE PATH  
PRINTCOMPUTING REGULARIZATION PATH USING THE LASSO  
ALPHASLASSO COEFSLASSO LASSOPATHX Y EPS FITINTERCEPTFALSE  
PRINTCOMPUTING REGULARIZATION PATH USING THE POSITIVE LASSO  
ALPHASPOSITIVELASSO COEFSPOSITIVELASSO LASSOPATH  
X Y EPS POSITIVETRUE FITINTERCEPTFALSE  
PRINTCOMPUTING REGULARIZATION PATH USING THE ELASTIC NET  
ALPHASENET COEFSNET ENETPATH  
X Y EPSEPS L1RATIO08 FITINTERCEPTFALSE  
PRINTCOMPUTING REGULARIZATION PATH USING THE POSITIVE ELASTIC NET  
ALPHASPOSITIVEENET COEFSPOSITIVEENET ENETPATH  
X Y EPSEPS L1RATIO08 POSITIVETRUE FITINTERCEPTFALSE  
DISPLAY RESULTS

PLTFigure1

COLORS CYCLEB R G C K

NEGLOGALPHASLASSO NPLOG10ALPHASLASSO  
NEGLOGALPHASENET NPLOG10ALPHASENET  
FORCOEFL COEFE C INZIPCOEFSLASSO COEFSNET COLORS  
L1 PLTPLOTNEGLOGALPHASLASSO COEFL CC  
L2 PLTPLOTNEGLOGALPHASENET COEFE LINESTYLE CC  
PLTXLABELLOGALPHA  
PLTYLABELCOEFFICIENTS  
PLTTITLELASSO AND ELASTICNET PATHS  
PLTLEGENDL11 L21 LASSO ELASTICNET LOCLOWER LEFT  
PLTAXISTIGHT

PLTFigure2

NEGLOGALPHASPOSITIVELASSO NPLOG10ALPHASPOSITIVELASSO  
FORCOEFL COEFPL C INZIPCOEFSLASSO COEFSPOSITIVELASSO COLORS  
L1 PLTPLOTNEGLOGALPHASLASSO COEFL CC  
L2 PLTPLOTNEGLOGALPHASPOSITIVELASSO COEFPL LINESTYLE CC  
PLTXLABELLOGALPHA  
PLTYLABELCOEFFICIENTS  
PLTTITLELASSO AND POSITIVE LASSO  
PLTLEGENDL11 L21 LASSO POSITIVE LASSO LOCLOWER LEFT  
PLTAXISTIGHT

PLTFigure3

NEGLOGALPHASPOSITIVEENET NPLOG10ALPHASPOSITIVEENET  
FORCOEFE COEFPE C INZIPCOEFSNET COEFSPOSITIVEENET COLORS  
L1 PLTPLOTNEGLOGALPHASENET COEFE CC  
L2 PLTPLOTNEGLOGALPHASPOSITIVEENET COEFPE LINESTYLE CC  
PLTXLABELLOGALPHA

518 GENERALIZED LINEAR MODELS 1189

SCIKITLEARN USER GUIDE RELEASE 0213  
PLTYLABELCOEFFICIENTS  
PLTTITLEELASTICNET AND POSITIVE ELASTICNET  
PLTLEGENDL11 L21 ELASTICNET POSITIVE ELASTICNET  
LOCLOWER LEFT  
PLTAXISTIGHT  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0232 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
51829 AUTOMATIC RELEVANCE DETERMINATION REGRESSION ARD  
FIT REGRESSION MODEL WITH BAYESIAN RIDGE REGRESSION  
SEEBAYESIAN RIDGE REGRESSION FOR MORE INFORMATION ON THE REGRESSOR  
COMPARED TO THE OLS ORDINARY LEAST SQUARES ESTIMATOR THE COEFFICIENT WEIGHTS ARE SLIGHTLY SHIFTED TOWARD ZEROS  
WHICH STABILISES THEM  
THE HISTOGRAM OF THE ESTIMATED WEIGHTS IS VERY PEAKED AS A SPARSITYINDUCING PRIOR IS IMPLIED ON THE WEIGHTS  
THE ESTIMATION OF THE MODEL IS DONE BY ITERATIVELY MAXIMIZING THE MARGINAL LOGLIKELIHOOD OF THE OBSERVATIONS  
WE ALSO PLOT PREDICTIONS AND UNCERTAINTIES FOR ARD FOR ONE DIMENSIONAL REGRESSION USING POLYNOMIAL FEATURE EXPAN  
SION NOTE THE UNCERTAINTY STARTS GOING UP ON THE RIGHT SIDE OF THE PLOT THIS IS BECAUSE THESE TEST SAMPLES ARE OUTSIDE  
OF THE RANGE OF THE TRAINING SAMPLES  
1190 CHAPTER 5 EXAMPLES





SCIKITLEARN USER GUIDE RELEASE 0213

- 518 GENERALIZED LINEAR MODELS 1193

SCIKITLEARN USER GUIDE RELEASE 0213

```
•
PRINTDOC
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SCIPY IMPORT STATS
FROM SKLEARNLINEARMODEL IMPORT ARDREGRESSION LINEARREGRESSION
```

```
GENERATING SIMULATED DATA WITH GAUSSIAN WEIGHTS
PARAMETERS OF THE EXAMPLE
NPRANDOMSEED0
NSAMPLES NFEATURES 100 100
CREATE GAUSSIAN DATA
X NPRANDOMRANDNNSAMPLES NFEATURES
CREATE WEIGHTS WITH A PRECISION LAMBDA OF 4
LAMBDA 4
W NPZEROSNFEATURES
ONLY KEEP 10 WEIGHTS OF INTEREST
RELEVANTFEATURES NPRANDOMRANDINT0 NFEATURES 10
FORIINRELEVANTFEATURES
WI STATSNORMMRVSLOC0 SCALE1 NPSQRTLAMBDA
CREATE NOISE WITH A PRECISION ALPHA OF 50
ALPHA 50
1194 CHAPTER 5 EXAMPLES
```



SCIKITLEARN USER GUIDE RELEASE 0213

NOISE STATS NORM RV SLOC0 SCALE1 NPSQRT ALPHA SIZE N SAMPLES

CREATE THE TARGET

Y NPDOTX W NOISE

FIT THE ARD REGRESSION

CLF ARDREGRESSION COMPUTE SCORE TRUE

CLFFITX Y

OLS LINEAR REGRESSION

OLS FITX Y

PLOT THE TRUE WEIGHTS THE ESTIMATED WEIGHTS THE HISTOGRAM OF THE WEIGHTS AND PREDICTIONS WITH STANDARD DEVIATIONS

PLTFigureFIGSIZE6 5

PLTTITLEWEIGHTS OF THE MODEL

PLTPLOTCLFCOEF COLOR DARKBLUE LINESTYLE LINEWIDTH2

LABELARD ESTIMATE

PLTPLOTOLS COEF COLOR YELLOWGREEN LINESTYLE LINEWIDTH2

LABELOLS ESTIMATE

PLTPLOTW COLOR ORANGE LINESTYLE LINEWIDTH2 LABELGROUND TRUTH

PLTXLABELFEATURES

PLTYLABELVALUES OF THE WEIGHTS

PLTLEGENDLOC1

PLTFigureFIGSIZE6 5

PLTTITLEHISTOGRAM OF THE WEIGHTS

PLTHISTCLFCOEF BINS NFEATURES COLOR NAVY LOG TRUE

PLTSCATTERCLFCOEF RELEVANTFEATURES NPFULL LEN RELEVANTFEATURES 5

COLOR GOLD MARKER O LABEL RELEVANTFEATURES

PLTYLABELFEATURES

PLTXLABELVALUES OF THE WEIGHTS

PLTLEGENDLOC1

PLTFigureFIGSIZE6 5

PLTTITLEMARGINAL LOG LIKELIHOOD

PLTPLOTCLFSCORES COLOR NAVY LINEWIDTH2

PLTYLABELSCORE

PLTXLABELITERATIONS

PLOTTING SOME PREDICTIONS FOR POLYNOMIAL REGRESSION

DEFFX NOISEAMOUNT

Y NPSQRTX NPSINX

NOISE NPRANDOMNORMAL0 1 LENX

RETURN Y NOISEAMOUNT NOISE

DEGREE 10

X NPLINSPACE0 10 100

Y FX NOISEAMOUNT1

CLFPOLY ARDREGRESSIONTHRESHOLDLAMBDA1E5

CLFPOLYFITNPVANDERX DEGREE Y

XPLOT NPLINSPACE0 11 25

YPLOT FXPLOT NOISEAMOUNT0

YMEAN YSTD CLFPOLYPREDICTNPVANDERXPLOT DEGREE RETURNSTD TRUE

518 GENERALIZED LINEAR MODELS 1195

SCIKITLEARN USER GUIDE RELEASE 0213  
PLTFIGUREFIGSIZE6 5  
PLTERRORBARXPLOT YMEAN YSTD COLORNAVY  
LABELPOLYNOMIAL ARD LINEWIDTH2  
PLTPLOTXPLOT YPLOT COLORGOLD LINEWIDTH2  
LABELGROUND TRUTH  
PLTYLABELOUTPUT Y  
PLTXLABELFEATURE X  
PLTLEGENDLOCLOWER LEFT  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0293 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
51830 BAYESIAN RIDGE REGRESSION  
COMPUTES A BAYESIAN RIDGE REGRESSION ON A SYNTHETIC DATASET  
SEEBAYESIAN RIDGE REGRESSION FOR MORE INFORMATION ON THE REGRESSOR  
COMPARED TO THE OLS ORDINARY LEAST SQUARES ESTIMATOR THE COEFFICIENT WEIGHTS ARE SLIGHTLY SHIFTED TOWARD ZEROS  
WHICH STABILISES THEM  
AS THE PRIOR ON THE WEIGHTS IS A GAUSSIAN PRIOR THE HISTOGRAM OF THE ESTIMATED WEIGHTS IS GAUSSIAN  
THE ESTIMATION OF THE MODEL IS DONE BY ITERATIVELY MAXIMIZING THE MARGINAL LOGLIKELIHOOD OF THE OBSERVATIONS  
WE ALSO PLOT PREDICTIONS AND UNCERTAINTIES FOR BAYESIAN RIDGE REGRESSION FOR ONE DIMENSIONAL REGRESSION USING POLY  
NOMIAL FEATURE EXPANSION NOTE THE UNCERTAINTY STARTS GOING UP ON THE RIGHT SIDE OF THE PLOT THIS IS BECAUSE THESE TEST  
SAMPLES ARE OUTSIDE OF THE RANGE OF THE TRAINING SAMPLES  
1196 CHAPTER 5 EXAMPLES







SCIKITLEARN USER GUIDE RELEASE 0213

```
•
PRINTDOC
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SCIPY IMPORT STATS
FROM SKLEARNLINEARMODEL IMPORT BAYESIANRIDGE LINEARREGRESSION
```

```
GENERATING SIMULATED DATA WITH GAUSSIAN WEIGHTS
NPRANDOMSEED0
NSAMPLES NFEATURES 100 100
X NPRANDOMRANDNNSAMPLES NFEATURES CREATE GAUSSIAN DATA
CREATE WEIGHTS WITH A PRECISION LAMBDA OF 4
LAMBDA 4
W NPZEROSNFEATURES
ONLY KEEP 10 WEIGHTS OF INTEREST
RELEVANTFEATURES NPRANDOMRANDINT0 NFEATURES 10
FORIINRELEVANTFEATURES
WI STATSNORMRVSLOC0 SCALE1 NPSQRTLAMBDA
CREATE NOISE WITH A PRECISION ALPHA OF 50
ALPHA 50
NOISE STATSNORMRVSLOC0 SCALE1 NPSQRTALPHA SIZENSAMPLES
CREATE THE TARGET
Y NPDOTX W NOISE
1200 CHAPTER 5 EXAMPLES
```

SCIKITLEARN USER GUIDE RELEASE 0213

FIT THE BAYESIAN RIDGE REGRESSION AND AN OLS FOR COMPARISON  
CLF BAYESIANRIDGECOMPUTESCORETRUE  
CLFFITX Y  
OLS LINEARREGRESSION  
OLSFITX Y

PLOT TRUE WEIGHTS ESTIMATED WEIGHTS HISTOGRAM OF THE WEIGHTS AND  
PREDICTIONS WITH STANDARD DEVIATIONS  
LW 2  
PLTFIGUREFIGSIZE6 5  
PLTTITLEWEIGHTS OF THE MODEL  
PLTPLOTCLFCOEF COLORLIGHTGREEN LINEWIDTHLW  
LABELBAYESIAN RIDGE ESTIMATE  
PLTPLOTW COLORGOLD LINEWIDTHLW LABELGROUND TRUTH  
PLTPLOTOLSCOEF COLORNNAVY LINESTYLE LABELOLS ESTIMATE  
PLTXLABELFEATURES  
PLTYLABELVALUES OF THE WEIGHTS  
PLTLEGENDLOCBEST PROPDICTSIZE12  
PLTFIGUREFIGSIZE6 5  
PLTTITLEHISTOGRAM OF THE WEIGHTS  
PLTHISTCLFCOEF BINSNFEATURES COLORGOLD LOGTRUE  
EDGECOLORBLACK  
PLTSCATTERCLFCOEFRELEVANTFEATURES NPFULLLENRELEVANTFEATURES 5  
COLORNNAVY LABELRELEVANT FEATURES  
PLTYLABELFEATURES  
PLTXLABELVALUES OF THE WEIGHTS  
PLTLEGENDLOCUPPER LEFT  
PLTFIGUREFIGSIZE6 5  
PLTTITLEMARGINAL LOGLIKELIHOOD  
PLTPLOTCLFSCORES COLORNNAVY LINEWIDTHLW  
PLTYLABELSCORE  
PLTXLABELITERATIONS  
PLOTING SOME PREDICTIONS FOR POLYNOMIAL REGRESSION  
DEFFX NOISEAMOUNT  
Y NPSQRTX NPSINX  
NOISE NPRANDOMNORMAL0 1 LENX  
RETURNY NOISEAMOUNT NOISE  
DEGREE 10  
X NPLINSPACE0 10 100  
Y FX NOISEAMOUNT01  
CLFPOLY BAYESIANRIDGE  
CLFPOLYFITNPVANDERX DEGREE Y  
XPLOT NPLINSPACE0 11 25  
YPLOT FXPLOT NOISEAMOUNT0  
YMEAN YSTD CLFPOLYPREDICTNPVANDERXPLOT DEGREE RETURNSTDTRUE  
PLTFIGUREFIGSIZE6 5  
PLTERRORBARXPLOT YMEAN YSTD COLORNNAVY  
518 GENERALIZED LINEAR MODELS 1201

SCIKITLEARN USER GUIDE RELEASE 0213  
LABELPOLYNOMIAL BAYESIAN RIDGE REGRESSION LINEWIDTHLW  
PLTPLOTXPLOT YPLOT COLORGOLD LINEWIDTHLW  
LABELGROUND TRUTH  
PLTYLABELOUTPUT Y  
PLTXLABELFEATURE X  
PLTLEGENDLOCLOWER LEFT  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0143 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
51831 LASSO MODEL SELECTION CROSSVALIDATION AIC BIC  
USE THE AKAIKE INFORMATION CRITERION AIC THE BAYES INFORMATION CRITERION BIC AND CROSSVALIDATION TO SELECT AN  
OPTIMAL VALUE OF THE REGULARIZATION PARAMETER ALPHA OF THE LASSO ESTIMATOR  
RESULTS OBTAINED WITH LASSOLARSIC ARE BASED ON AICBIC CRITERIA  
INFORMATIONCRITERION BASED MODEL SELECTION IS VERY FAST BUT IT RELIES ON A PROPER ESTIMATION OF DEGREES OF FREEDOM ARE  
DERIVED FOR LARGE SAMPLES ASYMPTOTIC RESULTS AND ASSUME THE MODEL IS CORRECT IE THAT THE DATA ARE ACTUALLY GENERATED  
BY THIS MODEL THEY ALSO TEND TO BREAK WHEN THE PROBLEM IS BADLY CONDITIONED MORE FEATURES THAN SAMPLES  
FOR CROSSVALIDATION WE USE 20FOLD WITH 2 ALGORITHMS TO COMPUTE THE LASSO PATH COORDINATE DESCENT AS IMPLEMENTED  
BY THE LASSOCV CLASS AND LARS LEAST ANGLE REGRESSION AS IMPLEMENTED BY THE LASSOLARSCV CLASS BOTH ALGORITHMS  
GIVE ROUGHLY THE SAME RESULTS THEY DIFFER WITH REGARDS TO THEIR EXECUTION SPEED AND SOURCES OF NUMERICAL ERRORS  
LARS COMPUTES A PATH SOLUTION ONLY FOR EACH KINK IN THE PATH AS A RESULT IT IS VERY EFFICIENT WHEN THERE ARE ONLY OF FEW  
KINKS WHICH IS THE CASE IF THERE ARE FEW FEATURES OR SAMPLES ALSO IT IS ABLE TO COMPUTE THE FULL PATH WITHOUT SETTING  
ANY META PARAMETER ON THE OPPOSITE COORDINATE DESCENT COMPUTE THE PATH POINTS ON A PRESPECIFIED GRID HERE WE USE  
THE DEFAULT THUS IT IS MORE EFFICIENT IF THE NUMBER OF GRID POINTS IS SMALLER THAN THE NUMBER OF KINKS IN THE PATH SUCH  
A STRATEGY CAN BE INTERESTING IF THE NUMBER OF FEATURES IS REALLY LARGE AND THERE ARE ENOUGH SAMPLES TO SELECT A LARGE  
AMOUNT IN TERMS OF NUMERICAL ERRORS FOR HEAVILY CORRELATED VARIABLES LARS WILL ACCUMULATE MORE ERRORS WHILE THE  
COORDINATE DESCENT ALGORITHM WILL ONLY SAMPLE THE PATH ON A GRID  
NOTE HOW THE OPTIMAL VALUE OF ALPHA VARIES FOR EACH FOLD THIS ILLUSTRATES WHY NESTEDCROSS VALIDATION IS NECESSARY  
WHEN TRYING TO EVALUATE THE PERFORMANCE OF A METHOD FOR WHICH A PARAMETER IS CHOSEN BY CROSSVALIDATION THIS CHOICE  
OF PARAMETER MAY NOT BE OPTIMAL FOR UNSEEN DATA  
1202 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213

- 518 GENERALIZED LINEAR MODELS 1203



SCIKITLEARN USER GUIDE RELEASE 0213

- 

OUT

COMPUTING REGULARIZATION PATH USING THE COORDINATE DESCENT LASSO

COMPUTING REGULARIZATION PATH USING THE LARS LASSO

PRINTDOC

AUTHOR OLIVIER GRISEL GAELE VAROQUAUX ALEXANDRE GRAMFORT

LICENSE BSD 3 CLAUSE

IMPORT TIME

IMPORT NUMPY AS NP

IMPORT MATPLOTLIBPYPLOT AS PLT

FROM SKLEARNLINEARMODEL IMPORT LASSOCV LASSOLARSCV LASSOLARSIC

FROM SKLEARN IMPORT DATASETS

THIS IS TO AVOID DIVISION BY ZERO WHILE DOING NPLOG10

EPSILON 1E4

DIABETES DATASETSLOADDIABETES

518 GENERALIZED LINEAR MODELS 1205

```
SCIKITLEARN USER GUIDE RELEASE 0213
X DIABETESDATA
Y DIABETESTARGET
RNG NPRANDOMRANDOMSTATE42
X NPCX RNGRANDNXSHAPE0 14 ADD SOME BAD FEATURES
NORMALIZE DATA AS DONE BY LARS TO ALLOW FOR COMPARISON
X NPSQRTNPSUMX 2 AXIS0

LASSOLARSIC LEAST ANGLE REGRESSION WITH BICAIC CRITERION
MODEL BIC LASSOLARSICCRITERION BIC
T1 TIMETIME
MODEL BIC FITX Y
TBIC TIMETIME T1
ALPHA BIC MODEL BIC ALPHA
MODEL AIC LASSOLARSICCRITERION AIC
MODEL AIC FITX Y
ALPHA AIC MODEL AIC ALPHA
DEF PLOT I C CRITERION MODEL NAME COLOR
ALPHA MODEL ALPHA EPSILON
ALPHAS MODEL ALPHAS EPSILON
CRITERION MODEL CRITERION
PLT PLOT NPLOG10 ALPHAS CRITERION COLOR COLOR
LINEWIDTH3 LABEL SCRITERION NAME
PLT AX V LINE NPLOG10 ALPHA COLOR COLOR LINEWIDTH3
LABEL ALPHA SESTIMATE NAME
PLT X LABEL LOG ALPHA
PLT Y LABEL CRITERION
PLT FIGURE
PLOT I C CRITERION MODEL AIC B
PLOT I C CRITERION MODEL BIC BIC R
PLT LEGEND
PLT TITLE INFORMATION CRITERION FOR MODEL SELECTION TRAINING TIME 3FS
TBIC

LASSOCV COORDINATE DESCENT
COMPUTE PATHS
PRINT COMPUTING REGULARIZATION PATH USING THE COORDINATE DESCENT LASSO
T1 TIMETIME
MODEL LASSOCVCV20 FITX Y
TLASSOCV TIMETIME T1
DISPLAY RESULTS
MLOGALPHAS NPLOG10 MODEL ALPHAS EPSILON
PLT FIGURE
YMIN YMAX 2300 3800
PLT PLOT MLOGALPHAS MODEL MSE PATH
PLT PLOT MLOGALPHAS MODEL MSE PATH MEAN AXIS1 K
LABEL AVERAGE ACROSS THE FOLDS LINEWIDTH2
1206 CHAPTER 5 EXAMPLES
```

SCIKITLEARN USER GUIDE RELEASE 0213  
PLTAXVLINENPLOG10MODELALPHA EPSILON LINESTYLE COLORK  
LABELALPHA CV ESTIMATE  
PLTLEGEND  
PLTXLABELLOGALPHA  
PLTYLABELMEAN SQUARE ERROR  
PLTTITLEMEAN SQUARE ERROR ON EACH FOLD COORDINATE DESCENT  
TRAIN TIME 2FS TLASSOCV  
PLTAXISTIGHT  
PLTYLIMYMIN YMAX

LASSOLARSCV LEAST ANGLE REGRESSION  
COMPUTE PATHS  
PRINTCOMPUTING REGULARIZATION PATH USING THE LARS LASSO  
T1 TIMETIME  
MODEL LASSOLARSCVCV20FITX Y  
TLASSOLARSCV TIMETIME T1  
DISPLAY RESULTS  
MLOGALPHAS NPLOG10MODELCVALPHAS EPSILON  
PLTFigure  
PLTPLOTMLOGALPHAS MODELMSEPATH  
PLTPLOTMLOGALPHAS MODELMSEPATHMEANAXIS1 K  
LABELAVERAGE ACROSS THE FOLDS LINEWIDTH2  
PLTAXVLINENPLOG10MODELALPHA LINESTYLE COLORK  
LABELALPHA CV  
PLTLEGEND  
PLTXLABELLOGALPHA  
PLTYLABELMEAN SQUARE ERROR  
PLTTITLEMEAN SQUARE ERROR ON EACH FOLD LARS TRAIN TIME 2FS  
TLASSOLARSCV  
PLTAXISTIGHT  
PLTYLIMYMIN YMAX  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0811 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
51832 MULTICLASS SPARSE LOGISITIC REGRESSION ON NEWGROUPS20  
COMPARISON OF MULTINOMIAL LOGISTIC L1 VS ONEVERSUSREST L1 LOGISTIC REGRESSION TO CLASSIFY DOCUMENTS FROM THE NEW  
GROUPS20 DATASET MULTINOMIAL LOGISTIC REGRESSION YIELDS MORE ACCURATE RESULTS AND IS FASTER TO TRAIN ON THE LARGER SCALE  
DATASET  
HERE WE USE THE L1 SPARSITY THAT TRIMS THE WEIGHTS OF NOT INFORMATIVE FEATURES TO ZERO THIS IS GOOD IF THE GOAL IS TO  
EXTRACT THE STRONGLY DISCRIMINATIVE VOCABULARY OF EACH CLASS IF THE GOAL IS TO GET THE BEST PREDICTIVE ACCURACY IT IS BETTER  
TO USE THE NON SPARSITYINDUCING L2 PENALTY INSTEAD  
518 GENERALIZED LINEAR MODELS 1207

SCIKITLEARN USER GUIDE RELEASE 0213

A MORE TRADITIONAL AND POSSIBLY BETTER WAY TO PREDICT ON A SPARSE SUBSET OF INPUT FEATURES WOULD BE TO USE UNIVARIATE  
FEATURE SELECTION FOLLOWED BY A TRADITIONAL L2PENALISED LOGISTIC REGRESSION MODEL

OUT

DATASET 20NEWSGROUP TRAINSAMPLES9000 NFEATURES130107 NCLASSES20

MODELONE VERSUS REST SOLVERSAGA NUMBER OF EPOCHS 1

MODELONE VERSUS REST SOLVERSAGA NUMBER OF EPOCHS 2

MODELONE VERSUS REST SOLVERSAGA NUMBER OF EPOCHS 4

TEST ACCURACY FOR MODEL OVR 07490

NONZERO COEFFICIENTS FOR MODEL OVR PER CLASS

031743104 036815852 04181174 046115889 024595141 041350581

031281945 027054655 058720899 032972861 04158116 03312658

041888599 041120001 059643217 031666244 034279478 028130692

035278655 024748861

RUN TIME 4 EPOCHS FOR MODEL OVR306

MODELMULTINOMIAL SOLVERSAGA NUMBER OF EPOCHS 1

MODELMULTINOMIAL SOLVERSAGA NUMBER OF EPOCHS 3

MODELMULTINOMIAL SOLVERSAGA NUMBER OF EPOCHS 7

TEST ACCURACY FOR MODEL MULTINOMIAL 07450

NONZERO COEFFICIENTS FOR MODEL MULTINOMIAL PER CLASS

013219888 011452112 013066169 013681047 012066991 015909982

013450468 009146318 007916561 012143851 013911627 010760374

018984374 012143851 017524038 022289346 011605832 007916561

007301682 015141384

RUN TIME 7 EPOCHS FOR MODEL MULTINOMIAL255

EXAMPLE RUN IN 11262 S

1208 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
IMPORT TIMEIT
IMPORT WARNINGS
IMPORT MATPLOTLIBPYPLOT AS PLT
IMPORT NUMPY AS NP
FROM SKLEARNDATASETS IMPORT FETCH20NEWSGROUPSVECTORIZED
FROM SKLEARNLINEARMODEL IMPORT LOGISTICREGRESSION
FROM SKLEARNMODELSELECTION IMPORT TRAINTESTSPLIT
FROM SKLEARNEXCEPTIONS IMPORT CONVERGENCEWARNING
PRINTDOC
  AUTHOR ARTHUR MENSCH
  WARNINGSFILTERWARNINGSIGNORE CATEGORYCONVERGENCEWARNING
MODULESKLEARN
T0 TIMEITDEFAULTTIMER
  WE USE SAGA SOLVER
SOLVER SAGA
  TURN DOWN FOR FASTER RUN TIME
NSAMPLES 10000
  MEMORIZED FETCHRCV1 FOR FASTER ACCESS
DATASET FETCH20NEWSGROUPSVECTORIZEDALL
X DATASETDATA
Y DATASETTARGET
X XNSAMPLES
Y YNSAMPLES
XTRAIN XTEST YTRAIN YTEST TRAINTESTSPLITX Y
RANDOMSTATE42
STRATIFY
TESTSIZE01
TRAINSAMPLES NFEATURES XTRAINSHAPE
NCLASSES NPUNIQUEYSHAPE0
PRINTDATASET 20NEWSGROUP TRAINSAMPLES I NFEATURES I NCLASSES I
  TRAINSAMPLES NFEATURES NCLASSES
MODELS OVR NAME ONE VERSUS REST ITERS 1 2 4
MULTINOMIAL NAME MULTINOMIAL ITERS 1 3 7
FORMODELINMODELS
  ADD INITIAL CHANCELEVEL VALUES FOR PLOTTING PURPOSE
ACCURACIES 1 NCLASSES
TIMES 0
DENSITIES 1
MODELPARAMS MODELSMODEL
518 GENERALIZED LINEAR MODELS 1209
```

SCIKITLEARN USER GUIDE RELEASE 0213  
SMALL NUMBER OF EPOCHS FOR FAST RUNTIME  
FORTHISMAXITER INMODELPARAMSITERS  
PRINTMODEL S SOLVER S NUMBER OF EPOCHS S  
MODELPARAMSNAME SOLVER THISMAXITER  
LR LOGISTICREGRESSIONSOLVERSOLVER  
MULTICLASSMODEL  
C1  
PENALTYL1  
FITINTERCEPTTRUE  
MAXITERTHISMAXITER  
RANDOMSTATE42  
  
T1 TIMEITDEFAULTTIMER  
LRFITXTRAIN YTRAIN  
TRAINTIME TIMEITDEFAULTTIMER T1  
YPRED LRPREDICTXTEST  
ACCURACY NPSUMYPRED YTEST YTESTSHAPE0  
DENSITY NPMEANLRCOEF 0 AXIS1 100  
ACCURACIESAPPENDACCURACY  
DENSITIESAPPENDDENSITY  
TIMESAPPENDTRAINTIME  
MODELSMODELTIMES TIMES  
MODELSMODELDENSITIES DENSITIES  
MODELSMODELACCURACIES ACCURACIES  
PRINTTEST ACCURACY FOR MODEL S4F MODEL ACCURACIES1  
PRINTNONZERO COEFFICIENTS FOR MODEL S  
PER CLASS NS MODEL DENSITIES1  
PRINTRUN TIME IEPOCHS FOR MODEL S  
2F MODELPARAMSITERS1 MODEL TIMES1  
FIG PLTFigure  
AX FIGADDSUBPLOT111  
FORMODELINMODELS  
NAME MODELSMODELNAME  
TIMES MODELSMODELTIMES  
ACCURACIES MODELSMODELACCURACIES  
AXPLOTTIMES ACCURACIES MARKERO  
LABELMODEL S NAME  
AXSETXLABELTRAIN TIME S  
AXSETYLABELTEST ACCURACY  
AXLEGEND  
FIGSUPTITLEMULTINOMIAL VS ONEVSREST LOGISTIC L1 N  
DATASET S 20NEWSGROUPS  
FIGTIGHTLAYOUT  
FIGSUBPLOTSADJUSTTOP085  
RUNTIME TIMEITDEFAULTTIMER T0  
PRINTEXAMPLE RUN IN 3FS RUNTIME  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 11263 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
1210 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213

51833 EARLY STOPPING OF STOCHASTIC GRADIENT DESCENT

STOCHASTIC GRADIENT DESCENT IS AN OPTIMIZATION TECHNIQUE WHICH MINIMIZES A LOSS FUNCTION IN A STOCHASTIC FASHION PERFORMING A GRADIENT DESCENT STEP SAMPLE BY SAMPLE IN PARTICULAR IT IS A VERY EFFICIENT METHOD TO FIT LINEAR MODELS AS A STOCHASTIC METHOD THE LOSS FUNCTION IS NOT NECESSARILY DECREASING AT EACH ITERATION AND CONVERGENCE IS ONLY GUARANTEED IN EXPECTATION FOR THIS REASON MONITORING THE CONVERGENCE ON THE LOSS FUNCTION CAN BE DIFFICULT ANOTHER APPROACH IS TO MONITOR CONVERGENCE ON A VALIDATION SCORE IN THIS CASE THE INPUT DATA IS SPLIT INTO A TRAINING SET AND A VALIDATION SET THE MODEL IS THEN FITTED ON THE TRAINING SET AND THE STOPPING CRITERION IS BASED ON THE PREDICTION SCORE COMPUTED ON THE VALIDATION SET THIS ENABLES US TO FIND THE LEAST NUMBER OF ITERATIONS WHICH IS SUFFICIENT TO BUILD A MODEL THAT GENERALIZES WELL TO UNSEEN DATA AND REDUCES THE CHANCE OF OVERFITTING THE TRAINING DATA THIS EARLY STOPPING STRATEGY IS ACTIVATED IF EARLYSTOPPINGTRUE OTHERWISE THE STOPPING CRITERION ONLY USES THE TRAINING LOSS ON THE ENTIRE INPUT DATA TO BETTER CONTROL THE EARLY STOPPING STRATEGY WE CAN SPECIFY A PARAMETER VALIDATIONFRACTION WHICH SET THE FRACTION OF THE INPUT DATASET THAT WE KEEP ASIDE TO COMPUTE THE VALIDATION SCORE THE OPTIMIZATION WILL CONTINUE UNTIL THE VALIDATION SCORE DID NOT IMPROVE BY AT LEAST TOL DURING THE LAST NITERNOCHANGE ITERATIONS THE ACTUAL NUMBER OF ITERATIONS IS AVAILABLE AT THE ATTRIBUTE NITER THIS EXAMPLE ILLUSTRATES HOW THE EARLY STOPPING CAN USED IN THE SKLEARNLINEARMODELSGDCCLASSIFIER MODEL TO ACHIEVE ALMOST THE SAME ACCURACY AS COMPARED TO A MODEL BUILT WITHOUT EARLY STOPPING THIS CAN SIGNIFICANTLY REDUCE TRAINING TIME NOTE THAT SCORES DIFFER BETWEEN THE STOPPING CRITERIA EVEN FROM EARLY ITERATIONS BECAUSE SOME OF THE TRAINING DATA IS HELD OUT WITH THE VALIDATION STOPPING CRITERION

OUT  
NO STOPPING CRITERION  
TRAINING LOSS  
VALIDATION SCORE

518 GENERALIZED LINEAR MODELS 1211

SCIKITLEARN USER GUIDE RELEASE 0213  
AUTHORS TOM DUPRE LA TOUR

```
LICENSE BSD 3 CLAUSE
IMPORT TIME
IMPORT SYS
IMPORT PANDAS AS PD
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARN IMPORT LINEARMODEL
FROM SKLEARNDATASETS IMPORT FETCHOPENML
FROM SKLEARNMODELSELECTION IMPORT TRAINTESTSPLIT
FROM SKLEARNUTILSTESTING IMPORT IGNOREWARNINGS
FROM SKLEARNEXCEPTIONS IMPORT CONVERGENCEWARNING
FROM SKLEARNUTILS IMPORT SHUFFLE
PRINTDOC
DEFOADMNISTNSAMPLESNONE CLASS00 CLASS18
LOAD MNIST SELECT TWO CLASSES SHUFFLE AND RETURN ONLY NSAMPLES
  LOAD DATA FROM HTTPOPENMLORGD554
MNIST  FETCHOPENMLMNIST784 VERSION1
  TAKE ONLY TWO CLASSES FOR BINARY CLASSIFICATION
MASK  NPLOGICALORMNISTTARGET  CLASS0 MNISTTARGET  CLASS1
X Y  SHUFFLEMNISTDATAMASK MNISTTARGETMASK RANDOMSTATE42
IFNSAMPLES IS NOTNONE
X Y  XNSAMPLES YNSAMPLES
RETURNX Y
IGNOREWARNINGS CATEGORYCONVERGENCEWARNING
DEFFITANDSCOREESTIMATOR MAXITER XTRAIN XTEST YTRAIN YTEST
FIT THE ESTIMATOR ON THE TRAIN SET AND SCORE IT ON BOTH SETS
ESTIMATORSETPARAMSMAXITERMAXITER
ESTIMATORSETPARAMSRANDOMSTATE0
START  TIMETIME
ESTIMATORFITXTRAIN YTRAIN
FITTIME  TIMETIME  START
NITER  ESTIMATORNITER
TRAINSCORE  ESTIMATORSCOREXTRAIN YTRAIN
TESTSCORE  ESTIMATORSCOREXTEST YTEST
RETURNFITTIME NITER TRAINSCORE TESTSCORE
  DEFINE THE ESTIMATORS TO COMPARE
ESTIMATORDICT
NO STOPPING CRITERION
LINEARMODELSGDCLASSIFIERTOL1E3 NITERNOCHANGE3
1212 CHAPTER 5 EXAMPLES
```

```
SCIKITLEARN USER GUIDE RELEASE 0213
TRAINING LOSS
LINEARMODELSGDCLASSIFIEREARLYSTOPPINGFALSE NITERNOCHANGE3
TOL01
VALIDATION SCORE
LINEARMODELSGDCLASSIFIEREARLYSTOPPINGTRUE NITERNOCHANGE3
TOL00001 VALIDATIONFRACTION02

LOAD THE DATASET
X Y LOADMNISTNSAMPLES10000
XTRAIN XTEST YTRAIN YTEST  TRAINTESTSPLITX Y TESTSIZE05
RANDOMSTATE0
RESULTS
FORESTIMATORNAME ESTIMATOR INESTIMATORORDICTITEMS
PRINTESTIMATORNAME  END
FORMAXITER INRANGE1 50
PRINT END
SYSTDOUTFLUSH
FITTIME NITER TRAINSCORE TESTSCORE  FITANDSCORE
ESTIMATOR MAXITER XTRAIN XTEST YTRAIN YTEST
RESULTSAPPENDESTIMATORNAME MAXITER FITTIME NITER
TRAINSCORE TESTSCORE
PRINT
  TRANSFORM THE RESULTS IN A PANDAS DATAFRAME FOR EASY PLOTTING
COLUMNS
STOPPING CRITERION MAXITER FIT TIME SEC NITER
TRAIN SCORE TEST SCORE

RESULTSDF PDDATAFRAMERESULTS COLUMNSCOLUMNS
  DEFINE WHAT TO PLOT XAXIS YAXIS
LINES  STOPPING CRITERION
PLOTLIST
MAXITER TRAIN SCORE
MAXITER TEST SCORE
MAXITER NITER
MAXITER FIT TIME SEC

NROWS 2
NCOLS INTNPCEILLENPLOTLIST 2
FIG AXES  PLTSUBPLOTSNROWSNROWS NCOLSNCOLS FIGSIZE6 NCOLS
4NROWS
AXES0 0GETSHAREDYAXESJOINAXES0 0 AXES0 1
FORAX XAXIS YAXIS INZIPAXESRAVEL PLOTLIST
FORCRITERION GROUPDF INRESULTSDFGROUPBYLINES
GROUPDFPLOTXXAXIS YYAXIS LABELCRITERION AXAX
AXSETTITLEYAXIS
AXLEGENDTITLELINES
FIGTIGHTLAYOUT
PLTSHOW
518 GENERALIZED LINEAR MODELS 1213
```

SCIKITLEARN USER GUIDE RELEASE 0213  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 45461 SECONDS  
519 MANIFOLD LEARNING  
EXAMPLES CONCERNING THE SKLEARNMANIFOLD MODULE  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
5191 SWISS ROLL REDUCTION WITH LLE  
AN ILLUSTRATION OF SWISS ROLL REDUCTION WITH LOCALLY LINEAR EMBEDDING  
OUT  
COMPUTING LLE EMBEDDING  
DONE RECONSTRUCTION ERROR 732714E08  
1214 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
AUTHOR FABIAN PEDREGOSA FABIANPEDREGOSAINRIA.FR
LICENSE BSD 3 CLAUSE C INRIA 2011
PRINTDOC
IMPORT MATPLOTLIB.PY.PLOT AS PLT
THIS IMPORT IS NEEDED TO MODIFY THE WAY FIGURE BEHAVES
FROM MPLTOOLKITSM.PLOT3D IMPORT AXES3D
AXES3D

LOCALLY LINEAR EMBEDDING OF THE SWISS ROLL
FROM SKLEARN IMPORT MANIFOLD.DATASETS
X COLOR DATASET.SAMPLES.GENERATOR.MAKESWISSROLL.N.SAMPLES1500
PRINTCOMPUTING LLE EMBEDDING
XR ERR MANIFOLD.LOCALLY.LINEAR.EMBEDDING.X.N.NEIGHBORS12
NCOMPONENTS2
PRINTDONE RECONSTRUCTION ERROR G ERR

PLOT RESULT
FIG PLT.FIGURE
AX FIG.ADDSUBPLOT211 PROJECTION3D
AXSCATTERX 0 X 1 X 2 C.COLOR CMAPPLT.CMSPECTRAL
AXSETTITLEORIGINAL DATA
AX FIG.ADDSUBPLOT212
AXSCATTERXR 0 XR 1 C.COLOR CMAPPLT.CMSPECTRAL
PLTAXISTIGHT
PLXTICKS PLTYTICKS
PLTTITLEPROJECTED DATA
PLTSHOW
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0182 SECONDS
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE
5192 MULTIDIMENSIONAL SCALING
AN ILLUSTRATION OF THE METRIC AND NONMETRIC MDS ON GENERATED NOISY DATA
THE RECONSTRUCTED POINTS USING THE METRIC MDS AND NON METRIC MDS ARE SLIGHTLY SHIFTED TO AVOID OVERLAPPING
519 MANIFOLD LEARNING 1215
```

```
SCIKITLEARN USER GUIDE RELEASE 0213
AUTHOR NELLE VAROQUAUX NELLEVAROQUAUXGMAILCOM
LICENSE BSD
PRINTDOC
IMPORT NUMPY AS NP
FROM MATPLOTLIB IMPORT PYPLOTASPLT
FROM MATPLOTLIBCOLLECTIONS IMPORT LINECOLLECTION
FROM SKLEARN IMPORT MANIFOLD
FROM SKLEARNMETRICS IMPORT EUCLIDEANDISTANCES
FROM SKLEARNDECOMPOSITION IMPORT PCA
NSAMPLES 20
SEED NPRANDOMRANDOMSTATESEED3
XTRUE SEEDRANDINT0 20 2 NSAMPLESASTYPENPFLOAT
XTRUE XTRUERESHAPENSAMPLES 2
CENTER THE DATA
XTRUE XTRUEMEAN
SIMILARITIES EUCLIDEANDISTANCESXTRUE
ADD NOISE TO THE SIMILARITIES
NOISE NPRANDOMRANDNSAMPLES NSAMPLES
NOISE NOISE NOISET
NOISENPARANGENOISESHAPED0 NPARANGENOISESHAPED0
1216 CHAPTER 5 EXAMPLES
```

SCIKITLEARN USER GUIDE RELEASE 0213  
SIMILARITIES NOISE  
MDS MANIFOLDMDSNCOMPONENTS2 MAXITER3000 EPS1E9 RANDOMSTATESEED  
DISSIMILARITYPRECOMPUTED NJOBS1  
POS MDSFITSIMILARITIESEMBEDDING  
NMDS MANIFOLDMDSNCOMPONENTS2 METRICFALSE MAXITER3000 EPS1E12  
DISSIMILARITYPRECOMPUTED RANDOMSTATESEED NJOBS1  
NINIT1  
NPOS NMDSFITTRANSFORMSIMILARITIES INITPOS  
RESCALE THE DATA  
POS NPSQRTXTRUE 2SUM NPSQRTPOS 2SUM  
NPOS NPSQRTXTRUE 2SUM NPSQRTNPOS 2SUM  
ROTATE THE DATA  
CLF PCANCOMPONENTS2  
XTRUE CLFFITTRANSFORMXTRUE  
POS CLFFITTRANSFORMPOS  
NPOS CLFFITTRANSFORMNPOS  
FIG PLTFigure1  
AX PLTAXES0 0 1 1  
S 100  
PLTSCATTERXTRUE 0 XTRUE 1 COLORNAVY SS LW0  
LABELTRUE POSITION  
PLTSCATTERPOS 0 POS 1 COLORTURQUOISE SS LW0 LABELMDS  
PLTSCATTERNPOS 0 NPOS 1 COLORDARKORANGE SS LW0 LABELNMDS  
PLTLEGENDSCATTERPOINTS1 LOCBEST SHADOWFALSE  
SIMILARITIES SIMILARITIESMAX SIMILARITIES 100  
SIMILARITIESNPISINFSIMILARITIES 0  
PLOT THE EDGES  
STARTIDX ENDIDX NPWHEREPOS  
A SEQUENCE OF LINE0LINE1LINE2 WHERE  
LINEN X0 Y0 X1 Y1 XM YM  
SEGMENTS XTRUEI XTRUEJ  
FORIINRANGELENPOS FORJINRANGELENPOS  
VALUES NPABSSIMILARITIES  
LC LINECOLLECTIONSEGMENTS  
ZORDER0 CMAPPLTCMBLUES  
NORMPLTNORMALIZE0 VALUESMAX  
LCSETARRAYSIMILARITIESFLATTEN  
LCSETLINEWIDTHHSNPFULLLLENSEGMENTS 05  
AXADDCOLLECTIONLC  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0063 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
519 MANIFOLD LEARNING 1217

SCIKITLEARN USER GUIDE RELEASE 0213

5193 TSNE THE EFFECT OF VARIOUS PERPLEXITY VALUES ON THE SHAPE

AN ILLUSTRATION OF TSNE ON THE TWO CONCENTRIC CIRCLES AND THE SCURVE DATASETS FOR DIFFERENT PERPLEXITY VALUES

WE OBSERVE A TENDENCY TOWARDS CLEARER SHAPES AS THE PREPLEXITY VALUE INCREASES

THE SIZE THE DISTANCE AND THE SHAPE OF CLUSTERS MAY VARY UPON INITIALIZATION PERPLEXITY VALUES AND DOES NOT ALWAYS CONVEY A MEANING

AS SHOWN BELOW TSNE FOR HIGHER PERPLEXITIES FINDS MEANINGFUL TOPOLOGY OF TWO CONCENTRIC CIRCLES HOWEVER THE SIZE AND THE DISTANCE OF THE CIRCLES VARIES SLIGHTLY FROM THE ORIGINAL CONTRARY TO THE TWO CIRCLES DATASET THE SHAPES VISUALLY DIVERGE FROM SCURVE TOPOLOGY ON THE SCURVE DATASET EVEN FOR LARGER PERPLEXITY VALUES

FOR FURTHER DETAILS “HOW TO USE TSNE EFFECTIVELY” [HTTPSDISTILLPUB2016MISREADTSNE](https://distill.pub/2016/misreadtsne) PROVIDES A GOOD DISCUSSION OF THE EFFECTS OF VARIOUS PARAMETERS AS WELL AS INTERACTIVE PLOTS TO EXPLORE THOSE EFFECTS

OUT

CIRCLES PERPLEXITY5 IN 089 SEC

CIRCLES PERPLEXITY30 IN 12 SEC

CIRCLES PERPLEXITY50 IN 13 SEC

CIRCLES PERPLEXITY100 IN 18 SEC

SCURVE PERPLEXITY5 IN 092 SEC

SCURVE PERPLEXITY30 IN 12 SEC

SCURVE PERPLEXITY50 IN 14 SEC

SCURVE PERPLEXITY100 IN 19 SEC

UNIFORM GRID PERPLEXITY5 IN 088 SEC

UNIFORM GRID PERPLEXITY30 IN 11 SEC

UNIFORM GRID PERPLEXITY50 IN 11 SEC

UNIFORM GRID PERPLEXITY100 IN 17 SEC

1218 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213  
AUTHOR NARINE KOKHLYAN NARINESLICECOM  
LICENSE BSD  
PRINTDOC  
IMPORT NUMPY AS NP  
IMPORT MATPLOTLIBPYPLOT AS PLT  
FROM MATPLOTLIBTICKER IMPORT NULLFORMATTER  
FROM SKLEARN IMPORT MANIFOLD DATASETS  
FROM TIME IMPORT TIME  
NSAMPLES 300  
NCOMPONENTS 2  
FIG SUBPLOTS PLTSUBPLOTS3 5 FIGSIZE15 8  
PERPLEXITIES 5 30 50 100  
X Y DATASETSMAKECIRCLESNSAMPLESNSAMPLES FACTOR5 NOISE05  
RED Y 0  
GREEN Y 1  
AX SUBPLOTS00  
AXSCATTERXRED 0 XRED 1 CR  
AXSCATTERXGREEN 0 XGREEN 1 CG  
AXXAXISSETMAJORFORMATTERNULLFORMATTER  
AXYAXISSETMAJORFORMATTERNULLFORMATTER  
PLTAXISTIGHT  
FORI PERPLEXITY INENUMERATEPERPLEXITIES  
AX SUBPLOTS0I 1  
T0 TIME  
TSNE MANIFOLDTSNENCOMPONENTSNCOMPONENTS INITRANDOM  
RANDOMSTATE0 PERPLEXITYPERPLEXITY  
Y TSNEFITTRANSFORMX  
T1 TIME  
PRINTCIRCLES PERPLEXITY DIN2GSEC PERPLEXITY T1 T0  
AXSETTITLEPERPLEXITY D PERPLEXITY  
AXSCATTERYRED 0 YRED 1 CR  
AXSCATTERYGREEN 0 YGREEN 1 CG  
AXXAXISSETMAJORFORMATTERNULLFORMATTER  
AXYAXISSETMAJORFORMATTERNULLFORMATTER  
AXAXISTIGHT  
ANOTHER EXAMPLE USING SCURVE  
X COLOR DATASETSSAMPLESGENERATORMAKESCURVENSAMPLES RANDOMSTATE0  
AX SUBPLOTS10  
AXSCATTERX 0 X 2 CCOLOR  
AXXAXISSETMAJORFORMATTERNULLFORMATTER  
AXYAXISSETMAJORFORMATTERNULLFORMATTER  
FORI PERPLEXITY INENUMERATEPERPLEXITIES  
AX SUBPLOTS1I 1  
T0 TIME  
TSNE MANIFOLDTSNENCOMPONENTSNCOMPONENTS INITRANDOM  
519 MANIFOLD LEARNING 1219

SCIKITLEARN USER GUIDE RELEASE 0213  
RANDOMSTATE0 PERPLEXITYPERPLEXITY  
Y TSNEFITTRANSFORMX  
T1 TIME  
PRINTSCURVE PERPLEXITY DIN2GSEC PERPLEXITY T1 TO  
AXSETTITLEPERPLEXITY D PERPLEXITY  
AXSCATTERY 0 Y 1 CCOLOR  
AXXAXISSETMAJORFORMATTERNULLFORMATTER  
AXYAXISSETMAJORFORMATTERNULLFORMATTER  
AXAXISTIGHT  
ANOTHER EXAMPLE USING A 2D UNIFORM GRID  
X NPLinspace0 1 INTNPSQRTNSAMPLES  
XX YY NPMESHGRIDX X  
X NPHSTACK  
XXRAVELRESHAPE1 1  
YYRAVELRESHAPE1 1  
  
COLOR XXRAVEL  
AX SUBPLOTS20  
AXSCATTERX 0 X 1 CCOLOR  
AXXAXISSETMAJORFORMATTERNULLFORMATTER  
AXYAXISSETMAJORFORMATTERNULLFORMATTER  
FORI PERPLEXITY INENUMERATEPERPLEXITIES  
AX SUBPLOTS2I 1  
T0 TIME  
TSNE MANIFOLDTSNENCOMPONENTSNCOMPONENTS INITRANDOM  
RANDOMSTATE0 PERPLEXITYPERPLEXITY  
Y TSNEFITTRANSFORMX  
T1 TIME  
PRINTUNIFORM GRID PERPLEXITY DIN2GSEC PERPLEXITY T1 TO  
AXSETTITLEPERPLEXITY D PERPLEXITY  
AXSCATTERY 0 Y 1 CCOLOR  
AXXAXISSETMAJORFORMATTERNULLFORMATTER  
AXYAXISSETMAJORFORMATTERNULLFORMATTER  
AXAXISTIGHT  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 15568 SECONDS  
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE  
5194 COMPARISON OF MANIFOLD LEARNING METHODS  
AN ILLUSTRATION OF DIMENSIONALITY REDUCTION ON THE SCURVE DATASET WITH VARIOUS MANIFOLD LEARNING METHODS  
FOR A DISCUSSION AND COMPARISON OF THESE ALGORITHMS SEE THE MANIFOLD MODULE PAGE  
FOR A SIMILAR EXAMPLE WHERE THE METHODS ARE APPLIED TO A SPHERE DATASET SEE MANIFOLD LEARNING METHODS ON A SEVERED  
SPHERE  
1220 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
NOTE THAT THE PURPOSE OF THE MDS IS TO FIND A LOWDIMENSIONAL REPRESENTATION OF THE DATA HERE 2D IN WHICH THE
DISTANCES RESPECT WELL THE DISTANCES IN THE ORIGINAL HIGHDIMENSIONAL SPACE UNLIKE OTHER MANIFOLDLEARNING ALGORITHMS
IT DOES NOT SEEKS AN ISOTROPIC REPRESENTATION OF THE DATA IN THE LOWDIMENSIONAL SPACE
OUT
STANDARD 0095 SEC
LTSA 02 SEC
HESSIAN 036 SEC
MODIFIED 02 SEC
ISOMAP 036 SEC
MDS 2 SEC
SPECTRALEMBEDDING 01 SEC
TSNE 6 SEC
AUTHOR JAKE VANDERPLAS VANDERPLASASTROWASHINGTONEDU
PRINTDOC
FROM TIME IMPORT TIME
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM MPLTOOLKITSMPLPLOT3D IMPORT AXES3D
FROM MATPLOTLIBTICKER IMPORT NULLFORMATTER
FROM SKLEARN IMPORT MANIFOLD DATASETS
NEXT LINE TO SILENCE PYFLAKES THIS IMPORT IS NEEDED
AXES3D
519 MANIFOLD LEARNING 1221
```

SCIKITLEARN USER GUIDE RELEASE 0213

NPOINTS 1000

X COLOR DATASETSSAMPLESGENERATORMAKESCURVENPOINTS RANDOMSTATEO

NNEIGHBORS 10

NCOMPONENTS 2

FIG PLTFIGUREFIGSIZE15 8

PLTSUPTITLEMANIFOLD LEARNING WITH IPOINTS INEIGHBORS

1000 NNEIGHBORS FONTSIZE14

AX FIGADDSUBPLOT251 PROJECTION3D

AXSCATTERX 0 X 1 X 2 CCOLOR CMAPPLTCMSPECTRAL

AXVIEWINIT4 72

METHODS STANDARD LTSA HESSIAN MODIFIED

LABELS LLE LTSA HESSIAN LLE MODIFIED LLE

FORI METHOD INENUMERATEMETHODS

T0 TIME

Y MANIFOLDLOCALLYLINEAREMBEDDINGNNEIGHBORS NCOMPONENTS

EIGENSOLVERAUTO

METHODMETHODFITTRANSFORMX

T1 TIME

PRINTS2GSEC METHODSI T1 T0

AX FIGADDSUBPLOT252 I

PLTSCATTERY 0 Y 1 CCOLOR CMAPPLTCMSPECTRAL

PLTTITLE S2GSEC LABELSI T1 T0

AXXAXISSETMAJORFORMATTERNULLFORMATTER

AXYAXISSETMAJORFORMATTERNULLFORMATTER

PLTAXISTIGHT

T0 TIME

Y MANIFOLDISOMAPNNEIGHBORS NCOMPONENTSFITTRANSFORMX

T1 TIME

PRINTISOMAP 2GSEC T1 T0

AX FIGADDSUBPLOT257

PLTSCATTERY 0 Y 1 CCOLOR CMAPPLTCMSPECTRAL

PLTTITLEISOMAP 2GSEC T1 T0

AXXAXISSETMAJORFORMATTERNULLFORMATTER

AXYAXISSETMAJORFORMATTERNULLFORMATTER

PLTAXISTIGHT

T0 TIME

MDS MANIFOLDMDSNCOMPONENTS MAXITER100 NINIT1

Y MDSFITTRANSFORMX

T1 TIME

PRINTMDS2GSEC T1 T0

AX FIGADDSUBPLOT258

PLTSCATTERY 0 Y 1 CCOLOR CMAPPLTCMSPECTRAL

PLTTITLEMDS 2GSEC T1 T0

AXXAXISSETMAJORFORMATTERNULLFORMATTER

AXYAXISSETMAJORFORMATTERNULLFORMATTER

PLTAXISTIGHT

T0 TIME

1222 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

SE MANIFOLDSPECTRALEMMBEDDINGNCOMPONENTSNCOMPONENTS

NNEIGHBORSNNEIGHBORS

Y SEFITTRANSFORMX

T1 TIME

PRINTSPECTRALEMMBEDDING 2GSEC T1 T0

AX FIGADDSUBPLOT259

PLTSCATTERY 0 Y 1 CCOLOR CMAPPLTCMSPECTRAL

PLTTITLESPECTRALEMMBEDDING 2GSEC T1 T0

AXXAXISSETMAJORFORMATTERNULLFORMATTER

AXYAXISSETMAJORFORMATTERNULLFORMATTER

PLTAXISTIGHT

T0 TIME

TSNE MANIFOLDTSNENCOMPONENTSNCOMPONENTS INITPCA RANDOMSTATEO

Y TSNEFITTRANSFORMX

T1 TIME

PRINTTSNE 2GSEC T1 T0

AX FIGADDSUBPLOT2 5 10

PLTSCATTERY 0 Y 1 CCOLOR CMAPPLTCMSPECTRAL

PLTTITLETSNE 2GSEC T1 T0

AXXAXISSETMAJORFORMATTERNULLFORMATTER

AXYAXISSETMAJORFORMATTERNULLFORMATTER

PLTAXISTIGHT

PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 9515 SECONDS

NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)

5195 MANIFOLD LEARNING METHODS ON A SEVERED SPHERE

AN APPLICATION OF THE DIFFERENT MANIFOLD LEARNING TECHNIQUES ON A SPHERICAL DATASET HERE ONE CAN SEE THE USE OF DIMENSIONALITY REDUCTION IN ORDER TO GAIN SOME INTUITION REGARDING THE MANIFOLD LEARNING METHODS REGARDING THE DATASET THE POLES ARE CUT FROM THE SPHERE AS WELL AS A THIN SLICE DOWN ITS SIDE THIS ENABLES THE MANIFOLD LEARNING TECHNIQUES TO ‘SPREAD IT OPEN’ WHILST PROJECTING IT ONTO TWO DIMENSIONS

FOR A SIMILAR EXAMPLE WHERE THE METHODS ARE APPLIED TO THE SCURVE DATASET SEE [COMPARISON OF MANIFOLD LEARNING METHODS](#)

NOTE THAT THE PURPOSE OF THE MDS IS TO FIND A LOWDIMENSIONAL REPRESENTATION OF THE DATA HERE 2D IN WHICH THE DISTANCES RESPECT WELL THE DISTANCES IN THE ORIGINAL HIGHDIMENSIONAL SPACE UNLIKE OTHER MANIFOLDLEARNING ALGORITHMS IT DOES NOT SEEKS AN ISOTROPIC REPRESENTATION OF THE DATA IN THE LOWDIMENSIONAL SPACE HERE THE MANIFOLD PROBLEM MATCHES FAIRLY THAT OF REPRESENTING A FLAT MAP OF THE EARTH AS WITH MAP PROJECTION

519 MANIFOLD LEARNING 1223

SCIKITLEARN USER GUIDE RELEASE 0213  
OUT  
STANDARD 0069 SEC  
LTSA 012 SEC  
HESSIAN 029 SEC  
MODIFIED 016 SEC  
ISO 027 SEC  
MDS 12 SEC  
SPECTRAL EMBEDDING 015 SEC  
TSNE 35 SEC  
AUTHOR JAQUES GROBLER JAQUESGROBLERINRIA.FR  
LICENSE BSD 3 CLAUSE  
PRINTDOC  
FROM TIME IMPORT TIME  
IMPORT NUMPY AS NP  
IMPORT MATPLOTLIBPY.PLOT AS PLT  
FROM MPLTOOLKITSMPLPLOT3D IMPORT AXES3D  
FROM MATPLOTLIBTICKER IMPORT NULLFORMATTER  
FROM SKLEARN IMPORT MANIFOLD  
FROM SKLEARNUTILS IMPORT CHECKRANDOMSTATE  
NEXT LINE TO SILENCE PYFLAKES  
AXES3D  
VARIABLES FOR MANIFOLD LEARNING  
1224 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
NNEIGHBORS 10  
NSAMPLES 1000  
CREATE OUR SPHERE  
RANDOMSTATE CHECKRANDOMSTATE0  
P RANDOMSTATERANDNSAMPLES 2NPPI 055  
T RANDOMSTATERANDNSAMPLES NPPI  
SEVER THE POLES FROM THE SPHERE  
INDICES T NPPI NPPI 8 T NPPI 8  
COLORS PINDICES  
X Y Z NPSINTINDICES NPCOSPINDICES  
NPSINTINDICES NPSINPINDICES  
NPCOSTINDICES  
PLOT OUR DATASET  
FIG PLTFIGUREFIGSIZE15 8  
PLTSUPTITLEMANIFOLD LEARNING WITH IPOINTS INEIGHBORS  
1000 NNEIGHBORS FONTSIZE14  
AX FIGADDSUBPLOT251 PROJECTION3D  
AXSCATTERX Y Z CPINDICES CMAPPLTCMRainbow  
AXVIEWINIT40 10  
SPHEREDATA NPARRAYX Y ZT  
PERFORM LOCALLY LINEAR EMBEDDING MANIFOLD LEARNING  
METHODS STANDARD LTSA HESSIAN MODIFIED  
LABELS LLE LTSA HESSIAN LLE MODIFIED LLE  
FORI METHOD INENUMERATEMETHODS  
T0 TIME  
TRANSDATA MANIFOLD  
LOCALLYLINEAREMBEDDINGNNEIGHBORS 2  
METHODMETHODFITTRANSFORMSPHEREDATAT  
T1 TIME  
PRINTS2GSEC METHODSI T1 T0  
AX FIGADDSUBPLOT252 I  
PLTSCATTERTRANSDATA0 TRANSDATA1 CCOLORS CMAPPLTCMRainbow  
PLTTITLE S2GSEC LABELSI T1 T0  
AXXAXISSETMAJORFORMATTERNULLFORMATTER  
AXYAXISSETMAJORFORMATTERNULLFORMATTER  
PLTAXISTIGHT  
PERFORM ISOMAP MANIFOLD LEARNING  
T0 TIME  
TRANSDATA MANIFOLDISOMAPNNEIGHBORS NCOMPONENTS2  
FITTRANSFORMSPHEREDATAT  
T1 TIME  
PRINTS2GSEC ISO T1 T0  
AX FIGADDSUBPLOT257  
PLTSCATTERTRANSDATA0 TRANSDATA1 CCOLORS CMAPPLTCMRainbow  
PLTTITLE S2GSEC ISOMAP T1 T0  
AXXAXISSETMAJORFORMATTERNULLFORMATTER  
AXYAXISSETMAJORFORMATTERNULLFORMATTER  
PLTAXISTIGHT  
519 MANIFOLD LEARNING 1225

SCIKITLEARN USER GUIDE RELEASE 0213  
PERFORM MULTIDIMENSIONAL SCALING  
T0 TIME  
MDS MANIFOLDMDS2 MAXITER100 NINIT1  
TRANSDATA MDSFITTRANSFORMSPHEREDATAT  
T1 TIME  
PRINTMDS2GSEC T1 T0  
AX FIGADDSUBPLOT258  
PLTSCATTERTRANSDATA0 TRANSDATA1 CCOLORS CMAPPLTCMRainbow  
PLTTITLEMDS 2GSEC T1 T0  
AXXAXISSETMAJORFORMATTERNULLFORMATTER  
AXYAXISSETMAJORFORMATTERNULLFORMATTER  
PLTAXISTIGHT  
PERFORM SPECTRAL EMBEDDING  
T0 TIME  
SE MANIFOLDSPECTRALEMBEDDINGNCOMPONENTS2  
NNEIGHBORSNNEIGHBORS  
TRANSDATA SEFITTRANSFORMSPHEREDATAT  
T1 TIME  
PRINTSPECTRAL EMBEDDING 2GSEC T1 T0  
AX FIGADDSUBPLOT259  
PLTSCATTERTRANSDATA0 TRANSDATA1 CCOLORS CMAPPLTCMRainbow  
PLTTITLESPECTRAL EMBEDDING 2GSEC T1 T0  
AXXAXISSETMAJORFORMATTERNULLFORMATTER  
AXYAXISSETMAJORFORMATTERNULLFORMATTER  
PLTAXISTIGHT  
PERFORM TDISTRIBUTED STOCHASTIC NEIGHBOR EMBEDDING  
T0 TIME  
TSNE MANIFOLDTSNENCOMPONENTS2 INITPCA RANDOMSTATE0  
TRANSDATA TSNEFITTRANSFORMSPHEREDATAT  
T1 TIME  
PRINTTSNE 2GSEC T1 T0  
AX FIGADDSUBPLOT2 5 10  
PLTSCATTERTRANSDATA0 TRANSDATA1 CCOLORS CMAPPLTCMRainbow  
PLTTITLETSNE 2GSEC T1 T0  
AXXAXISSETMAJORFORMATTERNULLFORMATTER  
AXYAXISSETMAJORFORMATTERNULLFORMATTER  
PLTAXISTIGHT  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 5954 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
5196 MANIFOLD LEARNING ON HANDWRITTEN DIGITS LOCALLY LINEAR EMBEDDING  
ISOMAP  
AN ILLUSTRATION OF VARIOUS EMBEDDINGS ON THE DIGITS DATASET  
1226 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213

THE RANDOMTREESEMBEDDING FROM THE SKLEARNENSEMBLE MODULE IS NOT TECHNICALLY A MANIFOLD EMBEDDING METHOD AS IT LEARN A HIGHDIMENSIONAL REPRESENTATION ON WHICH WE APPLY A DIMENSIONALITY REDUCTION METHOD HOWEVER IT IS OFTEN USEFUL TO CAST A DATASET INTO A REPRESENTATION IN WHICH THE CLASSES ARE LINEARLYSEPARABLE TSNE WILL BE INITIALIZED WITH THE EMBEDDING THAT IS GENERATED BY PCA IN THIS EXAMPLE WHICH IS NOT THE DEFAULT SETTING IT ENSURES GLOBAL STABILITY OF THE EMBEDDING IE THE EMBEDDING DOES NOT DEPEND ON RANDOM INITIALIZATION LINEAR DISCRIMINANT ANALYSIS FROM THE SKLEARNDISCRIMINANTANALYSIS MODULE AND NEIGHBORHOOD COMPONENTS ANALYSIS FROM THE SKLEARNNEIGHBORS MODULE ARE SUPERVISED DIMENSIONALITY REDUCTION METHOD IE THEY MAKE USE OF THE PROVIDED LABELS CONTRARY TO OTHER METHODS

- 

519 MANIFOLD LEARNING 1227



SCIKITLEARN USER GUIDE RELEASE 0213

- 519 MANIFOLD LEARNING 1229

SCIKITLEARN USER GUIDE RELEASE 0213

- 1230 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

- 

519 MANIFOLD LEARNING 1231



SCIKITLEARN USER GUIDE RELEASE 0213

- 519 MANIFOLD LEARNING 1233





SCIKITLEARN USER GUIDE RELEASE 0213

- 519 MANIFOLD LEARNING 1235



SCIKITLEARN USER GUIDE RELEASE 0213

- 519 MANIFOLD LEARNING 1237



SCIKITLEARN USER GUIDE RELEASE 0213

- 519 MANIFOLD LEARNING 1239

SCIKITLEARN USER GUIDE RELEASE 0213

- 

OUT

COMPUTING RANDOM PROJECTION

COMPUTING PCA PROJECTION

COMPUTING LINEAR DISCRIMINANT ANALYSIS PROJECTION

COMPUTING ISOMAP PROJECTION

DONE

COMPUTING LLE EMBEDDING

DONE RECONSTRUCTION ERROR 163544E06

COMPUTING MODIFIED LLE EMBEDDING

DONE RECONSTRUCTION ERROR 0360652

COMPUTING HESSIAN LLE EMBEDDING

DONE RECONSTRUCTION ERROR 0212801

COMPUTING LTSA EMBEDDING

DONE RECONSTRUCTION ERROR 0212808

COMPUTING MDS EMBEDDING

DONE STRESS 148085982692961

COMPUTING TOTALLY RANDOM TREES EMBEDDING

COMPUTING SPECTRAL EMBEDDING

COMPUTING TSNE EMBEDDING

COMPUTING NCA PROJECTION

1240 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
AUTHORS FABIAN PEDREGOSA FABIANPEDREGOSAINRIA.FR
OLIVIER GRISEL OLIVIERGRISELENSTA.ORG
MATHIEU BLONDEL MATHIEU.BLONDEL.ORG
GAEL VAROQUAUX
LICENSE BSD 3 CLAUSE C INRIA 2011
PRINTDOC
FROM TIME IMPORT TIME
IMPORT NUMPY AS NP
IMPORT MATPLOTLIB.PY.PLOT AS PLT
FROM MATPLOTLIB IMPORT OFFSETBOX
FROM SKLEARN IMPORT MANIFOLD.DATASETS.DECOMPOSITION.ENSEMBLE
DISCRIMINANT.ANALYSIS.RANDOM.PROJECTION.NEIGHBORS
DIGITS.DATASETS.LOAD.DIGITS.N.CLASS.6
X = DIGITS.DATA
Y = DIGITS.TARGET
N.SAMPLES = N.FEATURES
X.SHAPE
N.NEIGHBORS = 30

SCALE AND VISUALIZE THE EMBEDDING VECTORS
DEF.PLOT.EMBEDDING(X, TITLE=NONE,
XMIN=X.MAX, NPMIN=X.0, NPMAX=X.0,
X=X, X=XMIN, X=XMAX, X=XMIN,
PLT.FIGURE
AX = PLT.SUBPLOT(1,1,1)
FOR I IN RANGE(X.SHAPE[0])
    PLT.TEXT(X[I,0], X[I,1], STR(Y[I]))
    COLOR = PLT.CM.SET(1,Y[I],10)
    FONT.DICT['weight':'bold', 'size':9]
    IF HAS.ATTR('offsetbox', 'annotationbbox')
        ONLY PRINT THUMBNAILS WITH MATPLOTLIB 10
    SHOWNIMAGES = NP.ARRAY(1,1) JUST SOMETHING BIG
    FOR I IN RANGE(X.SHAPE[0])
        DIST = NPSUM(X[I], SHOWNIMAGES[2,1])
        IF NPMIN.DIST > 4E3
            DONT SHOW POINTS THAT ARE TOO CLOSE
        CONTINUE
    SHOWNIMAGES = NPR.SHOWNIMAGES(X[I])
    IMAGEBOX = OFFSETBOX.ANNOTATION.BBOX
    OFFSETBOX.OFFSET.IMAGE.DIGITS.IMAGES(I).CMAP(PLT.CM.GRAY)
    XI
    AX.ADD.ARTIST.IMAGEBOX
    PLT.XTICKS(PLT.YTICKS)
    IF TITLE IS NOT NONE
        PLT.TITLE(TITLE)

PLOT IMAGES OF THE DIGITS
NIMG.PERROW = 20
IMG = NP.ZEROS(10, NIMG.PERROW, 10, NIMG.PERROW)
FOR I IN RANGE(NIMG.PERROW)
519 MANIFOLD LEARNING 1241
```

SCIKITLEARN USER GUIDE RELEASE 0213

IX 10 I 1  
FORJINRANGENIMGPERROW  
IY 10 J 1  
IMGIXIX 8 IYIY 8 XI NIMGPERROW JRESHAPE8 8  
PLTIMSHOWIMG CMAPPLTCMBINARY  
PLTXTICKS  
PLTYTICKS  
PLTTITLEA SELECTION FROM THE 64DIMENSIONAL DIGITS DATASET

RANDOM 2D PROJECTION USING A RANDOM UNITARY MATRIX  
PRINTCOMPUTING RANDOM PROJECTION  
RP RANDOMPROJECTIONSPARSERANDOMPROJECTIONNCOMPONENTS2 RANDOMSTATE42  
XPROJECTED RPFITTRANSFORMX  
PLOTEMBEDDINGXPROJECTED RANDOM PROJECTION OF THE DIGITS

PROJECTION ON TO THE FIRST 2 PRINCIPAL COMPONENTS  
PRINTCOMPUTING PCA PROJECTION  
T0 TIME  
XPCA DECOMPOSITIONTRUNCATEDSVDNCOMPONENTS2FITTRANSFORMX  
PLOTEMBEDDINGXPCA  
PRINCIPAL COMPONENTS PROJECTION OF THE DIGITS TIME 2FS  
TIME T0

PROJECTION ON TO THE FIRST 2 LINEAR DISCRIMINANT COMPONENTS  
PRINTCOMPUTING LINEAR DISCRIMINANT ANALYSIS PROJECTION  
X2 XCOPY  
X2FLATXSHAPE1 1 001 MAKE X INVERTIBLE  
T0 TIME  
XLDA DISCRIMINANTANALYSISLINEARDISCRIMINANTANALYSISNCOMPONENTS2FIT  
↔TRANSFORMX2 Y  
PLOTEMBEDDINGXLDA  
LINEAR DISCRIMINANT PROJECTION OF THE DIGITS TIME 2FS  
TIME T0

ISOMAP PROJECTION OF THE DIGITS DATASET  
PRINTCOMPUTING ISOMAP PROJECTION  
T0 TIME  
XISO MANIFOLDISOMAPNNEIGHBORS NCOMPONENTS2FITTRANSFORMX  
PRINTDONE  
PLOTEMBEDDINGXISO  
ISOMAP PROJECTION OF THE DIGITS TIME 2FS  
TIME T0

LOCALLY LINEAR EMBEDDING OF THE DIGITS DATASET  
PRINTCOMPUTING LLE EMBEDDING  
CLF MANIFOLDLOCALLYLINEAREMBEDDINGNNEIGHBORS NCOMPONENTS2  
1242 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213

METHODSTANDARD

T0 TIME

XLLE CLFFITTRANSFORMX

PRINTDONE RECONSTRUCTION ERROR G CLFRECONSTRUCTIONERROR

PLOTEMBEDDINGXLLE

LOCALLY LINEAR EMBEDDING OF THE DIGITS TIME 2FS

TIME T0

MODIFIED LOCALLY LINEAR EMBEDDING OF THE DIGITS DATASET

PRINTCOMPUTING MODIFIED LLE EMBEDDING

CLF MANIFOLDLOCALLYLINEAREMBEDDINGNNEIGHBORS NCOMPONENTS2

METHODMODIFIED

T0 TIME

XMLLE CLFFITTRANSFORMX

PRINTDONE RECONSTRUCTION ERROR G CLFRECONSTRUCTIONERROR

PLOTEMBEDDINGXMLLE

MODIFIED LOCALLY LINEAR EMBEDDING OF THE DIGITS TIME 2FS

TIME T0

HLL EMBEDDING OF THE DIGITS DATASET

PRINTCOMPUTING HESSIAN LLE EMBEDDING

CLF MANIFOLDLOCALLYLINEAREMBEDDINGNNEIGHBORS NCOMPONENTS2

METHODHESSIAN

T0 TIME

XHLL CLFFITTRANSFORMX

PRINTDONE RECONSTRUCTION ERROR G CLFRECONSTRUCTIONERROR

PLOTEMBEDDINGXHLL

HESSIAN LOCALLY LINEAR EMBEDDING OF THE DIGITS TIME 2FS

TIME T0

L TSA EMBEDDING OF THE DIGITS DATASET

PRINTCOMPUTING L TSA EMBEDDING

CLF MANIFOLDLOCALLYLINEAREMBEDDINGNNEIGHBORS NCOMPONENTS2

METHODL TSA

T0 TIME

XL TSA CLFFITTRANSFORMX

PRINTDONE RECONSTRUCTION ERROR G CLFRECONSTRUCTIONERROR

PLOTEMBEDDINGXL TSA

LOCAL TANGENT SPACE ALIGNMENT OF THE DIGITS TIME 2FS

TIME T0

MDS EMBEDDING OF THE DIGITS DATASET

PRINTCOMPUTING MDS EMBEDDING

CLF MANIFOLDMDSNCOMPONENTS2 NINIT1 MAXITER100

T0 TIME

XMDS CLFFITTRANSFORMX

PRINTDONE STRESS F CLFSTRESS

PLOTEMBEDDINGXMDS

MDS EMBEDDING OF THE DIGITS TIME 2FS

TIME T0

519 MANIFOLD LEARNING 1243

SCIKITLEARN USER GUIDE RELEASE 0213

RANDOM TREES EMBEDDING OF THE DIGITS DATASET  
PRINTCOMPUTING TOTALLY RANDOM TREES EMBEDDING  
HASHER ENSEMBLERANDOMTREEEMBEDDINGNESTIMATORS200 RANDOMSTATE0  
MAXDEPTH5  
T0 TIME  
XTRANSFORMED HASHERFITTRANSFORMX  
PCA DECOMPOSITIONTRUNCATEDSVDNCOMPONENTS2  
XREDUCED PCAFITTRANSFORMXTRANSFORMED  
PLOT EMBEDDINGXREDUCED  
RANDOM FOREST EMBEDDING OF THE DIGITS TIME 2FS  
TIME T0

SPECTRAL EMBEDDING OF THE DIGITS DATASET  
PRINTCOMPUTING SPECTRAL EMBEDDING  
EMBEDDER MANIFOLDSPECTRALEMBEDDINGNCOMPONENTS2 RANDOMSTATE0  
EIGENSOLVERARPACK  
T0 TIME  
XSE EMBEDDERFITTRANSFORMX  
PLOT EMBEDDINGXSE  
SPECTRAL EMBEDDING OF THE DIGITS TIME 2FS  
TIME T0

TSNE EMBEDDING OF THE DIGITS DATASET  
PRINTCOMPUTING TSNE EMBEDDING  
TSNE MANIFOLDTSNENCOMPONENTS2 INITPCA RANDOMSTATE0  
T0 TIME  
XTSNE TSNEFITTRANSFORMX  
PLOT EMBEDDINGXTSNE  
TSNE EMBEDDING OF THE DIGITS TIME 2FS  
TIME T0

NCA PROJECTION OF THE DIGITS DATASET  
PRINTCOMPUTING NCA PROJECTION  
NCA NEIGHBORSNEIGHBORHOODCOMPONENTSANALYSISNCOMPONENTS2 RANDOMSTATE0  
T0 TIME  
XNCA NCAFITTRANSFORMX Y  
PLOT EMBEDDINGXNCA  
NCA EMBEDDING OF THE DIGITS TIME 2FS  
TIME T0  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 17641 SECONDS  
520 GAUSSIAN MIXTURE MODELS  
EXAMPLES CONCERNING THE SKLEARNMIXTURE MODULE  
1244 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE
5201 DENSITY ESTIMATION FOR A GAUSSIAN MIXTURE
PLOT THE DENSITY ESTIMATION OF A MIXTURE OF TWO GAUSSIANS DATA IS GENERATED FROM TWO GAUSSIANS WITH DIFFERENT CENTERS
AND COVARIANCE MATRICES
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM MATPLOTLIBCOLORS IMPORT LOGNORM
FROM SKLEARN IMPORT MIXTURE
NSAMPLES 300
GENERATE RANDOM SAMPLE TWO COMPONENTS
NPRANDOMSEED0
GENERATE SPHERICAL DATA CENTERED ON 20 20
SHIFTEDGAUSSIAN NPRANDOMRANDNNSAMPLES 2 NPARRAY20 20
GENERATE ZERO CENTERED STRETCHED GAUSSIAN DATA
C NPARRAY0 07 35 7
STRETCHEDGAUSSIAN NPDOTNPRANDOMRANDNNSAMPLES 2 C
520 GAUSSIAN MIXTURE MODELS 1245
```

SCIKITLEARN USER GUIDE RELEASE 0213

CONCATENATE THE TWO DATASETS INTO THE FINAL TRAINING SET

XTRAIN NPVSTACKSHIFTEDGAUSSIAN STRETCHEDGAUSSIAN

FIT A GAUSSIAN MIXTURE MODEL WITH TWO COMPONENTS

CLF MIXTUREGAUSSIANMIXTURENCOMPONENTS2 COVARIANCETYPEFULL

CLFFITXTRAIN

DISPLAY PREDICTED SCORES BY THE MODEL AS A CONTOUR PLOT

X NPLinspace20 30

Y NPLinspace20 40

X Y NPMESHGRIDX Y

XX NPARRAYXRavel YRAVELT

Z CLFScoresSAMPLESXX

Z ZRESHAPEXSHAPE

CS PLTCONTOURX Y Z NORMLOGNORMVMIN10 VMAX10000

LEVELSNPLOGSPACE0 3 10

CB PLTCOLORBARCS SHRINK08 EXTENDBOTH

PLTSCATTERXTRAIN 0 XTRAIN 1 8

PLTTITLENEGATIVE LOGLIKELIHOOD PREDICTED BY A GMM

PLTAXISTIGHT

PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0040 SECONDS

NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE

5202 GAUSSIAN MIXTURE MODEL ELLIPSOIDS

PLOT THE CONFIDENCE ELLIPSOIDS OF A MIXTURE OF TWO GAUSSIANS OBTAINED WITH EXPECTATION MAXIMISATION

GAUSSIANMIXTURE CLASS AND VARIATIONAL INFERENCE BAYESIANGAUSSIANMIXTURE CLASS MODELS WITH A

DIRICHLET PROCESS PRIOR

BOTH MODELS HAVE ACCESS TO FIVE COMPONENTS WITH WHICH TO FIT THE DATA NOTE THAT THE EXPECTATION MAXIMISATION

MODEL WILL NECESSARILY USE ALL FIVE COMPONENTS WHILE THE VARIATIONAL INFERENCE MODEL WILL EFFECTIVELY ONLY USE AS MANY

AS ARE NEEDED FOR A GOOD FIT HERE WE CAN SEE THAT THE EXPECTATION MAXIMISATION MODEL SPLITS SOME COMPONENTS

ARBITRARILY BECAUSE IT IS TRYING TO FIT TOO MANY COMPONENTS WHILE THE DIRICHLET PROCESS MODEL ADAPTS IT NUMBER OF STATE

AUTOMATICALLY

THIS EXAMPLE DOESN'T SHOW IT AS WE'RE IN A LOWDIMENSIONAL SPACE BUT ANOTHER ADVANTAGE OF THE DIRICHLET PROCESS

MODEL IS THAT IT CAN FIT FULL COVARIANCE MATRICES EFFECTIVELY EVEN WHEN THERE ARE LESS EXAMPLES PER CLUSTER THAN THERE ARE

DIMENSIONS IN THE DATA DUE TO REGULARIZATION PROPERTIES OF THE INFERENCE ALGORITHM

1246 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
IMPORT ITERTOOLS
IMPORT NUMPY AS NP
FROM SCIPY IMPORT LINALG
IMPORT MATPLOTLIBPYPLOT AS PLT
IMPORT MATPLOTLIB AS MPL
FROM SKLEARN IMPORT MIXTURE
COLORITER ITERTOOLSCYCLENNAVY C CORNFLOWERBLUE GOLD
DARKORANGE
DEFPLOTRRESULTSX Y MEANS COVARIANCES INDEX TITLE
SPLOT PLTSUBPLOT2 1 1 INDEX
FORI MEAN COVAR COLOR INENUMERATEZIP
MEANS COVARIANCES COLORITER
V W LINALGEIGHCOVAR
V 2NPSQRT2 NPSQRTV
U W0 LINALGNORMW0
AS THE DP WILL NOT USE EVERY COMPONENT IT HAS ACCESS TO
UNLESS IT NEEDS IT WE SHOULDNT PLOT THE REDUNDANT
COMPONENTS
IF NOTNPANYY I
CONTINUE
PLTSCATTERXY I 0 XY I 1 8 COLORCOLOR
520 GAUSSIAN MIXTURE MODELS 1247
```

SCIKITLEARN USER GUIDE RELEASE 0213

PLOT AN ELLIPSE TO SHOW THE GAUSSIAN COMPONENT

ANGLE NPARCTANU1 U0

ANGLE 180 ANGLE NPPI CONVERT TO DEGREES

ELL MPLPATCHESELLIPSEMEAN V0 V1 180 ANGLE COLORCOLOR

ELLSETCLIPBOXSPLOTBBOX

ELLSETALPHA05

SPLOTADDARTISTELL

PLTXLIM9 5

PLTYLIM3 6

PLXTTICKS

PLTYTICKS

PLTTITLETITLE

NUMBER OF SAMPLES PER COMPONENT

NSAMPLES 500

GENERATE RANDOM SAMPLE TWO COMPONENTS

NPRANDOMSEED0

C NPARRAY0 01 17 4

X NPRNPDOTNPRANDOMRANDNNSAMPLES 2 C

7NPRANDOMRANDNNSAMPLES 2 NPARRAY6 3

FIT A GAUSSIAN MIXTURE WITH EM USING FIVE COMPONENTS

GMM MIXTUREGAUSSIANMIXTURENCOMPONENTS5 COVARIANCETYPEFULLFITX

PLOTRESULTSX GMMMPREDICTX GMMMEANS GMMCovARIANCES 0

GAUSSIAN MIXTURE

FIT A DIRICHLET PROCESS GAUSSIAN MIXTURE USING FIVE COMPONENTS

DPGMM MIXTUREBAYESIANGAUSSIANMIXTURENCOMPONENTS5

COVARIANCETYPEFULLFITX

PLOTRESULTSX DPGMMPREDICTX DPGMMMEANS DPGMMCovARIANCES 1

BAYESIAN GAUSSIAN MIXTURE WITH A DIRICHLET PROCESS PRIOR

PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0138 SECONDS

NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE

5203 GAUSSIAN MIXTURE MODEL SELECTION

THIS EXAMPLE SHOWS THAT MODEL SELECTION CAN BE PERFORMED WITH GAUSSIAN MIXTURE MODELS USING INFORMATIONTHEORETIC CRITERIA BIC MODEL SELECTION CONCERNS BOTH THE COVARIANCE TYPE AND THE NUMBER OF COMPONENTS IN THE MODEL IN THAT CASE AIC ALSO PROVIDES THE RIGHT RESULT NOT SHOWN TO SAVE TIME BUT BIC IS BETTER SUITED IF THE PROBLEM IS TO IDENTIFY THE RIGHT MODEL UNLIKE BAYESIAN PROCEDURES SUCH INFERENCES ARE PRIORFREE

IN THAT CASE THE MODEL WITH 2 COMPONENTS AND FULL COVARIANCE WHICH CORRESPONDS TO THE TRUE GENERATIVE MODEL IS SELECTED

1248 CHAPTER 5 EXAMPLES

```

SCIKITLEARN USER GUIDE RELEASE 0213
IMPORT NUMPY AS NP
IMPORT ITERTOOLS
FROM SCIPY IMPORT LINALG
IMPORT MATPLOTLIBPYPLOT AS PLT
IMPORT MATPLOTLIB AS MPL
FROM SKLEARN IMPORT MIXTURE
PRINTDOC
  NUMBER OF SAMPLES PER COMPONENT
NSAMPLES 500
GENERATE RANDOM SAMPLE TWO COMPONENTS
NPRANDOMSEED0
C  NPARRAY0 01 17 4
X  NPRNPDOTNPRANDOMRANDNNSAMPLES 2 C
7NPRANDOMRANDNNSAMPLES 2  NPARRAY6 3
LOWESTBIC  NPINFTY
BIC
NCOMPONENTSRANGE  RANGE1 7
CVTYPES  SPHERICAL TIED DIAG FULL
FORCVTYPE INCVTYPES
FORNCOMPONENTS INNCOMPONENTSRANGE
520 GAUSSIAN MIXTURE MODELS 1249

```

SCIKITLEARN USER GUIDE RELEASE 0213  
FIT A GAUSSIAN MIXTURE WITH EM  
GMM MIXTUREGAUSSIANMIXTURENCOMPONENTSNCOMPONENTS  
COVARIANCETYPEPCVTYPE  
GMMFITX  
BICAPPENDGMMBICX  
IFBIC1 LOWESTBIC  
LOWESTBIC BIC1  
BESTGMM GMM  
BIC NPARRAYBIC  
COLORITER ITERTOOLSCYCLENAVY TURQUOISE CORNFLOWERBLUE  
DARKORANGE  
CLF BESTGMM  
BARS  
PLOT THE BIC SCORES  
PLTFIGUREFIGSIZE8 6  
SPL PLTSUBPLOT2 1 1  
FORI CVTYPE COLOR INENUMERATEZIPCVTYPES COLORITER  
XPOS NPARRAYNCOMPONENTSRANGE 2 1 2  
BARSAPPENDPLTBARXPOS BICI LENNCOMPONENTSRANGE  
I 1 LENNCOMPONENTSRANGE  
WIDTH2 COLORCOLOR  
PLTXTICKSNCOMPONENTSRANGE  
PLTYLIMBICMIN 101 01 BICMAX BICMAX  
PLTTITLEBIC SCORE PER MODEL  
XPOS NPMODBICARGMIN LENNCOMPONENTSRANGE 65  
2NPFLOORBICARGMIN LENNCOMPONENTSRANGE  
PLTTEXTXPOS BICMIN 097 03 BICMAX FONTSIZE14  
SPLSETXLABELNUMBER OF COMPONENTS  
SPLLEGENDB0 FORBINBARS CVTYPES  
PLOT THE WINNER  
SPLOT PLTSUBPLOT2 1 2  
Y CLFPREDICTX  
FORI MEAN COV COLOR INENUMERATEZIPCLFMEANS CLFCOVARIANCES  
COLORITER  
V W LINALGEIGHCOV  
IF NOTNPANY Y I  
CONTINUE  
PLTSCATTERXY I 0 XY I 1 8 COLORCOLOR  
PLOT AN ELLIPSE TO SHOW THE GAUSSIAN COMPONENT  
ANGLE NPARCTAN2W01 W00  
ANGLE 180 ANGLE NPPI CONVERT TO DEGREES  
V 2NPSQRT2 NPSQRTV  
ELL MPLPATCHESELLIPSEMEAN V0 V1 180 ANGLE COLORCOLOR  
ELLSETCLIPBOXSPLOTBBOX  
ELLSETALPHA5  
SPLOTADDARTISTELL  
PLTXTICKS  
PLTYTICKS  
PLTTITLESELECTED GMM FULL MODEL 2 COMPONENTS  
PLTSUBPLOTSADJUSTHSPACE35 BOTTOM02  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0207 SECONDS  
1250 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213

NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)

5204 GMM COVARIANCES

DEMONSTRATION OF SEVERAL COVARIANCES TYPES FOR GAUSSIAN MIXTURE MODELS

SEEGAUSSIAN MIXTURE MODELS FOR MORE INFORMATION ON THE ESTIMATOR

ALTHOUGH GMM ARE OFTEN USED FOR CLUSTERING WE CAN COMPARE THE OBTAINED CLUSTERS WITH THE ACTUAL CLASSES FROM THE DATASET WE INITIALIZE THE MEANS OF THE GAUSSIANS WITH THE MEANS OF THE CLASSES FROM THE TRAINING SET TO MAKE THIS COMPARISON VALID

WE PLOT PREDICTED LABELS ON BOTH TRAINING AND HELD OUT TEST DATA USING A VARIETY OF GMM COVARIANCE TYPES ON THE IRIS DATASET WE COMPARE GMMS WITH SPHERICAL DIAGONAL FULL AND TIED COVARIANCE MATRICES IN INCREASING ORDER OF PERFORMANCE ALTHOUGH ONE WOULD EXPECT FULL COVARIANCE TO PERFORM BEST IN GENERAL IT IS PRONE TO OVERFITTING ON SMALL DATASETS AND DOES NOT GENERALIZE WELL TO HELD OUT TEST DATA

ON THE PLOTS TRAIN DATA IS SHOWN AS DOTS WHILE TEST DATA IS SHOWN AS CROSSES THE IRIS DATASET IS FOURDIMENSIONAL ONLY THE FIRST TWO DIMENSIONS ARE SHOWN HERE AND THUS SOME POINTS ARE SEPARATED IN OTHER DIMENSIONS

520 GAUSSIAN MIXTURE MODELS 1251

SCIKITLEARN USER GUIDE RELEASE 0213  
AUTHOR RON WEISS RONWEISSGMAILCOM GAE VAROQUAUX  
MODIFIED BY THIERRY GUILLEMOT THIERRYGUILLEMOTWORKGMAILCOM  
LICENSE BSD 3 CLAUSE  
IMPORT MATPLOTLIB AS MPL  
IMPORT MATPLOTLIBPYPLOT AS PLT  
IMPORT NUMPY AS NP  
FROM SKLEARN IMPORT DATASETS  
FROM SKLEARNMIXTURE IMPORT GAUSSIANMIXTURE  
FROM SKLEARNMODELSELECTION IMPORT STRATIFIEDKFOLD  
PRINTDOC  
COLORS NAVY TURQUOISE DARKORANGE  
1252 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
DEFMAKEELLIPSESGMM AX  
FORN COLOR INENUMERATECOLORS  
IFGMMCOVARIANCETYPE FULL  
COVARIANCES GMMCOVARIANCESN2 2  
ELIFGMMCOVARIANCETYPE TIED  
COVARIANCES GMMCOVARIANCES2 2  
ELIFGMMCOVARIANCETYPE DIAG  
COVARIANCES NPDIAGGMMCOVARIANCESN2  
ELIFGMMCOVARIANCETYPE SPHERICAL  
COVARIANCES NPEYEGMMMEANSSHAPE1 GMMCOVARIANCESN  
V W NPLINALGEIGHCOVARIANCES  
U W0 NPLINALGNORMW0  
ANGLE NPARCTAN2U1 U0  
ANGLE 180 ANGLE NPPI CONVERT TO DEGREES  
V 2NPSQRT2 NPSQRTV  
ELL MPLPATCHESELLIPSEGMMMEANSN 2 V0 V1  
180 ANGLE COLORCOLOR  
ELLSETCLIPBOXAXBBOX  
ELLSETALPHA05  
AXADDARTISTELL  
AXSETASPECTEQUAL DATALIM  
IRIS DATASETSLOADIRIS  
BREAK UP THE DATASET INTO NONOVERLAPPING TRAINING 75 AND TESTING  
25 SETS  
SKF STRATIFIEDKFOLDNSPLITS4  
ONLY TAKE THE FIRST FOLD  
TRAININDEX TESTINDEX NEXTITERSKFSPLITIRISDATA IRISTARGET  
XTRAIN IRISDATATRAININDEX  
YTRAIN IRISTARGETTRAININDEX  
XTEST IRISDATATESTINDEX  
YTEST IRISTARGETTESTINDEX  
NCLASSES LENNPUNIQUEYTRAIN  
TRY GMMS USING DIFFERENT TYPES OF COVARIANCES  
ESTIMATORS COVTYPE GAUSSIANMIXTURENCOMPONENTSNCLASSES  
COVARIANCETYPECOVTYPE MAXITER20 RANDOMSTATE0  
FORCOVTYPE INSPHERICAL DIAG TIED FULL  
NESTIMATORS LENESTIMATORS  
PLTFIGUREFIGSIZE3 NESTIMATORS 2 6  
PLTSUBPLOTSADJUSTBOTTOM01 TOP095 HSPACE15 WSPACE05  
LEFT01 RIGHT99  
FORINDEX NAME ESTIMATOR INENUMERATEESTIMATORSITEMS  
SINCE WE HAVE CLASS LABELS FOR THE TRAINING DATA WE CAN  
INITIALIZE THE GMM PARAMETERS IN A SUPERVISED MANNER  
ESTIMATORMEANSINIT NPARRAYXTRAINYTRAIN IMEANAXISO  
FORIINRANGENCLASSES  
TRAIN THE OTHER PARAMETERS USING THE EM ALGORITHM  
ESTIMATORFITXTRAIN  
520 GAUSSIAN MIXTURE MODELS 1253

SCIKITLEARN USER GUIDE RELEASE 0213  
H PLTSUBPLOT2 NESTIMATORS 2 INDEX 1  
MAKEELLIPSESESTIMATOR H  
FORN COLOR INENUMERATECOLORS  
DATA IRISDATAIRISTARGET N  
PLTSCATTERDATA 0 DATA 1 S08 COLORCOLOR  
LABELIRISTARGETNAMESN  
PLOT THE TEST DATA WITH CROSSES  
FORN COLOR INENUMERATECOLORS  
DATA XTESTYTEST N  
PLTSCATTERDATA 0 DATA 1 MARKERX COLORCOLOR  
YTRAINPRED ESTIMATORPREDICTXTRAIN  
TRAINACCURACY NPMEANYTRAINPREDRAVEL YTRAINRAVEL 100  
PLTTEXT005 09 TRAIN ACCURACY 1F TRAINACCURACY  
TRANSFORMHTRANSAXES  
YTESTPRED ESTIMATORPREDICTXTEST  
TESTACCURACY NPMEANYTESTPREDRAVEL YTESTRAVEL 100  
PLTTEXT005 08 TEST ACCURACY 1F TESTACCURACY  
TRANSFORMHTRANSAXES  
PLTXTICKS  
PLTYTICKS  
PLTTITLENAME  
PLTLEGENDSCATTERPOINTS1 LOCLOWER RIGHT PROPDICTSIZE12  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0097 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
5205 GAUSSIAN MIXTURE MODEL SINE CURVE  
THIS EXAMPLE DEMONSTRATES THE BEHAVIOR OF GAUSSIAN MIXTURE MODELS FIT ON DATA THAT WAS NOT SAMPLED FROM A MIXTURE  
OF GAUSSIAN RANDOM VARIABLES THE DATASET IS FORMED BY 100 POINTS LOOSELY SPACED FOLLOWING A NOISY SINE CURVE THERE  
IS THEREFORE NO GROUND TRUTH VALUE FOR THE NUMBER OF GAUSSIAN COMPONENTS  
THE FIRST MODEL IS A CLASSICAL GAUSSIAN MIXTURE MODEL WITH 10 COMPONENTS FIT WITH THE EXPECTATIONMAXIMIZATION  
ALGORITHM  
THE SECOND MODEL IS A BAYESIAN GAUSSIAN MIXTURE MODEL WITH A DIRICHLET PROCESS PRIOR FIT WITH VARIATIONAL INFERENCE  
THE LOW VALUE OF THE CONCENTRATION PRIOR MAKES THE MODEL FAVOR A LOWER NUMBER OF ACTIVE COMPONENTS THIS MODELS  
“DECIDES” TO FOCUS ITS MODELING POWER ON THE BIG PICTURE OF THE STRUCTURE OF THE DATASET GROUPS OF POINTS WITH ALTERNATING  
DIRECTIONS MODELED BY NONDIAGONAL COVARIANCE MATRICES THOSE ALTERNATING DIRECTIONS ROUGHLY CAPTURE THE ALTERNATING  
NATURE OF THE ORIGINAL SINE SIGNAL  
THE THIRD MODEL IS ALSO A BAYESIAN GAUSSIAN MIXTURE MODEL WITH A DIRICHLET PROCESS PRIOR BUT THIS TIME THE VALUE OF THE  
CONCENTRATION PRIOR IS HIGHER GIVING THE MODEL MORE LIBERTY TO MODEL THE FINEGRAINED STRUCTURE OF THE DATA THE RESULT  
IS A MIXTURE WITH A LARGER NUMBER OF ACTIVE COMPONENTS THAT IS SIMILAR TO THE FIRST MODEL WHERE WE ARBITRARILY DECIDED  
TO FIX THE NUMBER OF COMPONENTS TO 10  
1254 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

WHICH MODEL IS THE BEST IS A MATTER OF SUBJECTIVE JUDGEMENT DO WE WANT TO FAVOR MODELS THAT ONLY CAPTURE THE BIG PICTURE TO SUMMARIZE AND EXPLAIN MOST OF THE STRUCTURE OF THE DATA WHILE IGNORING THE DETAILS OR DO WE PREFER MODELS THAT CLOSELY FOLLOW THE HIGH DENSITY REGIONS OF THE SIGNAL

THE LAST TWO PANELS SHOW HOW WE CAN SAMPLE FROM THE LAST TWO MODELS THE RESULTING SAMPLES DISTRIBUTIONS DO NOT LOOK EXACTLY LIKE THE ORIGINAL DATA DISTRIBUTION THE DIFFERENCE PRIMARILY STEMS FROM THE APPROXIMATION ERROR WE MADE BY USING A MODEL THAT ASSUMES THAT THE DATA WAS GENERATED BY A FINITE NUMBER OF GAUSSIAN COMPONENTS INSTEAD OF A CONTINUOUS NOISY SINE CURVE

IMPORT ITERTOOLS

IMPORT NUMPY AS NP

FROM SCIPY IMPORT LINALG

IMPORT MATPLOTLIBPYPLOT AS PLT

IMPORT MATPLOTLIB AS MPL

520 GAUSSIAN MIXTURE MODELS 1255

SCIKITLEARN USER GUIDE RELEASE 0213  
FROM SKLEARN IMPORT MIXTURE  
PRINTDOC  
COLORITER ITERTOOLSCYCLENAVY C CORNFLOWERBLUE GOLD  
DARKORANGE  
DEFPLOTRESULTSX Y MEANS COVARIANCES INDEX TITLE  
SPLOT PLTSUBPLOT5 1 1 INDEX  
FORI MEAN COVAR COLOR INENUMERATEZIP  
MEANS COVARIANCES COLORITER  
V W LINALGEIGHCOVAR  
V 2NPSQRT2 NPSQRTV  
U W0 LINALGNORMW0  
AS THE DP WILL NOT USE EVERY COMPONENT IT HAS ACCESS TO  
UNLESS IT NEEDS IT WE SHOULDNT PLOT THE REDUNDANT  
COMPONENTS  
IF NOTNPANY Y I  
CONTINUE  
PLTSCATTERXY I 0 XY I 1 8 COLORCOLOR  
PLOT AN ELLIPSE TO SHOW THE GAUSSIAN COMPONENT  
ANGLE NPARCTANU1 U0  
ANGLE 180 ANGLE NPPI CONVERT TO DEGREES  
ELL MPLPATCHESELLIPSEMEAN V0 V1 180 ANGLE COLORCOLOR  
ELLSETCLIPBOXSPLOTBBOX  
ELLSETALPHA05  
SPLOTADDARTISTELL  
PLTXLIM6 4 NPPI 6  
PLTYLIM5 5  
PLTTITLETITLE  
PLTXTICKS  
PLTYTICKS  
DEFPLOTSAMPLESX Y NCOMPONENTS INDEX TITLE  
PLTSUBPLOT5 1 4 INDEX  
FORI COLOR INZIPRANGENCOMPONENTS COLORITER  
AS THE DP WILL NOT USE EVERY COMPONENT IT HAS ACCESS TO  
UNLESS IT NEEDS IT WE SHOULDNT PLOT THE REDUNDANT  
COMPONENTS  
IF NOTNPANY Y I  
CONTINUE  
PLTSCATTERXY I 0 XY I 1 8 COLORCOLOR  
PLTXLIM6 4 NPPI 6  
PLTYLIM5 5  
PLTTITLETITLE  
PLTXTICKS  
PLTYTICKS  
PARAMETERS  
NSAMPLES 100  
1256 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
GENERATE RANDOM SAMPLE FOLLOWING A SINE CURVE
NPRANDOMSEED0
X NPZEROSNSAMPLES 2
STEP 4 NPPI NSAMPLES
FORIINRANGEXSHAPE0
X ISTEP 6
XI 0 X NPRANDOMNORMAL0 01
XI 1 3 NPSINX NPRANDOMNORMAL0 2
PLTFIGUREFIGSIZE10 10
PLTSUBPLOTSADJUSTBOTTOM04 TOP095 HSPACE2 WSPACE05
LEFT03 RIGHT97
FIT A GAUSSIAN MIXTURE WITH EM USING TEN COMPONENTS
GMM MIXTUREGAUSSIANMIXTURENCOMPONENTS10 COVARIANCETYPEFULL
MAXITER100FITX
PLOTRESULTSX GMPREDICTX GMMMEANS GMMCOVARIANCES 0
EXPECTATIONMAXIMIZATION
DPGMM MIXTUREBAYESIANGAUSSIANMIXTURE
NCOMPONENTS10 COVARIANCETYPEFULL WEIGHTCONCENTRATIONPRIOR1E2
WEIGHTCONCENTRATIONPRIORTYPEDIRICHLETPROCESS
MEANPRECISIONPRIOR1E2 COVARIANCEPRIOR1E0 NPEYE2
INITPARAMSRANDOM MAXITER100 RANDOMSTATE2FITX
PLOTRESULTSX DPGMMPREDICTX DPGMMMEANS DPGMMCovARIANCES 1
BAYESIAN GAUSSIAN MIXTURE MODELS WITH A DIRICHLET PROCESS PRIOR
RFOR GAMMA0001
XS YS DPGMMSAMPLENSAMPLES2000
PLOTSAMPLESXS YS DPGMMNCOMPONENTS 0
GAUSSIAN MIXTURE WITH A DIRICHLET PROCESS PRIOR
RFOR GAMMA0001 SAMPLED WITH 2000 SAMPLES
DPGMM MIXTUREBAYESIANGAUSSIANMIXTURE
NCOMPONENTS10 COVARIANCETYPEFULL WEIGHTCONCENTRATIONPRIOR1E2
WEIGHTCONCENTRATIONPRIORTYPEDIRICHLETPROCESS
MEANPRECISIONPRIOR1E2 COVARIANCEPRIOR1E0 NPEYE2
INITPARAMSKMEANS MAXITER100 RANDOMSTATE2FITX
PLOTRESULTSX DPGMMPREDICTX DPGMMMEANS DPGMMCovARIANCES 2
BAYESIAN GAUSSIAN MIXTURE MODELS WITH A DIRICHLET PROCESS PRIOR
RFOR GAMMA0100
XS YS DPGMMSAMPLENSAMPLES2000
PLOTSAMPLESXS YS DPGMMNCOMPONENTS 1
GAUSSIAN MIXTURE WITH A DIRICHLET PROCESS PRIOR
RFOR GAMMA0100 SAMPLED WITH 2000 SAMPLES
PLTSHOW
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0301 SECONDS
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE
520 GAUSSIAN MIXTURE MODELS 1257
```

SCIKITLEARN USER GUIDE RELEASE 0213

5206 CONCENTRATION PRIOR TYPE ANALYSIS OF VARIATION BAYESIAN GAUSSIAN MIXTURE

THIS EXAMPLE PLOTS THE ELLIPSOIDS OBTAINED FROM A TOY DATASET MIXTURE OF THREE GAUSSIANS FITTED BY THE BAYESIANGAUSSIANMIXTURE CLASS MODELS WITH A DIRICHLET DISTRIBUTION PRIOR WEIGHTCONCENTRATIONPRIORTYPEDIRICHLETDISTRIBUTION AND A DIRICHLET PROCESS PRIOR WEIGHTCONCENTRATIONPRIORTYPEDIRICHLETPROCESS ON EACH FIGURE WE PLOT THE RESULTS FOR THREE DIFFERENT VALUES OF THE WEIGHT CONCENTRATION PRIOR

THEBAYESIANGAUSSIANMIXTURE CLASS CAN ADAPT ITS NUMBER OF MIXTURE COMPONENTS AUTOMATICALLY THE PARAMETERWEIGHTCONCENTRATIONPRIOR HAS A DIRECT LINK WITH THE RESULTING NUMBER OF COMPONENTS WITH NONZERO WEIGHTS SPECIFYING A LOW VALUE FOR THE CONCENTRATION PRIOR WILL MAKE THE MODEL PUT MOST OF THE WEIGHT ON FEW COMPONENTS SET THE REMAINING COMPONENTS WEIGHTS VERY CLOSE TO ZERO HIGH VALUES OF THE CONCENTRATION PRIOR WILL ALLOW A LARGER NUMBER OF COMPONENTS TO BE ACTIVE IN THE MIXTURE

THE DIRICHLET PROCESS PRIOR ALLOWS TO DEFINE AN INFINITE NUMBER OF COMPONENTS AND AUTOMATICALLY SELECTS THE CORRECT NUMBER OF COMPONENTS IT ACTIVATES A COMPONENT ONLY IF IT IS NECESSARY

ON THE CONTRARY THE CLASSICAL FINITE MIXTURE MODEL WITH A DIRICHLET DISTRIBUTION PRIOR WILL FAVOR MORE UNIFORMLY WEIGHTED COMPONENTS AND THEREFORE TENDS TO DIVIDE NATURAL CLUSTERS INTO UNNECESSARY SUBCOMPONENTS

- 

1258 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213

```
•
AUTHOR THIERRY GUILLEMOT THIERRYGUILLEMOTWORKGMAILCOM
LICENSE BSD 3 CLAUSE
IMPORT NUMPY AS NP
IMPORT MATPLOTLIB AS MPL
IMPORT MATPLOTLIBPYPLOT AS PLT
IMPORT MATPLOTLIBGRIDSPEC AS GRIDSPEC
FROM SKLEARNMIXTURE IMPORT BAYESIANGAUSSIANMIXTURE
PRINTDOC
DEFPLOTELLIPSESAX WEIGHTS MEANS COVARs
FORNINRANGEMEANSSSHAPE0
EIGVALS EIGVECS NPLINALGEIGHCOVARSN
UNITEIGVEC EIGVECS0 NPLINALGNORMEIGVECS0
ANGLE NPARCTAN2UNITEIGVEC1 UNITEIGVEC0
ELLIPSE NEEDS DEGREES
ANGLE 180 ANGLE NPPI
EIGENVECTOR NORMALIZATION
EIGVALS 2 NPSQRT2 NPSQRTEIGVALS
ELL MPLPATCHESELLIPSEMEANSN EIGVALS0 EIGVALS1
180 ANGLE EDGECOLORBLACK
ELLSETCLIPBOXAXBBOX
ELLSETALPHAWEIGHTSN
ELLSETFACECOLOR56B4E9
AXADDARTISTELL
DEFPLOTRESULTSAX1 AX2 ESTIMATOR X Y TITLE PLOTTITLEFALSE
AX1SETTITLETITLE
AX1SCATTERX 0 X 1 S5 MARKERO COLORCOLORSY ALPHA08
AX1SETXLIM2 2
AX1SETYLIM3 3
520 GAUSSIAN MIXTURE MODELS 1259
```

SCIKITLEARN USER GUIDE RELEASE 0213  
AX1SETXTICKS  
AX1SETYTICKS  
PLOTELLIPSESAX1 ESTIMATORWEIGHTS ESTIMATORMEANS  
ESTIMATORCOVARIANCES  
AX2GETXAXISSETTICKPARAMSDIRECTIONOUT  
AX2YAXISGRIDTRUE ALPHA07  
FORK WINENUMERATEESTIMATORWEIGHTS  
AX2BARK W WIDTH09 COLOR56B4E9 ZORDER3  
ALIGNCENTER EDGECOLORBLACK  
AX2TEXTK W 0007 1F W100  
HORIZONTALALIGNMENTCENTER  
AX2SETXLIM6 2 NCOMPONENTS 4  
AX2SETYLIM0 11  
AX2TICKPARAMSAXISY WHICHBOTH LEFTFALSE  
RIGHTFALSE LABELLEFTFALSE  
AX2TICKPARAMSAXISX WHICHBOTH TOPFALSE  
IFPLOTTITLE  
AX1SETYLABELESTIMATED MIXTURES  
AX2SETYLABELWEIGHT OF EACH COMPONENT  
PARAMETERS OF THE DATASET  
RANDOMSTATE NCOMPONENTS NFEATURES 2 3 2  
COLORS NPARRAY0072B2 F0E442 D55E00  
COVARs NPARRAY7 0 0 1  
5 0 0 1  
5 0 0 1  
SAMPLES NPARRAY200 500 200  
MEANS NPARRAY0 70  
0 0  
0 70  
MEANPRECISIONPRIOR 08 TO MINIMIZE THE INFLUENCE OF THE PRIOR  
ESTIMATORS  
FINITE MIXTURE WITH A DIRICHLET DISTRIBUTION NPRIOR AND  
RGAMMA0 BAYESIANGAUSSIANMIXTURE  
WEIGHTCONCENTRATIONPRIORTYPEDIRICHLETDISTRIBUTION  
NCOMPONENTS2 NCOMPONENTS REGCOVAR0 INITPARAMSRANDOM  
MAXITER1500 MEANPRECISIONPRIOR8  
RANDOMSTATERANDOMSTATE 0001 1 1000  
INFINITE MIXTURE WITH A DIRICHLET PROCESS NPRIOR AND RGAMMA0  
BAYESIANGAUSSIANMIXTURE  
WEIGHTCONCENTRATIONPRIORTYPEDIRICHLETPROCESS  
NCOMPONENTS2 NCOMPONENTS REGCOVAR0 INITPARAMSRANDOM  
MAXITER1500 MEANPRECISIONPRIOR8  
RANDOMSTATERANDOMSTATE 1 1000 100000  
GENERATE DATA  
RNG NPRANDOMRANDOMSTATERANDOMSTATE  
X NPVSTACK  
RNGMULTIVARIATENORMALMEANSJ COVARsJ SAMPLESJ  
FORJINRANGENCOMPONENTS  
Y NPCONCATENATENPFULLSAMPLESJ J DTYPEINT  
FORJINRANGENCOMPONENTS  
PLOT RESULTS IN TWO DIFFERENT FIGURES  
1260 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
FORTITLE ESTIMATOR CONCENTRATIONSPRIOR INESTIMATORS  
PLTFIGUREFIGSIZE47 3 8  
PLTSUBPLOTSADJUSTBOTTOM04 TOP090 HSPACE05 WSPACE05  
LEFT03 RIGHT99  
GS GRIDSPECGRIDSPEC3 LENCONCENTRATIONSPRIOR  
FORK CONCENTRATION INENUMERATECONCENTRATIONSPRIOR  
ESTIMATORWEIGHTCONCENTRATIONPRIOR CONCENTRATION  
ESTIMATORFITX  
PLOTRESULTSPLOTSUBPLOTGS02 K PLTSUBPLOTGS2 K ESTIMATOR  
X Y R S1E TITLE CONCENTRATION  
PLOTTITLEK 0  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 6923 SECONDS  
521 MODEL SELECTION  
EXAMPLES RELATED TO THE SKLEARNMODELSELECTION MODULE  
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE  
5211 PLOTTING CROSSVALIDATED PREDICTIONS  
THIS EXAMPLE SHOWS HOW TO USE CROSSVALPREDICT TO VISUALIZE PREDICTION ERRORS  
521 MODEL SELECTION 1261

```

SCIKITLEARN USER GUIDE RELEASE 0213
FROM SKLEARN IMPORT DATASETS
FROM SKLEARNMODELSELECTION IMPORT CROSSVALPREDICT
FROM SKLEARN IMPORT LINEARMODEL
IMPORT MATPLOTLIBPYPLOT AS PLT
LR LINEARMODELLINEARREGRESSION
BOSTON DATASETSLOADBOSTON
Y BOSTONTARGET
CROSSVALPREDICT RETURNS AN ARRAY OF THE SAME SIZE AS Y WHERE EACH ENTRY
IS A PREDICTION OBTAINED BY CROSS VALIDATION
PREDICTED CROSSVALPREDICTLR BOSTONDATA Y CV10
FIG AX PLTSUBPLOTS
AXSCATTERY PREDICTED EDGECOLORS0 0 0
AXPLOTYMIN YMAX YMIN YMAX K LW4
AXSETXLABELMEASURED
AXSETYLABELPREDICTED
PLTSHOW
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0024 SECONDS
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE
1262 CHAPTER 5 EXAMPLES

```

```
SCIKITLEARN USER GUIDE RELEASE 0213
5212 PLOTTING VALIDATION CURVES
IN THIS PLOT YOU CAN SEE THE TRAINING SCORES AND VALIDATION SCORES OF AN SVM FOR DIFFERENT VALUES OF THE KERNEL PARAMETER
GAMMA FOR VERY LOW VALUES OF GAMMA YOU CAN SEE THAT BOTH THE TRAINING SCORE AND THE VALIDATION SCORE ARE LOW
THIS IS CALLED UNDERFITTING MEDIUM VALUES OF GAMMA WILL RESULT IN HIGH VALUES FOR BOTH SCORES IE THE CLASSIFIER IS
PERFORMING FAIRLY WELL IF GAMMA IS TOO HIGH THE CLASSIFIER WILL OVERFIT WHICH MEANS THAT THE TRAINING SCORE IS GOOD BUT
THE VALIDATION SCORE IS POOR
PRINTDOC
IMPORT MATPLOTLIBPYPLOT AS PLT
IMPORT NUMPY AS NP
FROM SKLEARNDATASETS IMPORT LOADDIGITS
FROM SKLEARN SVM IMPORT SVC
FROM SKLEARNMODELSELECTION IMPORT VALIDATIONCURVE
DIGITS LOADDIGITS
X Y DIGITSDATA DIGITSTARGET
PARAMRANGE NPLOGSPACE6 1 5
TRAJNSCORES TESTSCORES VALIDATIONCURVE
SVC X Y PARAMNAMEGAMMA PARAMRANGEPARAMRANGE
CV5 SCORINGACCURACY NJOBS1
TRAJNSCORESMEAN NPMEANTRAJNSCORES AXIS1
TRAJNSCORESSTD NPSTDTRAJNSCORES AXIS1
521 MODEL SELECTION 1263
```

SCIKITLEARN USER GUIDE RELEASE 0213  
TESTSCORESMEAN NPMEANTESTSCORES AXIS1  
TESTSCORESSTD NPSTDTESTSCORES AXIS1  
PLTTITLEVALIDATION CURVE WITH SVM  
PLTXLABELRGAMMA  
PLTYLABELSCORE  
PLTYLIM00 11  
LW 2  
PLTSEMILOGXPARAMRANGE TRAINSCORESMEAN LABELTRAINING SCORE  
COLORDARKORANGE LWLW  
PLTFILLBETWEENPARAMRANGE TRAINSCORESMEAN TRAINSCORESSTD  
TRAINSCORESMEAN TRAINSCORESSTD ALPHA02  
COLORDARKORANGE LWLW  
PLTSEMILOGXPARAMRANGE TESTSCORESMEAN LABELCROSSVALIDATION SCORE  
COLORNAVY LWLW  
PLTFILLBETWEENPARAMRANGE TESTSCORESMEAN TESTSCORESSTD  
TESTSCORESMEAN TESTSCORESSTD ALPHA02  
COLORNAVY LWLW  
PLTLEGENDLOCBEST  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 13416 SECONDS  
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE  
5213 UNDERFITTING VS OVERFITTING  
THIS EXAMPLE DEMONSTRATES THE PROBLEMS OF UNDERFITTING AND OVERFITTING AND HOW WE CAN USE LINEAR REGRESSION WITH  
POLYNOMIAL FEATURES TO APPROXIMATE NONLINEAR FUNCTIONS THE PLOT SHOWS THE FUNCTION THAT WE WANT TO APPROXIMATE  
WHICH IS A PART OF THE COSINE FUNCTION IN ADDITION THE SAMPLES FROM THE REAL FUNCTION AND THE APPROXIMATIONS OF  
DIFFERENT MODELS ARE DISPLAYED THE MODELS HAVE POLYNOMIAL FEATURES OF DIFFERENT DEGREES WE CAN SEE THAT A LINEAR  
FUNCTION POLYNOMIAL WITH DEGREE 1 IS NOT SUFFICIENT TO FIT THE TRAINING SAMPLES THIS IS CALLED UNDERFITTING A  
POLYNOMIAL OF DEGREE 4 APPROXIMATES THE TRUE FUNCTION ALMOST PERFECTLY HOWEVER FOR HIGHER DEGREES THE MODEL  
WILL OVERFIT THE TRAINING DATA IE IT LEARNS THE NOISE OF THE TRAINING DATA WE EVALUATE QUANTITATIVELY OVERFITTING  
UNDERFITTING BY USING CROSSVALIDATION WE CALCULATE THE MEAN SQUARED ERROR MSE ON THE VALIDATION SET THE HIGHER  
THE LESS LIKELY THE MODEL GENERALIZES CORRECTLY FROM THE TRAINING DATA  
1264 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
PRINTDOC
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARNPIPELINE IMPORT PIPELINE
FROM SKLEARNPREPROCESSING IMPORT POLYNOMIALFEATURES
FROM SKLEARNLINEARMODEL IMPORT LINEARREGRESSION
FROM SKLEARNMODELSELECTION IMPORT CROSSVALSCORE
DEFTRUEFUNX
RETURNNPCOS15 NPPIX
NPRANDOMSEED0
NSAMPLES 30
DEGREES 1 4 15
X NPSORTNPRANDOMRANDNSAMPLES
Y TRUEFUNX NPRANDOMRANDNNSAMPLES 01
PLTFIGUREFIGSIZE14 5
FORIIRANGELENDREES
AX PLTSUBPLOT1 LENDEGREES I 1
PLTSETPAX XTICKS YTICKS
POLYNOMIALFEATURES POLYNOMIALFEATURESDEGREEDEGREES
INCLUDEBIASFALSE
LINEARREGRESSION LINEARREGRESSION
PIPELINE PIPELINEPOLYNOMIALFEATURES POLYNOMIALFEATURES
LINEARREGRESSION LINEARREGRESSION
PIPELINEFITX NPNEWAXIS Y
EVALUATE THE MODELS USING CROSSVALIDATION
SCORES CROSSVALSCOREPIPELINE X NPNEWAXIS Y
SCORINGNEGMEANSQUAREDERROR CV10
XTEST NPLINSPACE0 1 100
PLTPLOTXTEST PIPELINEPREDICTXTEST NPNEWAXIS LABELMODEL
PLTPLOTXTEST TRUEFUNXTEST LABELTRUE FUNCTION
PLTSCATTERX Y EDGECOLORB S20 LABELSAMPLES
PLTXLABELX
PLTYLABELY
PLTXLIM0 1
PLTYLIM2 2
PLTLEGENDLOCBEST
PLTTITLEDEGREE NMSE 2E 2EFORMAT
DEGREES SCORESMEAN SCORESSTD
PLTSHOW
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0084 SECONDS
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE
521 MODEL SELECTION 1265
```

SCIKITLEARN USER GUIDE RELEASE 0213

5214 PARAMETER ESTIMATION USING GRID SEARCH WITH CROSSVALIDATION

THIS EXAMPLES SHOWS HOW A CLASSIFIER IS OPTIMIZED BY CROSSVALIDATION WHICH IS DONE USING THE SKLEARN  
MODELSELECTIONGRIDSEARCHCV OBJECT ON A DEVELOPMENT SET THAT COMPRISES ONLY HALF OF THE AVAILABLE LABELED  
DATA

THE PERFORMANCE OF THE SELECTED HYPERPARAMETERS AND TRAINED MODEL IS THEN MEASURED ON A DEDICATED EVALUATION SET  
THAT WAS NOT USED DURING THE MODEL SELECTION STEP

MORE DETAILS ON TOOLS AVAILABLE FOR MODEL SELECTION CAN BE FOUND IN THE SECTIONS ON CROSSVALIDATION EVALUATING  
ESTIMATOR PERFORMANCE ANDTUNING THE HYPERPARAMETERS OF AN ESTIMATOR

OUT

TUNING HYPERPARAMETERS FOR PRECISION

BEST PARAMETERS SET FOUND ON DEVELOPMENT SET

C 10 GAMMA 0001 KERNEL RBF

GRID SCORES ON DEVELOPMENT SET

0986 0016 FOR C 1 GAMMA 0001 KERNEL RBF

0959 0029 FOR C 1 GAMMA 00001 KERNEL RBF

0988 0017 FOR C 10 GAMMA 0001 KERNEL RBF

0982 0026 FOR C 10 GAMMA 00001 KERNEL RBF

0988 0017 FOR C 100 GAMMA 0001 KERNEL RBF

0982 0025 FOR C 100 GAMMA 00001 KERNEL RBF

0988 0017 FOR C 1000 GAMMA 0001 KERNEL RBF

0982 0025 FOR C 1000 GAMMA 00001 KERNEL RBF

0975 0014 FOR C 1 KERNEL LINEAR

0975 0014 FOR C 10 KERNEL LINEAR

0975 0014 FOR C 100 KERNEL LINEAR

0975 0014 FOR C 1000 KERNEL LINEAR

DETAILED CLASSIFICATION REPORT

THE MODEL IS TRAINED ON THE FULL DEVELOPMENT SET

THE SCORES ARE COMPUTED ON THE FULL EVALUATION SET

PRECISION RECALL F1SCORE SUPPORT

0 100 100 100 89

1 097 100 098 90

2 099 098 098 92

3 100 099 099 93

4 100 100 100 76

5 099 098 099 108

6 099 100 099 89

7 099 100 099 78

8 100 098 099 92

9 099 099 099 92

ACCURACY 099 899

MACRO AVG 099 099 099 899

WEIGHTED AVG 099 099 099 899

TUNING HYPERPARAMETERS FOR RECALL

1266 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213  
BEST PARAMETERS SET FOUND ON DEVELOPMENT SET  
C 10 GAMMA 0001 KERNEL RBF  
GRID SCORES ON DEVELOPMENT SET  
0986 0019 FOR C 1 GAMMA 0001 KERNEL RBF  
0957 0029 FOR C 1 GAMMA 00001 KERNEL RBF  
0987 0019 FOR C 10 GAMMA 0001 KERNEL RBF  
0981 0028 FOR C 10 GAMMA 00001 KERNEL RBF  
0987 0019 FOR C 100 GAMMA 0001 KERNEL RBF  
0981 0026 FOR C 100 GAMMA 00001 KERNEL RBF  
0987 0019 FOR C 1000 GAMMA 0001 KERNEL RBF  
0981 0026 FOR C 1000 GAMMA 00001 KERNEL RBF  
0972 0012 FOR C 1 KERNEL LINEAR  
0972 0012 FOR C 10 KERNEL LINEAR  
0972 0012 FOR C 100 KERNEL LINEAR  
0972 0012 FOR C 1000 KERNEL LINEAR  
DETAILED CLASSIFICATION REPORT  
THE MODEL IS TRAINED ON THE FULL DEVELOPMENT SET  
THE SCORES ARE COMPUTED ON THE FULL EVALUATION SET  
PRECISION RECALL F1SCORE SUPPORT  
0 100 100 100 89  
1 097 100 098 90  
2 099 098 098 92  
3 100 099 099 93  
4 100 100 100 76  
5 099 098 099 108  
6 099 100 099 89  
7 099 100 099 78  
8 100 098 099 92  
9 099 099 099 92  
ACCURACY 099 899  
MACRO AVG 099 099 099 899  
WEIGHTED AVG 099 099 099 899  
FROM SKLEARN IMPORT DATASETS  
FROM SKLEARNMODELSELECTION IMPORT TRAJNTESTSPLIT  
FROM SKLEARNMODELSELECTION IMPORT GRIDSEARCHCV  
FROM SKLEARNMETRICS IMPORT CLASSIFICATIONREPORT  
FROM SKLEARN SVM IMPORT SVC  
PRINTDOC  
LOADING THE DIGITS DATASET  
DIGITS DATASETSLOADDIGITS  
521 MODEL SELECTION 1267

```
SCIKITLEARN USER GUIDE RELEASE 0213
TO APPLY AN CLASSIFIER ON THIS DATA WE NEED TO FLATTEN THE IMAGE TO
TURN THE DATA IN A SAMPLES FEATURE MATRIX
NSAMPLES LENDIGITSIMAGES
X DIGITSIMAGESRESHAPENSAMPLES 1
Y DIGITSTARGET
SPLIT THE DATASET IN TWO EQUAL PARTS
XTRAIN XTEST YTRAIN YTEST TRAINTESTSPLIT
X Y TESTSIZE05 RANDOMSTATE0
SET THE PARAMETERS BY CROSSVALIDATION
TUNEDPARAMETERS KERNEL RBF GAMMA 1E3 1E4
C 1 10 100 1000
KERNEL LINEAR C 1 10 100 1000
SCORES PRECISION RECALL
FORSCOREINSCORES
PRINT TUNING HYPERPARAMETERS FOR S SCORE
PRINT
CLF GRIDSEARCHCVSVC TUNEDPARAMETERS CV5
SCORING SMACRO SCORE
CLFFITXTRAIN YTRAIN
PRINTBEST PARAMETERS SET FOUND ON DEVELOPMENT SET
PRINT
PRINTCLFBESTPARAMS
PRINT
PRINTGRID SCORES ON DEVELOPMENT SET
PRINT
MEANS CLFCVRESULTSMEANTESTSCORE
STDS CLFCVRESULTSSTDTESTSCORE
FORMEAN STD PARAMS INZIPMEANS STDS CLFCVRESULTSPARAMS
PRINT03F003F FORR
MEAN STD 2 PARAMS
PRINT
PRINTDETAILED CLASSIFICATION REPORT
PRINT
PRINTTHE MODEL IS TRAINED ON THE FULL DEVELOPMENT SET
PRINTTHE SCORES ARE COMPUTED ON THE FULL EVALUATION SET
PRINT
YTRUE YPRED YTEST CLFPREDICTXTEST
PRINTCLASSIFICATIONREPORTYTRUE YPRED
PRINT
NOTE THE PROBLEM IS TOO EASY THE HYPERPARAMETER PLATEAU IS TOO FLAT AND THE
OUTPUT MODEL IS THE SAME FOR PRECISION AND RECALL WITH TIES IN QUALITY
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 4344 SECONDS
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE
1268 CHAPTER 5 EXAMPLES
```

SCIKITLEARN USER GUIDE RELEASE 0213

5215 TRAIN ERROR VS TEST ERROR

ILLUSTRATION OF HOW THE PERFORMANCE OF AN ESTIMATOR ON UNSEEN DATA TEST DATA IS NOT THE SAME AS THE PERFORMANCE ON TRAINING DATA AS THE REGULARIZATION INCREASES THE PERFORMANCE ON TRAIN DECREASES WHILE THE PERFORMANCE ON TEST IS OPTIMAL WITHIN A RANGE OF VALUES OF THE REGULARIZATION PARAMETER THE EXAMPLE WITH AN ELASTICNET REGRESSION MODEL AND THE PERFORMANCE IS MEASURED USING THE EXPLAINED VARIANCE AKA R2

OUT

OPTIMAL REGULARIZATION PARAMETER 000013141473626117567

PRINTDOC

AUTHOR ALEXANDRE GRAMFORT ALEXANDREGRAMFORTINRIA.FR

LICENSE BSD 3 CLAUSE

IMPORT NUMPY AS NP

FROM SKLEARN IMPORT LINEARMODEL

GENERATE SAMPLE DATA

521 MODEL SELECTION 1269

```
SCIKITLEARN USER GUIDE RELEASE 0213
NSAMPLESTRAIN NSAMPLESTEST NFEATURES 75 150 500
NPRANDOMSEED0
COEF NPRANDOMRANDNNFEATURES
COEF50 00 ONLY THE TOP 10 FEATURES ARE IMPACTING THE MODEL
X NPRANDOMRANDNNSAMPLESTRAIN NSAMPLESTEST NFEATURES
Y NPDOTX COEF
  SPLIT TRAIN AND TEST DATA
XTRAIN XTEST XNSAMPLESTRAIN XNSAMPLESTRAIN
YTRAIN YTEST YNSAMPLESTRAIN YNSAMPLESTRAIN

  COMPUTE TRAIN AND TEST ERRORS
ALPHAS NPLOGSPACE5 1 60
ENET LINEARMODELELASTICNETL1RATIO07 MAXITER10000
TRAINERRORS LIST
TESTERRORS LIST
FORALPHAINALPHAS
ENETSETPARAMSALPHAALPHA
ENETFITXTRAIN YTRAIN
TRAINERRORSAPPENDENETSCOREXTRAIN YTRAIN
TESTERRORSAPPENDENETSCOREXTEST YTEST
IALPHAOPTIM NPARGMAXTESTERRORS
ALPHAOPTIM ALPHASIALPHAOPTIM
PRINTOPTIMAL REGULARIZATION PARAMETER S ALPHAOPTIM
  ESTIMATE THE COEF ON FULL DATA WITH OPTIMAL REGULARIZATION PARAMETER
ENETSETPARAMSALPHAALPHAOPTIM
COEF ENETFITX YCOEF

  PLOT RESULTS FUNCTIONS
IMPORT MATPLOTLIBPYPLOT AS PLT
PLTSUBPLOT2 1 1
PLTSEMILOGXALPHAS TRAINERRORS LABELTRAIN
PLTSEMILOGXALPHAS TESTERRORS LABELTEST
PLTVLINESALPHAOPTIM PLTYLIM0 NPMAXTESTERRORS COLORK
LINEWIDTH3 LABELOPTIMUM ON TEST
PLTLEGENDLOCLOWER LEFT
PLTYLIM0 12
PLTXLABELREGULARIZATION PARAMETER
PLTYLABELPERFORMANCE
  SHOW ESTIMATED COEF VS TRUE COEF
PLTSUBPLOT2 1 2
PLTPLOTCOEF LABELTRUE COEF
PLTPLOTCOEF LABELESTIMATED COEF
PLTLEGEND
PLTSUBPLOTSADJUST009 004 094 094 026 026
PLTSHOW
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 3507 SECONDS
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE
1270 CHAPTER 5 EXAMPLES
```

SCIKITLEARN USER GUIDE RELEASE 0213

5216 RECEIVER OPERATING CHARACTERISTIC ROC WITH CROSS VALIDATION

EXAMPLE OF RECEIVER OPERATING CHARACTERISTIC ROC METRIC TO EVALUATE CLASSIFIER OUTPUT QUALITY USING CROSSVALIDATION

ROC CURVES TYPICALLY FEATURE TRUE POSITIVE RATE ON THE Y AXIS AND FALSE POSITIVE RATE ON THE X AXIS THIS MEANS THAT THE TOP LEFT CORNER OF THE PLOT IS THE “IDEAL” POINT A FALSE POSITIVE RATE OF ZERO AND A TRUE POSITIVE RATE OF ONE THIS IS NOT VERY REALISTIC BUT IT DOES MEAN THAT A LARGER AREA UNDER THE CURVE AUC IS USUALLY BETTER

THE “STEEPNESS” OF ROC CURVES IS ALSO IMPORTANT SINCE IT IS IDEAL TO MAXIMIZE THE TRUE POSITIVE RATE WHILE MINIMIZING THE FALSE POSITIVE RATE

THIS EXAMPLE SHOWS THE ROC RESPONSE OF DIFFERENT DATASETS CREATED FROM KFOLD CROSSVALIDATION TAKING ALL OF THESE CURVES IT IS POSSIBLE TO CALCULATE THE MEAN AREA UNDER CURVE AND SEE THE VARIANCE OF THE CURVE WHEN THE TRAINING SET IS SPLIT INTO DIFFERENT SUBSETS THIS ROUGHLY SHOWS HOW THE CLASSIFIER OUTPUT IS AFFECTED BY CHANGES IN THE TRAINING DATA AND HOW DIFFERENT THE SPLITS GENERATED BY KFOLD CROSSVALIDATION ARE FROM ONE ANOTHER

NOTE

SEE ALSO SKLEARNMETRICSROCAUCSCORE SKLEARNMODELSELECTIONCROSSVALSCORE

RECEIVER OPERATING CHARACTERISTIC ROC

PRINTDOC

IMPORT NUMPY AS NP

FROM SCIPY IMPORT INTERP

521 MODEL SELECTION 1271

```
SCIKITLEARN USER GUIDE RELEASE 0213
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARN IMPORT SVM DATASETS
FROM SKLEARNMETRICS IMPORT ROCCURVE AUC
FROM SKLEARNMODELSELECTION IMPORT STRATIFIEDKFOLD
```

```
DATA IO AND GENERATION
IMPORT SOME DATA TO PLAY WITH
IRIS DATASETSLOADIRIS
X IRISDATA
Y IRISTARGET
X Y XY 2 YY 2
NSAMPLES NFEATURES XSHAPE
ADD NOISY FEATURES
RANDOMSTATE NPRANDOMRANDOMSTATE0
X NPCX RANDOMSTATETERANDNNSAMPLES 200 NFEATURES
```

```
CLASSIFICATION AND ROC ANALYSIS
RUN CLASSIFIER WITH CROSSVALIDATION AND PLOT ROC CURVES
CV STRATIFIEDKFOLDNSPLITS6
CLASSIFIER SVMKERNELLINEAR PROBABILITYTRUE
RANDOMSTATETERANDOMSTATE
TPRS
AUCS
MEANFPR NPLINSPACE0 1 100
I 0
FORTRAIN TEST INCVSPLITX Y
PROBAS CLASSIFIERFITXTRAIN YTRAINPREDICTPROBAXTEST
COMPUTE ROC CURVE AND AREA THE CURVE
FPR TPR THRESHOLDS ROCCURVEYTEST PROBAS 1
TPRSAPPENDINTERPMEANFPR FPR TPR
TPRS10 00
ROCAUC AUCFPR TPR
AUCSAPPENDROCAUC
PLTPLOTFPR TPR LW1 ALPHA03
LABELROC FOLD DAUC 02F I ROCAUC
I 1
PLTPLOT0 1 0 1 LINESTYLE LW2 COLORR
LABELCHANCE ALPHA8
MEANTPR NPMEANTPRS AXIS0
MEANTPR1 10
MEANAUC AUCMEANFPR MEANTPR
STDAUC NPSTDAUCS
PLTPLOTMEANFPR MEANTPR COLORB
LABELRMEAN ROC AUC 02FPM02F MEANAUC STDAUC
LW2 ALPHA8
STDTPR NPSTDTPRS AXIS0
TPRSUPPER NPMINIMUMMEANTPR STDTPR 1
1272 CHAPTER 5 EXAMPLES
```

SCIKITLEARN USER GUIDE RELEASE 0213  
TPRSLOWER NPMAXIMUMMEANTPR STDTPR 0  
PLTFILLBETWEENMEANFPR TPRSLOWER TPRSUPPER COLORGREY ALPHA2  
LABELRPM 1 STD DEV  
PLTXLIM005 105  
PLTYLIM005 105  
PLTXLABELFALSE POSITIVE RATE  
PLTYLABELTRUE POSITIVE RATE  
PLTTITLERECEIVER OPERATING CHARACTERISTIC EXAMPLE  
PLTLEGENDLOCLOWER RIGHT  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0245 SECONDS  
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE  
5217 COMPARING RANDOMIZED SEARCH AND GRID SEARCH FOR HYPERPARAMETER ESTIMATION  
COMPARE RANDOMIZED SEARCH AND GRID SEARCH FOR OPTIMIZING HYPERPARAMETERS OF A RANDOM FOREST ALL PARAMETERS THAT INFLUENCE THE LEARNING ARE SEARCHED SIMULTANEOUSLY EXCEPT FOR THE NUMBER OF ESTIMATORS WHICH POSES A TIME QUALITY TRADEOFF  
THE RANDOMIZED SEARCH AND THE GRID SEARCH EXPLORE EXACTLY THE SAME SPACE OF PARAMETERS THE RESULT IN PARAMETER SETTINGS IS QUITE SIMILAR WHILE THE RUN TIME FOR RANDOMIZED SEARCH IS DRASTICALLY LOWER  
THE PERFORMANCE IS SLIGHTLY WORSE FOR THE RANDOMIZED SEARCH THOUGH THIS IS MOST LIKELY A NOISE EFFECT AND WOULD NOT CARRY OVER TO A HELDOUT TEST SET  
NOTE THAT IN PRACTICE ONE WOULD NOT SEARCH OVER THIS MANY DIFFERENT PARAMETERS SIMULTANEOUSLY USING GRID SEARCH BUT PICK ONLY THE ONES DEEMED MOST IMPORTANT  
OUT  
RANDOMIZEDSEARCHCV TOOK 442 SECONDS FOR 20 CANDIDATES PARAMETER SETTINGS  
MODEL WITH RANK 1  
MEAN VALIDATION SCORE 0939 STD 0024  
PARAMETERS BOOTSTRAP FALSE CRITERION ENTROPY MAXDEPTH NONE MAX  
↪FEATURES 7 MINSAMPLESSPLIT 3  
MODEL WITH RANK 2  
MEAN VALIDATION SCORE 0933 STD 0022  
PARAMETERS BOOTSTRAP FALSE CRITERION GINI MAXDEPTH NONE MAXFEATURES  
↪ 6 MINSAMPLESSPLIT 6  
MODEL WITH RANK 3  
MEAN VALIDATION SCORE 0930 STD 0031  
PARAMETERS BOOTSTRAP TRUE CRITERION GINI MAXDEPTH NONE MAXFEATURES  
↪ 6 MINSAMPLESSPLIT 6  
GRIDSEARCHCV TOOK 1324 SECONDS FOR 72 CANDIDATE PARAMETER SETTINGS  
MODEL WITH RANK 1  
MEAN VALIDATION SCORE 0937 STD 0019  
PARAMETERS BOOTSTRAP FALSE CRITERION ENTROPY MAXDEPTH NONE MAX  
↪FEATURES 10 MINSAMPLESSPLIT 2  
521 MODEL SELECTION 1273

```
SCIKITLEARN USER GUIDE RELEASE 0213
MODEL WITH RANK 2
MEAN VALIDATION SCORE 0936 STD 0020
PARAMETERS BOOTSTRAP FALSE CRITERION GINI MAXDEPTH NONE MAXFEATURES
↩→ 10 MINSAMPLESSPLIT 2
MODEL WITH RANK 3
MEAN VALIDATION SCORE 0931 STD 0029
PARAMETERS BOOTSTRAP FALSE CRITERION ENTROPY MAXDEPTH NONE MAX
↩→FEATURES 10 MINSAMPLESSPLIT 3
PRINTDOC
IMPORT NUMPY AS NP
FROM TIME IMPORT TIME
FROM SCIPYSTATS IMPORT RANDINT ASSPRANDINT
FROM SKLEARNMODELSELECTION IMPORT GRIDSEARCHCV
FROM SKLEARNMODELSELECTION IMPORT RANDOMIZEDSEARCHCV
FROM SKLEARNDATASETS IMPORT LOADDIGITS
FROM SKLEARNENSEMBLE IMPORT RANDOMFORESTCLASSIFIER
GET SOME DATA
DIGITS LOADDIGITS
X Y DIGITSDATA DIGITSTARGET
BUILD A CLASSIFIER
CLF RANDOMFORESTCLASSIFIERNESTIMATORS20
UTILITY FUNCTION TO REPORT BEST SCORES
DEFREPORTRESULTS NTOP3
FORIINRANGE1 NTOP 1
CANDIDATES NPFLATNONZERORESULTSRANKTESTSCORE I
FORCANDIDATE INCANDIDATES
PRINTMODEL WITH RANK 0FORMATI
PRINTMEAN VALIDATION SCORE 03F STD 13FFORMAT
RESULTSMEANTESTSCORECANDIDATE
RESULTSSTDTESTSCORECANDIDATE
PRINTPARAMETERS 0FORMATRESULTSPARAMSCANDIDATE
PRINT
SPECIFY PARAMETERS AND DISTRIBUTIONS TO SAMPLE FROM
PARAMDIST MAXDEPTH 3 NONE
MAXFEATURES SPRANDINT1 11
MINSAMPLESSPLIT SPRANDINT2 11
BOOTSTRAP TRUE FALSE
CRITERION GINI ENTROPY
RUN RANDOMIZED SEARCH
NITERSEARCH 20
RANDOMSEARCH RANDOMIZEDSEARCHCVCLF PARAMDISTRIBUTIONSPARAMDIST
1274 CHAPTER 5 EXAMPLES
```



SCIKITLEARN USER GUIDE RELEASE 0213

NITERNITERSEARCH CV5 IIDFALSE

START TIME

RANDOMSEARCHFITX Y

PRINTRANDOMIZEDSEARCHCV TOOK 2FSECONDS FOR DCANDIDATES

PARAMETER SETTINGS TIME START NITERSEARCH

REPORTRANDOMSEARCHCVRESULTS

USE A FULL GRID OVER ALL PARAMETERS

PARAMGRID MAXDEPTH 3 NONE

MAXFEATURES 1 3 10

MINSAMPLESSPLIT 2 3 10

BOOTSTRAP TRUE FALSE

CRITERION GINI ENTROPY

RUN GRID SEARCH

GRIDSEARCH GRIDSEARCHCVCLF PARAMGRIDPARAMGRID CV5 IIDFALSE

START TIME

GRIDSEARCHFITX Y

PRINTGRIDSEARCHCV TOOK 2FSECONDS FOR DCANDIDATE PARAMETER SETTINGS

TIME START LENGRIDSEARCHCVRESULTSPARAMS

REPORTGRIDSEARCHCVRESULTS

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 17712 SECONDS

NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE

5218 NESTED VERSUS NONNESTED CROSSVALIDATION

THIS EXAMPLE COMPARES NONNESTED AND NESTED CROSSVALIDATION STRATEGIES ON A CLASSIFIER OF THE IRIS DATA SET NESTED CROSSVALIDATION CV IS OFTEN USED TO TRAIN A MODEL IN WHICH HYPERPARAMETERS ALSO NEED TO BE OPTIMIZED NESTED CV ESTIMATES THE GENERALIZATION ERROR OF THE UNDERLYING MODEL AND ITS HYPERPARAMETER SEARCH CHOOSING THE PARAMETERS THAT MAXIMIZE NONNESTED CV BIASES THE MODEL TO THE DATASET YIELDING AN OVERLYOPTIMISTIC SCORE

MODEL SELECTION WITHOUT NESTED CV USES THE SAME DATA TO TUNE MODEL PARAMETERS AND EVALUATE MODEL PERFORMANCE INFORMATION MAY THUS “LEAK” INTO THE MODEL AND OVERFIT THE DATA THE MAGNITUDE OF THIS EFFECT IS PRIMARILY DEPENDENT ON THE SIZE OF THE DATASET AND THE STABILITY OF THE MODEL SEE CAWLEY AND TALBOT1FOR AN ANALYSIS OF THESE ISSUES

TO AVOID THIS PROBLEM NESTED CV EFFECTIVELY USES A SERIES OF TRAINVALIDATIONTEST SET SPLITS IN THE INNER LOOP

HERE EXECUTED BY GRIDSEARCHCV THE SCORE IS APPROXIMATELY MAXIMIZED BY FITTING A MODEL TO EACH TRAINING

SET AND THEN DIRECTLY MAXIMIZED IN SELECTING HYPERPARAMETERS OVER THE VALIDATION SET IN THE OUTER LOOP HERE IN

CROSSVALSCORE GENERALIZATION ERROR IS ESTIMATED BY AVERAGING TEST SET SCORES OVER SEVERAL DATASET SPLITS

THE EXAMPLE BELOW USES A SUPPORT VECTOR CLASSIFIER WITH A NONLINEAR KERNEL TO BUILD A MODEL WITH OPTIMIZED HYPERPARAMETERS BY GRID SEARCH WE COMPARE THE PERFORMANCE OF NONNESTED AND NESTED CV STRATEGIES BY TAKING THE DIFFERENCE BETWEEN THEIR SCORES

SEE ALSO

•CROSSVALIDATION EVALUATING ESTIMATOR PERFORMANCE

1CAWLEY GC TALBOT NLC ON OVERFITTING IN MODEL SELECTION AND SUBSEQUENT SELECTION BIAS IN PERFORMANCE EVALUATION J MACH LEARN 201011 20792107

521 MODEL SELECTION 1275

SCIKITLEARN USER GUIDE RELEASE 0213

- TUNING THE HYPERPARAMETERS OF AN ESTIMATOR

REFERENCES

OUT

AVERAGE DIFFERENCE OF 0007581 WITH STD DEV OF 0007833

FROM SKLEARNDATASETS IMPORT LOADIRIS

FROM MATPLOTLIB IMPORT PYPLLOTASPLT

FROM SKLEARN SVM IMPORT SVC

FROM SKLEARNMODELSELECTION IMPORT GRIDSEARCHCV CROSSVALSCORE KFOLD

IMPORT NUMPY AS NP

PRINTDOC

NUMBER OF RANDOM TRIALS

NUMTRIALS 30

1276 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

LOAD THE DATASET

IRIS LOADIRIS

XIRIS IRISDATA

YIRIS IRISTARGET

SET UP POSSIBLE VALUES OF PARAMETERS TO OPTIMIZE OVER

PGRID C 1 10 100

GAMMA 01 1

WE WILL USE A SUPPORT VECTOR CLASSIFIER WITH RBF KERNEL

SVM SVCKERNELRBF

ARRAYS TO STORE SCORES

NONNESTEDSCORES NPZEROSNUMTRIALS

NESTEDSCORES NPZEROSNUMTRIALS

LOOP FOR EACH TRIAL

FORIINRANGENUMTRIALS

CHOOSE CROSSVALIDATION TECHNIQUES FOR THE INNER AND OUTER LOOPS

INDEPENDENTLY OF THE DATASET

EG GROUPKFOLD LEAVEONEOUT LEAVEONEGROUPOUT ETC

INNERCV KFOLDNSPLITS4 SHUFFLETRUE RANDOMSTATEI

OUTERCV KFOLDNSPLITS4 SHUFFLETRUE RANDOMSTATEI

NONNESTED PARAMETER SEARCH AND SCORING

CLF GRIDSEARCHCVESTIMATORSVM PARAMGRIDPGRID CVINNERCV

IIDFALSE

CLFFITXIRIS YIRIS

NONNESTEDSCORESI CLFBESTSCORE

NESTED CV WITH PARAMETER OPTIMIZATION

NESTEDSCORE CROSSVALSCORECLF XXIRIS YYIRIS CVOUTERCV

NESTEDSCORESI NESTEDSCOREMEAN

SCOREDIFFERENCE NONNESTEDSCORES NESTEDSCORES

PRINTAVERAGE DIFFERENCE OF 6F WITH STD DEV OF 6F

FORMATSCOREDIFFERENCEMEAN SCOREDIFFERENCESTD

PLOT SCORES ON EACH TRIAL FOR NESTED AND NONNESTED CV

PLTFigure

PLTSUBPLOT211

NONNESTEDSCORESLINE PLTPLOTNONNESTEDSCORES COLORR

NESTEDLINE PLTPLOTNESTEDSCORES COLORB

PLTYLABELSCORE FONTSIZE14

PLTLEGENDNONNESTEDSCORESLINE NESTEDLINE

NONNESTED CV NESTED CV

BBOXTOANCHOR0 4 5 0

PLTTITLENONNESTED AND NESTED CROSS VALIDATION ON IRIS DATASET

X5 Y11 FONTSIZE15

PLOT BAR CHART OF THE DIFFERENCE

PLTSUBPLOT212

DIFFERENCEPLOT PLTBARRANGENUMTRIALS SCOREDIFFERENCE

PLTXLABELINDIVIDUAL TRIAL

PLTLEGENDDIFFERENCEPLOT

521 MODEL SELECTION 1277

SCIKITLEARN USER GUIDE RELEASE 0213  
NONNESTED CV NESTED CV SCORE  
BBOXTOANCHOR0 1 8 0  
PLTYLABELSCORE DIFFERENCE FONTSIZE14  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 3447 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
5219 DEMONSTRATION OF MULTIMETRIC EVALUATION ON CROSSVALSCORE AND GRID  
SEARCHCV  
MULTIPLE METRIC PARAMETER SEARCH CAN BE DONE BY SETTING THE SCORING PARAMETER TO A LIST OF METRIC SCORER NAMES OR A  
DICT MAPPING THE SCORER NAMES TO THE SCORER CALLABLES  
THE SCORES OF ALL THE SCORERS ARE AVAILABLE IN THE CVRESULTS DICT AT KEYS ENDING IN SCORERNAME  
MEANTESTPRECISION RANKTESTPRECISION ETC  
THEBESTESTIMATOR BESTINDEX BESTSCORE ANDBESTPARAMS CORRESPOND TO THE SCORER KEY  
THAT IS SET TO THE REFIT ATTRIBUTE  
AUTHOR RAGHAV RV RVRAGHAV93GMAILCOM  
LICENSE BSD  
IMPORT NUMPY AS NP  
FROM MATPLOTLIB IMPORT PYPLOTASPLT  
FROM SKLEARNDATASETS IMPORT MAKEHASTIE102  
FROM SKLEARNMODELSELECTION IMPORT GRIDSEARCHCV  
FROM SKLEARNMETRICS IMPORT MAKESCORER  
FROM SKLEARNMETRICS IMPORT ACCURACYSORE  
FROM SKLEARNTREE IMPORT DECISIONTREECLASSIFIER  
PRINTDOC  
RUNNINGGRIDSEARCHCV USING MULTIPLE EVALUATION METRICS  
X Y MAKEHASTIE102NSAMPLES8000 RANDOMSTATE42  
THE SCORERS CAN BE EITHER BE ONE OF THE PREDEFINED METRIC STRINGS OR A SCORER  
CALLABLE LIKE THE ONE RETURNED BY MAKESCORER  
SCORING AUC ROCAUC ACCURACY MAKESCORERACCURACYSORE  
SETTING REFITAUC REFITS AN ESTIMATOR ON THE WHOLE DATASET WITH THE  
PARAMETER SETTING THAT HAS THE BEST CROSSVALIDATED AUC SCORE  
THAT ESTIMATOR IS MADE AVAILABLE AT GSBESTESTIMATOR ALONG WITH  
PARAMETERS LIKE GSBESTSCORE GSBESTPARAMS AND  
GSBESTINDEX  
GS GRIDSEARCHCVDECISIONTREECLASSIFIERRANDOMSTATE42  
PARAMGRIDMINSAMPLESSPLIT RANGE2 403 10  
SCORINGSCORING CV5 REFITAUC RETURNTRAINSCORETRUE  
1278 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
GSFITX Y  
RESULTS GSCVRESULTS  
PLOT THE RESULT  
PLTFIGUREFIGSIZE13 13  
PLTTITLEGRIDSEARCHCV EVALUATING USING MULTIPLE SCORERS SIMULTANEOUSLY  
FONTSIZE16  
PLTXLABELMINSAMPLESSPLIT  
PLTYLABELSCORE  
AX PLTGCA  
AXSETXLIM0 402  
AXSETYLIM073 1  
GET THE REGULAR NUMPY ARRAY FROM THE MASKEDARRAY  
XAXIS NPARRAYRESULTSPARAMMINSAMPLESSPLITDATA DTYPEFLOAT  
FORSCORER COLOR INZIPSORTEDSCORING G K  
FORSAMPLE STYLE INTRAIN TEST  
SAMPLESCOREMEAN RESULTSMEAN SS SAMPLE SCORER  
SAMPLESCORESTD RESULTSTD SS SAMPLE SCORER  
AXFILLBETWEENXAXIS SAMPLESCOREMEAN SAMPLESCORESTD  
SAMPLESCOREMEAN SAMPLESCORESTD  
ALPHA01 IFSAMPLE TEST ELSE0 COLORCOLOR  
AXPLOTXAXIS SAMPLESCOREMEAN STYLE COLORCOLOR  
ALPHA1 IFSAMPLE TEST ELSE07  
LABELSS SCORER SAMPLE  
BESTINDEX NPNONZERORESULTSRANKTEST S SCORER 100  
BESTSCORE RESULTSMEANTEST S SCORERBESTINDEX  
PLOT A DOTTED VERTICAL LINE AT THE BEST SCORE FOR THAT SCORER MARKED BY X  
AXPLOTXAXISBESTINDEX 2 0 BESTSCORE  
LINESTYLE COLORCOLOR MARKERX MARKEREDGEWIDTH3 MS8  
ANNOTATE THE BEST SCORE FOR THAT SCORER  
AXANNOTATE 02F BESTSCORE  
XAXISBESTINDEX BESTSCORE 0005  
PLTLEGENDLOCBEST  
PLTGRIDFALSE  
PLTSHOW  
521 MODEL SELECTION 1279

SCIKITLEARN USER GUIDE RELEASE 0213

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 20458 SECONDS

NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)

52110 BALANCE MODEL COMPLEXITY AND CROSSVALIDATED SCORE

THIS EXAMPLE BALANCES MODEL COMPLEXITY AND CROSSVALIDATED SCORE BY FINDING A DECENT ACCURACY WITHIN 1 STANDARD DEVIATION OF THE BEST ACCURACY SCORE WHILE MINIMISING THE NUMBER OF PCA COMPONENTS 1

THE FIGURE SHOWS THE TRADEOFF BETWEEN CROSSVALIDATED SCORE AND THE NUMBER OF PCA COMPONENTS THE BALANCED CASE IS WHEN NCOMPONENTS6 AND ACCURACY080 WHICH FALLS INTO THE RANGE WITHIN 1 STANDARD DEVIATION OF THE BEST ACCURACY SCORE

1280 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
1 HASTIE T TIBSHIRANI R,, FRIEDMAN J 2001 MODEL ASSESSMENT AND SELECTION THE ELEMENTS OF STATISTICAL  
LEARNING PP 219260 NEW YORK NY USA SPRINGER NEW YORK INC  
OUT  
THE BESTINDEX IS 2  
THE NCOMPONENTS SELECTED IS 6  
THE CORRESPONDING ACCURACY SCORE IS 080  
AUTHOR WENHAO ZHANG WENHAOZUCLAEDU  
PRINTDOC  
IMPORT NUMPY AS NP  
IMPORT MATPLOTLIBPYPLOT AS PLT  
FROM SKLEARNDATASETS IMPORT LOADDIGITS  
FROM SKLEARNDECOMPOSITION IMPORT PCA  
FROM SKLEARNMODELSELECTION IMPORT GRIDSEARCHCV  
FROM SKLEARNPIPELINE IMPORT PIPELINE  
FROM SKLEARN SVM IMPORT LINEARSVC  
521 MODEL SELECTION 1281

SCIKITLEARN USER GUIDE RELEASE 0213  
DEFLOWERBOUND CVRESULTS

CALCULATE THE LOWER BOUND WITHIN 1 STANDARD DEVIATION  
OF THE BEST MEANTESTSCORES  
PARAMETERS

CVRESULTS DICT OF NUMPYMASKED NDARRAYS  
SEE ATTRIBUTE CVRESULTS OF GRIDSEARCHCV  
RETURNS

FLOAT  
LOWER BOUND WITHIN 1 STANDARD DEVIATION OF THE  
BEST MEANTESTSCORE

BESTSCOREIDX NPARGMAXCVRESULTSMEANTESTSCORE  
RETURN CVRESULTSMEANTESTSCOREBESTSCOREIDX  
CVRESULTSSTDTESTSCOREBESTSCOREIDX  
DEFBESTLOWCOMPLEXITYCVRESULTS

BALANCE MODEL COMPLEXITY WITH CROSSVALIDATED SCORE  
PARAMETERS

CVRESULTS DICT OF NUMPYMASKED NDARRAYS  
SEE ATTRIBUTE CVRESULTS OF GRIDSEARCHCV  
RETURN

INT  
INDEX OF A MODEL THAT HAS THE FEWEST PCA COMPONENTS  
WHILE HAS ITS TEST SCORE WITHIN 1 STANDARD DEVIATION OF THE BEST  
MEANTESTSCORE

THRESHOLD LOWERBOUND CVRESULTS  
CANDIDATEIDX NPFLATNONZEROCVRESULTSMEANTESTSCORE THRESHOLD  
BESTIDX CANDIDATEIDXCVRRESULTSPARAMREDUCEDIMNCOMPONENTS  
CANDIDATEIDXARGMIN  
RETURNBESTIDX  
PIPE PIPELINE  
REDUCEDIM PCARANDOMSTATE42  
CLASSIFY LINEARSVCRANDOMSTATE42

PARAMGRID  
REDUCEDIMNCOMPONENTS 2 4 6 8

GRID GRIDSEARCHCVPIPE CV10 NJOBS1 PARAMGRIDPARAMGRID  
SCORINGACCURACY REFITBESTLOWCOMPLEXITY  
1282 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213  
DIGITS LOADDIGITS  
GRIDFITDIGITS DATA DIGITSTARGET  
NCOMPONENTS GRIDCVRESULTSPARAMREDUCEDIMNCOMPONENTS  
TESTSCORES GRIDCVRESULTSMEANTESTSCORE  
PLTFigure  
PLTBARNCOMPONENTS TESTSCORES WIDTH13 COLORB  
LOWER LOWERBOUNDGRIDCVRESULTS  
PLTAXHLINENPMAXTTESTSCORES LINESTYLE COLORY  
LABELBEST SCORE  
PLTAXHLINELOWER LINESTYLE COLOR5 LABELBEST SCORE 1 STD  
PLTTITLEBALANCE MODEL COMPLEXITY AND CROSSVALIDATED SCORE  
PLTXLABELNUMBER OF PCA COMPONENTS USED  
PLTYLABELDIGIT CLASSIFICATION ACCURACY  
PLTXTICKSNCOMPONENTSTOLIST  
PLTYLIM0 10  
PLTLEGENDLOCUPPER LEFT  
BESTINDEX GRIDBESTINDEX  
PRINTTHE BESTINDEX IS D BESTINDEX  
PRINTTHE NCOMPONENTS SELECTED IS D NCOMPONENTSBESTINDEX  
PRINTTHE CORRESPONDING ACCURACY SCORE IS 2F  
GRIDCVRESULTSMEANTESTSCOREBESTINDEX  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 16219 SECONDS  
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE  
52111 SAMPLE PIPELINE FOR TEXT FEATURE EXTRACTION AND EVALUATION  
THE DATASET USED IN THIS EXAMPLE IS THE 20 NEWSGROUPS DATASET WHICH WILL BE AUTOMATICALLY DOWNLOADED AND THEN CACHED  
AND REUSED FOR THE DOCUMENT CLASSIFICATION EXAMPLE  
YOU CAN ADJUST THE NUMBER OF CATEGORIES BY GIVING THEIR NAMES TO THE DATASET LOADER OR SETTING THEM TO NONE TO GET THE  
20 OF THEM  
HERE IS A SAMPLE OUTPUT OF A RUN ON A QUADCORE MACHINE  
LOADING 20 NEWSGROUPS DATASET FORCATEGORIES  
ALTATHEISM TALKRELIGIONMISC  
1427 DOCUMENTS  
2 CATEGORIES  
PERFORMING GRID SEARCH  
PIPELINE VECT TFIDF CLF  
PARAMETERS  
CLFALPHA 100000000000000001E05 99999999999999995E07  
CLFMAXITER 10 50 80  
CLFPENALTY L2 ELASTICNET  
TFIDFUSEIDF TRUEFALSE  
521 MODEL SELECTION 1283

SCIKITLEARN USER GUIDE RELEASE 0213

VECTMAXN 1 2

VECTMAXDF 05 075 10

VECTMAXFEATURES NONE 5000 10000 50000

DONEIN1737030S

BEST SCORE 0940

BEST PARAMETERS SET

CLFALPHA 99999999999999995E07

CLFMAXITER 50

CLFPENALTY ELASTICNET

TFIDFUSEIDF TRUE

VECTMAXN 2

VECTMAXDF 075

VECTMAXFEATURES 50000

AUTHOR OLIVIER GRISEL OLIVIERGRISELENSTAORG

PETER PRETTENHOFER PETERPRETTENHOFERGMAILCOM

MATHIEU BLONDEL MATHIEUMBLODELORG

LICENSE BSD 3 CLAUSE

FROM PPRINT IMPORT PPRINT

FROM TIME IMPORT TIME

IMPORT LOGGING

FROM SKLEARNDATASETS IMPORT FETCH20NEWSGROUPS

FROM SKLEARNFEATUREEXTRACTIONTEXT IMPORT COUNTVECTORIZER

FROM SKLEARNFEATUREEXTRACTIONTEXT IMPORT TFIDFTRANSFORMER

FROM SKLEARNLINEARMODEL IMPORT SGDCLASSIFIER

FROM SKLEARNMODELSELECTION IMPORT GRIDSEARCHCV

FROM SKLEARNPIPELINE IMPORT PIPELINE

PRINTDOC

DISPLAY PROGRESS LOGS ON STDOUT

LOGGINGBASICCONFIGLEVELLOGGINGINFO

FORMAT ASCTIMES LEVELNAMES MESSAGES

LOAD SOME CATEGORIES FROM THE TRAINING SET

CATEGORIES

ALTATHEISM

TALKRELIGIONMISC

UNCOMMENT THE FOLLOWING TO DO THE ANALYSIS ON ALL THE CATEGORIES

CATEGORIES NONE

PRINTLOADING 20 NEWSGROUPS DATASET FOR CATEGORIES

PRINTCATEGORIES

DATA FETCH20NEWSGROUPSSUBSETTRAIN CATEGORIESCATEGORIES

PRINTDDOCUMENTS LENDATAFILENAMES

PRINTDCATEGORIES LENDATATARGETNAMES

PRINT

DEFINE A PIPELINE COMBINING A TEXT FEATURE EXTRACTOR WITH A SIMPLE

CLASSIFIER

PIPELINE PIPELINE

1284 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
VECT COUNTVECTORIZER  
TFIDF TFIDFTRANSFORMER  
CLF SGDCLASSIFIERTOL1E3

UNCOMMENTING MORE PARAMETERS WILL GIVE BETTER EXPLORING POWER BUT WILL  
INCREASE PROCESSING TIME IN A COMBINATORIAL WAY  
PARAMETERS  
VECTMAXDF 05 075 10  
VECTMAXFEATURES NONE 5000 10000 50000  
VECTNGRAMRANGE 1 1 1 2 UNIGRAMS OR BIGRAMS  
TFIDFUSEIDF TRUE FALSE  
TFIDFNORM L1 L2  
CLFMAXITER 20  
CLFALPHA 000001 0000001  
CLFPENALTY L2 ELASTICNET  
CLFMAXITER 10 50 80

IFNAME MAIN  
MULTIPROCESSING REQUIRES THE FORK TO HAPPEN IN A MAIN PROTECTED  
BLOCK  
FIND THE BEST PARAMETERS FOR BOTH THE FEATURE EXTRACTION AND THE  
CLASSIFIER  
GRIDSEARCH GRIDSEARCHCVPipeline PARAMETERS CV5  
NJOBS1 VERBOSE1  
PRINTPERFORMING GRID SEARCH  
PRINTPIPELINE NAME FORNAME INPIPELINESTEPS  
PRINTPARAMETERS  
PPRINTPARAMETERS  
TO TIME  
GRIDSEARCHFITDATADATA DATATARGET  
PRINTDONE IN 03FS TIME TO  
PRINT  
PRINTBEST SCORE 03F GRIDSEARCHBESTSCORE  
PRINTBEST PARAMETERS SET  
BESTPARAMETERS GRIDSEARCHBESTESTIMATORGETPARAMS  
FORPARAMNAME INSORTEDPARAMETERSKEYS  
PRINTTSR PARAMNAME BESTPARAMETERSPARAMNAME  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0000 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
52112 CONFUSION MATRIX  
EXAMPLE OF CONFUSION MATRIX USAGE TO EVALUATE THE QUALITY OF THE OUTPUT OF A CLASSIFIER ON THE IRIS DATA SET THE DIAGONAL  
ELEMENTS REPRESENT THE NUMBER OF POINTS FOR WHICH THE PREDICTED LABEL IS EQUAL TO THE TRUE LABEL WHILE OFFDIAGONAL  
ELEMENTS ARE THOSE THAT ARE MISLABELED BY THE CLASSIFIER THE HIGHER THE DIAGONAL VALUES OF THE CONFUSION MATRIX THE  
BETTER INDICATING MANY CORRECT PREDICTIONS  
THE FIGURES SHOW THE CONFUSION MATRIX WITH AND WITHOUT NORMALIZATION BY CLASS SUPPORT SIZE NUMBER OF ELEMENTS IN  
521 MODEL SELECTION 1285

SCIKITLEARN USER GUIDE RELEASE 0213

EACH CLASS THIS KIND OF NORMALIZATION CAN BE INTERESTING IN CASE OF CLASS IMBALANCE TO HAVE A MORE VISUAL INTERPRETATION OF WHICH CLASS IS BEING MISCLASSIFIED

HERE THE RESULTS ARE NOT AS GOOD AS THEY COULD BE AS OUR CHOICE FOR THE REGULARIZATION PARAMETER C WAS NOT THE BEST IN REAL LIFE APPLICATIONS THIS PARAMETER IS USUALLY CHOSEN USING TUNING THE HYPERPARAMETERS OF AN ESTIMATOR

- 

1286 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

•

OUT  
CONFUSION MATRIX WITHOUT NORMALIZATION

13 0 0  
0 10 6  
0 0 9

NORMALIZED CONFUSION MATRIX

1 0 0  
0 062 038  
0 0 1

PRINTDOC  
IMPORT NUMPY AS NP  
IMPORT MATPLOTLIBPYPLOT AS PLT  
FROM SKLEARN IMPORT SVM DATASETS  
FROM SKLEARNMODELSELECTION IMPORT TRAINTESTSPLIT  
FROM SKLEARNMETRICS IMPORT CONFUSIONMATRIX  
FROM SKLEARNUTILSMULTICLASS IMPORT UNIQUELABELS  
IMPORT SOME DATA TO PLAY WITH  
521 MODEL SELECTION 1287

```
SCIKITLEARN USER GUIDE RELEASE 0213
IRIS DATASETSLOADIRIS
X IRISDATA
Y IRISTARGET
CLASSNAMES IRISTARGETNAMES
SPLIT THE DATA INTO A TRAINING SET AND A TEST SET
XTRAIN XTEST YTRAIN YTEST TRAINTESTSPLITX Y RANDOMSTATE0
RUN CLASSIFIER USING A MODEL THAT IS TOO REGULARIZED C TOO LOW TO SEE
THE IMPACT ON THE RESULTS
CLASSIFIER SVMKERNELLINEAR C001
YPRED CLASSIFIERFITXTRAIN YTRAINPREDICTXTEST
DEFPLOTCONFUSIONMATRIXYTRUE YPRED CLASSES
NORMALIZEFALSE
TITLENONE
CMAPPLTCMBLUES
```

THIS FUNCTION PRINTS AND PLOTS THE CONFUSION MATRIX  
NORMALIZATION CAN BE APPLIED BY SETTING NORMALIZETRUE

```
IF NOTTITLE
IFNORMALIZE
TITLE NORMALIZED CONFUSION MATRIX
ELSE
TITLE CONFUSION MATRIX WITHOUT NORMALIZATION
COMPUTE CONFUSION MATRIX
CM CONFUSIONMATRIXYTRUE YPRED
ONLY USE THE LABELS THAT APPEAR IN THE DATA
CLASSES CLASSESUNIQUELABELSYTRUE YPRED
IFNORMALIZE
CM CMASTYPEFLOAT CMSUMAXIS1 NPNEWAXIS
PRINTNORMALIZED CONFUSION MATRIX
ELSE
PRINTCONFUSION MATRIX WITHOUT NORMALIZATION
PRINTCM
FIG AX PLTSUBPLOTS
IM AXIMSHOWCM INTERPOLATIONNEAREST CMAPCMAP
AXFIGURECOLORBARIM AXAX
WE WANT TO SHOW ALL TICKS
AXSETXTICKSNPARANGECSHAPE1
YTICKSNPARANGECSHAPE0
AND LABEL THEM WITH THE RESPECTIVE LIST ENTRIES
XTICKLABELSCLASSES YTICKLABELSCLASSES
TITLETITLE
YLABELTRUE LABEL
XLABELPREDICTED LABEL
ROTATE THE TICK LABELS AND SET THEIR ALIGNMENT
PLTSETPAXGETXTICKLABELS ROTATION45 HARIGHT
ROTATIONMODEANCHOR
LOOP OVER DATA DIMENSIONS AND CREATE TEXT ANNOTATIONS
FMT 2F IFNORMALIZE ELSE
1288 CHAPTER 5 EXAMPLES
```

```
SCIKITLEARN USER GUIDE RELEASE 0213
THRESH CMMAX 2
FORIINRANGECMSHAPE0
FORJINRANGECMSHAPE1
AXTEXTJ I FORMATCMI J FMT
HACENTER VACENTER
COLORWHITE IFCMI J THRESH ELSEBLACK
FIGTIGHTLAYOUT
RETURNAX
NPSETPRINTOPTIONSPRECISION2
 PLOT NONNORMALIZED CONFUSION MATRIX
PLOTCONFUSIONMATRIXYTEST YPRED CLASSESCLASSNAMES
TITLECONFUSION MATRIX WITHOUT NORMALIZATION
 PLOT NORMALIZED CONFUSION MATRIX
PLOTCONFUSIONMATRIXYTEST YPRED CLASSESCLASSNAMES NORMALIZETRUE
TITLENORMALIZED CONFUSION MATRIX
PLTSHOW
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0217 SECONDS
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE
52113 VISUALIZING CROSSVALIDATION BEHAVIOR IN SCIKITLEARN
CHOOSING THE RIGHT CROSSVALIDATION OBJECT IS A CRUCIAL PART OF FITTING A MODEL PROPERLY THERE ARE MANY WAYS TO SPLIT
DATA INTO TRAINING AND TEST SETS IN ORDER TO AVOID MODEL OVERFITTING TO STANDARDIZE THE NUMBER OF GROUPS IN TEST SETS ETC
THIS EXAMPLE VISUALIZES THE BEHAVIOR OF SEVERAL COMMON SCIKITLEARN OBJECTS FOR COMPARISON
FROM SKLEARNMODELSELECTION IMPORT TIMESERIESSPLIT KFOLD SHUFFLESPLIT
STRATIFIEDKFOLD GROUPSHUFFLESPLIT
GROUPKFOLD STRATIFIEDSHUFFLESPLIT
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM MATPLOTLIBPATCHES IMPORT PATCH
NPRANDOMSEED1338
CMAPDATA PLTCMPAIED
CMAPCV PLTCMCOOLWARM
NSPLITS 4
VISUALIZE OUR DATA
FIRST WE MUST UNDERSTAND THE STRUCTURE OF OUR DATA IT HAS 100 RANDOMLY GENERATED INPUT DATAPOINTS 3 CLASSES SPLIT
UNEVENLY ACROSS DATAPOINTS AND 10 “GROUPS” SPLIT EVENLY ACROSS DATAPOINTS
AS WE’LL SEE SOME CROSSVALIDATION OBJECTS DO SPECIFIC THINGS WITH LABELED DATA OTHERS BEHAVE DIFFERENTLY WITH GROUPED
DATA AND OTHERS DO NOT USE THIS INFORMATION
TO BEGIN WE’LL VISUALIZE OUR DATA
521 MODEL SELECTION 1289
```

SCIKITLEARN USER GUIDE RELEASE 0213  
GENERATE THE CLASSGROUP DATA  
NPOINTS 100  
X NPRANDOMRANDN100 10  
PERCENTILESCLASSES 1 3 6  
Y NPHSTACKII INT100 PERC  
FORII PERC INENUMERATEPERCENTILESCLASSES  
EVENLY SPACED GROUPS REPEATED ONCE  
GROUPS NPHSTACKII 10FORIIINRANGE10  
DEFVISUALIZEGROUPSCLASSES GROUPS NAME  
VISUALIZE DATASET GROUPS  
FIG AX PLTSUBPLOTS  
AXSCATTERRANGELENGROUPS 5 LENGROUPS CGROUPS MARKER  
LW50 CMAPCMAPDATA  
AXSCATTERRANGELENGROUPS 35 LENGROUPS CCLASSES MARKER  
LW50 CMAPCMAPDATA  
AXSETYLIM1 5 YTICKS5 35  
YTICKLABELSDATA NGROUP DATA NCLASS XLABELSAMPLE INDEX  
VISUALIZEGROUPSY GROUPS NO GROUPS  
1290 CHAPTER 5 EXAMPLES



```
SCIKITLEARN USER GUIDE RELEASE 0213
DEFINE A FUNCTION TO VISUALIZE CROSSVALIDATION BEHAVIOR
WE'LL DEFINE A FUNCTION THAT LETS US VISUALIZE THE BEHAVIOR OF EACH CROSSVALIDATION OBJECT WE'LL PERFORM 4 SPLITS OF THE
DATA ON EACH SPLIT WE'LL VISUALIZE THE INDICES CHOSEN FOR THE TRAINING SET IN BLUE AND THE TEST SET IN RED
DEFPLOTCVINDICESCV X Y GROUP AX NSPLITS LW10
CREATE A SAMPLE PLOT FOR INDICES OF A CROSSVALIDATION OBJECT
GENERATE THE TRAININGTESTING VISUALIZATIONS FOR EACH CV SPLIT
FORII TR TT INENUMERATECVSPLITXX YY GROUPSGROUP
FILL IN INDICES WITH THE TRAININGTEST GROUPS
INDICES NPARRAYNPNAN LENX
INDICESTT 1
INDICESTR 0
VISUALIZE THE RESULTS
AXSCATTERRANGELENINDICES II 5 LENINDICES
CINDICES MARKER LWLW CMAPCMAPCV
VMIN2 VMAX12
PLOT THE DATA CLASSES AND GROUPS AT THE END
AXSCATTERRANGELENX II 15 LENX
CY MARKER LWLW CMAPCMAPDATA
AXSCATTERRANGELENX II 25 LENX
CGROUP MARKER LWLW CMAPCMAPDATA
FORMATTING
YTICKLABELS LISTRANGENSPLITS CLASS GROUP
AXSETYTICKSNPARANGENSPLITS2 5 YTICKLABELSYTICKLABELS
XLABELSAMPLE INDEX YLABELCV ITERATION
YLIMNSPLITS22 2 XLIM0 100
AXSETTITLEFORMATTYPEPCVNAME FONTSIZE15
RETURNAX
LET'S SEE HOW IT LOOKS FOR THE KFOLD CROSSVALIDATION OBJECT
FIG AX PLTSUBPLOTS
CV KFOLDNSPLITS
PLOTCVINDICESCV X Y GROUPS AX NSPLITS
521 MODEL SELECTION 1291
```

SCIKITLEARN USER GUIDE RELEASE 0213  
AS YOU CAN SEE BY DEFAULT THE KFOLD CROSSVALIDATION ITERATOR DOES NOT TAKE EITHER DATAPOINT CLASS OR GROUP INTO  
CONSIDERATION WE CAN CHANGE THIS BY USING THE STRATIFIEDKFOLD LIKE SO  
FIG AX PLTSUBPLOTS  
CV STRATIFIEDKFOLDNSPLITS  
PLOT CV INDICES CV X Y GROUPS AX NSPLITS  
1292 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

IN THIS CASE THE CROSSVALIDATION RETAINED THE SAME RATIO OF CLASSES ACROSS EACH CV SPLIT NEXT WE'LL VISUALIZE THIS BEHAVIOR FOR A NUMBER OF CV ITERATORS

VISUALIZE CROSSVALIDATION INDICES FOR MANY CV OBJECTS

LET'S VISUALLY COMPARE THE CROSS VALIDATION BEHAVIOR FOR MANY SCIKITLEARN CROSSVALIDATION OBJECTS BELOW WE WILL LOOP THROUGH SEVERAL COMMON CROSSVALIDATION OBJECTS VISUALIZING THE BEHAVIOR OF EACH

NOTE HOW SOME USE THE GROUPCLASS INFORMATION WHILE OTHERS DO NOT

CVS KFOLD GROUPKFOLD SHUFFLESPLIT STRATIFIEDKFOLD

GROUPSHUFFLESPLIT STRATIFIEDSHUFFLESPLIT TIMESERIESSPLIT

FORCVINCVS

THISCV CVNSPLITSNSPLITS

FIG AX PLTSUBPLOTSFIGSIZE6 3

PLOTVCINDICESTHISCV X Y GROUPS AX NSPLITS

AXLEGENDPATCHCOLORCMAPCV8 PATCHCOLORCMAPCV02

TESTING SET TRAINING SET LOC102 8

MAKE THE LEGEND FIT

PLTTIGHTLAYOUT

FIGSUBPLOTSADJUSTRIGHT7

PLTSHOW

521 MODEL SELECTION 1293

SCIKITLEARN USER GUIDE RELEASE 0213

- 
- 

1294 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

- 
- 

521 MODEL SELECTION 1295

SCIKITLEARN USER GUIDE RELEASE 0213

- 
- 

1296 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

•

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0401 SECONDS

NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE

52114 RECEIVER OPERATING CHARACTERISTIC ROC

EXAMPLE OF RECEIVER OPERATING CHARACTERISTIC ROC METRIC TO EVALUATE CLASSIFIER OUTPUT QUALITY

ROC CURVES TYPICALLY FEATURE TRUE POSITIVE RATE ON THE Y AXIS AND FALSE POSITIVE RATE ON THE X AXIS THIS MEANS THAT THE TOP LEFT CORNER OF THE PLOT IS THE “IDEAL” POINT A FALSE POSITIVE RATE OF ZERO AND A TRUE POSITIVE RATE OF ONE THIS IS NOT VERY REALISTIC BUT IT DOES MEAN THAT A LARGER AREA UNDER THE CURVE AUC IS USUALLY BETTER

THE “STEEPNESS” OF ROC CURVES IS ALSO IMPORTANT SINCE IT IS IDEAL TO MAXIMIZE THE TRUE POSITIVE RATE WHILE MINIMIZING THE FALSE POSITIVE RATE

MULTICLASS SETTINGS

ROC CURVES ARE TYPICALLY USED IN BINARY CLASSIFICATION TO STUDY THE OUTPUT OF A CLASSIFIER IN ORDER TO EXTEND ROC CURVE AND ROC AREA TO MULTICLASS OR MULTILABEL CLASSIFICATION IT IS NECESSARY TO BINARIZE THE OUTPUT ONE ROC CURVE CAN BE DRAWN PER LABEL BUT ONE CAN ALSO DRAW A ROC CURVE BY CONSIDERING EACH ELEMENT OF THE LABEL INDICATOR MATRIX AS A BINARY PREDICTION MICROAVERAGING

ANOTHER EVALUATION MEASURE FOR MULTICLASS CLASSIFICATION IS MACROAVERAGING WHICH GIVES EQUAL WEIGHT TO THE CLASSIFICATION OF EACH LABEL

NOTE

SEE ALSO SKLEARNMETRICSROCAUCSCORE RECEIVER OPERATING CHARACTERISTIC ROC WITH CROSS VALIDATION

PRINTDOC

IMPORT NUMPY AS NP

521 MODEL SELECTION 1297

```
SCIKITLEARN USER GUIDE RELEASE 0213
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM ITERTOOLS IMPORT CYCLE
FROM SKLEARN IMPORT SVM DATASETS
FROM SKLEARNMETRICS IMPORT ROCCURVE AUC
FROM SKLEARNMODELSELECTION IMPORT TRAINTESTSPLIT
FROM SKLEARNPREPROCESSING IMPORT LABELBINARIZE
FROM SKLEARNMULTICLASS IMPORT ONEVSRESTCLASSIFIER
FROM SCIPY IMPORT INTERP
IMPORT SOME DATA TO PLAY WITH
IRIS DATASETSLOADIRIS
X IRISDATA
Y IRISTARGET
BINARIZE THE OUTPUT
Y LABELBINARIZEY CLASSES0 1 2
NCLASSES YSHAPE1
ADD NOISY FEATURES TO MAKE THE PROBLEM HARDER
RANDOMSTATE NPRANDOMRANDOMSTATE0
NSAMPLES NFEATURES XSHAPE
X NPCX RANDOMSTATETERANDNNSAMPLES 200 NFEATURES
SHUFFLE AND SPLIT TRAINING AND TEST SETS
XTRAIN XTEST YTRAIN YTEST TRAINTESTSPLITX Y TESTSIZE5
RANDOMSTATE0
LEARN TO PREDICT EACH CLASS AGAINST THE OTHER
CLASSIFIER ONEVSRESTCLASSIFIERSVMSVCKERNELLINEAR PROBABILITYTRUE
RANDOMSTATETERANDOMSTATE
YSCORE CLASSIFIERFITXTRAIN YTRAINDECISIONFUNCTIONXTEST
COMPUTE ROC CURVE AND ROC AREA FOR EACH CLASS
FPR DICT
TPR DICT
ROCAUC DICT
FORIINRANGENCLASSES
FPRI TPRI ROCCURVEYTEST I YSCORE I
ROCAUCI AUCFPRI TPRI
COMPUTE MICROAVERAGE ROC CURVE AND ROC AREA
FPRMICRO TPRMICRO ROCCURVEYTESTRAVEL YSCORERAVEL
ROCAUCMICRO AUCFPRMICRO TPRMICRO
PLOT OF A ROC CURVE FOR A SPECIFIC CLASS
PLTFigure
LW 2
PLTPLOTFPR2 TPR2 COLORDARKORANGE
LWLW LABELROC CURVE AREA 02F ROCAUC2
PLTPLOT0 1 0 1 COLORNAVY LWLW LInestyle
PLTXLIM00 10
PLTYLIM00 105
PLTXLABELFALSE POSITIVE RATE
PLTYLABELTRUE POSITIVE RATE
PLTTITLERECEIVER OPERATING CHARACTERISTIC EXAMPLE
PLTLEGENDLOCLOWER RIGHT
1298 CHAPTER 5 EXAMPLES
```



SCIKITLEARN USER GUIDE RELEASE 0213  
PLTSHOW  
PLOT ROC CURVES FOR THE MULTICLASS PROBLEM  
  COMPUTE MACROAVERAGE ROC CURVE AND ROC AREA  
  FIRST AGGREGATE ALL FALSE POSITIVE RATES  
ALLFPR NPUNIQUENPCONCATENATEFPRI FORIINRANGENCLASSES  
  THEN INTERPOLATE ALL ROC CURVES AT THIS POINTS  
MEANTPR NPZEROSLIKEALLFPR  
FORIINRANGENCLASSES  
MEANTPR INTERPALLFPR FPRI TPRI  
  FINALLY AVERAGE IT AND COMPUTE AUC  
MEANTPR NCLASSES  
FPRMACRO ALLFPR  
TPRMACRO MEANTPR  
ROCAUCMACRO AUCFPRMACRO TPRMACRO  
  PLOT ALL ROC CURVES  
PLTFigure  
PLTPLOTFPRMICRO TPRMICRO  
LABELMICROAVERAGE ROC CURVE AREA 002F  
FORMATROCAUCMICRO  
521 MODEL SELECTION 1299

SCIKITLEARN USER GUIDE RELEASE 0213  
COLORDEEPPINK LINESTYLE LINEWIDTH4  
PLTPLOTFPRMACRO TPRMACRO  
LABELMACROAVERAGE ROC CURVE AREA 002F  
FORMATROCAUCMACRO  
COLORNAVY LINESTYLE LINEWIDTH4  
COLORS CYCLEAQUA DARKORANGE CORNFLOWERBLUE  
FORI COLOR INZIPRANGENCLASSES COLORS  
PLTPLOTFPRI TPRI COLORCOLOR LWLW  
LABELROC CURVE OF CLASS 0 AREA 102F  
FORMATI ROCAUCI  
PLTPLOT0 1 0 1 K LWLW  
PLTXLIM00 10  
PLTYLIM00 105  
PLTXLABELFALSE POSITIVE RATE  
PLTYLABELTRUE POSITIVE RATE  
PLTTITLESOME EXTENSION OF RECEIVER OPERATING CHARACTERISTIC TO MULTICLASS  
PLTLEGENDLOCLOWER RIGHT  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0145 SECONDS  
1300 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
52115 PLOTTING LEARNING CURVES

ON THE LEFT SIDE THE LEARNING CURVE OF A NAIVE BAYES CLASSIFIER IS SHOWN FOR THE DIGITS DATASET NOTE THAT THE TRAINING SCORE AND THE CROSSVALIDATION SCORE ARE BOTH NOT VERY GOOD AT THE END HOWEVER THE SHAPE OF THE CURVE CAN BE FOUND IN MORE COMPLEX DATASETS VERY OFTEN THE TRAINING SCORE IS VERY HIGH AT THE BEGINNING AND DECREASES AND THE CROSS VALIDATION SCORE IS VERY LOW AT THE BEGINNING AND INCREASES ON THE RIGHT SIDE WE SEE THE LEARNING CURVE OF AN SVM WITH RBF KERNEL WE CAN SEE CLEARLY THAT THE TRAINING SCORE IS STILL AROUND THE MAXIMUM AND THE VALIDATION SCORE COULD BE INCREASED WITH MORE TRAINING SAMPLES

•  
521 MODEL SELECTION 1301

SCIKITLEARN USER GUIDE RELEASE 0213

•

```
PRINTDOC
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARNNAIVEBAYES IMPORT GAUSSIANNB
FROM SKLEARN SVM IMPORT SVC
FROM SKLEARN DATASETS IMPORT LOADDIGITS
FROM SKLEARN MODELSELECTION IMPORT LEARNINGCURVE
FROM SKLEARN MODELSELECTION IMPORT SHUFFLESPLIT
DEFPLOTLEARNINGCURVEESTIMATOR TITLE X Y YLIMNONE CVNONE
NJOBSNONE TRAINSIZESNPLINSPACE1 10 5
```

GENERATE A SIMPLE PLOT OF THE TEST AND TRAINING LEARNING CURVE  
PARAMETERS

ESTIMATOR OBJECT TYPE THAT IMPLEMENTS THE FIT AND PREDICT METHODS  
AN OBJECT OF THAT TYPE WHICH IS CLONED FOR EACH VALIDATION  
TITLE STRING  
TITLE FOR THE CHART  
X ARRAYLIKE SHAPE NSAMPLES NFEATURES  
TRAINING VECTOR WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND  
NFEATURES IS THE NUMBER OF FEATURES  
1302 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

Y ARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NFEATURES OPTIONAL  
TARGET RELATIVE TO X FOR CLASSIFICATION OR REGRESSION  
NONE FOR UNSUPERVISED LEARNING

YLIM TUPLE SHAPE YMIN YMAX OPTIONAL  
DEFINES MINIMUM AND MAXIMUM YVALUES PLOTTED

CV INT CROSSVALIDATION GENERATOR OR AN ITERABLE OPTIONAL  
DETERMINES THE CROSSVALIDATION SPLITTING STRATEGY  
POSSIBLE INPUTS FOR CV ARE  
NONE TO USE THE DEFAULT 3FOLD CROSSVALIDATION  
INTEGER TO SPECIFY THE NUMBER OF FOLDS  
TERMCV SPLITTER  
AN ITERABLE YIELDING TRAIN TEST SPLITS AS ARRAYS OF INDICES  
FOR INTEGERNONE INPUTS IF Y IS BINARY OR MULTICLASS  
CLASSTRATIFIEDKFOLD USED IF THE ESTIMATOR IS NOT A CLASSIFIER  
OR IF Y IS NEITHER BINARY NOR MULTICLASS CLASSKFOLD IS USED  
REFER REFUSER GUIDE CROSSVALIDATION FOR THE VARIOUS  
CROSSVALIDATORS THAT CAN BE USED HERE

NJOBS INT OR NONE OPTIONAL DEFAULTNONE  
NUMBER OF JOBS TO RUN IN PARALLEL  
NONE MEANS 1 UNLESS IN A OBJJOBLIBPARALLELBACKEND CONTEXT  
1 MEANS USING ALL PROCESSORS SEE TERMGLOSSARY NJOBS  
FOR MORE DETAILS

TRAINSIZES ARRAYLIKE SHAPE NTICKS DTYPE FLOAT OR INT  
RELATIVE OR ABSOLUTE NUMBERS OF TRAINING EXAMPLES THAT WILL BE USED TO  
GENERATE THE LEARNING CURVE IF THE DTYPE IS FLOAT IT IS REGARDED AS A  
FRACTION OF THE MAXIMUM SIZE OF THE TRAINING SET THAT IS DETERMINED  
BY THE SELECTED VALIDATION METHOD IE IT HAS TO BE WITHIN 0 1  
OTHERWISE IT IS INTERPRETED AS ABSOLUTE SIZES OF THE TRAINING SETS  
NOTE THAT FOR CLASSIFICATION THE NUMBER OF SAMPLES USUALLY HAVE TO  
BE BIG ENOUGH TO CONTAIN AT LEAST ONE SAMPLE FROM EACH CLASS  
DEFAULT NPLINSPACE01 10 5

PLTFigure

PLTTITLEtitle

IFYLIMIS NOTNONE

PLTYLIM YLIM

PLTXLABELTRAINING EXAMPLES

PLTYLABELSCORE

TRAINSIZES TRAINSCORES TESTSCORES LEARNINGCURVE  
ESTIMATOR X Y CVCV NJOBSNJOBS TRAINSIZES TRAINSIZES  
TRAINSCORESMEAN NPMEANTRAINSCORES AXIS1  
TRAINSCORESSTD NPSTDTRAINSCORES AXIS1  
TESTSCORESMEAN NPMEANTESTSCORES AXIS1  
TESTSCORESSTD NPSTDTESTSCORES AXIS1

PLTGRID

PLTFILLBETWEEN TRAINSIZES TRAINSCORESMEAN TRAINSCORESSTD  
TRAINSCORESMEAN TRAINSCORESSTD ALPHA01

COLORR

PLTFILLBETWEEN TRAINSIZES TESTSCORESMEAN TESTSCORESSTD  
TESTSCORESMEAN TESTSCORESSTD ALPHA01 COLORG

521 MODEL SELECTION 1303

SCIKITLEARN USER GUIDE RELEASE 0213  
PLTPLOTTRAINIZES TRAINSCORESMEAN O COLORR  
LABELTRAINING SCORE  
PLTPLOTTRAINIZES TESTSCORESMEAN O COLORG  
LABELCROSSVALIDATION SCORE  
PLTLEGENDLOCBEST  
RETURNPLT  
DIGITS LOADDIGITS  
X Y DIGITSDATA DIGITSTARGET  
TITLE LEARNING CURVES NAIVE BAYES  
CROSS VALIDATION WITH 100 ITERATIONS TO GET SMOOTHER MEAN TEST AND TRAIN  
SCORE CURVES EACH TIME WITH 20 DATA RANDOMLY SELECTED AS A VALIDATION SET  
CV SHUFFLESPLITNSPLITS100 TESTSIZE02 RANDOMSTATE0  
ESTIMATOR GAUSSIANNB  
PLOTLEARNINGCURVEESTIMATOR TITLE X Y YLIM07 101 CVCV NJOBS4  
TITLE RLEARNING CURVES SVM RBF KERNEL GAMMA0001  
SVC IS MORE EXPENSIVE SO WE DO A LOWER NUMBER OF CV ITERATIONS  
CV SHUFFLESPLITNSPLITS10 TESTSIZE02 RANDOMSTATE0  
ESTIMATOR SVCGAMMA0001  
PLOTLEARNINGCURVEESTIMATOR TITLE X Y 07 101 CVCV NJOBS4  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 2856 SECONDS  
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE  
52116 PRECISIONRECALL  
EXAMPLE OF PRECISIONRECALL METRIC TO EVALUATE CLASSIFIER OUTPUT QUALITY  
PRECISIONRECALL IS A USEFUL MEASURE OF SUCCESS OF PREDICTION WHEN THE CLASSES ARE VERY IMBALANCED IN INFORMATION  
RETRIEVAL PRECISION IS A MEASURE OF RESULT RELEVANCY WHILE RECALL IS A MEASURE OF HOW MANY TRULY RELEVANT RESULTS ARE  
RETURNED  
THE PRECISIONRECALL CURVE SHOWS THE TRADEOFF BETWEEN PRECISION AND RECALL FOR DIFFERENT THRESHOLD A HIGH AREA UNDER  
THE CURVE REPRESENTS BOTH HIGH RECALL AND HIGH PRECISION WHERE HIGH PRECISION RELATES TO A LOW FALSE POSITIVE RATE AND  
HIGH RECALL RELATES TO A LOW FALSE NEGATIVE RATE HIGH SCORES FOR BOTH SHOW THAT THE CLASSIFIER IS RETURNING ACCURATE RESULTS  
HIGH PRECISION AS WELL AS RETURNING A MAJORITY OF ALL POSITIVE RESULTS HIGH RECALL  
A SYSTEM WITH HIGH RECALL BUT LOW PRECISION RETURNS MANY RESULTS BUT MOST OF ITS PREDICTED LABELS ARE INCORRECT WHEN  
COMPARED TO THE TRAINING LABELS A SYSTEM WITH HIGH PRECISION BUT LOW RECALL IS JUST THE OPPOSITE RETURNING VERY FEW  
RESULTS BUT MOST OF ITS PREDICTED LABELS ARE CORRECT WHEN COMPARED TO THE TRAINING LABELS AN IDEAL SYSTEM WITH HIGH  
PRECISION AND HIGH RECALL WILL RETURN MANY RESULTS WITH ALL RESULTS LABELED CORRECTLY  
PRECISION  $\frac{TP}{P}$  IS DEFINED AS THE NUMBER OF TRUE POSITIVES  $TP$  OVER THE NUMBER OF TRUE POSITIVES PLUS THE NUMBER OF FALSE  
POSITIVES  $P$   
 $\frac{TP}{P}$   
 $\frac{TP}{P}$   
1304 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

RECALL  $\frac{TP}{TP + FN}$  IS DEFINED AS THE NUMBER OF TRUE POSITIVES  $TP$  OVER THE NUMBER OF TRUE POSITIVES PLUS THE NUMBER OF FALSE NEGATIVES  $FN$

$\frac{1}{\frac{1}{P} + \frac{1}{R}}$

THESE QUANTITIES ARE ALSO RELATED TO THE  $F_1$  SCORE WHICH IS DEFINED AS THE HARMONIC MEAN OF PRECISION AND RECALL

$F_1 = 2 \times \frac{P \times R}{P + R}$

NOTE THAT THE PRECISION MAY NOT DECREASE WITH RECALL THE DEFINITION OF PRECISION  $P = \frac{TP}{TP + FP}$

$F_1$  SCORE SHOWS THAT LOWERING THE THRESHOLD OF A CLASSIFIER MAY INCREASE THE DENOMINATOR BY INCREASING THE NUMBER OF RESULTS RETURNED IF THE THRESHOLD WAS PREVIOUSLY SET TOO HIGH THE NEW RESULTS MAY ALL BE TRUE POSITIVES WHICH WILL INCREASE PRECISION IF THE PREVIOUS THRESHOLD WAS ABOUT RIGHT OR TOO LOW FURTHER LOWERING THE THRESHOLD WILL INTRODUCE FALSE POSITIVES DECREASING PRECISION RECALL IS DEFINED AS  $R = \frac{TP}{TP + FN}$

$F_1$  SCORE WHERE  $F_1$  DOES NOT DEPEND ON THE CLASSIFIER THRESHOLD THIS MEANS THAT LOWERING THE CLASSIFIER THRESHOLD MAY INCREASE RECALL BY INCREASING THE NUMBER OF TRUE POSITIVE RESULTS IT IS ALSO POSSIBLE THAT LOWERING THE THRESHOLD MAY LEAVE RECALL UNCHANGED WHILE THE PRECISION FLUCTUATES THE RELATIONSHIP BETWEEN RECALL AND PRECISION CAN BE OBSERVED IN THE STAIRSTEP AREA OF THE PLOT AT THE EDGES OF THESE STEPS A SMALL CHANGE IN THE THRESHOLD CONSIDERABLY REDUCES PRECISION WITH ONLY A MINOR GAIN IN RECALL AVERAGE PRECISION AP SUMMARIZES SUCH A PLOT AS THE WEIGHTED MEAN OF PRECISIONS ACHIEVED AT EACH THRESHOLD WITH THE INCREASE IN RECALL FROM THE PREVIOUS THRESHOLD USED AS THE WEIGHT

$AP = \frac{\sum_{i=1}^N P_i \Delta R_i}{R_N}$

WHERE  $P_i$  AND  $R_i$  ARE THE PRECISION AND RECALL AT THE  $i$ TH THRESHOLD A PAIR  $(P_i, R_i)$  IS REFERRED TO AS AN OPERATING POINT AP AND THE TRAPEZOIDAL AREA UNDER THE OPERATING POINTS  $SKLEARNMETRICS.AUC$  ARE COMMON WAYS TO SUMMARIZE A PRECISIONRECALL CURVE THAT LEAD TO DIFFERENT RESULTS READ MORE IN THE USER GUIDE PRECISIONRECALL CURVES ARE TYPICALLY USED IN BINARY CLASSIFICATION TO STUDY THE OUTPUT OF A CLASSIFIER IN ORDER TO EXTEND THE PRECISIONRECALL CURVE AND AVERAGE PRECISION TO MULTICLASS OR MULTILABEL CLASSIFICATION IT IS NECESSARY TO BINARIZE THE OUTPUT ONE CURVE CAN BE DRAWN PER LABEL BUT ONE CAN ALSO DRAW A PRECISIONRECALL CURVE BY CONSIDERING EACH ELEMENT OF THE LABEL INDICATOR MATRIX AS A BINARY PREDICTION MICROAVERAGING

NOTE  
SEE ALSO  $SKLEARNMETRICS.AVERAGEPRECISIONSCORE$   $SKLEARNMETRICS.RECALLSCORE$   $SKLEARNMETRICS.PRECISIONSCORE$   $SKLEARNMETRICS.F1SCORE$  IN BINARY CLASSIFICATION SETTINGS

```
CREATE SIMPLE DATA
TRY TO DIFFERENTIATE THE TWO FIRST CLASSES OF THE IRIS DATA
FROM SKLEARN IMPORT SVM DATASETS
FROM SKLEARNMODELSELECTION IMPORT TRAINTESTSPLIT
IMPORT NUMPY AS NP
IRIS = DATASETS.LOADIRIS
X = IRIS.DATA
Y = IRIS.TARGET
ADD NOISY FEATURES
RANDOMSTATE = NPROBANDRANDOMSTATE0
NSAMPLES NFEATURES X.SHAPE
521 MODEL SELECTION 1305
```

```
SCIKITLEARN USER GUIDE RELEASE 0213
X NPCX RANDOMSTATERANDNNSAMPLES 200 NFEATURES
LIMIT TO THE TWO FIRST CLASSES AND SPLIT INTO TRAINING AND TEST
XTRAIN XTEST YTRAIN YTEST TRAINTESTSPLITXY 2 YY 2
TESTSIZE5
RANDOMSTATERANDOMSTATE
CREATE A SIMPLE CLASSIFIER
CLASSIFIER SVMLINEARSVCRANDOMSTATERANDOMSTATE
CLASSIFIERFITXTRAIN YTRAIN
YSCORE CLASSIFIERDECISIONFUNCTIONXTEST
COMPUTE THE AVERAGE PRECISION SCORE
FROM SKLEARNMETRICS IMPORT AVERAGEPRECISIONSCORE
AVERAGEPRECISION AVERAGEPRECISIONSCOREYTEST YSCORE
PRINTAVERAGE PRECISIONRECALL SCORE 002FFORMAT
AVERAGEPRECISION
OUT
AVERAGE PRECISIONRECALL SCORE 088
PLOT THE PRECISIONRECALL CURVE
FROM SKLEARNMETRICS IMPORT PRECISIONRECALLCURVE
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM INSPECT IMPORT SIGNATURE
PRECISION RECALL PRECISIONRECALLCURVEYTEST YSCORE
IN MATPLOTLIB 15 PLTFILLBETWEEN DOES NOT HAVE A STEP ARGUMENT
STEPKWARGS STEP POST
IFSTEPINSIGNATUREPLTFILLBETWEENPARAMETERS
ELSE
PLTSTEPRECALL PRECISION COLORB ALPHA02
WHEREPOST
PLTFILLBETWEENRECALL PRECISION ALPHA02 COLORB STEPKWARGS
PLTXLABELRECALL
PLTYLABELPRECISION
PLTYLIM00 105
PLTXLIM00 10
PLTTITLE2CLASS PRECISIONRECALL CURVE AP002FFORMAT
AVERAGEPRECISION
1306 CHAPTER 5 EXAMPLES
```



```
SCIKITLEARN USER GUIDE RELEASE 0213
IN MULTILABEL SETTINGS
CREATE MULTILABEL DATA FIT AND PREDICT
WE CREATE A MULTILABEL DATASET TO ILLUSTRATE THE PRECISIONRECALL IN MULTILABEL SETTINGS
FROM SKLEARNPREPROCESSING IMPORT LABELBINARIZE
USE LABELBINARIZE TO BE MULTILABEL LIKE SETTINGS
Y LABELBINARIZEY CLASSES0 1 2
NCLASSES YSHAPE1
SPLIT INTO TRAINING AND TEST
XTRAIN XTEST YTRAIN YTEST TRAINTESTSPLITX Y TESTSIZE5
RANDOMSTATERRANDOMSTATE
WE USE ONEVSRESTCLASSIFIER FOR MULTILABEL PREDICTION
FROM SKLEARNMULTICLASS IMPORT ONEVSRESTCLASSIFIER
RUN CLASSIFIER
CLASSIFIER ONEVSRESTCLASSIFIERSVMLINEARSVCRANDOMSTATERRANDOMSTATE
CLASSIFIERFITXTRAIN YTRAIN
YSCORE CLASSIFIERDECISIONFUNCTIONXTEST
521 MODEL SELECTION 1307
```

SCIKITLEARN USER GUIDE RELEASE 0213  
THE AVERAGE PRECISION SCORE IN MULTILABEL SETTINGS  
FROM SKLEARNMETRICS IMPORT PRECISIONRECALLCURVE  
FROM SKLEARNMETRICS IMPORT AVERAGEPRECISIONSCORE  
FOR EACH CLASS  
PRECISION DICT  
RECALL DICT  
AVERAGEPRECISION DICT  
FORIINRANGENCLASSES  
PRECISIONI RECALLI PRECISIONRECALLCURVEYTEST I  
YSCORE I  
AVERAGEPRECISIONI AVERAGEPRECISIONSCOREYTEST I YSCORE I  
A MICROAVERAGE QUANTIFYING SCORE ON ALL CLASSES JOINTLY  
PRECISIONMICRO RECALLMICRO PRECISIONRECALLCURVEYTESTRAVEL  
YSCORERAVEL  
AVERAGEPRECISIONMICRO AVERAGEPRECISIONSCOREYTEST YSCORE  
AVERAGEMICRO  
PRINTAVERAGE PRECISION SCORE MICROAVERAGED OVER ALL CLASSES 002F  
FORMATAVERAGEPRECISIONMICRO  
OUT  
AVERAGE PRECISION SCORE MICROAVERAGED OVER ALL CLASSES 043  
PLOT THE MICROAVERAGED PRECISIONRECALL CURVE  
PLTFigure  
PLTSTEPRECALLMICRO PRECISIONMICRO COLORB ALPHA02  
WHEREPOST  
PLTFILLBETWEENRECALLMICRO PRECISIONMICRO ALPHA02 COLORB  
STEPKWARGS  
PLTXLABELRECALL  
PLTYLABELPRECISION  
PLTYLIM00 105  
PLTXLIM00 10  
PLTTITLE  
AVERAGE PRECISION SCORE MICROAVERAGED OVER ALL CLASSES AP002F  
FORMATAVERAGEPRECISIONMICRO  
1308 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
PLOT PRECISIONRECALL CURVE FOR EACH CLASS AND ISOF1 CURVES
FROM ITERTOOLS IMPORT CYCLE
  SETUP PLOT DETAILS
COLORS  CYCLENAVY TURQUOISE DARKORANGE CORNFLOWERBLUE TEAL
PLTFIGUREFIGSIZE7 8
FSCORES  NPLinspace02 08 NUM4
LINES
LABELS
FORFSCORE INFSCORES
X  NPLinspace001 1
Y  FSCORE X  2X  FSCORE
L  PLTPLOTXY 0 YY 0 COLORGRAY ALPHA02
PLTANNOTATEF1001FFORMATFSCORE XY09 Y45 002
LINESAPPENDL
LABELSAPPENDISOF1 CURVES
L  PLTPLOTRECALLMICRO PRECISIONMICRO COLORGOLD LW2
LINESAPPENDL
LABELSAPPENDMICROAVERAGE PRECISIONRECALL AREA 002F
FORMATAVERAGEPRECISIONMICRO
FORI COLOR INZIPRANGENCLASSES COLORS
521 MODEL SELECTION 1309
```

SCIKITLEARN USER GUIDE RELEASE 0213  
L PLTPLOTRECALLI PRECISIONI COLORCOLOR LW2  
LINESAPPENDL  
LABELSAPPENDPRECISIONRECALL FOR CLASS 0 AREA 102F  
FORMATI AVERAGEPRECISIONI  
FIG PLTGCF  
FIGSUBPLOTSADJUSTBOTTOM025  
PLTXLIM00 10  
PLTYLIM00 105  
PLTXLABELRECALL  
PLTYLABELPRECISION  
PLTTITLEEXTENSION OF PRECISIONRECALL CURVE TO MULTICLASS  
PLTLEGENDLINES LABELS LOC0 38 PROPDICTSIZE14  
PLTSHOW  
1310 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0064 SECONDS  
522 MULTIOUTPUT METHODS  
EXAMPLES CONCERNING THE SKLEARNMULTIOUTPUT MODULE  
522 MULTIOUTPUT METHODS 1311

SCIKITLEARN USER GUIDE RELEASE 0213

[NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)

5221 CLASSIFIER CHAIN

EXAMPLE OF USING CLASSIFIER CHAIN ON A MULTILABEL DATASET

FOR THIS EXAMPLE WE WILL USE THE YEAST DATASET WHICH CONTAINS 2417 DATAPPOINTS EACH WITH 103 FEATURES AND 14 POSSIBLE LABELS EACH DATA POINT HAS AT LEAST ONE LABEL AS A BASELINE WE FIRST TRAIN A LOGISTIC REGRESSION CLASSIFIER FOR EACH OF THE 14 LABELS TO EVALUATE THE PERFORMANCE OF THESE CLASSIFIERS WE PREDICT ON A HELDOUT TEST SET AND CALCULATE THE JACCARD SCORE FOR EACH SAMPLE

NEXT WE CREATE 10 CLASSIFIER CHAINS EACH CLASSIFIER CHAIN CONTAINS A LOGISTIC REGRESSION MODEL FOR EACH OF THE 14 LABELS THE MODELS IN EACH CHAIN ARE ORDERED RANDOMLY IN ADDITION TO THE 103 FEATURES IN THE DATASET EACH MODEL GETS THE PREDICTIONS OF THE PRECEDING MODELS IN THE CHAIN AS FEATURES NOTE THAT BY DEFAULT AT TRAINING TIME EACH MODEL GETS THE TRUE LABELS AS FEATURES THESE ADDITIONAL FEATURES ALLOW EACH CHAIN TO EXPLOIT CORRELATIONS AMONG THE CLASSES THE JACCARD SIMILARITY SCORE FOR EACH CHAIN TENDS TO BE GREATER THAN THAT OF THE SET INDEPENDENT LOGISTIC MODELS BECAUSE THE MODELS IN EACH CHAIN ARE ARRANGED RANDOMLY THERE IS SIGNIFICANT VARIATION IN PERFORMANCE AMONG THE CHAINS PRESUMABLY THERE IS AN OPTIMAL ORDERING OF THE CLASSES IN A CHAIN THAT WILL YIELD THE BEST PERFORMANCE HOWEVER WE DO NOT KNOW THAT ORDERING A PRIORI INSTEAD WE CAN CONSTRUCT AN VOTING ENSEMBLE OF CLASSIFIER CHAINS BY AVERAGING THE BINARY PREDICTIONS OF THE CHAINS AND APPLY A THRESHOLD OF 05 THE JACCARD SIMILARITY SCORE OF THE ENSEMBLE IS GREATER THAN THAT OF THE INDEPENDENT MODELS AND TENDS TO EXCEED THE SCORE OF EACH CHAIN IN THE ENSEMBLE ALTHOUGH THIS IS NOT GUARANTEED WITH RANDOMLY ORDERED CHAINS

AUTHOR ADAM KLECZEWSKI

LICENSE BSD 3 CLAUSE

IMPORT NUMPY AS NP

IMPORT MATPLOTLIBPYPLOT AS PLT

FROM SKLEARNDATASETS IMPORT FETCHOPENML

1312 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
FROM SKLEARNMULTIOUTPUT IMPORT CLASSIFIERCHAIN
FROM SKLEARNMODELSELECTION IMPORT TRAINTESTSPLIT
FROM SKLEARNMULTICLASS IMPORT ONEVSRESTCLASSIFIER
FROM SKLEARNMETRICS IMPORT JACCARDScore
FROM SKLEARNLINEARMODEL IMPORT LOGISTICREGRESSION
PRINTDOC
  LOAD A MULTILABEL DATASET FROM HTTPSWWWOPENMLORG40597
X Y FETCHOPENMLYEAST VERSION4 RETURNXYTRUE
Y Y TRUE
XTRAIN XTEST YTRAIN YTEST TRAINTESTSPLITX Y TESTSIZE2
RANDOMSTATE0
  FIT AN INDEPENDENT LOGISTIC REGRESSION MODEL FOR EACH CLASS USING THE
  ONEVSRESTCLASSIFIER WRAPPER
BASELR LOGISTICREGRESSIONSOLVERLBFGS
OVR ONEVSRESTCLASSIFIERBASELR
OVRFITXTRAIN YTRAIN
YPREDOVR OVRPREDICTXTEST
OVRJACCARDScore JACCARDScoreYTEST YPREDOVR AVERAGESAMPLES
  FIT AN ENSEMBLE OF LOGISTIC REGRESSION CLASSIFIER CHAINS AND TAKE THE
  TAKE THE AVERAGE PREDICTION OF ALL THE CHAINS
CHAINS CLASSIFIERCHAINBASELR ORDERRANDOM RANDOMSTATEI
FORIINRANGE10
FORCHAININCHAINS
CHAINFITXTRAIN YTRAIN
YPREDCHAINS NPARRAYCHAINPREDICTXTEST FORCHAININ
CHAINS
CHAINJACCARDScores JACCARDScoreYTEST YPREDCHAIN 5
AVERAGESAMPLES
FORYPREDCHAIN INYPREDCHAINS
YPREDENSEMBLE YPREDCHAINSMEANAXISO
ENSEMBLEJACCARDScore JACCARDScoreYTEST
YPREDENSEMBLE 5
AVERAGESAMPLES
MODELScores OVRJACCARDScore CHAINJACCARDScores
MODELScoresAPPENDENSEMBLEJACCARDScore
MODELNames INDEPENDENT
CHAIN 1
CHAIN 2
CHAIN 3
CHAIN 4
CHAIN 5
CHAIN 6
CHAIN 7
CHAIN 8
CHAIN 9
CHAIN 10
ENSEMBLE
XPOS NPARANGELENMODELNames
522 MULTIOUTPUT METHODS 1313
```

SCIKITLEARN USER GUIDE RELEASE 0213

PLOT THE JACCARD SIMILARITY SCORES FOR THE INDEPENDENT MODEL EACH OF THE CHAINS AND THE ENSEMBLE NOTE THAT THE VERTICAL AXIS ON THIS PLOT DOES NOT BEGIN AT 0

FIG AX PLTSUBPLOTSFIGSIZE7 4

AXGRIDTRUE

AXSETTITLECLASSIFIER CHAIN ENSEMBLE PERFORMANCE COMPARISON

AXSETXTICKSXPOS

AXSETXTICKLABELSMODELNAME ROTATIONVERTICAL

AXSETYLABELJACCARD SIMILARITY SCORE

AXSETYLIMMINMODELScores 9 MAXMODELScores 11

COLORS R B LENCHAINJACCARDSCORES G

AXBARXPOS MODELScores ALPHA05 COLORCOLORS

PLTTIGHTLAYOUT

PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 7947 SECONDS

523 NEAREST NEIGHBORS

EXAMPLES CONCERNING THE SKLEARNNEIGHBORS MODULE

NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)

5231 NEAREST NEIGHBORS REGRESSION

DEMONSTRATE THE RESOLUTION OF A REGRESSION PROBLEM USING A KNEAREST NEIGHBOR AND THE INTERPOLATION OF THE TARGET

USING BOTH BARYCENTER AND CONSTANT WEIGHTS

1314 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213  
PRINTDOC  
AUTHOR ALEXANDRE GRAMFORT ALEXANDREGRAMFORTINRIA  
FABIAN PEDREGOSA FABIANPEDREGOSAINRIA

LICENSE BSD 3 CLAUSE C INRIA

GENERATE SAMPLE DATA  
IMPORT NUMPY AS NP  
IMPORT MATPLOTLIBPYPLOT AS PLT  
FROM SKLEARN IMPORT NEIGHBORS  
NPRANDOMSEED0  
X NPSORT5 NPRANDOMRAND40 1 AXIS0  
T NPLinspace0 5 500 NPNEWAXIS  
Y NPSINXRAVEL  
ADD NOISE TO TARGETS  
Y5 1 05 NPRANDOMRAND8

FIT REGRESSION MODEL  
NNEIGHBORS 5  
523 NEAREST NEIGHBORS 1315

SCIKITLEARN USER GUIDE RELEASE 0213  
FOR I WEIGHTS INENUMERATEUNIFORM DISTANCE  
KNN NEIGHBORSKNEIGHBORSREGRESSORNNEIGHBORS WEIGHTSWEIGHTS  
Y KNNFITX YPREDICTT  
PLTSUBPLOT2 1 I 1  
PLTSCATTERX Y CK LABELDATA  
PLTPLOTT Y CG LABELPREDICTION  
PLTAXISTIGHT  
PLTLEGEND  
PLTTITLEKNEIGHBORSREGRESSOR K I WEIGHTS S NNEIGHBORS  
WEIGHTS  
PLTTIGHTLAYOUT  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0089 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
5232 OUTLIER DETECTION WITH LOCAL OUTLIER FACTOR LOF  
THE LOCAL OUTLIER FACTOR LOF ALGORITHM IS AN UNSUPERVISED ANOMALY DETECTION METHOD WHICH COMPUTES THE LOCAL  
DENSITY DEVIATION OF A GIVEN DATA POINT WITH RESPECT TO ITS NEIGHBORS IT CONSIDERS AS OUTLIERS THE SAMPLES THAT HAVE A  
SUBSTANTIALLY LOWER DENSITY THAN THEIR NEIGHBORS THIS EXAMPLE SHOWS HOW TO USE LOF FOR OUTLIER DETECTION WHICH IS  
THE DEFAULT USE CASE OF THIS ESTIMATOR IN SCIKITLEARN NOTE THAT WHEN LOF IS USED FOR OUTLIER DETECTION IT HAS NO PREDICT  
DECISIONFUNCTION AND SCORESAMPLES METHODS SEE USER GUIDE FOR DETAILS ON THE DIFFERENCE BETWEEN OUTLIER DETECTION  
AND NOVELTY DETECTION AND HOW TO USE LOF FOR NOVELTY DETECTION  
THE NUMBER OF NEIGHBORS CONSIDERED PARAMETER NNEIGHBORS IS TYPICALLY SET 1 GREATER THAN THE MINIMUM NUMBER OF  
SAMPLES A CLUSTER HAS TO CONTAIN SO THAT OTHER SAMPLES CAN BE LOCAL OUTLIERS RELATIVE TO THIS CLUSTER AND 2 SMALLER THAN  
THE MAXIMUM NUMBER OF CLOSE BY SAMPLES THAT CAN POTENTIALLY BE LOCAL OUTLIERS IN PRACTICE SUCH INFORMATIONS ARE  
GENERALLY NOT AVAILABLE AND TAKING NNEIGHBORS20 APPEARS TO WORK WELL IN GENERAL  
1316 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARNNEIGHBORS IMPORT LOCALOUTLIERFACTOR
PRINTDOC
NPRANDOMSEED42
GENERATE TRAIN DATA
XINLIERS 03 NPRANDOMRANDN100 2
XINLIERS NPRXINLIERS 2 XINLIERS 2
GENERATE SOME OUTLIERS
XOUTLIERS NPRANDOMUNIFORMLOW4 HIGH4 SIZE20 2
X NPRXINLIERS XOUTLIERS
NOUTLIERS LENXOUTLIERS
GROUNDTRUTH NPONESLENX DTYPEINT
GROUNDTRUTHNOUTLIERS 1
FIT THE MODEL FOR OUTLIER DETECTION DEFAULT
CLF LOCALOUTLIERFACTORNNNEIGHBORS20 CONTAMINATION01
USE FITPREDICT TO COMPUTE THE PREDICTED LABELS OF THE TRAINING SAMPLES
WHEN LOF IS USED FOR OUTLIER DETECTION THE ESTIMATOR HAS NO PREDICT
DECISIONFUNCTION AND SCORESAMPLES METHODS
YPRED CLFFITPREDICTX
523 NEAREST NEIGHBORS 1317
```

SCIKITLEARN USER GUIDE RELEASE 0213  
NERRORS YPRED GROUNDTRUTHSUM  
XSCORES CLFNegativeOutlierFactor  
PLTTITLELOCAL OUTLIER FACTOR LOF  
PLTSCATTERX 0 X 1 COLORK S3 LABELDATA POINTS  
PLOT CIRCLES WITH RADIUS PROPORTIONAL TO THE OUTLIER SCORES  
RADIUS XSCORESMAX XSCORES XSCORESMAX XSCORESMIN  
PLTSCATTERX 0 X 1 S1000 RADIUS EDGECOLORSR  
FACECOLORSNONE LABELOUTLIER SCORES  
PLTAXISTIGHT  
PLTXLIM5 5  
PLTYLIM5 5  
PLTXLABELPREDICTION ERRORS D NERRORS  
LEGEND PLTLEGENDLOCUPPER LEFT  
LEGENDLEGENDHANDLES0SIZES 10  
LEGENDLEGENDHANDLES1SIZES 20  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0017 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
5233 NEAREST NEIGHBORS CLASSIFICATION  
SAMPLE USAGE OF NEAREST NEIGHBORS CLASSIFICATION IT WILL PLOT THE DECISION BOUNDARIES FOR EACH CLASS  
1318 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

- 523 NEAREST NEIGHBORS 1319

```
SCIKITLEARN USER GUIDE RELEASE 0213
•
PRINTDOC
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM MATPLOTLIBCOLORS IMPORT LISTEDCOLORMAP
FROM SKLEARN IMPORT NEIGHBORS DATASETS
NNEIGHBORS 15
IMPORT SOME DATA TO PLAY WITH
IRIS DATASETSLOADIRIS
WE ONLY TAKE THE FIRST TWO FEATURES WE COULD AVOID THIS UGLY
SLICING BY USING A TWODIM DATASET
X IRISDATA 2
Y IRISTARGET
H 02 STEP SIZE IN THE MESH
CREATE COLOR MAPS
CMAPLIGHT LISTEDCOLORMAPFFAAAA AAFFAA AAAAFF
CMAPBOLD LISTEDCOLORMAPFF0000 00FF00 0000FF
FORWEIGHTS INUNIFORM DISTANCE
WE CREATE AN INSTANCE OF NEIGHBOURS CLASSIFIER AND FIT THE DATA
CLF NEIGHBORSKNEIGHBORSCLASSIFIERNNEIGHBORS WEIGHTSWEIGHTS
CLFFITX Y
1320 CHAPTER 5 EXAMPLES
```

```
SCIKITLEARN USER GUIDE RELEASE 0213
PLOT THE DECISION BOUNDARY FOR THAT WE WILL ASSIGN A COLOR TO EACH
POINT IN THE MESH XMIN XMAXXYMIN YMAX
XMIN XMAX X 0MIN 1 X 0MAX 1
YMIN YMAX X 1MIN 1 X 1MAX 1
XX YY NPMESHGRIDNPARANGEXMIN XMAX H
NPARANGEYMIN YMAX H
Z CLFPREDICTNPCXXRAVEL YYRAVEL
PUT THE RESULT INTO A COLOR PLOT
Z ZRESHAPEXXSHAPE
PLTFigure
PLTPCOLORMESHXX YY Z CMAPCMAPLIGHT
PLOT ALSO THE TRAINING POINTS
PLTSCATTERX 0 X 1 CY CMAPCMAPBOLD
EDGECOLORK S20
PLTXLIMXXMIN XXMAX
PLTYLIMYYMIN YYMAX
PLTTITLE3CLASS CLASSIFICATION K I WEIGHTS S
NNEIGHBORS WEIGHTS
PLTSHOW
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 1416 SECONDS
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE
5234 NEAREST CENTROID CLASSIFICATION
SAMPLE USAGE OF NEAREST CENTROID CLASSIFICATION IT WILL PLOT THE DECISION BOUNDARIES FOR EACH CLASS
523 NEAREST NEIGHBORS 1321
```





```
SCIKITLEARN USER GUIDE RELEASE 0213
•
OUT
NONE 08133333333333334
02 082
PRINTDOC
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM MATPLOTLIBCOLORS IMPORT LISTEDCOLORMAP
FROM SKLEARN IMPORT DATASETS
FROM SKLEARNNEIGHBORS IMPORT NEARESTCENTROID
NNEIGHBORS 15
IMPORT SOME DATA TO PLAY WITH
IRIS DATASETSLOADIRIS
WE ONLY TAKE THE FIRST TWO FEATURES WE COULD AVOID THIS UGLY
SLICING BY USING A TWODIM DATASET
X IRISDATA 2
Y IRISTARGET
523 NEAREST NEIGHBORS 1323
```

SCIKITLEARN USER GUIDE RELEASE 0213

H 02 STEP SIZE IN THE MESH

CREATE COLOR MAPS

CMAPLIGHT LISTEDCOLORMAPFFAAAA AAFFAA AAAAFF

CMAPBOLD LISTEDCOLORMAPFF0000 00FF00 0000FF

FORSHRINKAGE INNONE 2

WE CREATE AN INSTANCE OF NEIGHBOURS CLASSIFIER AND FIT THE DATA

CLF NEARESTCENTROIDSHRINKTHRESHOLDSHRINKAGE

CLFFITX Y

YPRED CLFPREDICTX

PRINTSHRINKAGE NPMEANY YPRED

PLOT THE DECISION BOUNDARY FOR THAT WE WILL ASSIGN A COLOR TO EACH

POINT IN THE MESH XMIN XMAXXYMIN YMAX

XMIN XMAX X 0MIN 1 X 0MAX 1

YMIN YMAX X 1MIN 1 X 1MAX 1

XX YY NPMESHGRIDNPARANGEXMIN XMAX H

NPARANGEYMIN YMAX H

Z CLFPREDICTNPCXXRAVEL YYRAVEL

PUT THE RESULT INTO A COLOR PLOT

Z ZRESHAPEXXSHAPE

PLTFigure

PLTPCOLORMESHXX YY Z CMAPCMAPLIGHT

PLOT ALSO THE TRAINING POINTS

PLTSCATTERX 0 X 1 CY CMAPCMAPBOLD

EDGECOLORK S20

PLTTITLE3CLASS CLASSIFICATION SHRINKTHRESHOLD R

SHRINKAGE

PLTAXISTIGHT

PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0051 SECONDS

NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)

5235 KERNEL DENSITY ESTIMATION

THIS EXAMPLE SHOWS HOW KERNEL DENSITY ESTIMATION KDE A POWERFUL NONPARAMETRIC DENSITY ESTIMATION TECHNIQUE

CAN BE USED TO LEARN A GENERATIVE MODEL FOR A DATASET WITH THIS GENERATIVE MODEL IN PLACE NEW SAMPLES CAN BE DRAWN

THESE NEW SAMPLES REFLECT THE UNDERLYING MODEL OF THE DATA

1324 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
OUT
BEST BANDWIDTH 379269019073225
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARNDATASETS IMPORT LOADDIGITS
FROM SKLEARNNEIGHBORS IMPORT KERNELDENSITY
FROM SKLEARNDECOMPOSITION IMPORT PCA
FROM SKLEARNMODELSELECTION IMPORT GRIDSEARCHCV
LOAD THE DATA
DIGITS LOADDIGITS
PROJECT THE 64DIMENSIONAL DATA TO A LOWER DIMENSION
PCA PCANCOMPONENTS15 WHITENFALSE
DATA PCAFITTRANSFORMDIGITSDATA
USE GRID SEARCH CROSSVALIDATION TO OPTIMIZE THE BANDWIDTH
PARAMS BANDWIDTH NPLOGSPACE1 1 20
523 NEAREST NEIGHBORS 1325
```

SCIKITLEARN USER GUIDE RELEASE 0213  
GRID GRIDSEARCHCVKERNELDENSITY PARAMS CV5 IIDFALSE  
GRIDFITDATA  
PRINTBEST BANDWIDTH 0FORMATGRIDBESTESTIMATORBANDWIDTH  
USE THE BEST ESTIMATOR TO COMPUTE THE KERNEL DENSITY ESTIMATE  
KDE GRIDBESTESTIMATOR  
SAMPLE 44 NEW POINTS FROM THE DATA  
NEWDATA KDESAMPLE44 RANDOMSTATE0  
NEWDATA PCAINVERSESTRANSFORMNEWDATA  
TURN DATA INTO A 4X11 GRID  
NEWDATA NEWDATAARESHAPE4 11 1  
REALDATA DIGITSDATA44RESHAPE4 11 1  
PLOT REAL DIGITS AND RESAMPLED DIGITS  
FIG AX PLTSUBPLOTS9 11 SUBPLOTKWIDICTXTICKS YTICKS  
FORJINRANGE11  
AX4 JSETVISIBLEFALSE  
FORIIRANGE4  
IM AXI JIMSHOWREALDATAI JRESHAPE8 8  
CMAPPLTCMBINARY INTERPOLATIONNEAREST  
IMSETCLIM0 16  
IM AXI 5 JIMSHOWNEWDATAI JRESHAPE8 8  
CMAPPLTCMBINARY INTERPOLATIONNEAREST  
IMSETCLIM0 16  
AX0 5SETTITLESELECTION FROM THE INPUT DATA  
AX5 5SETTITLENEW DIGITS DRAWN FROM THE KERNEL DENSITY MODEL  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 4482 SECONDS  
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE  
5236 NOVELTY DETECTION WITH LOCAL OUTLIER FACTOR LOF  
THE LOCAL OUTLIER FACTOR LOF ALGORITHM IS AN UNSUPERVISED ANOMALY DETECTION METHOD WHICH COMPUTES THE LOCAL  
DENSITY DEVIATION OF A GIVEN DATA POINT WITH RESPECT TO ITS NEIGHBORS IT CONSIDERS AS OUTLIERS THE SAMPLES THAT HAVE A  
SUBSTANTIALLY LOWER DENSITY THAN THEIR NEIGHBORS THIS EXAMPLE SHOWS HOW TO USE LOF FOR NOVELTY DETECTION NOTE  
THAT WHEN LOF IS USED FOR NOVELTY DETECTION YOU MUST NOT USE PREDICT DECISIONFUNCTION AND SCORESAMPLES ON THE  
TRAINING SET AS THIS WOULD LEAD TO WRONG RESULTS YOU MUST ONLY USE THESE METHODS ON NEW UNSEEN DATA WHICH ARE NOT  
IN THE TRAINING SET SEE USER GUIDE FOR DETAILS ON THE DIFFERENCE BETWEEN OUTLIER DETECTION AND NOVELTY DETECTION AND  
HOW TO USE LOF FOR OUTLIER DETECTION  
THE NUMBER OF NEIGHBORS CONSIDERED PARAMETER NNEIGHBORS IS TYPICALLY SET 1 GREATER THAN THE MINIMUM NUMBER  
OF SAMPLES A CLUSTER HAS TO CONTAIN SO THAT OTHER SAMPLES CAN BE LOCAL OUTLIERS RELATIVE TO THIS CLUSTER AND 2 SMALLER  
THAN THE MAXIMUM NUMBER OF CLOSE BY SAMPLES THAT CAN POTENTIALLY BE LOCAL OUTLIERS IN PRACTICE SUCH INFORMATIONS  
ARE GENERALLY NOT AVAILABLE AND TAKING NNEIGHBORS20 APPEARS TO WORK WELL IN GENERAL  
1326 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
IMPORT NUMPY AS NP
IMPORT MATPLOTLIB
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARNNEIGHBORS IMPORT LOCALOUTLIERFACTOR
PRINTDOC
NPRANDOMSEED42
XX YY NPMESHGRIDNPLINSPACE5 5 500 NPLINSPACE5 5 500
GENERATE NORMAL NOT ABNORMAL TRAINING OBSERVATIONS
X 03 NPRANDOMRANDN100 2
XTRAIN NPRX 2 X 2
GENERATE NEW NORMAL NOT ABNORMAL OBSERVATIONS
X 03 NPRANDOMRANDN20 2
XTEST NPRX 2 X 2
GENERATE SOME ABNORMAL NOVEL OBSERVATIONS
XOUTLIERS NPRANDOMUNIFORMLOW4 HIGH4 SIZE20 2
FIT THE MODEL FOR NOVELTY DETECTION NOVELTYTRUE
CLF LOCALOUTLIERFACTORNNEIGHBORS20 NOVELTYTRUE CONTAMINATION01
CLFFITXTRAIN
DO NOT USE PREDICT DECISIONFUNCTION AND SCORESAMPLES ON XTRAIN AS THIS
WOULD GIVE WRONG RESULTS BUT ONLY ON NEW UNSEEN DATA NOT USED IN XTRAIN
EG XTEST XOUTLIERS OR THE MESHGRID
YPREDTEST CLFPREDICTXTEST
523 NEAREST NEIGHBORS 1327
```

SCIKITLEARN USER GUIDE RELEASE 0213  
YPREDOUTLIERS CLFPREDICTXOUTLIERS  
NERRORTEST YPREDTESTYPREDTEST 1SIZE  
NERROROUTLIERS YPREDOUTLIERSYPREDOUTLIERS 1SIZE  
PLOT THE LEARNED FRONTIER THE POINTS AND THE NEAREST VECTORS TO THE PLANE  
Z CLFDECISIONFUNCTIONNPCXXRAVEL YYRAVEL  
Z ZRESHAPEXXSHAPE  
PLTTITLENOVELTY DETECTION WITH LOF  
PLTCONTOURFXX YY Z LEVELSNPLINSPACEZMIN 0 7 CMAPPLTCMPUBU  
A PLTCONTOURXX YY Z LEVELS0 LINEWIDTHS2 COLORSDARKRED  
PLTCONTOURFXX YY Z LEVELS0 ZMAX COLORSPALEVIOLETRED  
S 40  
B1 PLTSCATTERXTRAIN 0 XTRAIN 1 CWHITE SS EDGECOLORSK  
B2 PLTSCATTERXTEST 0 XTEST 1 CBLUEVIOLET SS  
EDGECOLORSK  
C PLTSCATTERXOUTLIERS 0 XOUTLIERS 1 CGOLD SS  
EDGECOLORSK  
PLTAXISTIGHT  
PLTXLIM5 5  
PLTYLIM5 5  
PLTLEGENDACOLLECTIONS0 B1 B2 C  
LEARNED FRONTIER TRAINING OBSERVATIONS  
NEW REGULAR OBSERVATIONS NEW ABNORMAL OBSERVATIONS  
LOCUPPER LEFT  
PROPMATPLOTLIBFONTMANAGERFONTPROPERTIESSIZE11  
PLTXLABEL  
ERRORS NOVEL REGULAR D40 ERRORS NOVEL ABNORMAL D40  
NERRORTEST NERROROUTLIERS  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0634 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
5237 COMPARING NEAREST NEIGHBORS WITH AND WITHOUT NEIGHBORHOOD COMPONENTS  
ANALYSIS  
AN EXAMPLE COMPARING NEAREST NEIGHBORS CLASSIFICATION WITH AND WITHOUT NEIGHBORHOOD COMPONENTS ANALYSIS  
IT WILL PLOT THE CLASS DECISION BOUNDARIES GIVEN BY A NEAREST NEIGHBORS CLASSIFIER WHEN USING THE EUCLIDEAN DISTANCE  
ON THE ORIGINAL FEATURES VERSUS USING THE EUCLIDEAN DISTANCE AFTER THE TRANSFORMATION LEARNED BY NEIGHBORHOOD COM  
PONENTS ANALYSIS THE LATTER AIMS TO FIND A LINEAR TRANSFORMATION THAT MAXIMISES THE STOCHASTIC NEAREST NEIGHBOR  
CLASSIFICATION ACCURACY ON THE TRAINING SET  
1328 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

- 523 NEAREST NEIGHBORS 1329

SCIKITLEARN USER GUIDE RELEASE 0213

```
•
LICENSE BSD 3 CLAUSE
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM MATPLOTLIBCOLORS IMPORT LISTEDCOLORMAP
FROM SKLEARN IMPORT DATASETS
FROM SKLEARNMODELSELECTION IMPORT TRAINTESTSPLIT
FROM SKLEARNPREPROCESSING IMPORT STANDARDSCALER
FROM SKLEARNNEIGHBORS IMPORT KNEIGHBORSCLASSIFIER
NEIGHBORHOODCOMPONENTSANALYSIS
FROM SKLEARNPIPELINE IMPORT PIPELINE
PRINTDOC
NNEIGHBORS 1
DATASET DATASETSLOADIRIS
X Y DATASETDATA DATASETTARGET
WE ONLY TAKE TWO FEATURES WE COULD AVOID THIS UGLY
SLICING BY USING A TWODIM DATASET
X X 0 2
XTRAIN XTEST YTRAIN YTEST
TRAINTESTSPLITX Y STRATIFY TESTSIZE07 RANDOMSTATE42
1330 CHAPTER 5 EXAMPLES
```



SCIKITLEARN USER GUIDE RELEASE 0213

H 01 STEP SIZE IN THE MESH

CREATE COLOR MAPS

CMAPLIGHT LISTEDCOLORMAPFFAAAA AAFFAA AAAAFF

CMAPBOLD LISTEDCOLORMAPFF0000 00FF00 0000FF

NAMES KNN NCA KNN

CLASSIFIERS PIPELINESCALER STANDARDSCALER

KNN KNEIGHBORSCLASSIFIERNNEIGHBORSNNEIGHBORS

PIPELINESCALER STANDARDSCALER

NCA NEIGHBORHOODCOMPONENTSANALYSIS

KNN KNEIGHBORSCLASSIFIERNNEIGHBORSNNEIGHBORS

XMIN XMAX X 0MIN 1 X 0MAX 1

YMIN YMAX X 1MIN 1 X 1MAX 1

XX YY NPMESHGRIDNPARANGEXMIN XMAX H

NPARANGYMIN YMAX H

FORNAME CLF INZIPNAMES CLASSIFIERS

CLFFITXTRAIN YTRAIN

SCORE CLFSCOREXTEST YTEST

PLOT THE DECISION BOUNDARY FOR THAT WE WILL ASSIGN A COLOR TO EACH

POINT IN THE MESH XMIN XMAXXYMIN YMAX

Z CLFPREDICTNPCXXRAVEL YYRAVEL

PUT THE RESULT INTO A COLOR PLOT

Z ZRESHAPEXXSHAPE

PLTFigure

PLTPCOLORMESHXX YY Z CMAPCMAPLIGHT ALPHA8

PLOT ALSO THE TRAINING AND TESTING POINTS

PLTSCATTERX 0 X 1 CY CMAPCMAPBOLD EDGECOLORK S20

PLTXLIMXXMIN XXMAX

PLTYLIMYYMIN YYMAX

PLTTITLE K FORMATNAME NNEIGHBORS

PLTTEXT09 01 2FFORMATSCORE SIZE15

HACENTER VACENTER TRANSFORMPLTGCATRANSAXES

PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 16006 SECONDS

NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE

5238 DIMENSIONALITY REDUCTION WITH NEIGHBORHOOD COMPONENTS ANALYSIS

SAMPLE USAGE OF NEIGHBORHOOD COMPONENTS ANALYSIS FOR DIMENSIONALITY REDUCTION

THIS EXAMPLE COMPARES DIFFERENT LINEAR DIMENSIONALITY REDUCTION METHODS APPLIED ON THE DIGITS DATA SET THE DATA

523 NEAREST NEIGHBORS 1331

SCIKITLEARN USER GUIDE RELEASE 0213

SET CONTAINS IMAGES OF DIGITS FROM 0 TO 9 WITH APPROXIMATELY 180 SAMPLES OF EACH CLASS EACH IMAGE IS OF DIMENSION 8X8 64 AND IS REDUCED TO A TWODIMENSIONAL DATA POINT

PRINCIPAL COMPONENT ANALYSIS PCA APPLIED TO THIS DATA IDENTIFIES THE COMBINATION OF ATTRIBUTES PRINCIPAL COMPONENTS OR DIRECTIONS IN THE FEATURE SPACE THAT ACCOUNT FOR THE MOST VARIANCE IN THE DATA HERE WE PLOT THE DIFFERENT SAMPLES ON THE 2 FIRST PRINCIPAL COMPONENTS

LINEAR DISCRIMINANT ANALYSIS LDA TRIES TO IDENTIFY ATTRIBUTES THAT ACCOUNT FOR THE MOST VARIANCE BETWEEN CLASSES IN PARTICULAR LDA IN CONTRAST TO PCA IS A SUPERVISED METHOD USING KNOWN CLASS LABELS

NEIGHBORHOOD COMPONENTS ANALYSIS NCA TRIES TO FIND A FEATURE SPACE SUCH THAT A STOCHASTIC NEAREST NEIGHBOR ALGORITHM WILL GIVE THE BEST ACCURACY LIKE LDA IT IS A SUPERVISED METHOD

ONE CAN SEE THAT NCA ENFORCES A CLUSTERING OF THE DATA THAT IS VISUALLY MEANINGFUL DESPITE THE LARGE REDUCTION IN DIMENSION

- 

1332 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

- 523 NEAREST NEIGHBORS 1333

SCIKITLEARN USER GUIDE RELEASE 0213

```
•
LICENSE BSD 3 CLAUSE
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARN IMPORT DATASETS
FROM SKLEARNMODELSELECTION IMPORT TRAINTESTSPLIT
FROM SKLEARNDECOMPOSITION IMPORT PCA
FROM SKLEARNDISCRIMINANTANALYSIS IMPORT LINEARDISCRIMINANTANALYSIS
FROM SKLEARNNEIGHBORS IMPORT KNEIGHBORSCLASSIFIER
NEIGHBORHOODCOMPONENTSANALYSIS
FROM SKLEARNPIPELINE IMPORT MAKEPIPELINE
FROM SKLEARNPREPROCESSING IMPORT STANDARDSCALER
PRINTDOC
NNEIGHBORS 3
RANDOMSTATE 0
LOAD DIGITS DATASET
DIGITS DATASETSLOADDIGITS
X Y DIGITSDATA DIGITSTARGET
SPLIT INTO TRAINTEST
XTRAIN XTEST YTRAIN YTEST
TRAINTESTSPLITX Y TESTSIZE05 STRATIFY
RANDOMSTATERRANDOMSTATE
1334 CHAPTER 5 EXAMPLES
```

```
SCIKITLEARN USER GUIDE RELEASE 0213
DIM LENX0
NCLASSES LENNPUNIQUEY
REDUCE DIMENSION TO 2 WITH PCA
PCA MAKEPIPELINESTANDARDSCALER
PCANCOMPONENTS2 RANDOMSTATERRANDOMSTATE
REDUCE DIMENSION TO 2 WITH LINEARDISCRIMINANTANALYSIS
LDA MAKEPIPELINESTANDARDSCALER
LINEARDISCRIMINANTANALYSISNCOMPONENTS2
REDUCE DIMENSION TO 2 WITH NEIGHBORHOODCOMPONENTANALYSIS
NCA MAKEPIPELINESTANDARDSCALER
NEIGHBORHOODCOMPONENTSANALYSISNCOMPONENTS2
RANDOMSTATERRANDOMSTATE
USE A NEAREST NEIGHBOR CLASSIFIER TO EVALUATE THE METHODS
KNN KNEIGHBORSCLASSIFIERNNEIGHBORSNNEIGHBORS
MAKE A LIST OF THE METHODS TO BE COMPARED
DIMREDUCTIONMETHODS PCA PCA LDA LDA NCA NCA
PLTFigure
FORI NAME MODEL INENUMERATEDDIMREDUCTIONMETHODS
PLTFigure
PLTSUBPLOT1 3 I 1 ASPECT1
FIT THE METHODS MODEL
MODELFITXTRAIN YTRAIN
FIT A NEAREST NEIGHBOR CLASSIFIER ON THE EMBEDDED TRAINING SET
KNNFITMODELTRANSFORMXTRAIN YTRAIN
COMPUTE THE NEAREST NEIGHBOR ACCURACY ON THE EMBEDDED TEST SET
ACCKNN KNNSCOREMODELTRANSFORMXTEST YTEST
EMBED THE DATA SET IN 2 DIMENSIONS USING THE FITTED MODEL
XEMBEDDED MODELTRANSFORMX
PLOT THE PROJECTED POINTS AND SHOW THE EVALUATION SCORE
PLTSCATTERXEMBEDDED 0 XEMBEDDED 1 CY S30 CMAPSET1
PLTTITLE KNN K NTEST ACCURACY 2FFORMATNAME
NNEIGHBORS
ACCKNN
PLTSHOW
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 2879 SECONDS
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE
5239 KERNEL DENSITY ESTIMATE OF SPECIES DISTRIBUTIONS
THIS SHOWS AN EXAMPLE OF A NEIGHBORSBASED QUERY IN PARTICULAR A KERNEL DENSITY ESTIMATE ON GEOSPATIAL DATA USING
A BALL TREE BUILT UPON THE HAVERSINE DISTANCE METRIC - IE DISTANCES OVER POINTS IN LATITUDELONGITUDE THE DATASET
523 NEAREST NEIGHBORS 1335
```

SCIKITLEARN USER GUIDE RELEASE 0213

IS PROVIDED BY PHILLIPS ET AL 2006 IF AVAILABLE THE EXAMPLE USES BASEMAP TO PLOT THE COAST LINES AND NATIONAL BOUNDARIES OF SOUTH AMERICA

THIS EXAMPLE DOES NOT PERFORM ANY LEARNING OVER THE DATA SEE SPECIES DISTRIBUTION MODELING FOR AN EXAMPLE OF CLASSIFICATION BASED ON THE ATTRIBUTES IN THIS DATASET IT SIMPLY SHOWS THE KERNEL DENSITY ESTIMATE OF OBSERVED DATA POINTS IN GEOSPATIAL COORDINATES

THE TWO SPECIES ARE

- “BRADYPUS VARIEGATUS” THE BROWNTHOATED SLOTH
- “MICRORYZOMYS MINUTUS” ALSO KNOWN AS THE FOREST SMALL RICE RAT A RODENT THAT LIVES IN PERU COLOMBIA ECUADOR PERU AND VENEZUELA

REFERENCES

- “MAXIMUM ENTROPY MODELING OF SPECIES GEOGRAPHIC DISTRIBUTIONS” S J PHILLIPS R P ANDERSON R E SCHAPIRE ECOLOGICAL MODELLING 190231259 2006

OUT

COMPUTING KDE IN SPHERICAL COORDINATES

PLOT COASTLINES FROM COVERAGE

COMPUTING KDE IN SPHERICAL COORDINATES

PLOT COASTLINES FROM COVERAGE

1336 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
AUTHOR JAKE VANDERPLAS JAKEVDPCSWASHINGTONEDU

```
LICENSE BSD 3 CLAUSE
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARNDATASETS IMPORT FETCHSPECIESDISTRIBUTIONS
FROM SKLEARNDATASETSSPECIESDISTRIBUTIONS IMPORT CONSTRUCTGRIDS
FROM SKLEARNNEIGHBORS IMPORT KERNELDENSITY
IF BASEMAP IS AVAILABLE WELL USE IT
OTHERWISE WELL IMPROVISE LATER
TRY
FROM MPLTOOLKITSBASEMAP IMPORT BASEMAP
BASEMAP TRUE
EXCEPTIMPORTERROR
BASEMAP FALSE
GET MATRICESARRAYS OF SPECIES IDS AND LOCATIONS
DATA FETCHSPECIESDISTRIBUTIONS
SPECIESNAMES BRADYPUS VARIEGATUS MICRORYZOMYS MINUTUS
XTRAIN NPVSTACKDATATRAINDD LAT
DATATRAINDD LONGT
YTRAIN NPARRAYDDECODEASCIISTARTSWITHMICRO
FORDINDATATRAINSPECIES DTYPEINT
XTRAIN NPPI 180 CONVERT LATLONG TO RADIANS
SET UP THE DATA GRID FOR THE CONTOUR PLOT
XGRID YGRID CONSTRUCTGRIDSDATA
X Y NPMESHGRIDXGRID5 YGRID51
LANDREFERENCE DATACOVERAGES65 5
LANDMASK LANDREFERENCE 9999RAVEL
XY NPVSTACKYRAVEL XRAVELT
XY XYLANDMASK
XY NPPI 180
PLOT MAP OF SOUTH AMERICA WITH DISTRIBUTIONS OF EACH SPECIES
FIG PLTFigure
FIGSUBPLOTSADJUSTLEFT005 RIGHT095 WSPACE005
FORIINRANGE2
PLTSUBPLOT1 2 I 1
CONSTRUCT A KERNEL DENSITY ESTIMATE OF THE DISTRIBUTION
PRINT COMPUTING KDE IN SPHERICAL COORDINATES
KDE KERNELDENSITYBANDWIDTH004 METRICHAVERSINE
KERNELGAUSSIAN ALGORITHMBALLTREE
KDEFITXTRAINYTRAIN I
EVALUATE ONLY ON THE LAND 9999 INDICATES OCEAN
Z NPFULLLANDMASKSHAPE0 9999 DTYPEINT
ZLANDMASK NPEXPKDESCORESAMPLESXY
Z ZRESHAPEXSHAPE
523 NEAREST NEIGHBORS 1337
```

SCIKITLEARN USER GUIDE RELEASE 0213  
PLOT CONTOURS OF THE DENSITY  
LEVELS NPLINSPACE0 ZMAX 25  
PLTCONTOURFX Y Z LEVELSLEVELS CMAPPLTCMREDS  
IFBASEMAP  
PRINT PLOT COASTLINES USING BASEMAP  
M BASEMAPPROJECTIONCYL LLCRNRLATYMIN  
URCRNRLATYMAX LLCRNRLONXMIN  
URCRNRLONXMAX RESOLUTIONC  
MDRAWCOASTLINES  
MDRAWCOUNTRIES  
ELSE  
PRINT PLOT COASTLINES FROM COVERAGE  
PLTCONTOURX Y LANDREFERENCE  
LEVELS9998 COLORSK  
LINESTYLESSOLID  
PLTXTICKS  
PLTYTICKS  
PLTTITLESPECIESNAMESI  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 5470 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
52310 NEIGHBORHOOD COMPONENTS ANALYSIS ILLUSTRATION  
AN EXAMPLE ILLUSTRATING THE GOAL OF LEARNING A DISTANCE METRIC THAT MAXIMIZES THE NEAREST NEIGHBORS CLASSIFICATION  
ACCURACY THE EXAMPLE IS SOLELY FOR ILLUSTRATION PURPOSES PLEASE REFER TO THE USER GUIDE FOR MORE INFORMATION  
1338 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213

- 523 NEAREST NEIGHBORS 1339

SCIKITLEARN USER GUIDE RELEASE 0213

```
•
LICENSE BSD 3 CLAUSE
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARNDATASETS IMPORT MAKECLASSIFICATION
FROM SKLEARNNEIGHBORS IMPORT NEIGHBORHOODCOMPONENTSANALYSIS
FROM MATPLOTLIB IMPORT CM
FROM SKLEARNUTILSFIXES IMPORT LOGSUMEXP
PRINTDOC
NNEIGHBORS 1
RANDOMSTATE 0
CREATE A TINY DATA SET OF 9 SAMPLES FROM 3 CLASSES
X Y MAKECLASSIFICATIONNNSAMPLES9 NFEATURES2 NINFORMATIVE2
NREDUNDANT0 NCLASSES3 NCLUSTERSPERCLASS1
CLASSEP10 RANDOMSTATERANDOMSTATE
PLOT THE POINTS IN THE ORIGINAL SPACE
PLTFigure
AX PLTGCA
DRAW THE GRAPH NODES
FORIINRANGEXSHAPE0
AXTEXTXI 0 XI 1 STRI VACENTER HACENTER
AXSCATTERXI 0 XI 1 S300 CCMSET1YI ALPHA04
1340 CHAPTER 5 EXAMPLES
```

```
SCIKITLEARN USER GUIDE RELEASE 0213
DEFPIX I
DIFFEMBEDDED XI X
DISTEMBEDDED NPEINSUMIJJI DIFFEMBEDDED
DIFFEMBEDDED
DISTEMBEDDEDI NPINF
  COMPUTE EXPONENTIATED DISTANCES USE THE LOGSUMEXP TRICK TO
  AVOID NUMERICAL INSTABILITIES
EXPDISTEMBEDDED NPEXPDISTEMBEDDED
LOGSUMEXPDISTEMBEDDED
RETURNEXPDISTEMBEDDED
DEFRELATEPOINTX I AX
PTI XI
FORJ PTJ INENUMERATEX
THICKNESS PIX I
IFI J
LINE PTI0 PTJ0 PTI1 PTJ1
AXPLOT LINE CCMSET1YJ
LINEWIDTH5 THICKNESSJ
  WE CONSIDER ONLY POINT 3
I 3
  PLOT BONDS LINKED TO SAMPLE I IN THE ORIGINAL SPACE
RELATEPOINTX I AX
AXSETTITLEORIGINAL POINTS
AXAXESGETXAXISSETVISIBLEFALSE
AXAXESGETYAXISSETVISIBLEFALSE
AXAXISEQUAL
  LEARN AN EMBEDDING WITH NEIGHBORHOODCOMPONENTSANALYSIS
NCA NEIGHBORHOODCOMPONENTSANALYSISMAXITER30 RANDOMSTATERANDOMSTATE
NCA NCAFITX Y
  PLOT THE POINTS AFTER TRANSFORMATION WITH NEIGHBORHOODCOMPONENTSANALYSIS
PLTFigure
AX2 PLTGCA
  GET THE EMBEDDING AND FIND THE NEW NEAREST NEIGHBORS
XEMBEDDED NCATransformX
RELATEPOINTXEMBEDDED I AX2
FORIINRANGELENX
AX2TEXTXEMBEDDEDI 0 XEMBEDDEDI 1 STRI
VACENTER HACENTER
AX2SCATTERXEMBEDDEDI 0 XEMBEDDEDI 1 S300 CCMSET1YI
ALPHA04
  MAKE AXES EQUAL SO THAT BOUNDARIES ARE DISPLAYED CORRECTLY AS CIRCLES
AX2SETTITLENCA EMBEDDING
AX2AXESGETXAXISSETVISIBLEFALSE
AX2AXESGETYAXISSETVISIBLEFALSE
523 NEAREST NEIGHBORS 1341
```

SCIKITLEARN USER GUIDE RELEASE 0213

AX2AXISEQUAL

PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0069 SECONDS

NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)

52311 SIMPLE 1D KERNEL DENSITY ESTIMATION

THIS EXAMPLE USES THE SKLEARNNEIGHBORSKERNELDENSITY CLASS TO DEMONSTRATE THE PRINCIPLES OF KERNEL DENSITY ESTIMATION IN ONE DIMENSION

THE FIRST PLOT SHOWS ONE OF THE PROBLEMS WITH USING HISTOGRAMS TO VISUALIZE THE DENSITY OF POINTS IN 1D INTUITIVELY A HISTOGRAM CAN BE THOUGHT OF AS A SCHEME IN WHICH A UNIT “BLOCK” IS STACKED ABOVE EACH POINT ON A REGULAR GRID AS THE TOP TWO PANELS SHOW HOWEVER THE CHOICE OF GRIDDING FOR THESE BLOCKS CAN LEAD TO WILDLY DIVERGENT IDEAS ABOUT THE UNDERLYING SHAPE OF THE DENSITY DISTRIBUTION IF WE INSTEAD CENTER EACH BLOCK ON THE POINT IT REPRESENTS WE GET THE ESTIMATE SHOWN IN THE BOTTOM LEFT PANEL THIS IS A KERNEL DENSITY ESTIMATION WITH A “TOP HAT” KERNEL THIS IDEA CAN BE GENERALIZED TO OTHER KERNEL SHAPES THE BOTTOMRIGHT PANEL OF THE FIRST FIGURE SHOWS A GAUSSIAN KERNEL DENSITY ESTIMATE OVER THE SAME DISTRIBUTION

SCIKITLEARN IMPLEMENTS EFFICIENT KERNEL DENSITY ESTIMATION USING EITHER A BALL TREE OR KD TREE STRUCTURE THROUGH THE SKLEARNNEIGHBORSKERNELDENSITY ESTIMATOR THE AVAILABLE KERNELS ARE SHOWN IN THE SECOND FIGURE OF THIS EXAMPLE

THE THIRD FIGURE COMPARES KERNEL DENSITY ESTIMATES FOR A DISTRIBUTION OF 100 SAMPLES IN 1 DIMENSION THOUGH THIS EXAMPLE USES 1D DISTRIBUTIONS KERNEL DENSITY ESTIMATION IS EASILY AND EFFICIENTLY EXTENSIBLE TO HIGHER DIMENSIONS AS WELL

1342 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

- 523 NEAREST NEIGHBORS 1343



SCIKITLEARN USER GUIDE RELEASE 0213

- AUTHOR JAKE VANDERPLAS JAKEVDPCSWASHINGTONEDU

```
IMPORT NUMPY AS NP
IMPORT MATPLOTLIB
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM DISTUTILSVERSION IMPORT LOOSEVERSION
FROM SCIPYSTATS IMPORT NORM
FROM SKLEARNNEIGHBORS IMPORT KERNELDENSITY
NORMED IS BEING DEPRECATED IN FAVOR OF DENSITY IN HISTOGRAMS
IFLOOSEVERSIONMATPLOTLIBVERSION 21
DENSITYPARAM DENSITY TRUE
ELSE
DENSITYPARAM NORMED TRUE
```

```

PLOT THE PROGRESSION OF HISTOGRAMS TO KERNELS
NPRANDOMSEED1
N 20
X NPCONCATENATENPRANDOMNORMAL0 1 INT03 N
NPRANDOMNORMAL5 1 INT07 N NPNEWAXIS
XPLOT NPLINSPACE5 10 1000 NPNEWAXIS
BINS NPLINSPACE5 10 10
FIG AX PLTSUBPLOTS2 2 SHAREXTRUE SHAREYTRUE
FIGSUBPLOTSADJUSTHSPACE005 WSPACE005
523 NEAREST NEIGHBORS 1345
```

SCIKITLEARN USER GUIDE RELEASE 0213

HISTOGRAM 1

AX0 0HISTX 0 BINSBINS FCAAAAFF DENSITYPARAM

AX0 0TEXT35 031 HISTOGRAM

HISTOGRAM 2

AX0 1HISTX 0 BINSBINS 075 FCAAAAFF DENSITYPARAM

AX0 1TEXT35 031 HISTOGRAM BINS SHIFTED

TOPHAT KDE

KDE KERNELDENSITYKERNELTOPHAT BANDWIDTH075FITX

LOGDENS KDESCORESAMPLESXPLOT

AX1 0FILLXPLOT 0 NPEXPLOGDENS FCAAAAFF

AX1 0TEXT35 031 TOPHAT KERNEL DENSITY

GAUSSIAN KDE

KDE KERNELDENSITYKERNELGAUSSIAN BANDWIDTH075FITX

LOGDENS KDESCORESAMPLESXPLOT

AX1 1FILLXPLOT 0 NPEXPLOGDENS FCAAAAFF

AX1 1TEXT35 031 GAUSSIAN KERNEL DENSITY

FORAXIINAXRAVEL

AXIPLOTX 0 NPFULLXSHAPE0 001 K

AXISETXLIM4 9

AXISETYLIM002 034

FORAXIINAX 0

AXISETYLABELNORMALIZED DENSITY

FORAXIINAX1

AXISETXLABELX

PLOT ALL AVAILABLE KERNELS

XPLOT NPLINSPACE6 6 1000 NONE

XSRC NPZEROS1 1

FIG AX PLTSUBPLOTS2 3 SHAREXTRUE SHAREYTRUE

FIGSUBPLOTSADJUSTLEFT005 RIGHT095 HSPACE005 WSPACE005

DEFFORMATFUNCX LOC

IFX 0

RETURN0

ELIFX 1

RETURNH

ELIFX 1

RETURNH

ELSE

RETURNH X

FORI KERNEL INENUMERATEGAUSSIAN TOPHAT EPANECHNIKOV

EXPONENTIAL LINEAR COSINE

AXI AXRAVELI

LOGDENS KERNELDENSITYKERNELKERNELFITXSRCSCORESAMPLESXPLOT

AXIFILLXPLOT 0 NPEXPLOGDENS K FCAAAAFF

AXITEXT26 095 KERNEL

AXIXAXISSETMAJORFORMATTERPLTFUNCFORMATTERFORMATFUNC

1346 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213  
AXIXAXISSETMAJORLOCATORPLTMULTIPLELOCATOR1  
AXIYAXISSETMAJORLOCATORPLTNULLOCALATOR  
AXISETYLIM0 105  
AXISETXLIM29 29  
AX0 1SETTITLEAVAILABLE KERNELS

PLOT A 1D DENSITY EXAMPLE  
N 100  
NPRANDOMSEED1  
X NPCONCATENATENPRANDOMNORMAL0 1 INT03 N  
NPRANDOMNORMAL5 1 INT07 N NPNEWAXIS  
XPLOT NPLINSPACES 10 1000 NPNEWAXIS  
TRUEDENS 03 NORM0 1PDFXPLOT 0  
07NORM5 1PDFXPLOT 0  
FIG AX PLTSUBPLOTS  
AXFILLXPLOT 0 TRUEDENS FCBLACK ALPHA02  
LABELINPUT DISTRIBUTION  
FORKERNELINGAUSSIAN TOPHAT EPANECHNIKOV  
KDE KERNELDENSITYKERNELKERNEL BANDWIDTH05FITX  
LOGDENS KDESCORESAMPLESXPLO  
AXPLOTXPLOT 0 NPEXPLOGDENS  
LABELKERNEL 0FORMATKERNEL  
AXTEXT6 038 N0 POINTSFORMATN  
AXLEGENDLOCUPPER LEFT  
AXPLOTX 0 0005 001 NPRANDOMRANDOMXSHAPE0 K  
AXSETXLIM4 9  
AXSETYLIM002 04  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0145 SECONDS  
524 NEURAL NETWORKS  
EXAMPLES CONCERNING THE SKLEARNNEURALNETWORK MODULE  
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE  
5241 VISUALIZATION OF MLP WEIGHTS ON MNIST  
SOMETIMES LOOKING AT THE LEARNED COEFFICIENTS OF A NEURAL NETWORK CAN PROVIDE INSIGHT INTO THE LEARNING BEHAVIOR FOR  
EXAMPLE IF WEIGHTS LOOK UNSTRUCTURED MAYBE SOME WERE NOT USED AT ALL OR IF VERY LARGE COEFFICIENTS EXIST MAYBE  
REGULARIZATION WAS TOO LOW OR THE LEARNING RATE TOO HIGH  
524 NEURAL NETWORKS 1347

SCIKITLEARN USER GUIDE RELEASE 0213

THIS EXAMPLE SHOWS HOW TO PLOT SOME OF THE FIRST LAYER WEIGHTS IN A MLPCLASSIFIER TRAINED ON THE MNIST DATASET

THE INPUT DATA CONSISTS OF 28X28 PIXEL HANDWRITTEN DIGITS LEADING TO 784 FEATURES IN THE DATASET THEREFORE THE FIRST LAYER WEIGHT MATRIX HAVE THE SHAPE 784 HIDDENLAYERSIZES0 WE CAN THEREFORE VISUALIZE A SINGLE COLUMN OF THE WEIGHT MATRIX AS A 28X28 PIXEL IMAGE

TO MAKE THE EXAMPLE RUN FASTER WE USE VERY FEW HIDDEN UNITS AND TRAIN ONLY FOR A VERY SHORT TIME TRAINING LONGER WOULD RESULT IN WEIGHTS WITH A MUCH SMOOTHER SPATIAL APPEARANCE

OUT

ITERATION 1 LOSS 032009978

ITERATION 2 LOSS 015347534

ITERATION 3 LOSS 011544755

ITERATION 4 LOSS 009279764

ITERATION 5 LOSS 007889367

ITERATION 6 LOSS 007170497

ITERATION 7 LOSS 006282111

ITERATION 8 LOSS 005530788

ITERATION 9 LOSS 004960484

ITERATION 10 LOSS 004645355

TRAINING SET SCORE 0986800

TEST SET SCORE 0970000

1348 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARNDATASETS IMPORT FETCHOPENML
FROM SKLEARNNEURALNETWORK IMPORT MLPCLASSIFIER
PRINTDOC
LOAD DATA FROM HTTPSWWWOPENMLORG554
X Y FETCHOPENMLMNIST784 VERSION1 RETURNXYTRUE
X X 255
RESCALE THE DATA USE THE TRADITIONAL TRAINTEST SPLIT
XTRAIN XTEST X60000 X60000
YTRAIN YTEST Y60000 Y60000
MLP MLPCLASSIFIERHIDDENLAYERSIZES100 100 MAXITER400 ALPHA1E4
SOLVERSGD VERBOSE10 TOL1E4 RANDOMSTATE1
MLP MLPCLASSIFIERHIDDENLAYERSIZES50 MAXITER10 ALPHA1E4
SOLVERSGD VERBOSE10 TOL1E4 RANDOMSTATE1
LEARNINGRATEINIT1
MLPFITXTRAIN YTRAIN
PRINTTRAINING SET SCORE F MLPSCOREXTRAIN YTRAIN
PRINTTEST SET SCORE F MLPSCOREXTEST YTEST
FIG AXES PLTSUBPLOTS4 4
USE GLOBAL MIN MAX TO ENSURE ALL WEIGHTS ARE SHOWN ON THE SAME SCALE
VMIN VMAX MLPCOEFS0MIN MLPCOEFS0MAX
FORCOEF AX INZIPMLPCOEFS0T AXESRAVEL
AXMATSHOWCOEFRESHAPE28 28 CMAPPLTCMGRAY VMIN5 VMIN
VMAX5 VMAX
AXSETXTICKS
AXSETYTTICKS
PLTSHOW
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 27631 SECONDS
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE
5242 RESTRICTED BOLTZMANN MACHINE FEATURES FOR DIGIT CLASSIFICATION
FOR GREYSCALE IMAGE DATA WHERE PIXEL VALUES CAN BE INTERPRETED AS DEGREES OF BLACKNESS ON A WHITE BACKGROUND LIKE
HANDWRITTEN DIGIT RECOGNITION THE BERNOULLI RESTRICTED BOLTZMANN MACHINE MODEL BERNOULLIRBM CAN PERFORM
EFFECTIVE NONLINEAR FEATURE EXTRACTION
IN ORDER TO LEARN GOOD LATENT REPRESENTATIONS FROM A SMALL DATASET WE ARTIFICIALLY GENERATE MORE LABELED DATA BY PER
TURBING THE TRAINING DATA WITH LINEAR SHIFTS OF 1 PIXEL IN EACH DIRECTION
THIS EXAMPLE SHOWS HOW TO BUILD A CLASSIFICATION PIPELINE WITH A BERNOULLIRBM FEATURE EXTRACTOR AND A
LOGISTICREGRESSION CLASSIFIER THE HYPERPARAMETERS OF THE ENTIRE MODEL LEARNING RATE HIDDEN LAYER SIZE REGU
LARIZATION WERE OPTIMIZED BY GRID SEARCH BUT THE SEARCH IS NOT REPRODUCED HERE BECAUSE OF RUNTIME CONSTRAINTS
LOGISTIC REGRESSION ON RAW PIXEL VALUES IS PRESENTED FOR COMPARISON THE EXAMPLE SHOWS THAT THE FEATURES EXTRACTED BY
THE BERNOULLIRBM HELP IMPROVE THE CLASSIFICATION ACCURACY
524 NEURAL NETWORKS 1349
```

SCIKITLEARN USER GUIDE RELEASE 0213

OUT

BERNOULLIRBM ITERATION 1 PSEUDOLIKELIHOOD	2539	TIME	015S
BERNOULLIRBM ITERATION 2 PSEUDOLIKELIHOOD	2372	TIME	034S
BERNOULLIRBM ITERATION 3 PSEUDOLIKELIHOOD	2272	TIME	033S
BERNOULLIRBM ITERATION 4 PSEUDOLIKELIHOOD	2186	TIME	032S
BERNOULLIRBM ITERATION 5 PSEUDOLIKELIHOOD	2166	TIME	035S
BERNOULLIRBM ITERATION 6 PSEUDOLIKELIHOOD	2100	TIME	037S
BERNOULLIRBM ITERATION 7 PSEUDOLIKELIHOOD	2075	TIME	035S
BERNOULLIRBM ITERATION 8 PSEUDOLIKELIHOOD	2052	TIME	034S
BERNOULLIRBM ITERATION 9 PSEUDOLIKELIHOOD	2038	TIME	031S
BERNOULLIRBM ITERATION 10 PSEUDOLIKELIHOOD	2023	TIME	035S
BERNOULLIRBM ITERATION 11 PSEUDOLIKELIHOOD	2002	TIME	034S
BERNOULLIRBM ITERATION 12 PSEUDOLIKELIHOOD	1993	TIME	035S
BERNOULLIRBM ITERATION 13 PSEUDOLIKELIHOOD	1971	TIME	033S
BERNOULLIRBM ITERATION 14 PSEUDOLIKELIHOOD	1969	TIME	033S
BERNOULLIRBM ITERATION 15 PSEUDOLIKELIHOOD	1961	TIME	033S
BERNOULLIRBM ITERATION 16 PSEUDOLIKELIHOOD	1957	TIME	032S
BERNOULLIRBM ITERATION 17 PSEUDOLIKELIHOOD	1936	TIME	033S
BERNOULLIRBM ITERATION 18 PSEUDOLIKELIHOOD	1922	TIME	036S
BERNOULLIRBM ITERATION 19 PSEUDOLIKELIHOOD	1931	TIME	035S
BERNOULLIRBM ITERATION 20 PSEUDOLIKELIHOOD	1921	TIME	034S

LOGISTIC REGRESSION USING RBM FEATURES

PRECISION RECALL F1SCORE SUPPORT

0 099 098 099 174

1 094 095 094 184

2 089 096 092 166

3 093 088 090 194

4 096 094 095 186

5 092 090 091 181

6 098 098 098 207

1350 CHAPTER 5 EXAMPLES

```

SCIKITLEARN USER GUIDE RELEASE 0213
7 092 099 096 154
8 087 084 086 182
9 091 091 091 169
ACCURACY 093 1797
MACRO AVG 093 093 093 1797
WEIGHTED AVG 093 093 093 1797
LOGISTIC REGRESSION USING RAW PIXEL FEATURES
PRECISION RECALL F1SCORE SUPPORT
0 090 091 091 174
1 060 058 059 184
2 075 085 080 166
3 078 078 078 194
4 081 084 083 186
5 076 077 077 181
6 091 087 089 207
7 085 088 087 154
8 067 057 062 182
9 075 077 076 169
ACCURACY 078 1797
MACRO AVG 078 078 078 1797
WEIGHTED AVG 078 078 078 1797
PRINTDOC
  AUTHORS YANN N DAUPHIN VLAD NICULAE GABRIEL SYNNAEVE
  LICENSE BSD
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SCIPYNDIMAGE IMPORT CONVOLVE
FROM SKLEARN IMPORT LINEARMODEL DATASETS METRICS
FROM SKLEARNMODELSELECTION IMPORT TRAINTESTSPLIT
FROM SKLEARNNEURALNETWORK IMPORT BERNOULLIRBM
FROM SKLEARNPIPELINE IMPORT PIPELINE
FROM SKLEARNBASE IMPORT CLONE

SETTING UP
DEFNUDGEDDATASET X Y

THIS PRODUCES A DATASET 5 TIMES BIGGER THAN THE ORIGINAL ONE
BY MOVING THE 8X8 IMAGES IN X AROUND BY 1PX TO LEFT RIGHT DOWN UP

DIRECTIONVECTORS
0 1 0
524 NEURAL NETWORKS 1351

```

SCIKITLEARN USER GUIDE RELEASE 0213

```
0 0 0
0 0 0
0 0 0
1 0 0
0 0 0
0 0 0
0 0 1
0 0 0
0 0 0
0 0 0
0 1 0
DEFSHIFTX W
RETURNCONVOLVEXRESHAPE8 8 MODECONSTANT WEIGHTSWRAVEL
X NPCONCATENATEX
NPAPPLYALONGAXISSHIFT 1 X VECTOR
FORVECTORINDIRECTIONVECTORS
Y NPCONCATENATEY FORINRANGE5 AXIS0
RETURNX Y
LOAD DATA
DIGITS DATASETSLOADDIGITS
X NPASARRAYDIGITSDATA FLOAT32
X Y NUDGEDDATASETX DIGITSTARGET
X X NPMINX 0 NPMAXX 0 00001 01 SCALING
XTRAIN XTEST YTRAIN YTEST TRAINTESTSPLIT
X Y TESTSIZE02 RANDOMSTATE0
MODELS WE WILL USE
LOGISTIC LINEARMODELLOGISTICREGRESSIONSOLVERNEWTONCG TOL1
MULTICLASSMULTINOMIAL
RBM BERNOULLIRBMRANDOMSTATE0 VERBOSETRUE
RBMFEATURESCCLASSIFIER PIPELINE
STEP SRBM RBM LOGISTIC LOGISTIC
```

```
TRAINING
HYPERPARAMETERS THESE WERE SET BY CROSSVALIDATION
USING A GRIDSEARCHCV HERE WE ARE NOT PERFORMING CROSSVALIDATION TO
SAVE TIME
RBMLEARNINGRATE 006
RBMNITER 20
MORE COMPONENTS TEND TO GIVE BETTER PREDICTION PERFORMANCE BUT LARGER
FITTING TIME
RBMNCOMPONENTS 100
LOGISTICC 6000
TRAINING RBMLOGISTIC PIPELINE
RBMFEATURESCCLASSIFIERFITXTRAIN YTRAIN
1352 CHAPTER 5 EXAMPLES
```

SCIKITLEARN USER GUIDE RELEASE 0213  
TRAINING THE LOGISTIC REGRESSION CLASSIFIER DIRECTLY ON THE PIXEL  
RAWPIXELCLASSIFIER CLONELOGISTIC  
RAWPIXELCLASSIFIERC 100  
RAWPIXELCLASSIFIERFITXTRAIN YTRAIN

EVALUATION  
YPRED RBMFEATURESClassifierPREDICTXTEST  
PRINTLOGISTIC REGRESSION USING RBM FEATURES NSN  
METRICSClassificationReportYTEST YPRED  
YPRED RAWPIXELCLASSIFIERPREDICTXTEST  
PRINTLOGISTIC REGRESSION USING RAW PIXEL FEATURES NSN  
METRICSClassificationReportYTEST YPRED

PLOTTING  
PLTFigureFIGSIZE42 4  
FORI COMP INENUMERATERBMCOMPONENTS  
PLTSUBPLOT10 10 | 1  
PLTImSHOWCOMPRESHAPE8 8 CMAPPLTCMGRAYR  
INTERPOLATIONNEAREST  
PLXTICKS  
PLTYTICKS  
PLTSUPTITLE100 COMPONENTS EXTRACTED BY RBM FONTSIZE16  
PLTSUBPLOTSADJUST008 002 092 085 008 023  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 11939 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
5243 VARYING REGULARIZATION IN MULTILAYER PERCEPTRON  
A COMPARISON OF DIFFERENT VALUES FOR REGULARIZATION PARAMETER 'ALPHA' ON SYNTHETIC DATASETS THE PLOT SHOWS THAT  
DIFFERENT ALPHAS YIELD DIFFERENT DECISION FUNCTIONS  
ALPHA IS A PARAMETER FOR REGULARIZATION TERM AKA PENALTY TERM THAT COMBATS OVERFITTING BY CONSTRAINING THE SIZE OF THE  
WEIGHTS INCREASING ALPHA MAY FIX HIGH VARIANCE A SIGN OF OVERFITTING BY ENCOURAGING SMALLER WEIGHTS RESULTING IN  
A DECISION BOUNDARY PLOT THAT APPEARS WITH LESSER CURVATURES SIMILARLY DECREASING ALPHA MAY FIX HIGH BIAS A SIGN OF  
UNDERFITTING BY ENCOURAGING LARGER WEIGHTS POTENTIALLY RESULTING IN A MORE COMPLICATED DECISION BOUNDARY  
524 NEURAL NETWORKS 1353

SCIKITLEARN USER GUIDE RELEASE 0213  
PRINTDOC  
AUTHOR ISSAM H LARADJI  
LICENSE BSD 3 CLAUSE  
IMPORT NUMPY AS NP  
FROM MATPLOTLIB IMPORT PYPLLOTASPLT  
FROM MATPLOTLIBCOLORS IMPORT LISTEDCOLORMAP  
FROM SKLEARNMODELSELECTION IMPORT TRAINTESTSPLIT  
FROM SKLEARNPREPROCESSING IMPORT STANDARDSCALER  
FROM SKLEARNDATASETS IMPORT MAKEMOONS MAKECIRCLES MAKECLASSIFICATION  
FROM SKLEARNNEURALNETWORK IMPORT MLPCLASSIFIER  
H 02 STEP SIZE IN THE MESH  
ALPHAS NPLOGSPACE5 3 5  
NAMES ALPHA STRI FORIINALPHAS  
CLASSIFIERS  
FORIINALPHAS  
CLASSIFIERSAPPENDMLPCLASSIFIERSOLVERLBFGS ALPHA1 RANDOMSTATE1  
HIDDENLAYERSIZES100 100  
X Y MAKECLASSIFICATIONNFEATURES2 NREDUNDANT0 NINFORMATIVE2  
RANDOMSTATE0 NCLUSTERSPERCLASS1  
RNG NPRANDOMRANDOMSTATE2  
X 2RNGUNIFORMSIZEXSHAPE  
LINEARLYSEPARABLE X Y  
DATASETS MAKEMOONSNOISE03 RANDOMSTATE0  
MAKECIRCLESNOISE02 FACTOR05 RANDOMSTATE1  
LINEARLYSEPARABLE  
FIGURE PLTFIGUREFIGSIZE17 9  
1354 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213

I 1

ITERATE OVER DATASETS

FORX YINDATASETS

PREPROCESS DATASET SPLIT INTO TRAINING AND TEST PART

X STANDARDSCALERFITTRANSFORMX

XTRAIN XTEST YTRAIN YTEST TRAINTESTSPLITX Y TESTSIZE4

XMIN XMAX X OMIN 5 X OMAX 5

YMIN YMAX X 1MIN 5 X 1MAX 5

XX YY NPMESHGRIDNPARANGEXMIN XMAX H

NPARANGEYMIN YMAX H

JUST PLOT THE DATASET FIRST

CM PLTCMRDBU

CMBRIGHT LISTEDCOLORMAPFF0000 0000FF

AX PLTSUBPLOTLENDATASETS LENCLASSIFIERS 1 I

PLOT THE TRAINING POINTS

AXSCATTERXTRAIN 0 XTRAIN 1 CYTRAIN CMAPCMBRIGHT

AND TESTING POINTS

AXSCATTERXTEST 0 XTEST 1 CYTEST CMAPCMBRIGHT ALPHA06

AXSETXLIMXXMIN XXMAX

AXSETYLIMYYMIN YYMAX

AXSETXTICKS

AXSETYTICKS

I 1

ITERATE OVER CLASSIFIERS

FORNAME CLF INZIPNAMES CLASSIFIERS

AX PLTSUBPLOTLENDATASETS LENCLASSIFIERS 1 I

CLFFITXTRAIN YTRAIN

SCORE CLFSCOREXTEST YTEST

PLOT THE DECISION BOUNDARY FOR THAT WE WILL ASSIGN A COLOR TO EACH

POINT IN THE MESH XMIN XMAXXYMIN YMAX

IFHASATTRCLF DECISIONFUNCTION

Z CLFDECISIONFUNCTIONNPCXXRAVEL YYRAVEL

ELSE

Z CLFPREDICTPROBANPCXXRAVEL YYRAVEL 1

PUT THE RESULT INTO A COLOR PLOT

Z ZRESHAPEXXSHAPE

AXCONTOURFXX YY Z CMAPCM ALPHA8

PLOT ALSO THE TRAINING POINTS

AXSCATTERXTRAIN 0 XTRAIN 1 CYTRAIN CMAPCMBRIGHT

EDGECOLORSBLACK S25

AND TESTING POINTS

AXSCATTERXTEST 0 XTEST 1 CYTEST CMAPCMBRIGHT

ALPHA06 EDGECOLORSBLACK S25

AXSETXLIMXXMIN XXMAX

AXSETYLIMYYMIN YYMAX

AXSETXTICKS

AXSETYTICKS

AXSETTITLENAME

AXTEXTXXMAX 3 YYMIN 3 2F SCORELSTRIPO

SIZE15 HORIZONTALALIGNMENTRIGHT

I 1

524 NEURAL NETWORKS 1355

SCIKITLEARN USER GUIDE RELEASE 0213  
FIGURESUBPLOTSADJUSTLEFT02 RIGHT98  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 7509 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
5244 COMPARE STOCHASTIC LEARNING STRATEGIES FOR MLPCLASSIFIER  
THIS EXAMPLE VISUALIZES SOME TRAINING LOSS CURVES FOR DIFFERENT STOCHASTIC LEARNING STRATEGIES INCLUDING SGD AND ADAM BECAUSE OF TIMECONSTRAINTS WE USE SEVERAL SMALL DATASETS FOR WHICH LBFGS MIGHT BE MORE SUITABLE THE GENERAL TREND SHOWN IN THESE EXAMPLES SEEMS TO CARRY OVER TO LARGER DATASETS HOWEVER  
NOTE THAT THOSE RESULTS CAN BE HIGHLY DEPENDENT ON THE VALUE OF LEARNINGRATEINIT  
OUT  
LEARNING ON DATASET IRIS  
TRAINING CONSTANT LEARNINGRATE  
TRAINING SET SCORE 0980000  
TRAINING SET LOSS 0096950  
TRAINING CONSTANT WITH MOMENTUM  
TRAINING SET SCORE 0980000  
TRAINING SET LOSS 0049530  
TRAINING CONSTANT WITH NESTEROVS MOMENTUM  
1356 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
TRAINING SET SCORE 0980000  
TRAINING SET LOSS 0049540  
TRAINING INVSCALING LEARNINGRATE  
TRAINING SET SCORE 0360000  
TRAINING SET LOSS 0978444  
TRAINING INVSCALING WITH MOMENTUM  
TRAINING SET SCORE 0860000  
TRAINING SET LOSS 0503452  
TRAINING INVSCALING WITH NESTEROVS MOMENTUM  
TRAINING SET SCORE 0860000  
TRAINING SET LOSS 0504185  
TRAINING ADAM  
TRAINING SET SCORE 0980000  
TRAINING SET LOSS 0045311  
LEARNING ON DATASET DIGITS  
TRAINING CONSTANT LEARNINGRATE  
TRAINING SET SCORE 0956038  
TRAINING SET LOSS 0243802  
TRAINING CONSTANT WITH MOMENTUM  
TRAINING SET SCORE 0992766  
TRAINING SET LOSS 0041297  
TRAINING CONSTANT WITH NESTEROVS MOMENTUM  
TRAINING SET SCORE 0993879  
TRAINING SET LOSS 0042898  
TRAINING INVSCALING LEARNINGRATE  
TRAINING SET SCORE 0638843  
TRAINING SET LOSS 1855465  
TRAINING INVSCALING WITH MOMENTUM  
TRAINING SET SCORE 0912632  
TRAINING SET LOSS 0290584  
TRAINING INVSCALING WITH NESTEROVS MOMENTUM  
TRAINING SET SCORE 0909293  
TRAINING SET LOSS 0318387  
TRAINING ADAM  
TRAINING SET SCORE 0991653  
TRAINING SET LOSS 0045934  
LEARNING ON DATASET CIRCLES  
TRAINING CONSTANT LEARNINGRATE  
TRAINING SET SCORE 0840000  
TRAINING SET LOSS 0601052  
TRAINING CONSTANT WITH MOMENTUM  
TRAINING SET SCORE 0940000  
TRAINING SET LOSS 0157334  
TRAINING CONSTANT WITH NESTEROVS MOMENTUM  
TRAINING SET SCORE 0940000  
TRAINING SET LOSS 0154453  
TRAINING INVSCALING LEARNINGRATE  
TRAINING SET SCORE 0500000  
TRAINING SET LOSS 0692470  
TRAINING INVSCALING WITH MOMENTUM  
TRAINING SET SCORE 0500000  
TRAINING SET LOSS 0689143  
TRAINING INVSCALING WITH NESTEROVS MOMENTUM  
TRAINING SET SCORE 0500000  
TRAINING SET LOSS 0689751  
TRAINING ADAM  
524 NEURAL NETWORKS 1357

SCIKITLEARN USER GUIDE RELEASE 0213  
TRAINING SET SCORE 0940000  
TRAINING SET LOSS 0150527  
LEARNING ON DATASET MOONS  
TRAINING CONSTANT LEARNINGRATE  
TRAINING SET SCORE 0850000  
TRAINING SET LOSS 0341523  
TRAINING CONSTANT WITH MOMENTUM  
TRAINING SET SCORE 0850000  
TRAINING SET LOSS 0336188  
TRAINING CONSTANT WITH NESTEROVS MOMENTUM  
TRAINING SET SCORE 0850000  
TRAINING SET LOSS 0335919  
TRAINING INVSCALING LEARNINGRATE  
TRAINING SET SCORE 0500000  
TRAINING SET LOSS 0689015  
TRAINING INVSCALING WITH MOMENTUM  
TRAINING SET SCORE 0830000  
TRAINING SET LOSS 0512595  
TRAINING INVSCALING WITH NESTEROVS MOMENTUM  
TRAINING SET SCORE 0830000  
TRAINING SET LOSS 0513034  
TRAINING ADAM  
TRAINING SET SCORE 0930000  
TRAINING SET LOSS 0170087  
PRINTDOC  
IMPORT WARNINGS  
IMPORT MATPLOTLIBPYPLOT AS PLT  
FROM SKLEARNNEURALNETWORK IMPORT MLPCLASSIFIER  
FROM SKLEARNPREPROCESSING IMPORT MINMAXSCALER  
FROM SKLEARN IMPORT DATASETS  
FROM SKLEARNEXCEPTIONS IMPORT CONVERGENCEWARNING  
DIFFERENT LEARNING RATE SCHEDULES AND MOMENTUM PARAMETERS  
PARAMS SOLVER SGD LEARNINGRATE CONSTANT MOMENTUM 0  
LEARNINGRATEINIT 02  
SOLVER SGD LEARNINGRATE CONSTANT MOMENTUM 9  
NESTEROVSMOMENTUM FALSE LEARNINGRATEINIT 02  
SOLVER SGD LEARNINGRATE CONSTANT MOMENTUM 9  
NESTEROVSMOMENTUM TRUE LEARNINGRATEINIT 02  
SOLVER SGD LEARNINGRATE INVSCALING MOMENTUM 0  
LEARNINGRATEINIT 02  
SOLVER SGD LEARNINGRATE INVSCALING MOMENTUM 9  
NESTEROVSMOMENTUM TRUE LEARNINGRATEINIT 02  
SOLVER SGD LEARNINGRATE INVSCALING MOMENTUM 9  
NESTEROVSMOMENTUM FALSE LEARNINGRATEINIT 02  
SOLVER ADAM LEARNINGRATEINIT 001  
LABELS CONSTANT LEARNINGRATE CONSTANT WITH MOMENTUM  
1358 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
CONSTANT WITH NESTEROVS MOMENTUM  
INVSCALING LEARNINGRATE INVSCALING WITH MOMENTUM  
INVSCALING WITH NESTEROVS MOMENTUM ADAM  
PLOTARGS C RED LIFESTYLE  
C GREEN LIFESTYLE  
C BLUE LIFESTYLE  
C RED LIFESTYLE  
C GREEN LIFESTYLE  
C BLUE LIFESTYLE  
C BLACK LIFESTYLE  
DEFPLOTONDATASET X Y AX NAME  
FOR EACH DATASET PLOT LEARNING FOR EACH LEARNING STRATEGY  
PRINTNLEARNING ON DATASET S NAME  
AXSETTITLENAME  
X MINMAXSCALERFITTRANSFORMX  
MLPS  
IFNAME DIGITS  
DIGITS IS LARGER BUT CONVERGES FAIRLY QUICKLY  
MAXITER 15  
ELSE  
MAXITER 400  
FORLABEL PARAM INZIPLABELS PARAMS  
PRINTTRAINING S LABEL  
MLP MLPCLASSIFIERVERBOSE0 RANDOMSTATE0  
MAXITERMAXITER PARAM  
SOME PARAMETER COMBINATIONS WILL NOT CONVERGE AS CAN BE SEEN ON THE  
PLOTS SO THEY ARE IGNORED HERE  
WITHWARNINGSCATCHWARNINGS  
WARNINGSFILTERWARNINGSIGNORE CATEGORYCONVERGENCEWARNING  
MODULESKLEARN  
MLPFIT X Y  
MLPSAPPENDMLP  
PRINTTRAINING SET SCORE F MLPSCORE X Y  
PRINTTRAINING SET LOSS F MLPLOSS  
FORMLP LABEL ARGS INZIPMLPS LABELS PLOTARGS  
AXPLOTMLPLOSSCURVE LABELLABEL ARGS  
FIG AXES PLTSUBPLOTS2 2 FIGSIZE15 10  
LOAD GENERATE SOME TOY DATASETS  
IRIS DATASETSLOADIRIS  
DIGITS DATASETSLOADDIGITS  
DATASETS IRISDATA IRISTARGET  
DIGITSDATA DIGITSTARGET  
DATASETSMAKECIRCLESNOISE02 FACTOR05 RANDOMSTATE1  
DATASETSMAKEMOONSNOISE03 RANDOMSTATE0  
FORAX DATA NAME INZIPAXESRAVEL DATASETS IRIS DIGITS  
CIRCLES MOONS  
PLOTONDATASET DATA AXAX NAMENAME  
524 NEURAL NETWORKS 1359

SCIKITLEARN USER GUIDE RELEASE 0213  
FIGLEGENDAXGETLINES LABELS NCOL3 LOCUPPER CENTER  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 4379 SECONDS  
525 PREPROCESSING  
EXAMPLES CONCERNING THE SKLEARNPREPROCESSING MODULE  
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE  
5251 USING FUNCTIONTRANSFORMER TO SELECT COLUMNS  
SHOWS HOW TO USE A FUNCTION TRANSFORMER IN A PIPELINE IF YOU KNOW YOUR DATASET'S FIRST PRINCIPLE COMPONENT IS IRRELEVANT  
FOR A CLASSIFICATION TASK YOU CAN USE THE FUNCTIONTRANSFORMER TO SELECT ALL BUT THE FIRST COLUMN OF THE PCA TRANSFORMED  
DATA  
•  
1360 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

```
•
IMPORT MATPLOTLIBPYPLOT AS PLT
IMPORT NUMPY AS NP
FROM SKLEARNMODELSELECTION IMPORT TRAINTESTSPLIT
FROM SKLEARNDECOMPOSITION IMPORT PCA
FROM SKLEARNPIPELINE IMPORT MAKEPIPELINE
FROM SKLEARNPREPROCESSING IMPORT FUNCTIONTRANSFORMER
DEFGENERATEVECTORSHIFT05 NOISE15
RETURNNPVRANGE1000 NPRANDOMRAND1000 SHIFT NOISE
DEFGENERATEDATASET
```

THIS DATASET IS TWO LINES WITH A SLOPE 1 WHERE ONE HAS  
A Y OFFSET OF 100

```
RETURNNPVSTACK
NPVSTACK
GENERATEVECTOR
GENERATEVECTOR 100
T
NPVSTACK
GENERATEVECTOR
GENERATEVECTOR
T
525 PREPROCESSING 1361
```

SCIKITLEARN USER GUIDE RELEASE 0213  
NPHSTACKNPZEROS1000 NPONES1000  
DEFALLBUTFIRSTCOLUMNX  
RETURNX 1  
DEFDROPFIRSTCOMPONENTX Y

CREATE A PIPELINE WITH PCA AND THE COLUMN SELECTOR AND USE IT TO  
TRANSFORM THE DATASET

PIPELINE MAKEPIPELINE  
PCA FUNCTIONTRANSFORMERALLBUTFIRSTCOLUMN

XTRAIN XTEST YTRAIN YTEST TRAINTESTSPLITX Y  
PIPELINEFITXTRAIN YTRAIN  
RETURNPIPELINETRANSFORMXTEST YTEST  
IFNAME MAIN  
X Y GENERATEDDATASET  
LW 0  
PLTFigure  
PLTSCATTERX 0 X 1 CY LWLW  
PLTFigure  
XTRANSFORMED YTRANSFORMED DROPFIRSTCOMPONENT GENERATEDDATASET  
PLTSCATTER  
XTRANSFORMED 0  
NPZEROSLENXTRANSFORMED  
CYTRANSFORMED  
LWLW  
S60

PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0028 SECONDS  
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE  
5252 USING KBINSDISCRETIZER TO DISCRETIZE CONTINUOUS FEATURES  
THE EXAMPLE COMPARES PREDICTION RESULT OF LINEAR REGRESSION LINEAR MODEL AND DECISION TREE TREE BASED MODEL WITH  
AND WITHOUT DISCRETIZATION OF REALVALUED FEATURES  
AS IS SHOWN IN THE RESULT BEFORE DISCRETIZATION LINEAR MODEL IS FAST TO BUILD AND RELATIVELY STRAIGHTFORWARD TO INTERPRET  
BUT CAN ONLY MODEL LINEAR RELATIONSHIPS WHILE DECISION TREE CAN BUILD A MUCH MORE COMPLEX MODEL OF THE DATA ONE  
WAY TO MAKE LINEAR MODEL MORE POWERFUL ON CONTINUOUS DATA IS TO USE DISCRETIZATION ALSO KNOWN AS BINNING IN THE  
EXAMPLE WE DISCRETIZE THE FEATURE AND ONEHOT ENCODE THE TRANSFORMED DATA NOTE THAT IF THE BINS ARE NOT REASONABLY  
WIDE THERE WOULD APPEAR TO BE A SUBSTANTIALLY INCREASED RISK OF OVERFITTING SO THE DISCRETIZER PARAMETERS SHOULD USUALLY  
BE TUNED UNDER CROSS VALIDATION  
AFTER DISCRETIZATION LINEAR REGRESSION AND DECISION TREE MAKE EXACTLY THE SAME PREDICTION AS FEATURES ARE CONSTANT  
WITHIN EACH BIN ANY MODEL MUST PREDICT THE SAME VALUE FOR ALL POINTS WITHIN A BIN COMPARED WITH THE RESULT BEFORE  
1362 CHAPTER 5 EXAMPLES



```
SCIKITLEARN USER GUIDE RELEASE 0213
DISCRETIZATION LINEAR MODEL BECOME MUCH MORE FLEXIBLE WHILE DECISION TREE GETS MUCH LESS FLEXIBLE NOTE THAT BIN
NING FEATURES GENERALLY HAS NO BENEFICIAL EFFECT FOR TREEBASED MODELS AS THESE MODELS CAN LEARN TO SPLIT UP THE DATA
ANYWHERE
AUTHOR ANDREAS MÜLLER
HANMIN QIN QINHANMIN2005SINACOM
LICENSE BSD 3 CLAUSE
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARNLINEARMODEL IMPORT LINEARREGRESSION
FROM SKLEARNPREPROCESSING IMPORT KBINSDISCRETIZER
FROM SKLEARNTREE IMPORT DECISIONTREEREgressor
PRINTDOC
CONSTRUCT THE DATASET
RND NPRANDOMRANDOMSTATE42
X RNDUNIFORM3 3 SIZE100
Y NPSINX RNDNORMALSIZELENX 3
X XRESHAPE1 1
TRANSFORM THE DATASET WITH KBINSDISCRETIZER
ENC KBINSDISCRETIZERNBINS10 ENCODEONEHOT
XBINNED ENCFITTRANSFORMX
PREDICT WITH ORIGINAL DATASET
FIG AX1 AX2 PLTSUBPLOTSNCOLS2 SHAREYTRUE FIGSIZE10 4
LINE NPLINSPACE3 3 1000 ENDPOINTFALSERESHAPE1 1
REG LINEARREGRESSIONFITX Y
AX1PLOTLINE REGPREDICTLINE LINEWIDTH2 COLORGREEN
LABELLINEAR REGRESSION
REG DECISIONTREEREgressorMINSAMPLESSPLIT3 RANDOMSTATE0FITX Y
AX1PLOTLINE REGPREDICTLINE LINEWIDTH2 COLORRED
LABELDECISION TREE
AX1PLOTX 0 Y O CK
AX1LEGENDLOCBEST
AX1SETYLABELREGRESSION OUTPUT
AX1SETXLABELINPUT FEATURE
AX1SETTITLERESULT BEFORE DISCRETIZATION
525 PREPROCESSING 1363
```

SCIKITLEARN USER GUIDE RELEASE 0213  
PREDICT WITH TRANSFORMED DATASET  
LINEBINNED ENCTRANSFORMLINE  
REG LINEARREGRESSIONFITXBINNED Y  
AX2PLOTLINE REGPREDICTLINEBINNED LINEWIDTH2 COLORGREEN  
LINESTYLE LABELLINEAR REGRESSION  
REG DECISIONTREEREGRESSORMINSAMPLESSPLIT3  
RANDOMSTATE0FITXBINNED Y  
AX2PLOTLINE REGPREDICTLINEBINNED LINEWIDTH2 COLORRED  
LINESTYLE LABELDECISION TREE  
AX2PLOTX 0 Y 0 CK  
AX2VLINESENCBINEDGES0 PLTGCAGETYLIM LINEWIDTH1 ALPHA2  
AX2LEGENDLOCBEST  
AX2SETXLABELINPUT FEATURE  
AX2SETTITLERESULT AFTER DISCRETIZATION  
PLTTIGHTLAYOUT  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0122 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
5253 DEMONSTRATING THE DIFFERENT STRATEGIES OF KBINSDISCRETIZER  
THIS EXAMPLE PRESENTS THE DIFFERENT STRATEGIES IMPLEMENTED IN KBINSDISCRETIZER  
• ‘UNIFORM’ THE DISCRETIZATION IS UNIFORM IN EACH FEATURE WHICH MEANS THAT THE BIN WIDTHS ARE CONSTANT IN EACH DIMENSION  
• ‘QUANTILE’ THE DISCRETIZATION IS DONE ON THE QUANTILED VALUES WHICH MEANS THAT EACH BIN HAS APPROXIMATELY THE SAME NUMBER OF SAMPLES  
• ‘KMEANS’ THE DISCRETIZATION IS BASED ON THE CENTROIDS OF A KMEANS CLUSTERING PROCEDURE  
THE PLOT SHOWS THE REGIONS WHERE THE DISCRETIZED ENCODING IS CONSTANT  
1364 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
AUTHOR TOM DUPRÉ LA TOUR
LICENSE BSD 3 CLAUSE
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARNPREPROCESSING IMPORT KBINSDISCRETIZER
FROM SKLEARNDATASETS IMPORT MAKEBLOBS
PRINTDOC
STRATEGIES UNIFORM QUANTILE KMEANS
NSAMPLES 200
CENTERS0 NPARRAY0 0 0 5 2 4 8 8
CENTERS1 NPARRAY0 0 3 1
CONSTRUCT THE DATASETS
RANDOMSTATE 42
XLIST
NPRANDOMRANDOMSTATERRANDOMSTATEUNIFORM3 3 SIZENSAMPLES 2
MAKEBLOBSNSAMPLESNSAMPLES 10 NSAMPLES 4 10
NSAMPLES 10 NSAMPLES 4 10
CLUSTERSTD05 CENTERSCENTERS0
RANDOMSTATERRANDOMSTATE0
MAKEBLOBSNSAMPLESNSAMPLES 5 NSAMPLES 4 5
CLUSTERSTD05 CENTERSCENTERS1
RANDOMSTATERRANDOMSTATE0

525 PREPROCESSING 1365
```

SCIKITLEARN USER GUIDE RELEASE 0213

FIGURE PLTFIGUREFIGSIZE14 9

I 1

FORDSCNT X INENUMERATEXLIST

AX PLTSUBPLOTLENXLIST LENSTRATEGIES 1 I

AXSCATTERX 0 X 1 EDGECOLORSK

IFDSCNT 0

AXSETTITLEINPUT DATA SIZE14

XX YY NPMESHGRID

NPLINSPACEX 0MIN X 0MAX 300

NPLINSPACEX 1MIN X 1MAX 300

GRID NPCXXRAVEL YYRAVEL

AXSETXLIMXXMIN XXMAX

AXSETYLIMYYMIN YYMAX

AXSETXTICKS

AXSETYTICKS

I 1

TRANSFORM THE DATASET WITH KBINSDISCRETIZER

FORSTRATEGY INSTRATEGIES

ENC KBINSDISCRETIZERNBINS4 ENCODEORDINAL STRATEGYSTRATEGY

ENCFITX

GRIDENCODED ENCTRANSFORMGRID

AX PLTSUBPLOTLENXLIST LENSTRATEGIES 1 I

HORIZONTAL STRIPES

HORIZONTAL GRIDENCODED 0RESHAPEXXSHAPE

AXCONTOURFXX YY HORIZONTAL ALPHA5

VERTICAL STRIPES

VERTICAL GRIDENCODED 1RESHAPEXXSHAPE

AXCONTOURFXX YY VERTICAL ALPHA5

AXSCATTERX 0 X 1 EDGECOLORSK

AXSETXLIMXXMIN XXMAX

AXSETYLIMYYMIN YYMAX

AXSETXTICKS

AXSETYTICKS

IFDSCNT 0

AXSETTITLESTRATEGY S STRATEGY SIZE14

I 1

PLTTIGHTLAYOUT

PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0890 SECONDS

NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)

1366 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

5254 IMPORTANCE OF FEATURE SCALING

FEATURE SCALING THROUGH STANDARDIZATION OR ZSCORE NORMALIZATION CAN BE AN IMPORTANT PREPROCESSING STEP FOR MANY MACHINE LEARNING ALGORITHMS STANDARDIZATION INVOLVES RESCALING THE FEATURES SUCH THAT THEY HAVE THE PROPERTIES OF A STANDARD NORMAL DISTRIBUTION WITH A MEAN OF ZERO AND A STANDARD DEVIATION OF ONE WHILE MANY ALGORITHMS SUCH AS SVM KNEAREST NEIGHBORS AND LOGISTIC REGRESSION REQUIRE FEATURES TO BE NORMALIZED INTUITIVELY WE CAN THINK OF PRINCIPLE COMPONENT ANALYSIS PCA AS BEING A PRIME EXAMPLE OF WHEN NORMALIZATION IS IMPORTANT IN PCA WE ARE INTERESTED IN THE COMPONENTS THAT MAXIMIZE THE VARIANCE IF ONE COMPONENT EG HUMAN HEIGHT VARIES LESS THAN ANOTHER EG WEIGHT BECAUSE OF THEIR RESPECTIVE SCALES METERS VS KILOS PCA MIGHT DETERMINE THAT THE DIRECTION OF MAXIMAL VARIANCE MORE CLOSELY CORRESPONDS WITH THE 'WEIGHT' AXIS IF THOSE FEATURES ARE NOT SCALED AS A CHANGE IN HEIGHT OF ONE METER CAN BE CONSIDERED MUCH MORE IMPORTANT THAN THE CHANGE IN WEIGHT OF ONE KILOGRAM THIS IS CLEARLY INCORRECT TO ILLUSTRATE THIS PCA IS PERFORMED COMPARING THE USE OF DATA WITH STANDARDSCALER APPLIED TO UNSCALED DATA THE RESULTS ARE VISUALIZED AND A CLEAR DIFFERENCE NOTED THE 1ST PRINCIPAL COMPONENT IN THE UNSCALED SET CAN BE SEEN IT CAN BE SEEN THAT FEATURE 13 DOMINATES THE DIRECTION BEING A WHOLE TWO ORDERS OF MAGNITUDE ABOVE THE OTHER FEATURES THIS IS CONTRASTED WHEN OBSERVING THE PRINCIPAL COMPONENT FOR THE SCALED VERSION OF THE DATA IN THE SCALED VERSION THE ORDERS OF MAGNITUDE ARE ROUGHLY THE SAME ACROSS ALL THE FEATURES THE DATASET USED IS THE WINE DATASET AVAILABLE AT UCI THIS DATASET HAS CONTINUOUS FEATURES THAT ARE HETEROGENEOUS IN SCALE DUE TO DIFFERING PROPERTIES THAT THEY MEASURE IE ALCOHOL CONTENT AND MALIC ACID THE TRANSFORMED DATA IS THEN USED TO TRAIN A NAIVE BAYES CLASSIFIER AND A CLEAR DIFFERENCE IN PREDICTION ACCURACIES IS OBSERVED WHEREIN THE DATASET WHICH IS SCALED BEFORE PCA VASTLY OUTPERFORMS THE UNSCALED VERSION

OUT  
525 PREPROCESSING 1367

SCIKITLEARN USER GUIDE RELEASE 0213  
PREDICTION ACCURACY FOR THE NORMAL TEST DATASET WITH PCA  
8148  
PREDICTION ACCURACY FOR THE STANDARDIZED TEST DATASET WITH PCA  
9815  
PC 1 WITHOUT SCALING  
176E03 836E04 155E04 531E03 202E02 102E03 153E03  
112E04 631E04 233E03 154E04 743E04 100E00  
PC 1 WITH SCALING  
013 026 001 023 016 039 042 028 033 011 03 038  
028  
FROM SKLEARNMODELSELECTION IMPORT TRAINTESTSPLIT  
FROM SKLEARNPREPROCESSING IMPORT STANDARDSCALER  
FROM SKLEARNDECOMPOSITION IMPORT PCA  
FROM SKLEARNNAIVEBAYES IMPORT GAUSSIANNB  
FROM SKLEARN IMPORT METRICS  
IMPORT MATPLOTLIBPYPLOT AS PLT  
FROM SKLEARNDATASETS IMPORT LOADWINE  
FROM SKLEARNPIPELINE IMPORT MAKEPIPELINE  
PRINTDOC  
CODE SOURCE TYLER LANIGAN TYLERLANIGANGMAILCOM  
SEBASTIAN RASCHKA MAILSEBASTIANRASCHKACOM  
LICENSE BSD 3 CLAUSE  
RANDOMSTATE 42  
FIGSIZE 10 7  
FEATURES TARGET LOADWINERETURNXYTRUE  
MAKE A TRAINTEST SPLIT USING 30 TEST SIZE  
XTRAIN XTEST YTRAIN YTEST TRAINTESTSPLITFEATURES TARGET  
TESTSIZE030  
RANDOMSTATERRANDOMSTATE  
FIT TO DATA AND PREDICT USING PIPELINED GNB AND PCA  
UNSCALEDCLF MAKEPIPELINEPCANCOMPONENTS2 GAUSSIANNB  
UNSCALEDCLFFITXTRAIN YTRAIN  
PREDTEST UNSCALEDCLFPREDICTXTEST  
FIT TO DATA AND PREDICT USING PIPELINED SCALING GNB AND PCA  
STDCLF MAKEPIPELINESTANDARDSCALER PCANCOMPONENTS2 GAUSSIANNB  
STDCLFFITXTRAIN YTRAIN  
PREDTESTSTD STDCLFPREDICTXTEST  
SHOW PREDICTION ACCURACIES IN SCALED AND UNSCALED DATA  
1368 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
 PRINTNPREDICTION ACCURACY FOR THE NORMAL TEST DATASET WITH PCA  
 PRINT2 NFORMATMETRICSACCURACYSCOREYTEST PREDTEST  
 PRINTNPREDICTION ACCURACY FOR THE STANDARDIZED TEST DATASET WITH PCA  
 PRINT2 NFORMATMETRICSACCURACYSCOREYTEST PREDTESTSTD  
 EXTRACT PCA FROM PIPELINE  
 PCA UNSCALEDCLFNAMEDSTEPSPCA  
 PCASTD STDCLFNAMEDSTEPSPCA  
 SHOW FIRST PRINCIPAL COMPONENTS  
 PRINTNPC 1 WITHOUT SCALING N PCACOMPONENTS0  
 PRINTNPC 1 WITH SCALING N PCASTDCOMPONENTS0  
 USE PCA WITHOUT AND WITH SCALE ON XTRAIN DATA FOR VISUALIZATION  
 XTRAINTRANSFORMED PCATransformXTRAIN  
 SCALER STDCLFNAMEDSTEPSSANDARDSCALER  
 XTRAINSTDTRANSFORMED PCASTDTRANSFORMSCALERTRANSFORMXTRAIN  
 VISUALIZE STANDARDIZED VS UNTOUCHED DATASET WITH PCA PERFORMED  
 FIG AX1 AX2 PLTSUBPLOTSNCOLS2 FIGSIZEFIGSIZE  
 FORL C M INZIPRANGE0 3 BLUE RED GREEN S O  
 AX1SCATTERXTRAINTRANSFORMEDYTRAIN L 0  
 XTRAINTRANSFORMEDYTRAIN L 1  
 COLORC  
 LABELCLASS S L  
 ALPHA05  
 MARKERM  
  
 FORL C M INZIPRANGE0 3 BLUE RED GREEN S O  
 AX2SCATTERXTRAINSTDTRANSFORMEDYTRAIN L 0  
 XTRAINSTDTRANSFORMEDYTRAIN L 1  
 COLORC  
 LABELCLASS S L  
 ALPHA05  
 MARKERM  
  
 AX1SETTITLETRAINING DATASET AFTER PCA  
 AX2SETTITLESTANDARDIZED TRAINING DATASET AFTER PCA  
 FORAXINAX1 AX2  
 AXSETXLABEL1ST PRINCIPAL COMPONENT  
 AXSETYLABEL2ND PRINCIPAL COMPONENT  
 AXLEGENDLOCUPPER RIGHT  
 AXGRID  
 PLTTIGHTLAYOUT  
 PLTSHOW  
 TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0109 SECONDS  
 525 PREPROCESSING 1369

SCIKITLEARN USER GUIDE RELEASE 0213

[NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)

5255 MAP DATA TO A NORMAL DISTRIBUTION

THIS EXAMPLE DEMONSTRATES THE USE OF THE BOXCOX AND YEOJOHNSON TRANSFORMS THROUGH PREPROCESSING POWERTRANSFORMER TO MAP DATA FROM VARIOUS DISTRIBUTIONS TO A NORMAL DISTRIBUTION

THE POWER TRANSFORM IS USEFUL AS A TRANSFORMATION IN MODELING PROBLEMS WHERE HOMOSCEDASTICITY AND NORMALITY ARE DESIRED BELOW ARE EXAMPLES OF BOXCOX AND YEOJOHNSON APPLIED TO SIX DIFFERENT PROBABILITY DISTRIBUTIONS LOGNORMAL CHISQUARED WEIBULL GAUSSIAN UNIFORM AND BIMODAL

NOTE THAT THE TRANSFORMATIONS SUCCESSFULLY MAP THE DATA TO A NORMAL DISTRIBUTION WHEN APPLIED TO CERTAIN DATASETS BUT ARE INEFFECTIVE WITH OTHERS THIS HIGHLIGHTS THE IMPORTANCE OF VISUALIZING THE DATA BEFORE AND AFTER TRANSFORMATION ALSO NOTE THAT EVEN THOUGH BOXCOX SEEMS TO PERFORM BETTER THAN YEOJOHNSON FOR LOGNORMAL AND CHISQUARED DISTRIBUTIONS KEEP IN MIND THAT BOXCOX DOES NOT SUPPORT INPUTS WITH NEGATIVE VALUES

FOR COMPARISON WE ALSO ADD THE OUTPUT FROM PREPROCESSINGQUANTILETRANSFORMER IT CAN FORCE ANY ARBITRARY DISTRIBUTION INTO A GAUSSIAN PROVIDED THAT THERE ARE ENOUGH TRAINING SAMPLES THOUSANDS BECAUSE IT IS A NONPARAMETRIC METHOD IT IS HARDER TO INTERPRET THAN THE PARAMETRIC ONES BOXCOX AND YEOJOHNSON

ON “SMALL” DATASETS LESS THAN A FEW HUNDRED POINTS THE QUANTILE TRANSFORMER IS PRONE TO OVERFITTING THE USE OF THE POWER TRANSFORM IS THEN RECOMMENDED

1370 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213  
AUTHOR ERIC CHANG ERICCHANG2017UNORTHWESTERNEDU  
NICOLAS HUG CONTACTNICOLASHUGCOM  
LICENSE BSD 3 CLAUSE  
IMPORT NUMPY AS NP  
525 PREPROCESSING 1371

```
SCIKITLEARN USER GUIDE RELEASE 0213
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARNPREPROCESSING IMPORT POWERTRANSFORMER
FROM SKLEARNPREPROCESSING IMPORT QUANTILETRANSFORMER
FROM SKLEARNMODELSELECTION IMPORT TRAINTESTSPLIT
PRINTDOC
NSAMPLES 1000
FONTSIZE 6
BINS 30
RNG NPRANDOMRANDOMSTATE304
BC POWERTRANSFORMERMETHODBOXCOX
YJ POWERTRANSFORMERMETHODYEOJOHNSON
NQUANTILES IS SET TO THE TRAINING SET SIZE RATHER THAN THE DEFAULT VALUE
TO AVOID A WARNING BEING RAISED BY THIS EXAMPLE
QT QUANTILETRANSFORMERNQUANTILES500 OUTPUTDISTRIBUTIONNORMAL
RANDOMSTATERNG
SIZE NSAMPLES 1
LOGNORMAL DISTRIBUTION
XLOGNORMAL RNGLOGNORMALSIZESIZE
CHISQUARED DISTRIBUTION
DF 3
XCHISQ RNGCHISQUAREDfdf SIZESIZE
WEIBULL DISTRIBUTION
A 50
XWEIBULL RNGWEIBULLAA SIZESIZE
GAUSSIAN DISTRIBUTION
LOC 100
XGAUSSIAN RNGNORMALLOCLOC SIZESIZE
UNIFORM DISTRIBUTION
XUNIFORM RNGUNIFORMLOW0 HIGH1 SIZESIZE
BIMODAL DISTRIBUTION
LOCA LOCB 100 105
XA XB RNGNORMALLOCLOCA SIZESIZE RNGNORMALLOCLOCB SIZESIZE
XBIMODAL NPConcatenateXA XB AXIS0
CREATE PLOTS
DISTRIBUTIONS
LOGNORMAL XLOGNORMAL
CHISQUARED XCHISQ
WEIBULL XWEIBULL
GAUSSIAN XGAUSSIAN
UNIFORM XUNIFORM
BIMODAL XBIMODAL
```

SCIKITLEARN USER GUIDE RELEASE 0213  
COLORS FIREBRICK DARKORANGE GOLDENROD  
SEAGREEN ROYALBLUE DARKORCHID  
FIG AXES PLTSUBPLOTSNROWS8 NCOLS3 FIGSIZEPLTFIGASPECT2  
AXES AXESFLATTEN  
AXESIDX 0 3 6 9 1 4 7 10 2 5 8 11 12 15 18 21  
13 16 19 22 14 17 20 23  
AXESLIST AXESI AXESJ AXESK AXESL  
FOR I J K L IN AXESIDX  
FORDISTRIBUTION COLOR AXES INZIPDISTRIBUTIONS COLORS AXESLIST  
NAME X DISTRIBUTION  
XTRAIN XTEST TRAINTESTSPLITX TESTSIZE5  
PERFORM POWER TRANSFORMS AND QUANTILE TRANSFORM  
XTRANSBC BCFITXTRAINTRANSFORMXTEST  
LMBDABC ROUNDCLAMBDA0 2  
XTRANSY YJFITXTRAINTRANSFORMXTEST  
LMBDAYJ ROUNDYJLAMBDA0 2  
XTRANSQT QTFITXTRAINTRANSFORMXTEST  
AXORIGINAL AXBC AXJ AXQT AXES  
AXORIGINALHISTXTRAIN COLORCOLOR BINSBINS  
AXORIGINALSETTITLENAME FONTSIZEFONTSIZE  
AXORIGINALTICKPARAMSAXISBOTH WHICHMAJOR LABELSIZEFONTSIZE  
FORAX XTRANS METHNAME LMBDA INZIP  
AXBC AXJ AXQT  
XTRANSBC XTRANSY XTRANSQT  
BOXCOX YEOJOHNSON QUANTILE TRANSFORM  
LMBDABC LMBDAYJ NONE  
AXHISTXTRANS COLORCOLOR BINSBINS  
TITLE AFTER FORMATMETHNAME  
IFLMBDAIS NOTNONE  
TITLE RNLAMBDA FORMATLMBDA  
AXSETTITLETITLE FONTSIZEFONTSIZE  
AXTICKPARAMSAXISBOTH WHICHMAJOR LABELSIZEFONTSIZE  
AXSETXLIM35 35  
PLTTIGHTLAYOUT  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 1536 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
5256 FEATURE DISCRETIZATION  
A DEMONSTRATION OF FEATURE DISCRETIZATION ON SYNTHETIC CLASSIFICATION DATASETS FEATURE DISCRETIZATION DECOMPOSES EACH  
FEATURE INTO A SET OF BINS HERE EQUALLY DISTRIBUTED IN WIDTH THE DISCRETE VALUES ARE THEN ONEHOT ENCODED AND GIVEN TO  
A LINEAR CLASSIFIER THIS PREPROCESSING ENABLES A NONLINEAR BEHAVIOR EVEN THOUGH THE CLASSIFIER IS LINEAR  
525 PREPROCESSING 1373

SCIKITLEARN USER GUIDE RELEASE 0213

ON THIS EXAMPLE THE FIRST TWO ROWS REPRESENT LINEARLY NONSEPARABLE DATASETS MOONS AND CONCENTRIC CIRCLES WHILE THE THIRD IS APPROXIMATELY LINEARLY SEPARABLE ON THE TWO LINEARLY NONSEPARABLE DATASETS FEATURE DISCRETIZATION LARGELY INCREASES THE PERFORMANCE OF LINEAR CLASSIFIERS ON THE LINEARLY SEPARABLE DATASET FEATURE DISCRETIZATION DECREASES THE PERFORMANCE OF LINEAR CLASSIFIERS TWO NONLINEAR CLASSIFIERS ARE ALSO SHOWN FOR COMPARISON

THIS EXAMPLE SHOULD BE TAKEN WITH A GRAIN OF SALT AS THE INTUITION CONVEYED DOES NOT NECESSARILY CARRY OVER TO REAL DATASETS PARTICULARLY IN HIGHDIMENSIONAL SPACES DATA CAN MORE EASILY BE SEPARATED LINEARLY MOREOVER USING FEATURE DISCRETIZATION AND ONEHOT ENCODING INCREASES THE NUMBER OF FEATURES WHICH EASILY LEAD TO OVERFITTING WHEN THE NUMBER OF SAMPLES IS SMALL

THE PLOTS SHOW TRAINING POINTS IN SOLID COLORS AND TESTING POINTS SEMITRANSSPARENT THE LOWER RIGHT SHOWS THE CLASSIFICATION ACCURACY ON THE TEST SET

OUT

DATASET 0

LOGISTICREGRESSION 086

LINEARSVC 086

KBINSDISCRETIZER LOGISTICREGRESSION 094

KBINSDISCRETIZER LINEARSVC 092

GRADIENTBOOSTINGCLASSIFIER 090

SVC 094

DATASET 1

LOGISTICREGRESSION 040

LINEARSVC 040

KBINSDISCRETIZER LOGISTICREGRESSION 088

KBINSDISCRETIZER LINEARSVC 086

GRADIENTBOOSTINGCLASSIFIER 080

SVC 084

DATASET 2

LOGISTICREGRESSION 096

LINEARSVC 098

KBINSDISCRETIZER LOGISTICREGRESSION 094

KBINSDISCRETIZER LINEARSVC 084

GRADIENTBOOSTINGCLASSIFIER 094

1374 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
SVC 098  
CODE SOURCE TOM DUPRÉ LA TOUR  
ADAPTED FROM PLOTCLASSIFIERCOMPARISON BY GAËL VAROQUAUX AND ANDREAS MÜLLER

LICENSE BSD 3 CLAUSE  
IMPORT NUMPY AS NP  
IMPORT MATPLOTLIBPYPLOT AS PLT  
FROM MATPLOTLIBCOLORS IMPORT LISTEDCOLORMAP  
FROM SKLEARNMODELSELECTION IMPORT TRAINTESTSPLIT  
FROM SKLEARNPREPROCESSING IMPORT STANDARDSCALER  
FROM SKLEARNDATASETS IMPORT MAKEMOONS MAKECIRCLES MAKECLASSIFICATION  
FROM SKLEARNLINEARMODEL IMPORT LOGISTICREGRESSION  
FROM SKLEARNMODELSELECTION IMPORT GRIDSEARCHCV  
FROM SKLEARNPIPELINE IMPORT MAKEPIPELINE  
FROM SKLEARNPREPROCESSING IMPORT KBINSDISCRETIZER  
FROM SKLEARN SVM IMPORT SVC LINEARSVC  
FROM SKLEARNENSEMBLE IMPORT GRADIENTBOOSTINGCLASSIFIER  
FROM SKLEARNUTILSTESTING IMPORT IGNOREWARNINGS  
FROM SKLEARNEXCEPTIONS IMPORT CONVERGENCEWARNING  
PRINTDOC  
H 02 STEP SIZE IN THE MESH  
DEFGETNAMEESTIMATOR  
NAME ESTIMATORCLASSNAME  
IFNAME PIPELINE  
NAME GETNAMEEST1 FORESTINESTIMATORSTEPS  
NAME JOINNAME  
RETURNNAME  
LIST OF ESTIMATOR PARAMGRID WHERE PARAMGRID IS USED IN GRIDSEARCHCV  
CLASSIFIERS  
LOGISTICREGRESSIONSOLVERLBFGS RANDOMSTATE0  
C NPLOGSPACE2 7 10  
  
LINEARSVCRANDOMSTATE0  
C NPLOGSPACE2 7 10  
  
MAKEPIPELINE  
KBINSDISCRETIZERENCONEHOT  
LOGISTICREGRESSIONSOLVERLBFGS RANDOMSTATE0  
KBINSDISCRETIZERNBINS NPARANGE2 10  
LOGISTICREGRESSIONC NPLOGSPACE2 7 10  
  
MAKEPIPELINE  
KBINSDISCRETIZERENCONEHOT LINEARSVCRANDOMSTATE0  
KBINSDISCRETIZERNBINS NPARANGE2 10  
LINEARSVCC NPLOGSPACE2 7 10  
525 PREPROCESSING 1375

SCIKITLEARN USER GUIDE RELEASE 0213

GRADIENTBOOSTINGCLASSIFIERNESTIMATORS50 RANDOMSTATE0  
LEARNINGRATE NPLOGSPACE4 0 10

SVCRANDOMSTATE0 GAMMASCALE  
C NPLOGSPACE2 7 10

NAMES GETNAMEE FORE GINCLASSIFIERS  
NSAMPLES 100  
DATASETS  
MAKEMOONSNSAMPLESNSAMPLES NOISE02 RANDOMSTATE0  
MAKECIRCLESNSAMPLESNSAMPLES NOISE02 FACTOR05 RANDOMSTATE1  
MAKECLASSIFICATIONNSAMPLESNSAMPLES NFEATURES2 NREDUNDANT0  
NINFORMATIVE2 RANDOMSTATE2  
NCLUSTERSPERCLASS1

FIG AXES PLTSUBPLOTSNROWSLENDATASETS NCOLSLENCLASSIFIERS 1  
FIGSIZE21 9  
CM PLTCMPIYG  
CMBRIGHT LISTEDCOLORMAPB30065 178000  
ITERATE OVER DATASETS  
FORDSCNT X Y INENUMERATEDATASETS  
PRINTNDATASET DN DSCNT  
PREPROCESS DATASET SPLIT INTO TRAINING AND TEST PART  
X STANDARDSCALERFITTRANSFORMX  
XTRAIN XTEST YTRAIN YTEST TRAINTESTSPLIT  
X Y TESTSIZE5 RANDOMSTATE42  
CREATE THE GRID FOR BACKGROUND COLORS  
XMIN XMAX X 0MIN 5 X 0MAX 5  
YMIN YMAX X 1MIN 5 X 1MAX 5  
XX YY NPMESHGRID  
NPARANGEXMIN XMAX H NPARANGEYMIN YMAX H  
PLOT THE DATASET FIRST  
AX AXESDSCNT 0  
IFDSCNT 0  
AXSETTITLEINPUT DATA  
PLOT THE TRAINING POINTS  
AXSCATTERXTRAIN 0 XTRAIN 1 CYTRAIN CMAPCMBRIGHT  
EDGECOLORSK  
AND TESTING POINTS  
AXSCATTERXTEST 0 XTEST 1 CYTEST CMAPCMBRIGHT ALPHA06  
EDGECOLORSK  
AXSETXLIMXXMIN XXMAX  
AXSETYLIMYYMIN YYMAX  
AXSETXTICKS  
AXSETYTICKS  
ITERATE OVER CLASSIFIERS  
FORESTIDX NAME ESTIMATOR PARAMGRID IN  
1376 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
ENUMERATEZIPNAMES CLASSIFIERS  
AX AXESDSCNT ESTIDX 1  
CLF GRIDSEARCHCVESTIMATORESTIMATOR PARAMGRIDPARAMGRID CV5  
IIDFALSE  
WITHIGNOREWARNINGSCATEGORYCONVERGENCEWARNING  
CLFFITXTRAIN YTRAIN  
SCORE CLFScoreXTEST YTEST  
PRINTS2F NAME SCORE  
PLOT THE DECISION BOUNDARY FOR THAT WE WILL ASSIGN A COLOR TO EACH  
POINT IN THE MESH XMIN XMAX YMIN YMAX  
IFHASATTRCLF DECISIONFUNCTION  
Z CLFDECISIONFUNCTIONNPCXXRAVEL YYRAVEL  
ELSE  
Z CLFPREDICTPROBANPCXXRAVEL YYRAVEL 1  
PUT THE RESULT INTO A COLOR PLOT  
Z ZRESHAPEXXSHAPE  
AXCONTOURFXX YY Z CMAPCM ALPHA8  
PLOT THE TRAINING POINTS  
AXSCATTERXTRAIN 0 XTRAIN 1 CYTRAIN CMAPCMBRIGHT  
EDGECOLORSK  
AND TESTING POINTS  
AXSCATTERXTEST 0 XTEST 1 CYTEST CMAPCMBRIGHT  
EDGECOLORSK ALPHA06  
AXSETXLMXXMIN XXMAX  
AXSETYLIMYYMIN YYMAX  
AXSETXTICKS  
AXSETYTICKS  
IFDSCNT 0  
AXSETTITLENAMEREPLACE N  
AXTEXT095 006 2F SCORELSTRIP0 SIZE15  
BBOXDICTBOXSTYLEROUND ALPHA08 FACECOLORWHITE  
TRANSFORMAXTRANSAXES HORIZONTALALIGNMENTRIGHT  
PLTTIGHTLAYOUT  
ADD SUPTITLES ABOVE THE FIGURE  
PLTSUBPLOTSADJUSTTOP090  
SUPTITLES  
LINEAR CLASSIFIERS  
FEATURE DISCRETIZATION AND LINEAR CLASSIFIERS  
NONLINEAR CLASSIFIERS  
  
FORI SUPTITLE INZIP1 3 5 SUPTITLES  
AX AXES0 I  
AXTEXT105 125 SUPTITLE TRANSFORMAXTRANSAXES  
HORIZONTALALIGNMENTCENTER SIZEXLARGE  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 16451 SECONDS  
525 PREPROCESSING 1377

SCIKITLEARN USER GUIDE RELEASE 0213

NOTE [CLICK HERE](#) TO DOWNLOAD THE FULL EXAMPLE CODE

5257 COMPARE THE EFFECT OF DIFFERENT SCALERS ON DATA WITH OUTLIERS

FEATURE 0 MEDIAN INCOME IN A BLOCK AND FEATURE 5 NUMBER OF HOUSEHOLDS OF THE CALIFORNIA HOUSING DATASET HAVE VERY DIFFERENT SCALES AND CONTAIN SOME VERY LARGE OUTLIERS THESE TWO CHARACTERISTICS LEAD TO DIFFICULTIES TO VISUALIZE THE DATA AND MORE IMPORTANTLY THEY CAN DEGRADE THE PREDICTIVE PERFORMANCE OF MANY MACHINE LEARNING ALGORITHMS UNSCALED DATA CAN ALSO SLOW DOWN OR EVEN PREVENT THE CONVERGENCE OF MANY GRADIENTBASED ESTIMATORS INDEED MANY ESTIMATORS ARE DESIGNED WITH THE ASSUMPTION THAT EACH FEATURE TAKES VALUES CLOSE TO ZERO OR MORE IM

PORTANTLY THAT ALL FEATURES VARY ON COMPARABLE SCALES IN PARTICULAR METRICBASED AND GRADIENTBASED ESTIMATORS OFTEN ASSUME APPROXIMATELY STANDARDIZED DATA CENTERED FEATURES WITH UNIT VARIANCES A NOTABLE EXCEPTION ARE DECISION TREEBASED ESTIMATORS THAT ARE ROBUST TO ARBITRARY SCALING OF THE DATA

THIS EXAMPLE USES DIFFERENT SCALERS TRANSFORMERS AND NORMALIZERS TO BRING THE DATA WITHIN A PREDEFINED RANGE SCALERS ARE LINEAR OR MORE PRECISELY AFFINE TRANSFORMERS AND DIFFER FROM EACH OTHER IN THE WAY TO ESTIMATE THE PARAMETERS USED TO SHIFT AND SCALE EACH FEATURE QUANTILETRANSFORMER PROVIDES NONLINEAR TRANSFORMATIONS IN WHICH DISTANCES BETWEEN MARGINAL OUTLIERS AND INLIERS ARE SHRUNK POWERTRANSFORMER PROVIDES NONLINEAR TRANSFORMATIONS IN WHICH DATA IS MAPPED TO A NORMAL DISTRIBUTION TO STABILIZE VARIANCE AND MINIMIZE SKEWNESS UNLIKE THE PREVIOUS TRANSFORMATIONS NORMALIZATION REFERS TO A PER SAMPLE TRANSFORMATION INSTEAD OF A PER FEATURE TRANSFORMATION

THE FOLLOWING CODE IS A BIT VERBOSE FEEL FREE TO JUMP DIRECTLY TO THE ANALYSIS OF THE RESULTS

AUTHOR RAGHAV RV RVRAGHAV93GMAILCOM

GUILLAUME LEMAITRE GLEMAITRE58GMAILCOM

THOMAS UNTERTHINER

LICENSE BSD 3 CLAUSE

IMPORT NUMPY AS NP

IMPORT MATPLOTLIB AS MPL

FROM MATPLOTLIB IMPORT PYPLOTASPLT

FROM MATPLOTLIB IMPORT CM

FROM SKLEARNPREPROCESSING IMPORT MINMAXSCALER

FROM SKLEARNPREPROCESSING IMPORT MINMAXSCALE

FROM SKLEARNPREPROCESSING IMPORT MAXABSSCALER

FROM SKLEARNPREPROCESSING IMPORT STANDARDSCALER

FROM SKLEARNPREPROCESSING IMPORT ROBUSTSCALER

FROM SKLEARNPREPROCESSING IMPORT NORMALIZER

FROM SKLEARNPREPROCESSING IMPORT QUANTILETRANSFORMER

FROM SKLEARNPREPROCESSING IMPORT POWERTRANSFORMER

FROM SKLEARNDATASETS IMPORT FETCHCALIFORNIAHOUSING

PRINTDOC

DATASET FETCHCALIFORNIAHOUSING

XFULL YFULL DATASETDATA DATASETTARGET

TAKE ONLY 2 FEATURES TO MAKE VISUALIZATION EASIER

1378 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213  
FEATURE OF 0 HAS A LONG TAIL DISTRIBUTION  
FEATURE 5 HAS A FEW BUT VERY LARGE OUTLIERS  
X XFULL 0 5  
DISTRIBUTIONS  
UNSCALED DATA X  
DATA AFTER STANDARD SCALING  
STANDARDSCALERFITTRANSFORMX  
DATA AFTER MINMAX SCALING  
MINMAXSCALERFITTRANSFORMX  
DATA AFTER MAXABS SCALING  
MAXABSSCALERFITTRANSFORMX  
DATA AFTER ROBUST SCALING  
ROBUSTSCALERQUANTILERANGE25 75FITTRANSFORMX  
DATA AFTER POWER TRANSFORMATION YEOJOHNSON  
POWERTRANSFORMERMETHODYEOJOHNSONFITTRANSFORMX  
DATA AFTER POWER TRANSFORMATION BOXCOX  
POWERTRANSFORMERMETHODBOXCOXFITTRANSFORMX  
DATA AFTER QUANTILE TRANSFORMATION GAUSSIAN PDF  
QUANTILETRANSFORMEROUTPUTDISTRIBUTIONNORMAL  
FITTRANSFORMX  
DATA AFTER QUANTILE TRANSFORMATION UNIFORM PDF  
QUANTILETRANSFORMEROUTPUTDISTRIBUTIONUNIFORM  
FITTRANSFORMX  
DATA AFTER SAMPLEWISE L2 NORMALIZING  
NORMALIZERFITTRANSFORMX  
  
SCALE THE OUTPUT BETWEEN 0 AND 1 FOR THE COLORBAR  
Y MINMAXSCALEYFULL  
PLASMA DOES NOT EXIST IN MATPLOTLIB 15  
CMAP GETATTRCM PLASMAR CMHOTR  
DEFCREATEAXESTITLE FIGSIZE16 6  
FIG PLTFIGUREFIGSIZEFIGSIZE  
FIGSUPTITLETITLE  
DEFINE THE AXIS FOR THE FIRST PLOT  
LEFT WIDTH 01 022  
BOTTOM HEIGHT 01 07  
BOTTOMH HEIGHT 015  
LEFTH LEFT WIDTH 002  
RECTSCATTER LEFT BOTTOM WIDTH HEIGHT  
RECTHISTX LEFT BOTTOMH WIDTH 01  
RECTHISTY LEFTH BOTTOM 005 HEIGHT  
AXSCATTER PLTAXESRECTSCATTER  
AXHISTX PLTAXESRECTHISTX  
AXHISTY PLTAXESRECTHISTY  
DEFINE THE AXIS FOR THE ZOOMEDIN PLOT  
LEFT WIDTH LEFT 02  
LEFTH LEFT WIDTH 002  
RECTSCATTER LEFT BOTTOM WIDTH HEIGHT  
525 PREPROCESSING 1379

SCIKITLEARN USER GUIDE RELEASE 0213  
RECTHISTX LEFT BOTTOMH WIDTH 01  
RECTHISTY LEFTH BOTTOM 005 HEIGHT  
AXSCATTERZOOM PLTAXESRECTSCATTER  
AXHISTXZOOM PLTAXESRECTHISTX  
AXHISTYZOOM PLTAXESRECTHISTY  
DEFINE THE AXIS FOR THE COLORBAR  
LEFT WIDTH WIDTH LEFT 013 001  
RECTCOLORBAR LEFT BOTTOM WIDTH HEIGHT  
AXCOLORBAR PLTAXESRECTCOLORBAR  
RETURNAXSCATTER AXHISTY AXHISTX  
AXSCATTERZOOM AXHISTYZOOM AXHISTXZOOM  
AXCOLORBAR  
DEFPLOTDISTRIBUTIONAXES X Y HISTNBINS50 TITLE  
X0LABEL X1LABEL  
AX HISTX1 HISTX0 AXES  
AXSETTITLETITLE  
AXSETXLABELX0LABEL  
AXSETYLABELX1LABEL  
THE SCATTER PLOT  
COLORS CMAPY  
AXSCATTERX 0 X 1 ALPHA05 MARKERO S5 LW0 CCOLORS  
REMOVING THE TOP AND THE RIGHT SPINE FOR AESTHETICS  
MAKE NICE AXIS LAYOUT  
AXSPINESTOPSETVISIBLEFALSE  
AXSPINESRIGHTSETVISIBLEFALSE  
AXGETXAXISTICKBOTTOM  
AXGETYAXISTICKLEFT  
AXSPINESLEFTSETPOSITIONOUTWARD 10  
AXSPINESBOTTOMSETPOSITIONOUTWARD 10  
HISTOGRAM FOR AXIS X1 FEATURE 5  
HISTX1SETYLIMAXGETYLIM  
HISTX1HISTX 1 BINSHISTNBINS ORIENTATIONHORIZONTAL  
COLORGREY ECGREY  
HISTX1AXISOFF  
HISTOGRAM FOR AXIS X0 FEATURE 0  
HISTX0SETXLIMAXGETXLIM  
HISTX0HISTX 0 BINSHISTNBINS ORIENTATIONVERTICAL  
COLORGREY ECGREY  
HISTX0AXISOFF  
TWO PLOTS WILL BE SHOWN FOR EACH SCALERNORMALIZERTRANSFORMER THE LEFT FIGURE WILL SHOW A SCATTER PLOT OF THE FULL DATA  
SET WHILE THE RIGHT FIGURE WILL EXCLUDE THE EXTREME VALUES CONSIDERING ONLY 99 OF THE DATA SET EXCLUDING MARGINAL  
OUTLIERS IN ADDITION THE MARGINAL DISTRIBUTIONS FOR EACH FEATURE WILL BE SHOWN ON THE SIDE OF THE SCATTER PLOT  
DEMAKEPLOTITEMIDX  
TITLE X DISTRIBUTIONSITEMIDX  
AXZOOMOUT AXZOOMIN AXCOLORBAR CREATEAXESTITLE  
1380 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
 AXARR AXZOOMOUT AXZOOMIN  
 PLOTDISTRIBUTIONAXARR0 X Y HISTNBINS200  
 X0LABELMEDIAN INCOME  
 X1LABELNUMBER OF HOUSEHOLDS  
 TITLEFULL DATA  
 ZOOMIN  
 ZOOMINPERCENTILERANGE 0 99  
 CUTOFFSX0 NPPERCENTILEX 0 ZOOMINPERCENTILERANGE  
 CUTOFFSX1 NPPERCENTILEX 1 ZOOMINPERCENTILERANGE  
 NONOUTLIERSMASK  
 NPALLX CUTOFFSX00 CUTOFFSX10 AXIS1  
 NPALLX CUTOFFSX01 CUTOFFSX11 AXIS1  
 PLOTDISTRIBUTIONAXARR1 XNONOUTLIERSMASK YNONOUTLIERSMASK  
 HISTNBINS50  
 X0LABELMEDIAN INCOME  
 X1LABELNUMBER OF HOUSEHOLDS  
 TITLEZOOMIN  
 NORM MPLCOLORSNORMALIZEYFULLMIN YFULLMAX  
 MPLCOLORBARCOLORBARBASEAXCOLORBAR CMAPCMAP  
 NORMNORM ORIENTATIONVERTICAL  
 LABELCOLOR MAPPING FOR VALUES OF Y  
 ORIGINAL DATA  
 EACH TRANSFORMATION IS PLOTTED SHOWING TWO TRANSFORMED FEATURES WITH THE LEFT PLOT SHOWING THE ENTIRE DATASET AND THE RIGHT ZOOMEDIN TO SHOW THE DATASET WITHOUT THE MARGINAL OUTLIERS A LARGE MAJORITY OF THE SAMPLES ARE COMPACTED TO A SPECIFIC RANGE 0 10 FOR THE MEDIAN INCOME AND 0 6 FOR THE NUMBER OF HOUSEHOLDS NOTE THAT THERE ARE SOME MARGINAL OUTLIERS SOME BLOCKS HAVE MORE THAN 1200 HOUSEHOLDS THEREFORE A SPECIFIC PREPROCESSING CAN BE VERY BENEFICIAL DEPENDING OF THE APPLICATION IN THE FOLLOWING WE PRESENT SOME INSIGHTS AND BEHAVIORS OF THOSE PREPROCESSING METHODS IN THE PRESENCE OF MARGINAL OUTLIERS  
 MAKEPLOT0  
 525 PREPROCESSING 1381

SCIKITLEARN USER GUIDE RELEASE 0213

STANDARDSCALER

STANDARDSCALER REMOVES THE MEAN AND SCALES THE DATA TO UNIT VARIANCE HOWEVER THE OUTLIERS HAVE AN INFLUENCE WHEN COMPUTING THE EMPIRICAL MEAN AND STANDARD DEVIATION WHICH SHRINK THE RANGE OF THE FEATURE VALUES AS SHOWN IN THE LEFT FIGURE BELOW NOTE IN PARTICULAR THAT BECAUSE THE OUTLIERS ON EACH FEATURE HAVE DIFFERENT MAGNITUDES THE SPREAD OF THE TRANSFORMED DATA ON EACH FEATURE IS VERY DIFFERENT MOST OF THE DATA LIE IN THE 2 4 RANGE FOR THE TRANSFORMED MEDIAN INCOME FEATURE WHILE THE SAME DATA IS SQUEEZED IN THE SMALLER 02 02 RANGE FOR THE TRANSFORMED NUMBER OF HOUSEHOLDS

STANDARDSCALER THEREFORE CANNOT GUARANTEE BALANCED FEATURE SCALES IN THE PRESENCE OF OUTLIERS

MAKEPLOT1

MINMAXSCALER

MINMAXSCALER RESCALES THE DATA SET SUCH THAT ALL FEATURE VALUES ARE IN THE RANGE 0 1 AS SHOWN IN THE RIGHT PANEL BELOW HOWEVER THIS SCALING COMPRESS ALL INLIERS IN THE NARROW RANGE 0 0005 FOR THE TRANSFORMED NUMBER OF HOUSEHOLDS

ASSTANDARDSCALER MINMAXSCALER IS VERY SENSITIVE TO THE PRESENCE OF OUTLIERS

MAKEPLOT2

1382 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

MAXABSSCALER

MAXABSSCALER DIFFERS FROM THE PREVIOUS SCALER SUCH THAT THE ABSOLUTE VALUES ARE MAPPED IN THE RANGE 0 1 ON POSITIVE ONLY DATA THIS SCALER BEHAVES SIMILARLY TO MINMAXSCALER AND THEREFORE ALSO SUFFERS FROM THE PRESENCE OF LARGE OUTLIERS

MAKEPLOT3

ROBUSTSCALER

UNLIKE THE PREVIOUS SCALERS THE CENTERING AND SCALING STATISTICS OF THIS SCALER ARE BASED ON PERCENTILES AND ARE THEREFORE NOT INFLUENCED BY A FEW NUMBER OF VERY LARGE MARGINAL OUTLIERS CONSEQUENTLY THE RESULTING RANGE OF THE TRANSFORMED FEATURE VALUES IS LARGER THAN FOR THE PREVIOUS SCALERS AND MORE IMPORTANTLY ARE APPROXIMATELY SIMILAR FOR BOTH FEATURES MOST OF THE TRANSFORMED VALUES LIE IN A 2 3 RANGE AS SEEN IN THE ZOOMEDIN FIGURE NOTE THAT THE OUTLIERS THEMSELVES ARE STILL PRESENT IN THE TRANSFORMED DATA IF A SEPARATE OUTLIER CLIPPING IS DESIRABLE A NONLINEAR TRANSFORMATION IS REQUIRED SEE BELOW

MAKEPLOT4

POWERTRANSFORMER

POWERTRANSFORMER APPLIES A POWER TRANSFORMATION TO EACH FEATURE TO MAKE THE DATA MORE GAUSSIANLIKE CUR  
525 PREPROCESSING 1383

SCIKITLEARN USER GUIDE RELEASE 0213

RENTLYPOWERTRANSFORMER IMPLEMENTS THE YEOJOHNSON AND BOXCOX TRANSFORMS THE POWER TRANSFORM FINDS THE OPTIMAL SCALING FACTOR TO STABILIZE VARIANCE AND MIMIMIZE SKEWNESS THROUGH MAXIMUM LIKELIHOOD ESTIMATION BY DEFAULTPOWERTRANSFORMER ALSO APPLIES ZEROMEAN UNIT VARIANCE NORMALIZATION TO THE TRANSFORMED OUTPUT NOTE THAT BOXCOX CAN ONLY BE APPLIED TO STRICTLY POSITIVE DATA INCOME AND NUMBER OF HOUSEHOLDS HAPPEN TO BE STRICTLY POSITIVE BUT IF NEGATIVE VALUES ARE PRESENT THE YEOJOHNSON TRANSFORMED IS TO BE PREFERRED

MAKEPLOT5

MAKEPLOT6

- 
- 

QUANTILETRANSFORMER GAUSSIAN OUTPUT

QUANTILETRANSFORMER HAS AN ADDITIONAL OUTPUTDISTRIBUTION PARAMETER ALLOWING TO MATCH A GAUSSIAN DISTRIBUTION INSTEAD OF A UNIFORM DISTRIBUTION NOTE THAT THIS NONPARAMETETRIC TRANSFORMER INTRODUCES SATURATION ARTI FACTS FOR EXTREME VALUES

MAKEPLOT7

1384 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

QUANTILETRANSFORMER UNIFORM OUTPUT

QUANTILETRANSFORMER APPLIES A NONLINEAR TRANSFORMATION SUCH THAT THE PROBABILITY DENSITY FUNCTION OF EACH FEATURE WILL BE MAPPED TO A UNIFORM DISTRIBUTION IN THIS CASE ALL THE DATA WILL BE MAPPED IN THE RANGE 0 1 EVEN THE OUTLIERS WHICH CANNOT BE DISTINGUISHED ANYMORE FROM THE INLIERS

ASROBUSTSCALER QUANTILETRANSFORMER IS ROBUST TO OUTLIERS IN THE SENSE THAT ADDING OR REMOVING OUTLIERS IN THE TRAINING SET WILL YIELD APPROXIMATELY THE SAME TRANSFORMATION ON HELD OUT DATA BUT CONTRARY TO ROBUSTSCALER QUANTILETRANSFORMER WILL ALSO AUTOMATICALLY COLLAPSE ANY OUTLIER BY SETTING THEM TO THE A PRIORI DEFINED RANGE BOUNDARIES 0 AND 1

MAKEPLOT8

NORMALIZER

THENORMALIZER RESCALES THE VECTOR FOR EACH SAMPLE TO HAVE UNIT NORM INDEPENDENTLY OF THE DISTRIBUTION OF THE SAMPLES IT CAN BE SEEN ON BOTH FIGURES BELOW WHERE ALL SAMPLES ARE MAPPED ONTO THE UNIT CIRCLE IN OUR EXAMPLE THE TWO SELECTED FEATURES HAVE ONLY POSITIVE VALUES THEREFORE THE TRANSFORMED DATA ONLY LIE IN THE POSITIVE QUADRANT THIS WOULD NOT BE THE CASE IF SOME ORIGINAL FEATURES HAD A MIX OF POSITIVE AND NEGATIVE VALUES

525 PREPROCESSING 1385

SCIKITLEARN USER GUIDE RELEASE 0213  
MAKEPLOT9  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 4599 SECONDS  
526 SEMI SUPERVISED CLASSIFICATION  
EXAMPLES CONCERNING THE SKLEARNSEMISUPERVISED MODULE  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
5261 DECISION BOUNDARY OF LABEL PROPAGATION VERSUS SVM ON THE IRIS DATASET  
COMPARISON FOR DECISION BOUNDARY GENERATED ON IRIS DATASET BETWEEN LABEL PROPAGATION AND SVM  
THIS DEMONSTRATES LABEL PROPAGATION LEARNING A GOOD BOUNDARY EVEN WITH A SMALL AMOUNT OF LABELED DATA  
1386 CHAPTER 5 EXAMPLES



```
SCIKITLEARN USER GUIDE RELEASE 0213
PRINTDOC
AUTHORS CLAY WOOLAM CLAYWOOLAMORG
LICENSE BSD
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARN IMPORT DATASETS
FROM SKLEARN IMPORT SVM
FROM SKLEARNSEMISSUPERVISED IMPORT LABELPROPAGATION
RNG NPRANDOMRANDOMSTATE0
IRIS DATASETSLOADIRIS
X IRISDATA 2
Y IRISTARGET
STEP SIZE IN THE MESH
H 02
Y30 NPCOPYY
Y30RNGRANDLENY 03 1
Y50 NPCOPYY
Y50RNGRANDLENY 05 1
WE CREATE AN INSTANCE OF SVM AND FIT OUT DATA WE DO NOT SCALE OUR
526 SEMI SUPERVISED CLASSIFICATION 1387
```

SCIKITLEARN USER GUIDE RELEASE 0213

DATA SINCE WE WANT TO PLOT THE SUPPORT VECTORS

LS30 LABELPROPAGATIONLABELSPREADINGFITX Y30

Y30

LS50 LABELPROPAGATIONLABELSPREADINGFITX Y50

Y50

LS100 LABELPROPAGATIONLABELSPREADINGFITX Y Y

RBFSVC SVM SVCKERNELRBF GAMMA5FITX Y Y

CREATE A MESH TO PLOT IN

XMIN XMAX X 0MIN 1 X 0MAX 1

YMIN YMAX X 1MIN 1 X 1MAX 1

XX YY NPMESHGRIDNPARANGEXMIN XMAX H

NPARANGEYMIN YMAX H

TITLE FOR THE PLOTS

TITLES LABEL SPREADING 30 DATA

LABEL SPREADING 50 DATA

LABEL SPREADING 100 DATA

SVC WITH RBF KERNEL

COLORMAP 1 1 1 1 0 0 0 9 1 1 0 0 2 8 6 0

FORI CLF YTRAIN INENUMERATELS30 LS50 LS100 RBFSVC

PLOT THE DECISION BOUNDARY FOR THAT WE WILL ASSIGN A COLOR TO EACH

POINT IN THE MESH XMIN XMAXXYMIN YMAX

PLTSUBPLOT2 2 | 1

Z CLFPREDICTNPCXXRAVEL YYRAVEL

PUT THE RESULT INTO A COLOR PLOT

Z ZRESHAPEXXSHAPE

PLTCONTOURFXX YY Z CMAPPLTCMPAIED

PLTAXISOFF

PLOT ALSO THE TRAINING POINTS

COLORS COLORMAPY FORYINYTRAIN

PLTSCATTERX 0 X 1 CCOLORS EDGECOLORSBLACK

PLTTITLETITLES

PLTSUPTITLEUNLABELED POINTS ARE COLORED WHITE Y01

PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0915 SECONDS

NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)

5262 LABEL PROPAGATION LEARNING A COMPLEX STRUCTURE

EXAMPLE OF LABELPROPAGATION LEARNING A COMPLEX INTERNAL STRUCTURE TO DEMONSTRATE “MANIFOLD LEARNING” THE OUTER

CIRCLE SHOULD BE LABELED “RED” AND THE INNER CIRCLE “BLUE” BECAUSE BOTH LABEL GROUPS LIE INSIDE THEIR OWN DISTINCT

SHAPE WE CAN SEE THAT THE LABELS PROPAGATE CORRECTLY AROUND THE CIRCLE

1388 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
PRINTDOC
AUTHORS CLAY WOOLAM CLAYWOOLAMORG
ANDREAS MUELLER AMUELLERAISUNIBONNDE
LICENSE BSD
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARNSEMISUPERVISED IMPORT LABELPROPAGATION
FROM SKLEARNDATASETS IMPORT MAKECIRCLES
GENERATE RING WITH INNER BOX
NSAMPLES 200
X Y MAKECIRCLESNSAMPLESNSAMPLES SHUFFLEFALSE
OUTER INNER 0 1
LABELS NPFULLNSAMPLES 1
LABELS0 OUTER
LABELS1 INNER

LEARN WITH LABELSPREADING
LABELSPREAD LABELPROPAGATIONLABELSPREADINGKERNELKNN ALPHA08
LABELSPREADFITX LABELS

PLOT OUTPUT LABELS
OUTPUTLABELS LABELSPREADTRANSDUCTION
PLTFIGUREFIGSIZE85 4
PLTSUBPLOT1 2 1
PLTSCATTERXLABELS OUTER 0 XLABELS OUTER 1 COLORNAVY
MARKERS LW0 LABELOUTER LABELED S10
PLTSCATTERXLABELS INNER 0 XLABELS INNER 1 COLORC
MARKERS LW0 LABELINNER LABELED S10
PLTSCATTERXLABELS 1 0 XLABELS 1 1 COLORDARKORANGE
MARKER LABELUNLABELED
PLTLEGENDSCATTERPOINTS1 SHADOWFALSE LOCUPPER RIGHT
PLTTITLERAW DATA 2 CLASSESOUTER AND INNER
526 SEMI SUPERVISED CLASSIFICATION 1389
```

SCIKITLEARN USER GUIDE RELEASE 0213

PLTSUBPLOT1 2 2

OUTPUTLABELARRAY NPASARRAYOUTPUTLABELS

OUTERNUMBERS NPWHEREOUTPUTLABELARRAY OUTER0

INNERNUMBERS NPWHEREOUTPUTLABELARRAY INNER0

PLTSCATTERXOUTERNUMBERS 0 XOUTERNUMBERS 1 COLORNAVY

MARKERS LW0 S10 LABEL OUTER LEARNED

PLTSCATTERXINNERNUMBERS 0 XINNERNUMBERS 1 COLORC

MARKERS LW0 S10 LABEL INNER LEARNED

PLTLEGENDSCATTERPOINTS1 SHADOWFALSE LOCUPPER RIGHT

PLTTITLELABELS LEARNED WITH LABEL SPREADING KNN

PLTSUBPLOTSADJUSTLEFT007 BOTTOM007 RIGHT093 TOP092

PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0031 SECONDS

NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)

5263 LABEL PROPAGATION DIGITS DEMONSTRATING PERFORMANCE

THIS EXAMPLE DEMONSTRATES THE POWER OF SEMISUPERVISED LEARNING BY TRAINING A LABEL SPREADING MODEL TO CLASSIFY HANDWRITTEN DIGITS WITH SETS OF VERY FEW LABELS

THE HANDWRITTEN DIGIT DATASET HAS 1797 TOTAL POINTS THE MODEL WILL BE TRAINED USING ALL POINTS BUT ONLY 30 WILL BE LABELED RESULTS IN THE FORM OF A CONFUSION MATRIX AND A SERIES OF METRICS OVER EACH CLASS WILL BE VERY GOOD

AT THE END THE TOP 10 MOST UNCERTAIN PREDICTIONS WILL BE SHOWN

1390 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
 OUT  
 LABEL SPREADING MODEL 40 LABELED 300 UNLABELED POINTS 340 TOTAL  
 PRECISION RECALL F1SCORE SUPPORT  
 0 100 100 100 27  
 1 082 100 090 37  
 2 100 086 092 28  
 3 100 080 089 35  
 4 092 100 096 24  
 5 074 094 083 34  
 6 089 096 092 25  
 7 094 089 091 35  
 8 100 068 081 31  
 9 081 088 084 24  
 ACCURACY 090 300  
 MACRO AVG 091 090 090 300  
 WEIGHTED AVG 091 090 090 300  
 CONFUSION MATRIX  
 27 0 0 0 0 0 0 0 0 0  
 0 37 0 0 0 0 0 0 0  
 0 1 24 0 0 0 2 1 0 0  
 0 0 0 28 0 5 0 1 0 1  
 0 0 0 0 24 0 0 0 0 0  
 0 0 0 0 0 32 0 0 0 2  
 526 SEMI SUPERVISED CLASSIFICATION 1391

```
SCIKITLEARN USER GUIDE RELEASE 0213
0 0 0 0 0 1 24 0 0 0
0 0 0 0 1 3 0 31 0 0
0 7 0 0 0 0 1 0 21 2
0 0 0 0 1 2 0 0 0 21
PRINTDOC
AUTHORS CLAY WOOLAM CLAYWOOLAMORG
LICENSE BSD
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SCIPY IMPORT STATS
FROM SKLEARN IMPORT DATASETS
FROM SKLEARNSEMISSUPERVISED IMPORT LABELPROPAGATION
FROM SKLEARNMETRICS IMPORT CONFUSIONMATRIX CLASSIFICATIONREPORT
DIGITS DATASETSLOADDIGITS
RNG NPRANDOMRANDOMSTATE2
INDICES NPARANGELENDIGITSDATA
RNGSHUFFLEINDICES
X DIGITSDATAINDICES340
Y DIGITSTARGETINDICES340
IMAGES DIGITSIMAGESINDICES340
NTOTALSAMPLES LENY
NLABELEDPOINTS 40
INDICES NPARANGENTOTALSAMPLES
UNLABELEDSET INDICESNLABELEDPOINTS

SHUFFLE EVERYTHING AROUND
YTRAIN NPCOPYY
YTRAINUNLABELEDSET 1

LEARN WITH LABELSPREADING
LPMODEL LABELPROPAGATIONLABELSPREADINGGAMMA25 MAXITER20
LPMODELFITX YTRAIN
PREDICTEDLABELS LPMODELTRANSDUCTIONUNLABELEDSET
TRUELABELS YUNLABELEDSET
CM CONFUSIONMATRIXTRUELABELS PREDICTEDLABELS LABELSLPMODELCLASSES
PRINTLABEL SPREADING MODEL DLABELED DUNLABELED POINTS DTOTAL
NLABELEDPOINTS NTOTALSAMPLES NLABELEDPOINTS NTOTALSAMPLES
1392 CHAPTER 5 EXAMPLES
```

```
SCIKITLEARN USER GUIDE RELEASE 0213
PRINTCLASSIFICATIONREPORTTRUELABELS PREDICTEDLABELS
PRINTCONFUSION MATRIX
PRINTCM

CALCULATE UNCERTAINTY VALUES FOR EACH TRANSDUCED DISTRIBUTION
PREDENTROPIES STATSDISTRIBUTIONSENTROPYLPMODELLABELDISTRIBUTIONST

PICK THE TOP 10 MOST UNCERTAIN LABELS
UNCERTAINTYINDEX NPARGSORTPREDENTROPIES10

PLOT
F PLTFIGUREFIGSIZE7 5
FORINDEX IMAGEINDEX INENUMERATEUNCERTAINTYINDEX
IMAGE IMAGESIMAGEINDEX
SUB FADDSUBPLOT2 5 INDEX 1
SUBIMSHOWIMAGE CMAPPLTCMGRAYR
PLTXTICKS
PLTYTICKS
SUBSETTITLEPREDICT INTRUEI
LPMODELTRANSDUCTIONIMAGEINDEX YIMAGEINDEX
FSUPTITLELEARNING WITH SMALL AMOUNT OF LABELED DATA
PLTSHOW
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0225 SECONDS
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE
5264 LABEL PROPAGATION DIGITS ACTIVE LEARNING
DEMONSTRATES AN ACTIVE LEARNING TECHNIQUE TO LEARN HANDWRITTEN DIGITS USING LABEL PROPAGATION
WE START BY TRAINING A LABEL PROPAGATION MODEL WITH ONLY 10 LABELED POINTS THEN WE SELECT THE TOP FIVE MOST UNCERTAIN
POINTS TO LABEL NEXT WE TRAIN WITH 15 LABELED POINTS ORIGINAL 10 5 NEW ONES WE REPEAT THIS PROCESS FOUR TIMES
TO HAVE A MODEL TRAINED WITH 30 LABELED EXAMPLES NOTE YOU CAN INCREASE THIS TO LABEL MORE THAN 30 BY CHANGING
MAXITERATIONS LABELING MORE THAN 30 CAN BE USEFUL TO GET A SENSE FOR THE SPEED OF CONVERGENCE OF THIS ACTIVE
LEARNING TECHNIQUE
A PLOT WILL APPEAR SHOWING THE TOP 5 MOST UNCERTAIN DIGITS FOR EACH ITERATION OF TRAINING THESE MAY OR MAY NOT CONTAIN
MISTAKES BUT WE WILL TRAIN THE NEXT MODEL WITH THEIR TRUE LABELS
526 SEMI SUPERVISED CLASSIFICATION 1393
```

SCIKITLEARN USER GUIDE RELEASE 0213  
OUT  
ITERATION 0  
LABEL SPREADING MODEL 40 LABELED 290 UNLABELED 330 TOTAL  
PRECISION RECALL F1SCORE SUPPORT  
0 100 100 100 22  
1 078 069 073 26  
2 093 093 093 29  
3 100 089 094 27  
4 092 096 094 23  
5 096 070 081 33  
6 097 097 097 35  
7 094 091 092 33  
8 062 089 074 28  
9 073 079 076 34  
ACCURACY 087 290  
MACRO AVG 089 087 087 290  
WEIGHTED AVG 088 087 087 290  
CONFUSION MATRIX  
22 0 0 0 0 0 0 0 0  
0 18 2 0 0 0 1 0 5 0  
0 0 27 0 0 0 0 0 2 0  
0 0 0 24 0 0 0 0 3 0  
1394 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213

0 1 0 0 22 0 0 0 0  
0 0 0 0 0 23 0 0 0 10  
0 1 0 0 0 0 34 0 0 0  
0 0 0 0 0 0 0 30 3 0  
0 3 0 0 0 0 0 0 25 0  
0 0 0 0 2 1 0 2 2 27

ITERATION 1

LABEL SPREADING MODEL 45 LABELED 285 UNLABELED 330 TOTAL

PRECISION RECALL F1SCORE SUPPORT

0 100 100 100 22  
1 079 100 088 22  
2 100 093 096 29  
3 100 100 100 26  
4 092 096 094 23  
5 096 070 081 33  
6 100 097 099 35  
7 094 091 092 33  
8 077 086 081 28  
9 073 079 076 34

ACCURACY 090 285

MACRO AVG 091 091 091 285

WEIGHTED AVG 091 090 090 285

CONFUSION MATRIX

22 0 0 0 0 0 0 0 0 0  
0 22 0 0 0 0 0 0 0 0  
0 0 27 0 0 0 0 0 2 0  
0 0 0 26 0 0 0 0 0 0  
0 1 0 0 22 0 0 0 0 0  
0 0 0 0 0 23 0 0 0 10  
0 1 0 0 0 0 34 0 0 0  
0 0 0 0 0 0 0 30 3 0  
0 4 0 0 0 0 0 0 24 0  
0 0 0 0 2 1 0 2 2 27

ITERATION 2

LABEL SPREADING MODEL 50 LABELED 280 UNLABELED 330 TOTAL

PRECISION RECALL F1SCORE SUPPORT

0 100 100 100 22  
1 085 100 092 22  
2 100 100 100 28  
3 100 100 100 26  
4 087 100 093 20  
5 096 070 081 33  
6 100 097 099 35  
7 094 100 097 32  
8 092 086 089 28  
9 073 079 076 34

ACCURACY 092 280

MACRO AVG 093 093 093 280

WEIGHTED AVG 093 092 092 280

CONFUSION MATRIX

22 0 0 0 0 0 0 0 0 0  
0 22 0 0 0 0 0 0 0 0

526 SEMI SUPERVISED CLASSIFICATION 1395

SCIKITLEARN USER GUIDE RELEASE 0213

0 0 28 0 0 0 0 0 0  
0 0 0 26 0 0 0 0 0  
0 0 0 0 20 0 0 0 0  
0 0 0 0 0 23 0 0 0 10  
0 1 0 0 0 0 34 0 0 0  
0 0 0 0 0 0 0 32 0 0  
0 3 0 0 1 0 0 0 24 0  
0 0 0 0 2 1 0 2 2 27

ITERATION 3

LABEL SPREADING MODEL 55 LABELED 275 UNLABELED 330 TOTAL

PRECISION RECALL F1SCORE SUPPORT

0 100 100 100 22  
1 085 100 092 22  
2 100 100 100 27  
3 100 100 100 26  
4 087 100 093 20  
5 096 087 092 31  
6 100 097 099 35  
7 100 100 100 31  
8 092 086 089 28  
9 088 085 086 33

ACCURACY 095 275

MACRO AVG 095 095 095 275

WEIGHTED AVG 095 095 095 275

CONFUSION MATRIX

22 0 0 0 0 0 0 0 0 0  
0 22 0 0 0 0 0 0 0 0  
0 0 27 0 0 0 0 0 0 0  
0 0 0 26 0 0 0 0 0 0  
0 0 0 0 20 0 0 0 0 0  
0 0 0 0 0 27 0 0 0 4  
0 1 0 0 0 0 34 0 0 0  
0 0 0 0 0 0 0 31 0 0  
0 3 0 0 1 0 0 0 24 0  
0 0 0 0 2 1 0 0 2 28

ITERATION 4

LABEL SPREADING MODEL 60 LABELED 270 UNLABELED 330 TOTAL

PRECISION RECALL F1SCORE SUPPORT

0 100 100 100 22  
1 096 100 098 22  
2 100 096 098 27  
3 096 100 098 25  
4 086 100 093 19  
5 096 087 092 31  
6 100 097 099 35  
7 100 100 100 31  
8 092 096 094 25  
9 088 085 086 33

ACCURACY 096 270

MACRO AVG 095 096 096 270

WEIGHTED AVG 096 096 096 270

CONFUSION MATRIX

1396 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
22 0 0 0 0 0 0 0 0 0
0 22 0 0 0 0 0 0 0 0
0 0 26 1 0 0 0 0 0 0
0 0 0 25 0 0 0 0 0 0
0 0 0 0 19 0 0 0 0 0
0 0 0 0 0 27 0 0 0 4
0 1 0 0 0 0 34 0 0 0
0 0 0 0 0 0 0 31 0 0
0 0 0 0 1 0 0 0 24 0
0 0 0 0 2 1 0 0 2 28
PRINTDOC
AUTHORS CLAY WOOLAM CLAYWOOLAMORG
LICENSE BSD
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SCIPY IMPORT STATS
FROM SKLEARN IMPORT DATASETS
FROM SKLEARNSEMISSUPERVISED IMPORT LABELPROPAGATION
FROM SKLEARNMETRICS IMPORT CLASSIFICATIONREPORT CONFUSIONMATRIX
DIGITS DATASETSLOADDIGITS
RNG NPRANDOMRANDOMSTATE0
INDICES NPARANGELENDIGITSDATA
RNGSHUFFLEINDICES
X DIGITSDATAINDICES330
Y DIGITSTARGETINDICES330
IMAGES DIGITSIMAGESINDICES330
NTOTALSAMPLES LENY
NLABELEDPOINTS 40
MAXITERATIONS 5
UNLABELEDINDICES NPARANGENTOTALSAMPLESNLABELEDPOINTS
F PLTFigure
FORIINRANGEMAXITERATIONS
IFLENUNLABELEDINDICES 0
PRINTNO UNLABELED ITEMS LEFT TO LABEL
BREAK
YTRAIN NPCOPY
YTRAINUNLABELEDINDICES 1
LPModel LABELPROPAGATIONLABELSPREADINGGAMMA0.25 MAXITER20
LPModelFITX YTRAIN
PREDICTEDLABELS LPModelTRANSDUCTIONUNLABELEDINDICES
TRUELABELS YUNLABELEDINDICES
526 SEMI SUPERVISED CLASSIFICATION 1397
```

SCIKITLEARN USER GUIDE RELEASE 0213  
CM CONFUSIONMATRIXTRUELABELS PREDICTEDLABELS  
LABELSLPMODELCLASSES  
PRINTITERATION I S I 70  
PRINTLABEL SPREADING MODEL DLABELED DUNLABELED DTOTAL  
NLABELEDPOINTS NTOTALSAMPLES NLABELEDPOINTS  
NTOTALSAMPLES  
PRINTCLASSIFICATIONREPORTTRUELABELS PREDICTEDLABELS  
PRINTCONFUSION MATRIX  
PRINTCM  
COMPUTE THE ENTROPIES OF TRANSDUCED LABEL DISTRIBUTIONS  
PREDEENTROPIES STATSDISTRIBUTIONSENTROPY  
LPMODELLABELDISTRIBUTIONST  
SELECT UP TO 5 DIGIT EXAMPLES THAT THE CLASSIFIER IS MOST UNCERTAIN ABOUT  
UNCERTAINTYINDEX NPARGSORTPREDEENTROPIES1  
UNCERTAINTYINDEX UNCERTAINTYINDEX  
NPIN1DUNCERTAINTYINDEX UNLABELEDINDICES5  
KEEP TRACK OF INDICES THAT WE GET LABELS FOR  
DELETEINDICES NPARRAY DTYPEINT  
FOR MORE THAN 5 ITERATIONS VISUALIZE THE GAIN ONLY ON THE FIRST 5  
IFI 5  
FTEXT05 1 I 1 183  
MODELDNFFIT WITH NDLABELS  
I 1 I 5 10 SIZE10  
FORINDEX IMAGEINDEX INENUMERATEUNCERTAINTYINDEX  
IMAGE IMAGESIMAGEINDEX  
FOR MORE THAN 5 ITERATIONS VISUALIZE THE GAIN ONLY ON THE FIRST 5  
IFI 5  
SUB FADDSUBPLOT5 5 INDEX 1 5 I  
SUBIMSHOWIMAGE CMAPPLTCMGRAYR INTERPOLATIONNONE  
SUBSETTITLEPREDICT INTRUEI  
LPMODELTRANSDUCTIONIMAGEINDEX YIMAGEINDEX SIZE10  
SUBAXISOFF  
LABELING 5 POINTS REMOTE FROM LABELED SET  
DELETEINDEX NPWHEREUNLABELEDINDICES IMAGEINDEX  
DELETEINDICES NPCONCATENATEDELETEINDICES DELETEINDEX  
UNLABELEDINDICES NPDELETEUNLABELEDINDICES DELETEINDICES  
NLABELEDPOINTS LENUNCERTAINTYINDEX  
FSUPTITLEACTIVE LEARNING WITH LABEL PROPAGATION NROWS SHOW 5 MOST  
UNCERTAIN LABELS TO LEARN WITH THE NEXT MODEL Y115  
PLTSUBPLOTSADJUSTLEFT02 BOTTOM003 RIGHT09 TOP09 WSPACE02  
HSPACE085  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0575 SECONDS  
1398 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
527 SUPPORT VECTOR MACHINES  
EXAMPLES CONCERNING THE SKLEARN SVM MODULE  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
5271 NONLINEAR SVM  
PERFORM BINARY CLASSIFICATION USING NONLINEAR SVC WITH RBF KERNEL THE TARGET TO PREDICT IS A XOR OF THE INPUTS  
THE COLOR MAP ILLUSTRATES THE DECISION FUNCTION LEARNED BY THE SVC  
PRINTDOC  
IMPORT NUMPY AS NP  
IMPORT MATPLOTLIBPYPLOT AS PLT  
FROM SKLEARN IMPORT SVM  
XX YY NPMESHGRIDNPLINSPACE3 3 500  
NPLINSPACE3 3 500  
NPRANDOMSEED0  
X NPRANDOMRANDN300 2  
Y NPLOGICALXORX 0 0 X 1 0  
527 SUPPORT VECTOR MACHINES 1399

SCIKITLEARN USER GUIDE RELEASE 0213

FIT THE MODEL

CLF SVMNUSVCGAMMAAUTO

CLFFITX Y

PLOT THE DECISION FUNCTION FOR EACH DATAPOINT ON THE GRID

Z CLFDECISIONFUNCTIONNPCXXRAVEL YYRAVEL

Z ZRESHAPEXXSHAPE

PLTIMSHOWZ INTERPOLATIONNEAREST

EXTENTXXMIN XXMAX YYMIN YYMAX ASPECTAUTO

ORIGINLOWER CMAPPLTCMPUORR

CONTOURS PLTCONTOURXX YY Z LEVELS0 LINEWIDTHS2

LINESTYLES DASHED

PLTSCATTERX 0 X 1 S30 CY CMAPPLTCMPAIRE

EDGECOLORSK

PLXTICKS

PLTYTICKS

PLTAXIS3 3 3 3

PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 1073 SECONDS

NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)

5272 SVM MAXIMUM MARGIN SEPARATING HYPERPLANE

PLOT THE MAXIMUM MARGIN SEPARATING HYPERPLANE WITHIN A TWOCCLASS SEPARABLE DATASET USING A SUPPORT VECTOR MACHINE

CLASSIFIER WITH LINEAR KERNEL

1400 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
PRINTDOC
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARN IMPORT SVM
FROM SKLEARNDATASETS IMPORT MAKEBLOBS
    WE CREATE 40 SEPARABLE POINTS
X Y MAKEBLOBSNSAMPLES40 CENTERS2 RANDOMSTATE6
FIT THE MODEL DONT REGULARIZE FOR ILLUSTRATION PURPOSES
CLF SVMKERNELLINEAR C1000
CLFFITX Y
PLTSCATTERX 0 X 1 CY S30 CMAPPLTCMPAIED
    PLOT THE DECISION FUNCTION
AX PLTGCA
XLIM AXGETXLIM
YLIM AXGETYLIM
    CREATE GRID TO EVALUATE MODEL
XX NPLINSPACEXLIM0 XLIM1 30
YY NPLINSPACEYLIM0 YLIM1 30
YY XX NPMESHGRIDYY XX
527 SUPPORT VECTOR MACHINES 1401
```

```
SCIKITLEARN USER GUIDE RELEASE 0213
XY NPVSTACKXXRAVEL YYRAVELT
Z CLFDECISIONFUNCTIONXYRESHAPEXXSHAPE
PLOT DECISION BOUNDARY AND MARGINS
AXCONTOURXX YY Z COLORSK LEVELS1 0 1 ALPHA05
LINESTYLES
PLOT SUPPORT VECTORS
AXSCATTERCLFSUPPORTVECTORS 0 CLFSUPPORTVECTORS 1 S100
LINEWIDTH1 FACECOLORSNONE EDGECOLORSK
PLTSHOW
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0021 SECONDS
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE
5273 SVM WITH CUSTOM KERNEL
SIMPLE USAGE OF SUPPORT VECTOR MACHINES TO CLASSIFY A SAMPLE IT WILL PLOT THE DECISION SURFACE AND THE SUPPORT VECTORS
PRINTDOC
IMPORT NUMPY AS NP
1402 CHAPTER 5 EXAMPLES
```



```
SCIKITLEARN USER GUIDE RELEASE 0213
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARN IMPORT SVM DATASETS
IMPORT SOME DATA TO PLAY WITH
IRIS DATASETSLOADIRIS
X IRISDATA 2 WE ONLY TAKE THE FIRST TWO FEATURES WE COULD
AVOID THIS UGLY SLICING BY USING A TWODIM DATASET
Y IRISTARGET
DEFFMYKERNELX Y

WE CREATE A CUSTOM KERNEL
2 0
KX Y X YT
0 1

M NPARRAY2 0 0 10
RETURNNPDOTNPDOTX M YT
H 02 STEP SIZE IN THE MESH
WE CREATE AN INSTANCE OF SVM AND FIT OUT DATA
CLF SVMSVCKERNELMYKERNEL
CLFFITX Y
PLOT THE DECISION BOUNDARY FOR THAT WE WILL ASSIGN A COLOR TO EACH
POINT IN THE MESH XMIN XMAXXYMIN YMAX
XMIN XMAX X 0MIN 1 X 0MAX 1
YMIN YMAX X 1MIN 1 X 1MAX 1
XX YY NPMESHGRIDNPARANGEXMIN XMAX H NPARANGEYMIN YMAX H
Z CLFPREDICTNPCXXRAVEL YYRAVEL
PUT THE RESULT INTO A COLOR PLOT
Z ZRESHAPEXXSHAPE
PLTPCOLORMESHXX YY Z CMAPPLTCMPAIED
PLOT ALSO THE TRAINING POINTS
PLTSCATTERX 0 X 1 CY CMAPPLTCMPAIED EDGECOLORSK
PLTTITLE3CLASS CLASSIFICATION USING SUPPORT VECTOR MACHINE WITH CUSTOM
KERNEL
PLTAXISTIGHT
PLTSHOW
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0110 SECONDS
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE
5274 SVM WEIGHTED SAMPLES
PLOT DECISION FUNCTION OF A WEIGHTED DATASET WHERE THE SIZE OF POINTS IS PROPORTIONAL TO ITS WEIGHT
THE SAMPLE WEIGHTING RESCALES THE C PARAMETER WHICH MEANS THAT THE CLASSIFIER PUTS MORE EMPHASIS ON GETTING THESE
527 SUPPORT VECTOR MACHINES 1403
```

```
SCIKITLEARN USER GUIDE RELEASE 0213
POINTS RIGHT THE EFFECT MIGHT OFTEN BE SUBTLE TO EMPHASIZE THE EFFECT HERE WE PARTICULARLY WEIGHT OUTLIERS MAKING
THE DEFORMATION OF THE DECISION BOUNDARY VERY VISIBLE
PRINTDOC
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARN IMPORT SVM
DEFPLOTTDECISIONFUNCTIONCLASSIFIER SAMPLEWEIGHT AXIS TITLE
PLOT THE DECISION FUNCTION
XX YY NPMESHGRIDNPLINSPACE4 5 500 NPLINSPACE4 5 500
Z CLASSIFIERDECISIONFUNCTIONNPCXXRAVEL YYRAVEL
Z ZRESHAPEXXSHAPE
PLOT THE LINE THE POINTS AND THE NEAREST VECTORS TO THE PLANE
AXISCONTOURFXX YY Z ALPHA075 CMAPPLTCMBONE
AXISSCATTERX 0 X 1 CY S100 SAMPLEWEIGHT ALPHA09
CMAPPLTCMBONE EDGECOLORSBLACK
AXISAXISOFF
AXISSETTITLETITLE
WE CREATE 20 POINTS
NPRANDOMSEED0
X NPRNPRANDOMRANDN10 2 1 1 NPRANDOMRANDN10 2
Y 1 10 1 10
SAMPLEWEIGHTLASTTEN ABSNPRANDOMRANDNLENX
SAMPLEWEIGHTCONSTANT NPONESLENX
AND BIGGER WEIGHTS TO SOME OUTLIERS
SAMPLEWEIGHTLASTTEN15 5
SAMPLEWEIGHTLASTTEN9 15
FOR REFERENCE FIRST FIT WITHOUT SAMPLE WEIGHTS
FIT THE MODEL
1404 CHAPTER 5 EXAMPLES
```

SCIKITLEARN USER GUIDE RELEASE 0213

```
CLFWEIGHTS SVMSCVGAMMA1
CLFWEIGHTSFITX Y SAMPLEWEIGHTSAMPLEWEIGHTLASTTEN
CLFNOWEIGHTS SVMSCVGAMMA1
CLFNOWEIGHTSFITX Y
FIG AXES PLTSUBPLOTS1 2 FIGSIZE14 6
PLOTDECISIONFUNCTIONCLFNOWEIGHTS SAMPLEWEIGHTCONSTANT AXES0
CONSTANT WEIGHTS
PLOTDECISIONFUNCTIONCLFWEIGHTS SAMPLEWEIGHTLASTTEN AXES1
MODIFIED WEIGHTS
PLTSHOW
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0349 SECONDS
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE
5275 SVM SEPARATING HYPERPLANE FOR UNBALANCED CLASSES
FIND THE OPTIMAL SEPARATING HYPERPLANE USING AN SVC FOR CLASSES THAT ARE UNBALANCED
WE FIRST FIND THE SEPARATING PLANE WITH A PLAIN SVC AND THEN PLOT DASHED THE SEPARATING HYPERPLANE WITH AUTOMATICALLY
CORRECTION FOR UNBALANCED CLASSES
NOTE THIS EXAMPLE WILL ALSO WORK BY REPLACING SVCKERNELLINEAR WITH
SGDCLASSIFIERLOSSHINGE SETTING THE LOSS PARAMETER OF THE SGDCLASSIFIER EQUAL TOHINGE WILL
YIELD BEHAVIOUR SUCH AS THAT OF A SVC WITH A LINEAR KERNEL
FOR EXAMPLE TRY INSTEAD OF THE SVC
CLF SGDCLASSIFIERNITER100 ALPHA001
527 SUPPORT VECTOR MACHINES 1405
```

```
SCIKITLEARN USER GUIDE RELEASE 0213
PRINTDOC
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARN IMPORT SVM
FROM SKLEARNDATASETS IMPORT MAKEBLOBS
    WE CREATE TWO CLUSTERS OF RANDOM POINTS
NSAMPLES1 1000
NSAMPLES2 100
CENTERS 00 00 20 20
CLUSTERSSTD 15 05
X Y MAKEBLOBSNSAMPLESNSAMPLES1 NSAMPLES2
CENTERSCENTERS
CLUSTERSTDCLUSTERSSTD
RANDOMSTATE0 SHUFFLEFALSE
    FIT THE MODEL AND GET THE SEPARATING HYPERPLANE
CLF SVMKERNELLINEAR C10
CLFFITX Y
    FIT THE MODEL AND GET THE SEPARATING HYPERPLANE USING WEIGHTED CLASSES
WCLF SVMKERNELLINEAR CLASSWEIGHT1 10
WCLFFITX Y
    PLOT THE SAMPLES
1406 CHAPTER 5 EXAMPLES
```

SCIKITLEARN USER GUIDE RELEASE 0213  
PLTSCATTERX 0 X 1 CY CMAPPLTCMPAIREDCOLORSK  
PLOT THE DECISION FUNCTIONS FOR BOTH CLASSIFIERS  
AX PLTGCA  
XLIM AXGETXLIM  
YLIM AXGETYLIM  
CREATE GRID TO EVALUATE MODEL  
XX NPLINSPACEXLIM0 XLIM1 30  
YY NPLINSPACEYLIM0 YLIM1 30  
YY XX NPMESHGRIDYY XX  
XY NPVSTACKXXRAVEL YYRAVELT  
GET THE SEPARATING HYPERPLANE  
Z CLFDECISIONFUNCTIONXYRESHAPEXXSHAPE  
PLOT DECISION BOUNDARY AND MARGINS  
A AXCONTOURXX YY Z COLORSK LEVELS0 ALPHA05 LIFESTYLES  
GET THE SEPARATING HYPERPLANE FOR WEIGHTED CLASSES  
Z WCLFDECISIONFUNCTIONXYRESHAPEXXSHAPE  
PLOT DECISION BOUNDARY AND MARGINS FOR WEIGHTED CLASSES  
B AXCONTOURXX YY Z COLORSR LEVELS0 ALPHA05 LIFESTYLES  
PLTLEGENDACOLLECTIONS0 BCOLLECTIONS0 NON WEIGHTED WEIGHTED  
LOCUPPER RIGHT  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0034 SECONDS  
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE  
5276 SVMKERNELS  
THREE DIFFERENT TYPES OF SVMKERNELS ARE DISPLAYED BELOW THE POLYNOMIAL AND RBF ARE ESPECIALLY USEFUL WHEN THE  
DATAPOINTS ARE NOT LINEARLY SEPARABLE  
527 SUPPORT VECTOR MACHINES 1407

SCIKITLEARN USER GUIDE RELEASE 0213

- 
- 

1408 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

```
•
PRINTDOC
CODE SOURCE GAËL VAROQUAUX
LICENSE BSD 3 CLAUSE
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARN IMPORT SVM
OUR DATASET AND TARGETS
X NPC4 7
15 1
14 9
13 12
11 2
12 4
5 12
15 21
1 1

13 8
12 5
2 2
5 24
2 23
0 27
13 21T
Y 0 8 1 8
FIGURE NUMBER
FIGNUM 1
FIT THE MODEL
FORKERNELINLINEAR POLY RBF
CLF SVM SVCKERNELKERNEL GAMMA2
CLFFITX Y
527 SUPPORT VECTOR MACHINES 1409
```

SCIKITLEARN USER GUIDE RELEASE 0213

PLOT THE LINE THE POINTS AND THE NEAREST VECTORS TO THE PLANE

PLTFIGUREFIGNUM FIGSIZE4 3

PLTCLF

PLTSCATTERCLFSUPPORTVECTORS 0 CLFSUPPORTVECTORS 1 S80

FACECOLORSNONE ZORDER10 EDGECOLORSK

PLTSCATTERX 0 X 1 CY ZORDER10 CMAPPLTCMPAIED

EDGECOLORSK

PLTAXISTIGHT

XMIN 3

XMAX 3

YMIN 3

YMAX 3

XX YY NPMGRIDXMINXMAX200J YMINYMAX200J

Z CLFDECISIONFUNCTIONNPCXXRAVEL YYRAVEL

PUT THE RESULT INTO A COLOR PLOT

Z ZRESHAPEXXSHAPE

PLTFIGUREFIGNUM FIGSIZE4 3

PLTPCOLORMESHXX YY Z 0 CMAPPLTCMPAIED

PLTCONTOURXX YY Z COLORSK K K LINSTYLES

LEVELS5 0 5

PLTXLIMXMIN XMAX

PLTYLIMYMIN YMAX

PLTXTICKS

PLTYTICKS

FIGNUM FIGNUM 1

PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0096 SECONDS

NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)

5277 SVMANOVA SVM WITH UNIVARIATE FEATURE SELECTION

THIS EXAMPLE SHOWS HOW TO PERFORM UNIVARIATE FEATURE SELECTION BEFORE RUNNING A SVC SUPPORT VECTOR CLASSIFIER TO IMPROVE THE CLASSIFICATION SCORES WE USE THE IRIS DATASET 4 FEATURES AND ADD 36 NONINFORMATIVE FEATURES WE CAN FIND THAT OUR MODEL ACHIEVES BEST PERFORMANCE WHEN WE SELECT AROUND 10 OF FEATURES

1410 CHAPTER 5 EXAMPLES



```
SCIKITLEARN USER GUIDE RELEASE 0213
PRINTDOC
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARNDATASETS IMPORT LOADIRIS
FROM SKLEARNFEATURESELECTION IMPORT SELECTPERCENTILE CHI2
FROM SKLEARNMODELSELECTION IMPORT CROSSVALSCORE
FROM SKLEARNPIPELINE IMPORT PIPELINE
FROM SKLEARNPREPROCESSING IMPORT STANDARDSCALER
FROM SKLEARN SVM IMPORT SVC

IMPORT SOME DATA TO PLAY WITH
X Y LOADIRISRETURNXYTRUE
ADD NONINFORMATIVE FEATURES
NPRANDOMSEED0
X NPHSTACKX 2 NPRANDOMRANDOMXSHAPE0 36

CREATE A FEATURESELECTION TRANSFORM A SCALER AND AN INSTANCE OF SVM THAT WE
COMBINE TOGETHER TO HAVE AN FULLBLOWN ESTIMATOR
CLF PIPELINEANOVA SELECTPERCENTILECHI2
SCALER STANDARDSCALER
SVC SVCGAMMAAUTO
527 SUPPORT VECTOR MACHINES 1411
```

SCIKITLEARN USER GUIDE RELEASE 0213

PLOT THE CROSSVALIDATION SCORE AS A FUNCTION OF PERCENTILE OF FEATURES  
SCOREMEANS LIST  
SCORESTDS LIST  
PERCENTILES 1 3 6 10 15 20 30 40 60 80 100  
FORPERCENTILE INPERCENTILES  
CLFSETPARAMSANOVAPERCENTILEPERCENTILE  
THISSCORES CROSSVALSCORECLF X Y CV5  
SCOREMEANSAPPENDTHISSCORESMEAN  
SCORESTDSAPPENDTHISSCORESSTD  
PLTERRORBARPERCENTILES SCOREMEANS NPARRAYSCORESTDS  
PLTTITLE  
PERFORMANCE OF THE SVMANOVA VARYING THE PERCENTILE OF FEATURES SELECTED  
PLTXTICKSNPLINSPACE0 100 11 ENDPOINTTRUE  
PLTXLABELPERCENTILE  
PLTYLABELACCURACY SCORE  
PLTAXISTIGHT  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0199 SECONDS  
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE  
5278 SUPPORT VECTOR REGRESSION SVR USING LINEAR AND NONLINEAR KERNELS  
TOY EXAMPLE OF 1D REGRESSION USING LINEAR POLYNOMIAL AND RBF KERNELS  
1412 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
PRINTDOC  
IMPORT NUMPY AS NP  
FROM SKLEARN SVM IMPORT SVR  
IMPORT MATPLOTLIBPY PLOT AS PLT

GENERATE SAMPLE DATA  
X NPSORT5 NPRANDOMRAND40 1 AXIS0  
Y NPSINXRAVEL

ADD NOISE TO TARGETS  
Y5 3 05 NPRANDOMRAND8

FIT REGRESSION MODEL  
SVRRBF SVRKERNELRBF C100 GAMMA01 EPSILON1  
SVRLIN SVRKERNELLINEAR C100 GAMMAAUTO  
SVRPOLY SVRKERNELPOLY C100 GAMMAAUTO DEGREE3 EPSILON1  
COEF01

LOOK AT THE RESULTS  
LW 2  
SVRS SVRRBF SVRLIN SVRPOLY  
KERNELLABEL RBF LINEAR POLYNOMIAL  
MODEL COLOR M C G  
527 SUPPORT VECTOR MACHINES 1413

SCIKITLEARN USER GUIDE RELEASE 0213

FIG AXES PLTSUBPLOTSNROWS1 NCOLS3 FIGSIZE15 10 SHAREYTRUE

FORIX SVR INENUMERATESVRS

AXESIXPLOTX SVRFITX YPREDICTX COLORMODELCOLORIX LWLW

LABEL MODELFORMATKERNELLABELIX

AXESIXSCATTERXSVRSUPPORT YSVRSUPPORT FACECOLORNONE

EDGECOLORMODELCOLORIX S50

LABEL SUPPORT VECTORSFORMATKERNELLABELIX

AXESIXSCATTERXNPSETDIFF1DNPARANGELENX SVRSUPPORT

YNPSETDIFF1DNPARANGELENX SVRSUPPORT

FACECOLORNONE EDGECOLORK S50

LABELOTHER TRAINING DATA

AXESIXLEGENDLOCUPPER CENTER BBOXTOANCHOR05 11

NCOL1 FANCYBOXTRUE SHADOWTRUE

FIGTEXT05 004 DATA HACENTER VACENTER

FIGTEXT006 05 TARGET HACENTER VACENTER ROTATIONVERTICAL

FIGSUPTITLESUPPORT VECTOR REGRESSION FONTSIZE14

PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 3104 SECONDS

NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE

5279 SVM MARGINS EXAMPLE

THE PLOTS BELOW ILLUSTRATE THE EFFECT THE PARAMETER CHAS ON THE SEPARATION LINE A LARGE VALUE OF CBASICALLY TELLS OUR MODEL THAT WE DO NOT HAVE THAT MUCH FAITH IN OUR DATA’S DISTRIBUTION AND WILL ONLY CONSIDER POINTS CLOSE TO LINE OF SEPARATION

A SMALL VALUE OF CINCLUDES MOREALL THE OBSERVATIONS ALLOWING THE MARGINS TO BE CALCULATED USING ALL THE DATA IN THE AREA

- 

1414 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

•

```
PRINTDOC
CODE SOURCE GAËL VAROQUAUX
MODIFIED FOR DOCUMENTATION BY JAQUES GROBLER
LICENSE BSD 3 CLAUSE
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARN IMPORT SVM
WE CREATE 40 SEPARABLE POINTS
NPRANDOMSEED0
X NPRNPRANDOMRANDN20 2 2 2 NPRANDOMRANDN20 2 2 2
Y 0 20 1 20
FIGURE NUMBER
FIGNUM 1
FIT THE MODEL
FORNAME PENALTY INUNREG 1 REG 005
CLF SVM SVCKERNELLINEAR CPENALTY
CLFFITX Y
GET THE SEPARATING HYPERPLANE
W CLFCOEF0
A W0 W1
XX NPLINSPACE5 5
YY AXX CLFINTERCEPT0 W1
PLOT THE PARALLELS TO THE SEPARATING HYPERPLANE THAT PASS THROUGH THE
SUPPORT VECTORS MARGIN AWAY FROM HYPERPLANE IN DIRECTION
PERPENDICULAR TO HYPERPLANE THIS IS SQRT1A2 AWAY VERTICALLY IN
2D
MARGIN 1 NPSQRTNPSUMCLFCOEF 2
YYDOWN YY NPSQRT1 A 2MARGIN
YYUP YY NPSQRT1 A 2MARGIN
527 SUPPORT VECTOR MACHINES 1415
```

SCIKITLEARN USER GUIDE RELEASE 0213

PLOT THE LINE THE POINTS AND THE NEAREST VECTORS TO THE PLANE

PLTFIGUREFIGNUM FIGSIZE4 3

PLTCLF

PLTPLOTXX YY K

PLTPLOTXX YYDOWN K

PLTPLOTXX YYUP K

PLTSCATTERCLFSUPPORTVECTORS 0 CLFSUPPORTVECTORS 1 S80

FACECOLORSNONE ZORDER10 EDGECOLORSK

PLTSCATTERX 0 X 1 CY ZORDER10 CMAPPLTCMPAIED

EDGECOLORSK

PLTAXISTIGHT

XMIN 48

XMAX 42

YMIN 6

YMAX 6

XX YY NPMGRIDXMINXMAX200J YMINYMAX200J

Z CLFPREDICTNPCXXRAVEL YYRAVEL

PUT THE RESULT INTO A COLOR PLOT

Z ZRESHAPEXXSHAPE

PLTFIGUREFIGNUM FIGSIZE4 3

PLTPCOLORMESHXX YY Z CMAPPLTCMPAIED

PLTXLIMXMIN XMAX

PLTYLIMYMIN YMAX

PLXTTICKS

PLTYTICKS

FIGNUM FIGNUM 1

PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0061 SECONDS

NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)

52710 ONECLASS SVM WITH NONLINEAR KERNEL RBF

AN EXAMPLE USING A ONECLASS SVM FOR NOVELTY DETECTION

ONECLASS SVM IS AN UNSUPERVISED ALGORITHM THAT LEARNS A DECISION FUNCTION FOR NOVELTY DETECTION CLASSIFYING NEW DATA AS SIMILAR OR DIFFERENT TO THE TRAINING SET

1416 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
PRINTDOC
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
IMPORT MATPLOTLIBFONTMANAGER
FROM SKLEARN IMPORT SVM
XX YY NPMESHGRIDNPLINSPACE5 5 500 NPLINSPACE5 5 500
    GENERATE TRAIN DATA
X 03 NPRANDOMRANDN100 2
XTRAIN NPRX 2 X 2
    GENERATE SOME REGULAR NOVEL OBSERVATIONS
X 03 NPRANDOMRANDN20 2
XTEST NPRX 2 X 2
    GENERATE SOME ABNORMAL NOVEL OBSERVATIONS
XOUTLIERS NPRANDOMUNIFORMLOW4 HIGH4 SIZE20 2
    FIT THE MODEL
CLF SVMONECLASSSVMNU01 KERNELRBF GAMMA01
CLFFITXTRAIN
YPREDTRAIN CLFPREDICTXTRAIN
YPREDTEST CLFPREDICTXTEST
YPREDOUTLIERS CLFPREDICTXOUTLIERS
NERRORTRAIN YPREDTRAINYPREDTRAIN 1SIZE
NERRORTEST YPREDTESTYPREDTEST 1SIZE
NERROROUTLIERS YPREDOUTLIERSYPREDOUTLIERS 1SIZE
527 SUPPORT VECTOR MACHINES 1417
```

SCIKITLEARN USER GUIDE RELEASE 0213

PLOT THE LINE THE POINTS AND THE NEAREST VECTORS TO THE PLANE

Z CLFDECISIONFUNCTIONNPCXXRAVEL YYRAVEL

Z ZRESHAPEXXSHAPE

PLTTITLENOVELTY DETECTION

PLTCONTOURFXX YY Z LEVELSNPLINSPACEZMIN 0 7 CMAPPLTCMPUBU

A PLTCONTOURXX YY Z LEVELS0 LINEWIDTHS2 COLORSDARKRED

PLTCONTOURFXX YY Z LEVELS0 ZMAX COLORSPALEVIOLETRED

S 40

B1 PLTSCATTERXTRAIN 0 XTRAIN 1 CWHITE SS EDGECOLORSK

B2 PLTSCATTERXTEST 0 XTEST 1 CBLUEVIOLET SS

EDGECOLORSK

C PLTSCATTERXOUTLIERS 0 XOUTLIERS 1 CGOLD SS

EDGECOLORSK

PLTAXISTIGHT

PLTXLIM5 5

PLTYLIM5 5

PLTLEGENDACOLLECTIONS0 B1 B2 C

LEARNED FRONTIER TRAINING OBSERVATIONS

NEW REGULAR OBSERVATIONS NEW ABNORMAL OBSERVATIONS

LOCUPPER LEFT

PROPMATPLOTLIBFONTMANAGERFONTPROPERTIESSIZE11

PLTXLABEL

ERROR TRAIN D200 ERRORS NOVEL REGULAR D40

ERRORS NOVEL ABNORMAL D40

NERRORTRAIN NERRORTEST NERROROUTLIERS

PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0196 SECONDS

NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE

52711 PLOT DIFFERENT SVM CLASSIFIERS IN THE IRIS DATASET

COMPARISON OF DIFFERENT LINEAR SVM CLASSIFIERS ON A 2D PROJECTION OF THE IRIS DATASET WE ONLY CONSIDER THE FIRST 2

FEATURES OF THIS DATASET

- SEPAL LENGTH
- SEPAL WIDTH

THIS EXAMPLE SHOWS HOW TO PLOT THE DECISION SURFACE FOR FOUR SVM CLASSIFIERS WITH DIFFERENT KERNELS

THE LINEAR MODELS LINEARSVC ANDSVCKERNELLINEAR YIELD SLIGHTLY DIFFERENT DECISION BOUNDARIES

THIS CAN BE A CONSEQUENCE OF THE FOLLOWING DIFFERENCES

- LINEARSVC MINIMIZES THE SQUARED HINGE LOSS WHILE SVC MINIMIZES THE REGULAR HINGE LOSS
- LINEARSVC USES THE ONEVSALL ALSO KNOWN AS ONEVSREST MULTICLASS REDUCTION WHILE SVC USES THE ONE

VSONE MULTICLASS REDUCTION

BOTH LINEAR MODELS HAVE LINEAR DECISION BOUNDARIES INTERSECTING HYPERPLANES WHILE THE NONLINEAR KERNEL MODELS

POLYNOMIAL OR GAUSSIAN RBF HAVE MORE FLEXIBLE NONLINEAR DECISION BOUNDARIES WITH SHAPES THAT DEPEND ON THE KIND

OF KERNEL AND ITS PARAMETERS

1418 CHAPTER 5 EXAMPLES



```
SCIKITLEARN USER GUIDE RELEASE 0213
NOTE WHILE PLOTTING THE DECISION FUNCTION OF CLASSIFIERS FOR TOY 2D DATASETS CAN HELP GET AN INTUITIVE UNDERSTANDING
OF THEIR RESPECTIVE EXPRESSIVE POWER BE AWARE THAT THOSE INTUITIONS DON'T ALWAYS GENERALIZE TO MORE REALISTIC HIGH
DIMENSIONAL PROBLEMS
PRINTDOC
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARN IMPORT SVM DATASETS
DEFMAKEMESHGRIDX Y H02
CREATE A MESH OF POINTS TO PLOT IN
PARAMETERS

X DATA TO BASE XAXIS MESHGRID ON
Y DATA TO BASE YAXIS MESHGRID ON
H STEPSIZE FOR MESHGRID OPTIONAL
RETURNS

XX YY NDARRAY

527 SUPPORT VECTOR MACHINES 1419
```

SCIKITLEARN USER GUIDE RELEASE 0213

XMIN XMAX XMIN 1 XMAX 1  
YMIN YMAX YMIN 1 YMAX 1  
XX YY NPMESHGRIDNPARANGEXMIN XMAX H  
NPARANGEYMIN YMAX H  
RETURNXX YY  
DEFPLOTCONTOURSAX CLF XX YY PARAMS  
PLOT THE DECISION BOUNDARIES FOR A CLASSIFIER  
PARAMETERS

AX MATPLOTLIB AXES OBJECT  
CLF A CLASSIFIER  
XX MESHGRID NDARRAY  
YY MESHGRID NDARRAY  
PARAMS DICTIONARY OF PARAMS TO PASS TO CONTOURF OPTIONAL

Z CLFPREDICTNPCXXRAVEL YYRAVEL  
Z ZRESHAPEXXSHAPE  
OUT AXCONTOURFXX YY Z PARAMS  
RETURNOUT  
IMPORT SOME DATA TO PLAY WITH  
IRIS DATASETSLOADIRIS  
TAKE THE FIRST TWO FEATURES WE COULD AVOID THIS BY USING A TWODIM DATASET  
X IRISDATA 2  
Y IRISTARGET  
WE CREATE AN INSTANCE OF SVM AND FIT OUT DATA WE DO NOT SCALE OUR  
DATA SINCE WE WANT TO PLOT THE SUPPORT VECTORS  
C 10 SVM REGULARIZATION PARAMETER  
MODELS SVMKVCKERNELLINER CC  
SVMLINERSVCCC MAXITER10000  
SVMKVCKERNELRBF GAMMA07 CC  
SVMKVCKERNELPOLY DEGREE3 GAMMAAUTO CC  
MODELS CLFFITX Y FORCLFINMODELS  
TITLE FOR THE PLOTS  
TITLES SVC WITH LINEAR KERNEL  
LINEAR SVC LINEAR KERNEL  
SVC WITH RBF KERNEL  
SVC WITH POLYNOMIAL DEGREE 3 KERNEL  
SETUP 2X2 GRID FOR PLOTTING  
FIG SUB PLTSUBPLOTS2 2  
PLTSUBPLOTSADJUSTWSPACE04 HSPACE04  
X0 X1 X 0 X 1  
XX YY MAKEMESHGRIDX0 X1  
FORCLF TITLE AX INZIPMODELS TITLES SUBFLATTEN  
PLOTCONTOURSAX CLF XX YY  
CMAPPLTCMCOOLWARM ALPHA08  
AXSCATTERX0 X1 CY CMAPPLTCMCOOLWARM S20 EDGECOLORSK  
AXSETXLIMXXMIN XXMAX  
AXSETYLIMYYMIN YYMAX  
1420 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

AXSETXLABELSEPAL LENGTH

AXSETYLABELSEPAL WIDTH

AXSETXTICKS

AXSETYTICKS

AXSETTITLETITLE

PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0481 SECONDS

NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE

52712 SCALING THE REGULARIZATION PARAMETER FOR SVCS

THE FOLLOWING EXAMPLE ILLUSTRATES THE EFFECT OF SCALING THE REGULARIZATION PARAMETER WHEN USING SUPPORT VECTOR MA  
CHINES FORCLASSIFICATION FOR SVC CLASSIFICATION WE ARE INTERESTED IN A RISK MINIMIZATION FOR THE EQUATION

$\sum$

$\sum_{i=1}^n \frac{1}{n} \sum_{j=1}^m$

$\sum_{j=1}^m$

WHERE

- $C$  IS USED TO SET THE AMOUNT OF REGULARIZATION
- $\lambda$  IS ALOSS FUNCTION OF OUR SAMPLES AND OUR MODEL PARAMETERS
- 

IS APENALTY FUNCTION OF OUR MODEL PARAMETERS

IF WE CONSIDER THE LOSS FUNCTION TO BE THE INDIVIDUAL ERROR PER SAMPLE THEN THE DATAFIT TERM OR THE SUM OF THE ERROR FOR EACH SAMPLE WILL INCREASE AS WE ADD MORE SAMPLES THE PENALIZATION TERM HOWEVER WILL NOT INCREASE

WHEN USING FOR EXAMPLE CROSS VALIDATION TO SET THE AMOUNT OF REGULARIZATION WITH C THERE WILL BE A DIFFERENT AMOUNT OF SAMPLES BETWEEN THE MAIN PROBLEM AND THE SMALLER PROBLEMS WITHIN THE FOLDS OF THE CROSS VALIDATION

SINCE OUR LOSS FUNCTION IS DEPENDENT ON THE AMOUNT OF SAMPLES THE LATTER WILL INFLUENCE THE SELECTED VALUE OF C THE QUESTION THAT ARISES IS HOW DO WE OPTIMALLY ADJUST C TO ACCOUNT FOR THE DIFFERENT

AMOUNT OF TRAINING SAMPLES

THE FIGURES BELOW ARE USED TO ILLUSTRATE THE EFFECT OF SCALING OUR CTO COMPENSATE FOR THE CHANGE IN THE NUMBER OF SAMPLES IN THE CASE OF USING AN L1PENALTY AS WELL AS THE L2PENALTY

L1PENALTY CASE

IN THEL1CASE THEORY SAYS THAT PREDICTION CONSISTENCY IE THAT UNDER GIVEN HYPOTHESIS THE ESTIMATOR LEARNED PREDICTS AS WELL AS A MODEL KNOWING THE TRUE DISTRIBUTION IS NOT POSSIBLE BECAUSE OF THE BIAS OF THE L1 IT DOES SAY HOWEVER THAT MODEL CONSISTENCY IN TERMS OF FINDING THE RIGHT SET OF NONZERO PARAMETERS AS WELL AS THEIR SIGNS CAN BE ACHIEVED BY SCALINGC1

L2PENALTY CASE

THE THEORY SAYS THAT IN ORDER TO ACHIEVE PREDICTION CONSISTENCY THE PENALTY PARAMETER SHOULD BE KEPT CONSTANT AS THE NUMBER OF SAMPLES GROW

527 SUPPORT VECTOR MACHINES 1421

SCIKITLEARN USER GUIDE RELEASE 0213

SIMULATIONS

THE TWO FIGURES BELOW PLOT THE VALUES OF  $C$  ON THE X-AXIS AND THE CORRESPONDING CROSSVALIDATION SCORES ON THE Y-AXIS FOR SEVERAL DIFFERENT FRACTIONS OF A GENERATED DATASET. IN THE  $L_1$  PENALTY CASE, THE CROSSVALIDATION ERROR CORRELATES BEST WITH THE TEST ERROR WHEN SCALING  $C$  WITH THE NUMBER OF SAMPLES  $N$ , WHICH CAN BE SEEN IN THE FIRST FIGURE. FOR THE  $L_2$  PENALTY CASE, THE BEST RESULT COMES FROM THE CASE WHERE  $C$  IS NOT SCALED.

NOTE

TWO SEPARATE DATASETS ARE USED FOR THE TWO DIFFERENT PLOTS. THE REASON BEHIND THIS IS THE  $L_1$  CASE WORKS BETTER ON SPARSE DATA, WHILE  $L_2$  IS BETTER SUITED TO THE NONSPARSE CASE.

SCIKITLEARN USER GUIDE RELEASE 0213

- 

527 SUPPORT VECTOR MACHINES 1423

SCIKITLEARN USER GUIDE RELEASE 0213

•

PRINTDOC

AUTHOR ANDREAS MUELLER AMUELLERAISUNIBONNDE

JAQUES GROBLER JAQUESGROBLERINRIAFR

LICENSE BSD 3 CLAUSE

IMPORT NUMPY AS NP

IMPORT MATPLOTLIBPYPLOT AS PLT

FROM SKLEARN SVM IMPORT LINEARSVC

FROM SKLEARNMODELSELECTION IMPORT SHUFFLESPLIT

1424 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
FROM SKLEARNMODELSELECTION IMPORT GRIDSEARCHCV
FROM SKLEARNUTILS IMPORT CHECKRANDOMSTATE
FROM SKLEARN IMPORT DATASETS
RND CHECKRANDOMSTATE1
SET UP DATASET
NSAMPLES 100
NFEATURES 300
L1 DATA ONLY 5 INFORMATIVE FEATURES
X1 Y1 DATASETSMAKECLASSIFICATIONNSAMPLESNSAMPLES
NFEATURESNFEATURES NINFORMATIVE5
RANDOMSTATE1
L2 DATA NON SPARSE BUT LESS FEATURES
Y2 NPSIGN5 RNDRANDNSAMPLES
X2 RNDRANDNNSAMPLES NFEATURES 5 Y2 NPNEWAXIS
X2 5 RNDRANDNNSAMPLES NFEATURES 5
CLFSETS LINEARSVCPENALTYL1 LOSSSSQUAREDHINGE DUALFALSE
TOL1E3
NPLOGSPACE23 13 10 X1 Y1
LINEARSVCPENALTYL2 LOSSSSQUAREDHINGE DUALTRUE
TOL1E4
NPLOGSPACE45 2 10 X2 Y2
COLORS NAVY CYAN DARKORANGE
LW 2
FORCLF CS X Y INCLFSETS
SET UP THE PLOT FOR EACH REGRESSOR
FIG AXES PLTSUBPLOTSNROWS2 SHAREYTRUE FIGSIZE9 10
FORK TRAINSIZE INENUMERATENPLINSPACE03 07 31
PARAMGRID DICTCCS
TO GET NICE CURVE WE NEED A LARGE NUMBER OF ITERATIONS TO
REDUCE THE VARIANCE
GRID GRIDSEARCHCVCLF REFITFALSE PARAMGRIDPARAMGRID
CVSHUFFLESPLITTRAINSIZE TRAINSIZE
TESTSIZE3
NSPLITS250 RANDOMSTATE1
GRIDFITX Y
SCORES GRIDCVRESULTSMEANTESTSCORE
SCALES 1 NO SCALING
NSAMPLES TRAINSIZE 1NSAMPLES

FORAX SCALER NAME INZIPAXES SCALES
AXSETXLABELC
AXSETYLABELCV SCORE
GRIDCS CS FLOATSCALER SCALE THE CS
AXSEMILOGXGRIDCS SCORES LABELFRACTION 2F
TRAINSIZE COLORCOLORSK LWLW
AXSETTITLESCALING S PENALTY S LOSS
NAME CLFPENALTY CLFLOSS
527 SUPPORT VECTOR MACHINES 1425
```

SCIKITLEARN USER GUIDE RELEASE 0213

PLTLEGENDLOCBEST

PLTSHOW

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 14575 SECONDS

NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE

52713 RBF SVM PARAMETERS

THIS EXAMPLE ILLUSTRATES THE EFFECT OF THE PARAMETERS GAMMA AND C OF THE RADIAL BASIS FUNCTION RBF KERNEL SVM. INTUITIVELY, THE GAMMA PARAMETER DEFINES HOW FAR THE INFLUENCE OF A SINGLE TRAINING EXAMPLE REACHES. WITH LOW VALUES MEANING 'FAR' AND HIGH VALUES MEANING 'CLOSE', THE GAMMA PARAMETERS CAN BE SEEN AS THE INVERSE OF THE RADIUS OF INFLUENCE OF SAMPLES SELECTED BY THE MODEL AS SUPPORT VECTORS.

THE C PARAMETER TRADES OFF CORRECT CLASSIFICATION OF TRAINING EXAMPLES AGAINST MAXIMIZATION OF THE DECISION FUNCTION'S MARGIN. FOR LARGER VALUES OF C, A SMALLER MARGIN WILL BE ACCEPTED IF THE DECISION FUNCTION IS BETTER AT CLASSIFYING ALL TRAINING POINTS CORRECTLY. A LOWER C WILL ENCOURAGE A LARGER MARGIN. THEREFORE, A SIMPLER DECISION FUNCTION AT THE COST OF TRAINING ACCURACY. IN OTHER WORDS, "C" BEHAVES AS A REGULARIZATION PARAMETER IN THE SVM.

THE FIRST PLOT IS A VISUALIZATION OF THE DECISION FUNCTION FOR A VARIETY OF PARAMETER VALUES ON A SIMPLIFIED CLASSIFICATION PROBLEM INVOLVING ONLY 2 INPUT FEATURES AND 2 POSSIBLE TARGET CLASSES. BINARY CLASSIFICATION. NOTE THAT THIS KIND OF PLOT IS NOT POSSIBLE TO DO FOR PROBLEMS WITH MORE FEATURES OR TARGET CLASSES.

THE SECOND PLOT IS A HEATMAP OF THE CLASSIFIER'S CROSS-VALIDATION ACCURACY AS A FUNCTION OF C AND GAMMA. FOR THIS EXAMPLE, WE EXPLORE A RELATIVELY LARGE GRID FOR ILLUSTRATION PURPOSES. IN PRACTICE, A LOGARITHMIC GRID FROM 10<sup>-3</sup> TO 10<sup>3</sup> IS USUALLY SUFFICIENT. IF THE BEST PARAMETERS LIE ON THE BOUNDARIES OF THE GRID, IT CAN BE EXTENDED IN THAT DIRECTION IN A SUBSEQUENT SEARCH.

NOTE THAT THE HEAT MAP PLOT HAS A SPECIAL COLORBAR WITH A MIDPOINT VALUE CLOSE TO THE SCORE VALUES OF THE BEST PERFORMING MODELS, SO AS TO MAKE IT EASY TO TELL THEM APART IN THE BLINK OF AN EYE.

THE BEHAVIOR OF THE MODEL IS VERY SENSITIVE TO THE GAMMA PARAMETER. IF GAMMA IS TOO LARGE, THE RADIUS OF THE AREA OF INFLUENCE OF THE SUPPORT VECTORS ONLY INCLUDES THE SUPPORT VECTOR ITSELF, AND NO AMOUNT OF REGULARIZATION. WITH C, IT WILL BE ABLE TO PREVENT OVERFITTING.

WHEN GAMMA IS VERY SMALL, THE MODEL IS TOO CONSTRAINED AND CANNOT CAPTURE THE COMPLEXITY OR "SHAPE" OF THE DATA. THE REGION OF INFLUENCE OF ANY SELECTED SUPPORT VECTOR WOULD INCLUDE THE WHOLE TRAINING SET. THE RESULTING MODEL WILL BEHAVE SIMILARLY TO A LINEAR MODEL WITH A SET OF HYPERPLANES THAT SEPARATE THE CENTERS OF HIGH DENSITY OF ANY PAIR OF TWO CLASSES.

FOR INTERMEDIATE VALUES, WE CAN SEE ON THE SECOND PLOT THAT GOOD MODELS CAN BE FOUND ON A DIAGONAL OF C AND GAMMA. SMOOTH MODELS. LOWER GAMMA VALUES CAN BE MADE MORE COMPLEX BY INCREASING THE IMPORTANCE OF CLASSIFYING EACH POINT CORRECTLY. LARGER C VALUES HENCE THE DIAGONAL OF GOOD PERFORMING MODELS.

FINALLY, ONE CAN ALSO OBSERVE THAT FOR SOME INTERMEDIATE VALUES OF GAMMA, WE GET EQUALLY PERFORMING MODELS WHEN C BECOMES VERY LARGE. IT IS NOT NECESSARY TO REGULARIZE BY ENFORCING A LARGER MARGIN. THE RADIUS OF THE RBF KERNEL ALONE ACTS AS A GOOD STRUCTURAL REGULARIZER. IN PRACTICE, THOUGH IT MIGHT STILL BE INTERESTING TO SIMPLIFY THE DECISION FUNCTION WITH A LOWER VALUE OF C, SO AS TO FAVOR MODELS THAT USE LESS MEMORY AND THAT ARE FASTER TO PREDICT. WE SHOULD ALSO NOTE THAT SMALL DIFFERENCES IN SCORES RESULT FROM THE RANDOM SPLITS OF THE CROSS-VALIDATION PROCEDURE. THOSE SPURIOUS VARIATIONS CAN BE SMOOTHED OUT BY INCREASING THE NUMBER OF CV ITERATIONS (NSPLITS) AT THE EXPENSE OF COMPUTE TIME. INCREASING THE VALUE NUMBER OF C RANGE AND GAMMA RANGE STEPS WILL INCREASE THE RESOLUTION OF

THE HYPERPARAMETER HEAT MAP

1426 CHAPTER 5 EXAMPLES





SCIKITLEARN USER GUIDE RELEASE 0213

•

OUT

THE BEST PARAMETERS ARE C 10 GAMMA 01 WITH A SCORE OF 097

PRINTDOC

IMPORT NUMPY AS NP

IMPORT MATPLOTLIBPYPLOT AS PLT

FROM MATPLOTLIBCOLORS IMPORT NORMALIZE

FROM SKLEARN SVM IMPORT SVC

FROM SKLEARNPREPROCESSING IMPORT STANDARDSCALER

FROM SKLEARN DATASETS IMPORT LOADIRIS

FROM SKLEARNMODELSELECTION IMPORT STRATIFIEDSHUFFLESPLIT

FROM SKLEARNMODELSELECTION IMPORT GRIDSEARCHCV

UTILITY FUNCTION TO MOVE THE MIDPOINT OF A COLORMAP TO BE AROUND

THE VALUES OF INTEREST

CLASS MIDPOINTNORMALIZE NORMALIZE

1428 CHAPTER 5 EXAMPLES

```
SCIKITLEARN USER GUIDE RELEASE 0213
DEFINITSELF VMINNONE VMAXNONE MIDPOINTNONE CLIPFALSE
SELF MIDPOINT MIDPOINT
NORMALIZEINITSELF VMIN VMAX CLIP
DEFCALLSELF VALUE CLIPNONE
X Y SELFVMIN SELF MIDPOINT SELFVMAX 0 05 1
RETURN NPMAMASKEDARRAY N PINTERPVALUE X Y
```

LOAD AND PREPARE DATA SET

```
DATASET FOR GRID SEARCH
IRIS LOADIRIS
X IRISDATA
Y IRISTARGET
DATASET FOR DECISION FUNCTION VISUALIZATION WE ONLY KEEP THE FIRST TWO
FEATURES IN X AND SUBSAMPLE THE DATASET TO KEEP ONLY 2 CLASSES AND
MAKE IT A BINARY CLASSIFICATION PROBLEM
X2D X 2
X2D X2DY 0
Y2D YY 0
Y2D 1
IT IS USUALLY A GOOD IDEA TO SCALE THE DATA FOR SVM TRAINING
WE ARE CHEATING A BIT IN THIS EXAMPLE IN SCALING ALL OF THE DATA
INSTEAD OF FITTING THE TRANSFORMATION ON THE TRAINING SET AND
JUST APPLYING IT ON THE TEST SET
SCALER STANDARDSCALER
X SCALERFITTRANSFORMX
X2D SCALERFITTRANSFORMX2D
```

TRAIN CLASSIFIERS

```
FOR AN INITIAL SEARCH A LOGARITHMIC GRID WITH BASIS
10 IS OFTEN HELPFUL USING A BASIS OF 2 A FINER
TUNING CAN BE ACHIEVED BUT AT A MUCH HIGHER COST
CRANGE NPLOGSPACE2 10 13
GAMMARANGE NPLOGSPACE9 3 13
PARAMGRID DICTGAMMAGAMMARANGE CCRANGE
CV STRATIFIEDSHUFFLESPLITSPLITS5 TESTSIZE02 RANDOMSTATE42
GRID GRIDSEARCHCV SVC PARAMGRIDPARAMGRID CVCV
GRIDFITX Y
PRINTTHE BEST PARAMETERS ARE SWITCH A SCORE OF 02F
GRIDBESTPARAMS GRIDBESTSCORE
NOW WE NEED TO FIT A CLASSIFIER FOR ALL PARAMETERS IN THE 2D VERSION
WE USE A SMALLER SET OF PARAMETERS HERE BECAUSE IT TAKES A WHILE TO TRAIN
C2DRANGE 1E2 1 1E2
GAMMA2DRANGE 1E1 1 1E1
CLASSIFIERS
527 SUPPORT VECTOR MACHINES 1429
```

SCIKITLEARN USER GUIDE RELEASE 0213  
FORCINC2DRANGE  
FORGAMMAINGAMMA2DRANGE  
CLF SVCCC GAMMAGAMMA  
CLFFITX2D Y2D  
CLASSIFIERSAPPENDC GAMMA CLF

VISUALIZATION

DRAW VISUALIZATION OF PARAMETER EFFECTS  
PLTFIGUREFIGSIZE8 6  
XX YY NPMESHGRIDNPLINSPACE3 3 200 NPLINSPACE3 3 200  
FORK C GAMMA CLF INENUMERATECLASSIFIERS  
EVALUATE DECISION FUNCTION IN A GRID  
Z CLFDECISIONFUNCTIONNPCXXRAVEL YYRAVEL  
Z ZRESHAPEXXSHAPE  
VISUALIZE DECISION FUNCTION FOR THESE PARAMETERS  
PLTSUBPLOTLENC2DRANGE LENGAMMA2DRANGE K 1  
PLTTITLEGAMMA10 D C10D NPLOG10GAMMA NPLOG10C  
SIZEMEDIUM  
VISUALIZE PARAMETERS EFFECT ON DECISION FUNCTION  
PLTPCOLORMESHXX YY Z CMAPPLTCMRDBU  
PLTSCATTERX2D 0 X2D 1 CY2D CMAPPLTCMRDBUR  
EDGECOLORSK  
PLXTICKS  
PLTYTICKS  
PLTAXISTIGHT  
SCORES GRIDCVRESULTSMEANTESTSCORERESHAPELENCRANGE  
LENGAMMARANGE  
DRAW HEATMAP OF THE VALIDATION ACCURACY AS A FUNCTION OF GAMMA AND C

THE SCORE ARE ENCODED AS COLORS WITH THE HOT COLORMAP WHICH VARIES FROM DARK  
RED TO BRIGHT YELLOW AS THE MOST INTERESTING SCORES ARE ALL LOCATED IN THE  
092 TO 097 RANGE WE USE A CUSTOM NORMALIZER TO SET THE MIDPOINT TO 092 SO  
AS TO MAKE IT EASIER TO VISUALIZE THE SMALL VARIATIONS OF SCORE VALUES IN THE  
INTERESTING RANGE WHILE NOT BRUTALLY COLLAPSING ALL THE LOW SCORE VALUES TO  
THE SAME COLOR  
PLTFIGUREFIGSIZE8 6  
PLTSUBPLOTSADJUSTLEFT2 RIGHT095 BOTTOM015 TOP095  
PLTIMSHOWSCORES INTERPOLATIONNEAREST CMAPPLTCMHOT  
NORMMIDPOINTNORMALIZEVMIN02 MIDPOINT092  
PLTXLABELGAMMA  
PLTYLABELC  
PLTCOLORBAR  
PLXTICKSNPARANGELENGAMMARANGE GAMMARANGE ROTATION45  
PLTYTICKSNPARANGELENCRANGE CRANGE  
PLTTITLEVALIDATION ACCURACY  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 3597 SECONDS  
1430 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

528 WORKING WITH TEXT DOCUMENTS

EXAMPLES CONCERNING THE SKLEARNFEATUREEXTRACTIONTEXT MODULE

NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)

5281 FEATUREHASHER AND DICTVECTORIZER COMPARISON

COMPARES FEATUREHASHER AND DICTVECTORIZER BY USING BOTH TO VECTORIZE TEXT DOCUMENTS

THE EXAMPLE DEMONSTRATES SYNTAX AND SPEED ONLY IT DOESN'T ACTUALLY DO ANYTHING USEFUL WITH THE EXTRACTED VECTORS

SEE THE EXAMPLE SCRIPTS DOCUMENTCLASSIFICATION20NEWSGROUPSCLUSTERINGPY FOR ACTUAL LEARNING ON TEXT DOCUMENTS

A DISCREPANCY BETWEEN THE NUMBER OF TERMS REPORTED FOR DICTVECTORIZER AND FOR FEATUREHASHER IS TO BE EXPECTED DUE TO HASH COLLISIONS

OUT

USAGE HOMECIRCLECIPROJECTEXAMPLETEXTPLOTHASHINGVSDICTVECTORIZERPY N

↪FEATURESFORHASHING

THE DEFAULT NUMBER OF FEATURES IS 2 18

LOADING 20 NEWSGROUPS TRAINING DATA

3803 DOCUMENTS 6245MB

DICTVECTORIZER

DONE IN 0979095S AT 6378MBS

FOUND 47928 UNIQUE TERMS

FEATUREHASHER ON FREQUENCY DICTS

DONE IN 0816599S AT 7647MBS

FOUND 43873 UNIQUE TERMS

FEATUREHASHER ON RAW TOKENS

DONE IN 0935579S AT 6675MBS

FOUND 43873 UNIQUE TERMS

AUTHOR LARS BUITINCK

LICENSE BSD 3 CLAUSE

FROM COLLECTIONS IMPORT DEFAULTDICT

IMPORT RE

IMPORT SYS

FROM TIME IMPORT TIME

IMPORT NUMPY AS NP

FROM SKLEARNDATASETS IMPORT FETCH20NEWSGROUPS

FROM SKLEARNFEATUREEXTRACTION IMPORT DICTVECTORIZER FEATUREHASHER

528 WORKING WITH TEXT DOCUMENTS 1431

SCIKITLEARN USER GUIDE RELEASE 0213  
DEFNNONZEROCOLUMNSX  
RETURNS THE NUMBER OF NONZERO COLUMNS IN A CSR MATRIX X  
RETURNLENNPUNIQUEXNONZERO1  
DEFTOKENSDOC  
EXTRACT TOKENS FROM DOC  
THIS USES A SIMPLE REGEX TO BREAK STRINGS INTO TOKENS FOR A MORE  
PRINCIPLED APPROACH SEE COUNTVECTORIZER OR TFIDFVECTORIZER  
  
RETURNOKLOWER FORTOKINREFINDALLRW DOC  
DEFTOKENFREQSDOC  
EXTRACT A DICT MAPPING TOKENS FROM DOC TO THEIR FREQUENCIES  
FREQ DEFAULTDICTINT  
FORTOKINTOKENSDOC  
FREQTOK 1  
RETURNFREQ  
CATEGORIES  
ALTATHEISM  
COMPGRAPHICS  
COMPSYSIBMPCHARDWARE  
MISCFORSALE  
RECAUTOS  
SCISPACE  
TALKRELIGIONMISC  
  
UNCOMMENT THE FOLLOWING LINE TO USE A LARGER SET 11K DOCUMENTS  
CATEGORIES NONE  
PRINTDOC  
PRINTUSAGE SNFEATURESFORHASHING SYSARGV0  
PRINT THE DEFAULT NUMBER OF FEATURES IS 2 18  
PRINT  
TRY  
NFEATURES INTSYSARGV1  
EXCEPTINDEXERROR  
NFEATURES 2 18  
EXCEPTVALUEERROR  
PRINTNOT A VALID NUMBER OF FEATURES R SYSARGV1  
SYSEXIT1  
PRINTLOADING 20 NEWSGROUPS TRAINING DATA  
RAWDATA FETCH20NEWSGROUPSSUBSETTRAIN CATEGORIESCATEGORIESDATA  
DATASIZEMB SUMLENSCODEUTF8 FORSINRAWDATA 1E6  
PRINTDDOCUMENTS 03FMB LENRAWDATA DATASIZEMB  
PRINT  
PRINTDICTVECTORIZER  
TO TIME  
VECTORIZER DICTVECTORIZER  
VECTORIZERFITTRANSFORMTOKENFREQSD FORDINRAWDATA  
1432 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213

DURATION TIME TO

PRINTDONE IN FS AT03FMBS DURATION DATASIZEMB DURATION

PRINTFOUND DUNIQUE TERMS LENVECTORIZERGETFEATURENAMES

PRINT

PRINTFEATUREHASHER ON FREQUENCY DICTS

TO TIME

HASHER FEATUREHASHERNFEATURESNFEATURES

X HASHERTRANSFORMTOKENFREQSD FORDINRAWDATA

DURATION TIME TO

PRINTDONE IN FS AT03FMBS DURATION DATASIZEMB DURATION

PRINTFOUND DUNIQUE TERMS NNONZEROCOLUMNSX

PRINT

PRINTFEATUREHASHER ON RAW TOKENS

TO TIME

HASHER FEATUREHASHERNFEATURESNFEATURES INPUTYPESTRING

X HASHERTRANSFORMTOKENSD FORDINRAWDATA

DURATION TIME TO

PRINTDONE IN FS AT03FMBS DURATION DATASIZEMB DURATION

PRINTFOUND DUNIQUE TERMS NNONZEROCOLUMNSX

TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 3036 SECONDS

NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE

5282 CLUSTERING TEXT DOCUMENTS USING KMEANS

THIS IS AN EXAMPLE SHOWING HOW THE SCIKITLEARN CAN BE USED TO CLUSTER DOCUMENTS BY TOPICS USING A BAGOFWORDS APPROACH THIS EXAMPLE USES A SCIPYSPARSE MATRIX TO STORE THE FEATURES INSTEAD OF STANDARD NUMPY ARRAYS

TWO FEATURE EXTRACTION METHODS CAN BE USED IN THIS EXAMPLE

- TFIDFVECTORIZER USES A INMEMORY VOCABULARY A PYTHON DICT TO MAP THE MOST FREQUENT WORDS TO FEATURES INDICES AND HENCE COMPUTE A WORD OCCURRENCE FREQUENCY SPARSE MATRIX THE WORD FREQUENCIES ARE THEN REWEIGHTED USING THE INVERSE DOCUMENT FREQUENCY IDF VECTOR COLLECTED FEATUREWISE OVER THE CORPUS
- HASHINGVECTORIZER HASHES WORD OCCURRENCES TO A FIXED DIMENSIONAL SPACE POSSIBLY WITH COLLISIONS THE WORD COUNT VECTORS ARE THEN NORMALIZED TO EACH HAVE L2NORM EQUAL TO ONE PROJECTED TO THE EUCLIDEAN UNITBALL WHICH SEEMS TO BE IMPORTANT FOR KMEANS TO WORK IN HIGH DIMENSIONAL SPACE

HASHINGVECTORIZER DOES NOT PROVIDE IDF WEIGHTING AS THIS IS A STATELESS MODEL THE FIT METHOD DOES NOTHING

WHEN IDF WEIGHTING IS NEEDED IT CAN BE ADDED BY PIPELINING ITS OUTPUT TO A TFIDFTRANSFORMER INSTANCE

TWO ALGORITHMS ARE DEMOED ORDINARY KMEANS AND ITS MORE SCALABLE COUSIN MINIBATCH KMEANS

ADDITIONALLY LATENT SEMANTIC ANALYSIS CAN ALSO BE USED TO REDUCE DIMENSIONALITY AND DISCOVER LATENT PATTERNS IN THE DATA

IT CAN BE NOTED THAT KMEANS AND MINIBATCH KMEANS ARE VERY SENSITIVE TO FEATURE SCALING AND THAT IN THIS CASE THE IDF WEIGHTING HELPS IMPROVE THE QUALITY OF THE CLUSTERING BY QUITE A LOT AS MEASURED AGAINST THE “GROUND TRUTH” PROVIDED BY THE CLASS LABEL ASSIGNMENTS OF THE 20 NEWSGROUPS DATASET

THIS IMPROVEMENT IS NOT VISIBLE IN THE SILHOUETTE COEFFICIENT WHICH IS SMALL FOR BOTH AS THIS MEASURE SEEM TO SUFFER FROM THE PHENOMENON CALLED “CONCENTRATION OF MEASURE” OR “CURSE OF DIMENSIONALITY” FOR HIGH DIMENSIONAL DATASETS

528 WORKING WITH TEXT DOCUMENTS 1433

SCIKITLEARN USER GUIDE RELEASE 0213

SUCH AS TEXT DATA OTHER MEASURES SUCH AS VMEASURE AND ADJUSTED RAND INDEX ARE INFORMATION THEORETIC BASED EVALUATION SCORES AS THEY ARE ONLY BASED ON CLUSTER ASSIGNMENTS RATHER THAN DISTANCES HENCE NOT AFFECTED BY THE CURSE OF DIMENSIONALITY

NOTE AS KMEANS IS OPTIMIZING A NONCONVEX OBJECTIVE FUNCTION IT WILL LIKELY END UP IN A LOCAL OPTIMUM SEVERAL RUNS WITH INDEPENDENT RANDOM INIT MIGHT BE NECESSARY TO GET A GOOD CONVERGENCE

OUT

USAGE PLOTDOCUMENTCLUSTERINGPY OPTIONS

OPTIONS

H HELP SHOW THIS HELP MESSAGE AND EXIT

LSANCOMPONENTS PREPROCESS DOCUMENTS WITH LATENT SEMANTIC ANALYSIS

NOMINIBATCH USE ORDINARY KMEANS ALGORITHM IN BATCH MODE

NOIDF DISABLE INVERSE DOCUMENT FREQUENCY FEATURE WEIGHTING

USEHASHING USE A HASHING FEATURE VECTORIZER

NFEATURESNFEATURES

MAXIMUM NUMBER OF FEATURES DIMENSIONS TO EXTRACT FROM TEXT

VERBOSE PRINT PROGRESS REPORTS INSIDE KMEANS ALGORITHM

LOADING 20 NEWSGROUPS DATASET FOR CATEGORIES

ALTATHEISM TALKRELIGIONMISC COMPGraphics SCISPACE

3387 DOCUMENTS

4 CATEGORIES

EXTRACTING FEATURES FROM THE TRAINING DATASET USING A SPARSE VECTORIZER

DONE IN 0691398S

NSAMPLES 3387 NFEATURES 10000

CLUSTERING SPARSE DATA WITH MINIBATCHKMEANSBATCHSIZE1000 INITSIZE1000 N

↪CLUSTERS4 NINIT1

VERBOSEFALSE

DONE IN 0057S

HOMOGENEITY 0359

COMPLETENESS 0440

VMEASURE 0396

ADJUSTED RANDINDEX 0253

SILHOUETTE COEFFICIENT 0007

TOP TERMS PER CLUSTER

CLUSTER 0 HENRY ALASKA TORONTO MOON ZOO SPENCER AURORA SPACE NSMCA ZOOLOGY

CLUSTER 1 COM Graphics UNIVERSITY POSTING HOST NNTP KNOW UK ARTICLE CS

CLUSTER 2 GOD COM SANDVIK PEOPLE KEITH MORALITY SGI KENT LIVESEY JESUS

CLUSTER 3 SPACE NASA ACCESS GOV DIGEX PAT SHUTTLE HST ORBIT NET

AUTHOR PETER PRETTENHOFER PETERPRETTENHOFERGMAILCOM

LARS BUITINCK

LICENSE BSD 3 CLAUSE

FROM SKLEARNDATASETS IMPORT FETCH20NEWSGROUPS

FROM SKLEARNDECOMPOSITION IMPORT TRUNCATEDSVD

FROM SKLEARNFEATUREEXTRACTIONTEXT IMPORT TFIDFVECTORIZER

FROM SKLEARNFEATUREEXTRACTIONTEXT IMPORT HASHINGVECTORIZER

1434 CHAPTER 5 EXAMPLES



```
SCIKITLEARN USER GUIDE RELEASE 0213
FROM SKLEARNFEATUREEXTRACTIONTEXT IMPORT TFIDFTRANSFORMER
FROM SKLEARNPIPELINE IMPORT MAKEPIPELINE
FROM SKLEARNPREPROCESSING IMPORT NORMALIZER
FROM SKLEARN IMPORT METRICS
FROM SKLEARNCLUSTER IMPORT KMEANS MINIBATCHKMEANS
IMPORT LOGGING
FROM OPTPARSE IMPORT OPTIONPARSER
IMPORT SYS
FROM TIME IMPORT TIME
IMPORT NUMPY AS NP
    DISPLAY PROGRESS LOGS ON STDOUT
LOGGINGBASICCONFIGLEVELLOGGINGINFO
FORMAT ASCTIMES LEVELNAMES MESSAGES
    PARSE COMMANDLINE ARGUMENTS
OP  OPTIONPARSER
OPADDOPTIONLSA
DESTNCOMPONENTS TYPEINT
HELPPREPROCESS DOCUMENTS WITH LATENT SEMANTIC ANALYSIS
OPADDOPTIONNMINIBATCH
ACTIONSTOREFALSE DESTMINIBATCH DEFAULTTRUE
HELPUSE ORDINARY KMEANS ALGORITHM IN BATCH MODE
OPADDOPTIONNOIDF
ACTIONSTOREFALSE DESTUSEIDF DEFAULTTRUE
HELPPDISABLE INVERSE DOCUMENT FREQUENCY FEATURE WEIGHTING
OPADDOPTIONUSEHASHING
ACTIONSTORETRUE DEFAULTFALSE
HELPUSE A HASHING FEATURE VECTORIZER
OPADDOPTIONNFEATURES TYPEINT DEFAULT10000
HELPMAXIMUM NUMBER OF FEATURES DIMENSIONS
    TO EXTRACT FROM TEXT
OPADDOPTIONVERBOSE
ACTIONSTORETRUE DESTVERBOSE DEFAULTFALSE
HELPPPRINT PROGRESS REPORTS INSIDE KMEANS ALGORITHM
PRINTDOC
OPPRINTHELP
DEFISINTERACTIVE
RETURN NOT HASATTRSYSMODULESMAIN FILE
    WORKAROUND FOR JUPYTER NOTEBOOK AND IPYTHON CONSOLE
ARGV  IFISINTERACTIVE ELSESYSARGV1
OPTS ARGS  OPPARSEARGSARGV
IFLENARGS  0
OPERRORTHIS SCRIPT TAKES NO ARGUMENTS
SYSEXIT1

    LOAD SOME CATEGORIES FROM THE TRAINING SET
528 WORKING WITH TEXT DOCUMENTS 1435
```

```
SCIKITLEARN USER GUIDE RELEASE 0213
CATEGORIES
ALTATHEISM
TALKRELIGIONMISC
COMPGRAPHICS
SCISPACE

UNCOMMENT THE FOLLOWING TO DO THE ANALYSIS ON ALL THE CATEGORIES
CATEGORIES NONE
PRINTLOADING 20 NEWSGROUPS DATASET FOR CATEGORIES
PRINTCATEGORIES
DATASET FETCH20NEWSGROUPSSUBSETALL CATEGORIESCATEGORIES
SHUFFLETRUE RANDOMSTATE42
PRINTDOCUMENTS LENDATASETDATA
PRINTDCATEGORIES LENDATASETTARGETNAMES
PRINT
LABELS DATASETTARGET
TRUEK NPUNIQUELABELSSHAPE0
PRINTEXTRACTING FEATURES FROM THE TRAINING DATASET
USING A SPARSE VECTORIZER
TO TIME
IFOPTSUSEHASHING
IFOPTSUSEIDF
PERFORM AN IDF NORMALIZATION ON THE OUTPUT OF HASHINGVECTORIZER
HASHER HASHINGVECTORIZERNFEATURESOPTSNFEATURES
STOPWORDSENGLISH ALTERNATESIGNFALSE
NORMNONE BINARYFALSE
VECTORIZER MAKEPIPELINEHASHER TFIDFTRANSFORMER
ELSE
VECTORIZER HASHINGVECTORIZERNFEATURESOPTSNFEATURES
STOPWORDSENGLISH
ALTERNATESIGNFALSE NORML2
BINARYFALSE
ELSE
VECTORIZER TFIDFVECTORIZERMAXDF05 MAXFEATURESOPTSNFEATURES
MINDF2 STOPWORDSENGLISH
USEIDFOPTSUSEIDF
X VECTORIZERFITTRANSFORMDATASETDATA
PRINTDONE IN FS TIME TO
PRINTNSAMPLES D NFEATURES D XSHAPE
PRINT
IFOPTSNCOMPONENTS
PRINTPERFORMING DIMENSIONALITY REDUCTION USING LSA
TO TIME
VECTORIZER RESULTS ARE NORMALIZED WHICH MAKES KMEANS BEHAVE AS
SPHERICAL KMEANS FOR BETTER RESULTS SINCE LSASVD RESULTS ARE
NOT NORMALIZED WE HAVE TO REDO THE NORMALIZATION
SVD TRUNCATEDSVDOPTSNCOMPONENTS
NORMALIZER NORMALIZERCOPYFALSE
LSA MAKEPIPELINESVD NORMALIZER
X LSAFITTRANSFORMX
1436 CHAPTER 5 EXAMPLES
```

```
SCIKITLEARN USER GUIDE RELEASE 0213
PRINTDONE IN FS  TIME  TO
EXPLAINEDVARIANCE  SVDEXPLAINEDVARIANCERATIOSUM
PRINTEXPLAINED VARIANCE OF THE SVD STEP FORMAT
INTEXPLAINEDVARIANCE 100
PRINT

DO THE ACTUAL CLUSTERING
IFOPTSMINIBATCH
KM  MINIBATCHKMEANSNCLUSTERSTRUEK INITKMEANS NINIT1
INITSIZE1000 BATCHSIZE1000 VERBOSEOPTSVERBOSE
ELSE
KM  KMEANSNCLUSTERSTRUEK INITKMEANS MAXITER100 NINIT1
VERBOSEOPTSVERBOSE
PRINTCLUSTERING SPARSE DATA WITH S  KM
TO  TIME
KMFITX
PRINTDONE IN 03FS  TIME  TO
PRINT
PRINTHOMOGENEITY 03F  METRICSHOMOGENEITYSCORELABELS KMLABELS
PRINTCOMPLETENESS 03F  METRICSCOMPLETENESSSCORELABELS KMLABELS
PRINTVMEASURE 03F  METRICSVMEASURESCORELABELS KMLABELS
PRINTADJUSTED RANDINDEX 3F
METRICSSADJUSTEDRANDSCORELABELS KMLABELS
PRINTSILHOUETTE COEFFICIENT 03F
METRICSSILHOUETTESCOREX KMLABELS SAMPLESIZE1000
PRINT
IF NOTOPTSUSEHASHING
PRINTTOP TERMS PER CLUSTER
IFOPTSNCOMPONENTS
ORIGINALSPACECENTROIDS  SVDINVERSETRANSFORMKMCLUSTERCENTERS
ORDERCENTROIDS  ORIGINALSPACECENTROIDSARGSORT 1
ELSE
ORDERCENTROIDS  KMCLUSTERCENTERSARGSORT 1
TERMS  VECTORIZERGETFEATURENAMES
FORIINRANGETRUEK
PRINTCLUSTER D  I END
FORINDINORDERCENTROIDSI 10
PRINTS  TERMSIND END
PRINT
TOTAL RUNNING TIME OF THE SCRIPT  0 MINUTES 1137 SECONDS
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE
528 WORKING WITH TEXT DOCUMENTS 1437
```

SCIKITLEARN USER GUIDE RELEASE 0213

5283 CLASSIFICATION OF TEXT DOCUMENTS USING SPARSE FEATURES

THIS IS AN EXAMPLE SHOWING HOW SCIKITLEARN CAN BE USED TO CLASSIFY DOCUMENTS BY TOPICS USING A BAGOFWORDS APPROACH THIS EXAMPLE USES A SCIPYSPARSE MATRIX TO STORE THE FEATURES AND DEMONSTRATES VARIOUS CLASSIFIERS THAT CAN EFFICIENTLY HANDLE SPARSE MATRICES

THE DATASET USED IN THIS EXAMPLE IS THE 20 NEWSGROUPS DATASET IT WILL BE AUTOMATICALLY DOWNLOADED THEN CACHED

THE BAR PLOT INDICATES THE ACCURACY TRAINING TIME NORMALIZED AND TEST TIME NORMALIZED OF EACH CLASSIFIER

OUT

USAGE PLOTDOCUMENTCLASSIFICATION20NEWSGROUPSPY OPTIONS

OPTIONS

H HELP SHOW THIS HELP MESSAGE AND EXIT

REPORT PRINT A DETAILED CLASSIFICATION REPORT

CHI2SELECTSELECTCHI2

SELECT SOME NUMBER OF FEATURES USING A CHISQUARED

TEST

CONFUSIONMATRIX PRINT THE CONFUSION MATRIX

TOP10 PRINT TEN MOST DISCRIMINATIVE TERMS PER CLASS FOR

EVERY CLASSIFIER

ALLCATEGORIES WHETHER TO USE ALL CATEGORIES OR NOT

USEHASHING USE A HASHING VECTORIZER

NFEATURESNFEATURES

NFEATURES WHEN USING THE HASHING VECTORIZER

FILTERED REMOVE NEWSGROUP INFORMATION THAT IS EASILY OVERFIT

HEADERS SIGNATURES AND QUOTING

1438 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
LOADING 20 NEWSGROUPS DATASET FOR CATEGORIES  
ALTATHEISM TALKRELIGIONMISC COMPGRAPHICS SCISPACE  
DATA LOADED  
2034 DOCUMENTS 3980MB TRAINING SET  
1353 DOCUMENTS 2867MB TEST SET  
4 CATEGORIES  
EXTRACTING FEATURES FROM THE TRAINING DATA USING A SPARSE VECTORIZER  
DONE IN 0412178S AT 9655MBS  
NSAMPLES 2034 NFEATURES 33809  
EXTRACTING FEATURES FROM THE TEST DATA USING THE SAME VECTORIZER  
DONE IN 0351330S AT 8162MBS  
NSAMPLES 1353 NFEATURES 33809

RIDGE CLASSIFIER

TRAINING  
RIDGECLASSIFIERSOLVERSAG TOL001  
TRAIN TIME 0132S  
TEST TIME 0001S  
ACCURACY 0896  
DIMENSIONALITY 33809  
DENSITY 1000000

PERCEPTRON

TRAINING  
PERCEPTRONMAXITER50  
TRAIN TIME 0017S  
TEST TIME 0002S  
ACCURACY 0888  
DIMENSIONALITY 33809  
DENSITY 0255302

PASSIVEAGGRESSIVE

TRAINING  
PASSIVEAGGRESSIVECLASSIFIERMAXITER50  
TRAIN TIME 0031S  
TEST TIME 0002S  
ACCURACY 0904  
DIMENSIONALITY 33809  
DENSITY 0694674

KNN

TRAINING  
KNEIGHBORSCLASSIFIERNNEIGHBORS10  
TRAIN TIME 0002S  
TEST TIME 0317S  
528 WORKING WITH TEXT DOCUMENTS 1439

SCIKITLEARN USER GUIDE RELEASE 0213  
ACCURACY 0858

RANDOM FOREST

TRAINING  
RANDOMFORESTCLASSIFIERNESTIMATORS100  
TRAIN TIME 1671S  
TEST TIME 0071S  
ACCURACY 0840

L2 PENALTY

TRAINING  
LINEARSVCDUALFALSE TOL0001  
TRAIN TIME 0145S  
TEST TIME 0002S  
ACCURACY 0900  
DIMENSIONALITY 33809  
DENSITY 1000000

TRAINING  
SGDCLASSIFIERMAXITER50  
TRAIN TIME 0030S  
TEST TIME 0002S  
ACCURACY 0902  
DIMENSIONALITY 33809  
DENSITY 0579380

L1 PENALTY

TRAINING  
LINEARSVCDUALFALSE PENALTYL1 TOL0001  
TRAIN TIME 0301S  
TEST TIME 0002S  
ACCURACY 0873  
DIMENSIONALITY 33809  
DENSITY 0005553

TRAINING  
SGDCLASSIFIERMAXITER50 PENALTYL1  
TRAIN TIME 0093S  
TEST TIME 0002S  
ACCURACY 0887  
DIMENSIONALITY 33809  
DENSITY 0022901

ELASTICNET PENALTY

SCIKITLEARN USER GUIDE RELEASE 0213  
TRAINING  
SGDCLASSIFIERMAXITER50 PENALTYELASTICNET  
TRAIN TIME 0252S  
TEST TIME 0002S  
ACCURACY 0899  
DIMENSIONALITY 33809  
DENSITY 0187472

NEARESTCENTROID AKA ROCCHIO CLASSIFIER

TRAINING  
NEARESTCENTROID  
TRAIN TIME 0004S  
TEST TIME 0002S  
ACCURACY 0855

NAIVE BAYES

TRAINING  
MULTINOMIALNBALPHA001  
TRAIN TIME 0003S  
TEST TIME 0001S  
ACCURACY 0899  
DIMENSIONALITY 33809  
DENSITY 1000000

TRAINING  
BERNOULLINBALPHA001  
TRAIN TIME 0004S  
TEST TIME 0003S  
ACCURACY 0884  
DIMENSIONALITY 33809  
DENSITY 1000000

TRAINING  
COMPLEMENTNBALPHA01  
TRAIN TIME 0004S  
TEST TIME 0001S  
ACCURACY 0911  
DIMENSIONALITY 33809  
DENSITY 1000000

LINEARSVC WITH L1BASED FEATURE SELECTION

TRAINING  
PIPELINESTEPSFEATURESELECTION  
SELECTFROMMODELESTIMATORLINEARSVC  
DUALFALSE PENALTYL1  
TOL0001  
CLASSIFICATION LINEARSVC  
528 WORKING WITH TEXT DOCUMENTS 1441

SCIKITLEARN USER GUIDE RELEASE 0213  
TRAIN TIME 0252S  
TEST TIME 0002S  
ACCURACY 0880  
AUTHOR PETER PRETTENHOFER PETERPRETTENHOFERGMAILCOM  
OLIVIER GRISEL OLIVIERGRISELENSTAORG  
MATHIEU BLONDEL MATHIEUMBLONDELORG  
LARS BUITINCK  
LICENSE BSD 3 CLAUSE  
IMPORT LOGGING  
IMPORT NUMPY AS NP  
FROM OPTPARSE IMPORT OPTIONPARSER  
IMPORT SYS  
FROM TIME IMPORT TIME  
IMPORT MATPLOTLIBPYPLOT AS PLT  
FROM SKLEARNDATASETS IMPORT FETCH20NEWSGROUPS  
FROM SKLEARNFEATUREEXTRACTIONTEXT IMPORT TFIDFVECTORIZER  
FROM SKLEARNFEATUREEXTRACTIONTEXT IMPORT HASHINGVECTORIZER  
FROM SKLEARNFEATURESELECTION IMPORT SELECTFROMMODEL  
FROM SKLEARNFEATURESELECTION IMPORT SELECTKBEST CHI2  
FROM SKLEARNLINEARMODEL IMPORT RIDGECLASSIFIER  
FROM SKLEARNPIPELINE IMPORT PIPELINE  
FROM SKLEARNNSVM IMPORT LINEARSVC  
FROM SKLEARNLINEARMODEL IMPORT SGDCLASSIFIER  
FROM SKLEARNLINEARMODEL IMPORT PERCEPTRON  
FROM SKLEARNLINEARMODEL IMPORT PASSIVEAGGRESSIVECLASSIFIER  
FROM SKLEARNNAIVEBAYES IMPORT BERNOULLINB COMPLEMENTNB MULTINOMIALNB  
FROM SKLEARNNEIGHBORS IMPORT KNEIGHBORSCLASSIFIER  
FROM SKLEARNNEIGHBORS IMPORT NEARESTCENTROID  
FROM SKLEARNENSEMBLE IMPORT RANDOMFORESTCLASSIFIER  
FROM SKLEARNUTILSEXTMATH IMPORT DENSITY  
FROM SKLEARN IMPORT METRICS  
  DISPLAY PROGRESS LOGS ON STDOUT  
LOGGINGBASICCONFIGLEVELLOGGINGINFO  
FORMAT ASCTIMES LEVELNAMES MESSAGES  
  PARSE COMMANDLINE ARGUMENTS  
OP OPTIONPARSER  
OPADDOPTIONREPORT  
ACTIONSTORETRUE DESTPRINTREPORT  
HELPPRINT A DETAILED CLASSIFICATION REPORT  
OPADDOPTIONCHI2SELECT  
ACTIONSTORE TYPEINT DESTSELECTCHI2  
HELPSELECT SOME NUMBER OF FEATURES USING A CHISQUARED TEST  
OPADDOPTIONCONFUSIONMATRIX  
ACTIONSTORETRUE DESTPRINTCM  
HELPPRINT THE CONFUSION MATRIX  
OPADDOPTIONTOP10  
ACTIONSTORETRUE DESTPRINTTOP10  
1442 CHAPTER 5 EXAMPLES



SCIKITLEARN USER GUIDE RELEASE 0213  
HELPPRINT TEN MOST DISCRIMINATIVE TERMS PER CLASS  
FOR EVERY CLASSIFIER  
OPADDOPTIONALLCATEGORIES  
ACTIONSTORETRUE DESTALLCATEGORIES  
HELPPWHETHER TO USE ALL CATEGORIES OR NOT  
OPADDOPTIONUSEHASHING  
ACTIONSTORETRUE  
HELPUSE A HASHING VECTORIZER  
OPADDOPTIONNFEATURES  
ACTIONSTORE TYPEINT DEFAULT2 16  
HELPNFEATURES WHEN USING THE HASHING VECTORIZER  
OPADDOPTIONFILTERED  
ACTIONSTORETRUE  
HELPREMOVE NEWSGROUP INFORMATION THAT IS EASILY OVERFIT  
HEADERS SIGNATURES AND QUOTING  
DEFISINTERACTIVE  
RETURN NOT HASATTRSYSMODULESMAIN FILE  
WORKAROUND FOR JUPYTER NOTEBOOK AND IPYTHON CONSOLE  
ARGV IFISINTERACTIVE ELSESYSARGV1  
OPTS ARGS OPPARSEARGSARGV  
IFLENARGS 0  
OPERRORTHIS SCRIPT TAKES NO ARGUMENTS  
SYSEXIT1  
PRINTDOC  
OPPRINTHELP  
PRINT  
  
LOAD SOME CATEGORIES FROM THE TRAINING SET  
IFOPTSALLCATEGORIES  
CATEGORIES NONE  
ELSE  
CATEGORIES  
ALTATHEISM  
TALKRELIGIONMISC  
COMPGRAPHICS  
SCISPACE  
  
IFOPTSFILTERED  
REMOVE HEADERS FOOTERS QUOTES  
ELSE  
REMOVE  
PRINTLOADING 20 NEWSGROUPS DATASET FOR CATEGORIES  
PRINTCATEGORIES IFCATEGORIES ELSEALL  
DATATRAIN FETCH20NEWSGROUPSSUBSETTRAIN CATEGORIESCATEGORIES  
SHUFFLETRUE RANDOMSTATE42  
REMOVEREMOVE  
DATATEST FETCH20NEWSGROUPSSUBSETTEST CATEGORIESCATEGORIES  
528 WORKING WITH TEXT DOCUMENTS 1443

SCIKITLEARN USER GUIDE RELEASE 0213  
SHUFFLETRUE RANDOMSTATE42  
REMOVEREMOVE  
PRINTDATA LOADED  
ORDER OF LABELS IN TARGETNAMES CAN BE DIFFERENT FROM CATEGORIES  
TARGETNAMES DATATRAINTARGETNAMES  
DEFSIZEEMBDOCS  
RETURNSUMLENSENCODEUTF8 FORSINDOCS 1E6  
DATATRAINSIZEMB SIZEMBDATATRAINDATA  
DATATESTSIZEMB SIZEMBDATATESTDATA  
PRINTDDOCUMENTS 03FMB TRAINING SET  
LENDATATRAINDATA DATATRAINSIZEMB  
PRINTDDOCUMENTS 03FMB TEST SET  
LENDATATESTDATA DATATESTSIZEMB  
PRINTDCATEGORIES LENTARGETNAMES  
PRINT  
SPLIT A TRAINING SET AND A TEST SET  
YTRAIN YTEST DATATRAINTARGET DATATESTTARGET  
PRINTEXTRACTING FEATURES FROM THE TRAINING DATA USING A SPARSE VECTORIZER  
TO TIME  
IFOPTSUSEHASHING  
VECTORIZER HASHINGVECTORIZERSTOPWORDSENGLISH ALTERNATESIGNFALSE  
NFEATURESOPTSNFEATURES  
XTRAIN VECTORIZERTRANSFORMDATATRAINDATA  
ELSE  
VECTORIZER TFIDFVECTORIZERSUBLINEARTFTRUE MAXDF05  
STOPWORDSENGLISH  
XTRAIN VECTORIZERFITTRANSFORMDATATRAINDATA  
DURATION TIME TO  
PRINTDONE IN FS AT03FMBS DURATION DATATRAINSIZEMB DURATION  
PRINTNSAMPLES D NFEATURES D XTRAINSHAPE  
PRINT  
PRINTEXTRACTING FEATURES FROM THE TEST DATA USING THE SAME VECTORIZER  
TO TIME  
XTEST VECTORIZERTRANSFORMDATATESTDATA  
DURATION TIME TO  
PRINTDONE IN FS AT03FMBS DURATION DATATESTSIZEMB DURATION  
PRINTNSAMPLES D NFEATURES D XTESTSHAPE  
PRINT  
MAPPING FROM INTEGER FEATURE NAME TO ORIGINAL TOKEN STRING  
IFOPTSUSEHASHING  
FEATURENAMES NONE  
ELSE  
FEATURENAMES VECTORIZERGETFEATURENAMES  
IFOPTSSELECTCHI2  
PRINTEXTRACTING DBEST FEATURES BY A CHISQUARED TEST  
OPTSSELECTCHI2  
TO TIME  
1444 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
CH2 SELECTKBESTCHI2 KOPTSSELECTCHI2  
XTRAIN CH2FITTRANSFORMXTRAIN YTRAIN  
XTEST CH2TRANSFORMXTEST  
IFFEATURENAMES  
KEEP SELECTED FEATURE NAMES  
FEATURENAMES FEATURENAMESI FORI  
INCH2GETSUPPORTINDICESTRUE  
PRINTDONE IN FS TIME TO  
PRINT  
IFFEATURENAMES  
FEATURENAMES NPASARRAYFEATURENAMES  
DEFTRIMS  
TRIM STRING TO FIT ON TERMINAL ASSUMING 80COLUMN DISPLAY  
RETURNSIFLENS 80 ELSE577  
  
BENCHMARK CLASSIFIERS  
DEFBENCHMARKCLF  
PRINT80  
PRINTTRAINING  
PRINTCLF  
TO TIME  
CLFFITXTRAIN YTRAIN  
TRAINTIME TIME TO  
PRINTTRAIN TIME 03FS TRAINTIME  
TO TIME  
PRED CLFPREDICTXTEST  
TESTTIME TIME TO  
PRINTTEST TIME 03FS TESTTIME  
SCORE METRICSACCURACYSCOREYTEST PRED  
PRINTACCURACY 03F SCORE  
IFHASATTRCLF COEF  
PRINTDIMENSIONALITY D CLFCOEFSHAPE1  
PRINTDENSITY F DENSITYCLFCOEF  
IFOPTSPRINTTOP10 ANDFEATURENAMES IS NOTNONE  
PRINTTOP 10 KEYWORDS PER CLASS  
FORI LABEL INENUMERATETARGETNAMES  
TOP10 NPARGSORTCLFCOEFI10  
PRINTTRIMSS LABEL JOINFEATURENAMESTOP10  
PRINT  
IFOPTSPRINTREPORT  
PRINTCLASSIFICATION REPORT  
PRINTMETRICSCCLASSIFICATIONREPORTYTEST PRED  
TARGETNAMESTARGETNAMES  
IFOPTSPRINTCM  
PRINTCONFUSION MATRIX  
PRINTMETRICSCONFUSIONMATRIXYTEST PRED  
528 WORKING WITH TEXT DOCUMENTS 1445

SCIKITLEARN USER GUIDE RELEASE 0213  
PRINT  
CLFDESCR STRCLFSPLITO  
RETURNCLFDESCR SCORE TRAINTIME TESTTIME  
RESULTS  
FORCLF NAME IN  
RIDGECLASSIFIERTOL1E2 SOLVERSAG RIDGE CLASSIFIER  
PERCEPTRONMAXITER50 TOL1E3 PERCEPTRON  
PASSIVEAGGRESSIVECLASSIFIERMAXITER50 TOL1E3  
PASSIVEAGGRESSIVE  
KNEIGHBORSCLASSIFIERNNEIGHBORS10 KNN  
RANDOMFORESTCLASSIFIERNESTIMATORS100 RANDOM FOREST  
PRINT80  
PRINTNAME  
RESULTSAPPENDBENCHMARKCLF  
FORPENALTY INL2 L1  
PRINT80  
PRINTSPENALTY PENALTYUPPER  
  TRAIN LIBLINEAR MODEL  
RESULTSAPPENDBENCHMARKLINEARSVCPENALTYPENALTY DUALFALSE  
TOL1E3  
  TRAIN SGD MODEL  
RESULTSAPPENDBENCHMARKSGDCLASSIFIERALPHA0001 MAXITER50  
PENALTYPENALTY  
  TRAIN SGD WITH ELASTIC NET PENALTY  
PRINT80  
PRINTELASTICNET PENALTY  
RESULTSAPPENDBENCHMARKSGDCLASSIFIERALPHA0001 MAXITER50  
PENALTYELASTICNET  
  TRAIN NEARESTCENTROID WITHOUT THRESHOLD  
PRINT80  
PRINTNEARESTCENTROID AKA ROCCHIO CLASSIFIER  
RESULTSAPPENDBENCHMARKNEARESTCENTROID  
  TRAIN SPARSE NAIVE BAYES CLASSIFIERS  
PRINT80  
PRINTNAIVE BAYES  
RESULTSAPPENDBENCHMARKMULTINOMIALNBALPHA01  
RESULTSAPPENDBENCHMARKBERNOULLINBALPHA01  
RESULTSAPPENDBENCHMARKCOMPLEMENTNBALPHA1  
PRINT80  
PRINTLINEARSVC WITH L1BASED FEATURE SELECTION  
  THE SMALLER C THE STRONGER THE REGULARIZATION  
  THE MORE REGULARIZATION THE MORE SPARSITY  
RESULTSAPPENDBENCHMARKPIPELINE  
FEATURESELECTION SELECTFROMMODELLINEARSVCPENALTYL1 DUALFALSE  
TOL1E3  
CLASSIFICATION LINEARSVCPENALTYL2  
  MAKE SOME PLOTS  
INDICES NPARANGELENRESULTS  
1446 CHAPTER 5 EXAMPLES

SCIKITLEARN USER GUIDE RELEASE 0213  
RESULTS XI FORXINRESULTS FORIINRANGE4  
CLFNAMES SCORE TRAININGTIME TESTTIME RESULTS  
TRAININGTIME NPARRAYTRAININGTIME NPMAXTRAININGTIME  
TESTTIME NPARRAYTESTTIME NPMAXTESTTIME  
PLTFIGUREFIGSIZE12 8  
PLTTITLESCORE  
PLTBARHINDICES SCORE 2 LABELSCORE COLORNAVY  
PLTBARHINDICES 3 TRAININGTIME 2 LABELTRAINING TIME  
COLORC  
PLTBARHINDICES 6 TESTTIME 2 LABELTEST TIME COLORDARKORANGE  
PLTYTICKS  
PLTLEGENDLOCBEST  
PLTSUBPLOTSADJUSTLEFT25  
PLTSUBPLOTSADJUSTTOP95  
PLTSUBPLOTSADJUSTBOTTOM05  
FORI CINZIPINDICES CLFNAMES  
PLTTEXT3 I C  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 4728 SECONDS  
529 DECISION TREES  
EXAMPLES CONCERNING THE SKLEARN TREE MODULE  
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE  
5291 DECISION TREE REGRESSION  
A 1D REGRESSION WITH DECISION TREE  
THE DECISION TREES IS USED TO FIT A SINE CURVE WITH ADDITION NOISY OBSERVATION AS A RESULT IT LEARNS LOCAL LINEAR REGRESSIONS APPROXIMATING THE SINE CURVE  
WE CAN SEE THAT IF THE MAXIMUM DEPTH OF THE TREE CONTROLLED BY THE MAXDEPTH PARAMETER IS SET TOO HIGH THE DECISION TREES LEARN TOO FINE DETAILS OF THE TRAINING DATA AND LEARN FROM THE NOISE IE THEY OVERFIT  
529 DECISION TREES 1447

```
SCIKITLEARN USER GUIDE RELEASE 0213
PRINTDOC
IMPORT THE NECESSARY MODULES AND LIBRARIES
IMPORT NUMPY AS NP
FROM SKLEARNTREE IMPORT DECISIONTREEREgressor
IMPORT MATPLOTLIBPYplot AS PLT
CREATE A RANDOM DATASET
RNG NPRANDOMRANDOMSTATE1
X NPSORT5 RNGRAND80 1 AXIS0
Y NPSINXRavel
Y5 3 05 RNGRAND16
FIT REGRESSION MODEL
REGR1 DECISIONTREEREgressorMAXDEPTH2
REGR2 DECISIONTREEREgressorMAXDEPTH5
REGR1FITX Y
REGR2FITX Y
PREDICT
XTEST NPARANGE00 50 001 NPNEWAXIS
Y1 REGR1PREDICTXTEST
Y2 REGR2PREDICTXTEST
PLOT THE RESULTS
PLTFigure
1448 CHAPTER 5 EXAMPLES
```

SCIKITLEARN USER GUIDE RELEASE 0213  
PLTSCATTERX Y S20 EDGECOLORBLACK  
CDARKORANGE LABELDATA  
PLTPLOTXTEST Y1 COLORCORNFLOWERBLUE  
LABELMAXDEPTH2 LINEWIDTH2  
PLTPLOTXTEST Y2 COLORYELLOWGREEN LABELMAXDEPTH5 LINEWIDTH2  
PLTXLABELDATA  
PLTYLABELTARGET  
PLTTITLEDECISION TREE REGRESSION  
PLTLEGEND  
PLTSHOW  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0124 SECONDS  
NOTE [CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE](#)  
5292 MULTIOUTPUT DECISION TREE REGRESSION  
AN EXAMPLE TO ILLUSTRATE MULTIOUTPUT REGRESSION WITH DECISION TREE  
THEDECISION TREES IS USED TO PREDICT SIMULTANEOUSLY THE NOISY X AND Y OBSERVATIONS OF A CIRCLE GIVEN A SINGLE UNDERLYING  
FEATURE AS A RESULT IT LEARNS LOCAL LINEAR REGRESSIONS APPROXIMATING THE CIRCLE  
WE CAN SEE THAT IF THE MAXIMUM DEPTH OF THE TREE CONTROLLED BY THE MAXDEPTH PARAMETER IS SET TOO HIGH THE  
DECISION TREES LEARN TOO FINE DETAILS OF THE TRAINING DATA AND LEARN FROM THE NOISE IE THEY OVERFIT  
529 DECISION TREES 1449

```
SCIKITLEARN USER GUIDE RELEASE 0213
PRINTDOC
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARNTREE IMPORT DECISIONTREEREgressor
CREATE A RANDOM DATASET
RNG NPRANDOMRANDOMSTATE1
X NPSORT200 RNGRAND100 1 100 AXIS0
Y NPARRAYNPPI NPSINXRAVEL NPPI NPCOSXRAVELT
Y5 05 RNGRAND20 2
FIT REGRESSION MODEL
REGR1 DECISIONTREEREgressorMAXDEPTH2
REGR2 DECISIONTREEREgressorMAXDEPTH5
REGR3 DECISIONTREEREgressorMAXDEPTH8
REGR1FITX Y
REGR2FITX Y
REGR3FITX Y
PREDICT
XTEST NPARANGE1000 1000 001 NPNEWAXIS
Y1 REGR1PREDICTXTEST
Y2 REGR2PREDICTXTEST
Y3 REGR3PREDICTXTEST
1450 CHAPTER 5 EXAMPLES
```



```
SCIKITLEARN USER GUIDE RELEASE 0213
PLOT THE RESULTS
PLTFigure
S 25
PLTScatterY 0 Y 1 CNAVY SS
EDGEcolorBLACK LABELDATA
PLTScatterY1 0 Y1 1 CCORNFLOWERBLUE SS
EDGEcolorBLACK LABELMAXDEPTH2
PLTScatterY2 0 Y2 1 CRED SS
EDGEcolorBLACK LABELMAXDEPTH5
PLTScatterY3 0 Y3 1 CORANGE SS
EDGEcolorBLACK LABELMAXDEPTH8
PLTXLIM6 6
PLTYLIM6 6
PLTXLABELTARGET 1
PLTYLABELTARGET 2
PLTTITLEMULTIOUTPUT DECISION TREE REGRESSION
PLTLEGENDLOCBEST
PLTSHOW
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0109 SECONDS
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE
5293 PLOT THE DECISION SURFACE OF A DECISION TREE ON THE IRIS DATASET
PLOT THE DECISION SURFACE OF A DECISION TREE TRAINED ON PAIRS OF FEATURES OF THE IRIS DATASET
SEEDecision Tree FOR MORE INFORMATION ON THE ESTIMATOR
FOR EACH PAIR OF IRIS FEATURES THE DECISION TREE LEARNS DECISION BOUNDARIES MADE OF COMBINATIONS OF SIMPLE THRESHOLDING
RULES INFERRED FROM THE TRAINING SAMPLES
WE ALSO SHOW THE TREE STRUCTURE OF A MODEL BUILT ON ALL OF THE FEATURES
529 DECISION TREES 1451
```



SCIKITLEARN USER GUIDE RELEASE 0213

```
•
PRINTDOC
IMPORT NUMPY AS NP
IMPORT MATPLOTLIBPYPLOT AS PLT
FROM SKLEARNDATASETS IMPORT LOADIRIS
FROM SKLEARNTREE IMPORT DECISIONTREECLASSIFIER PLOTTREE
PARAMETERS
NCLASSES 3
PLOTCOLORS RYB
PLOTSTEP 002
LOAD DATA
IRIS LOADIRIS
FORPAIRIDX PAIR INENUMERATE0 1 0 2 0 3
1 2 1 3 2 3
WE ONLY TAKE THE TWO CORRESPONDING FEATURES
X IRISDATA PAIR
Y IRISTARGET
TRAIN
CLF DECISIONTREECLASSIFIERFITX Y
PLOT THE DECISION BOUNDARY
PLTSUBPLOT2 3 PAIRIDX 1
529 DECISION TREES 1453
```

```
SCIKITLEARN USER GUIDE RELEASE 0213
XMIN XMAX X 0MIN 1 X 0MAX 1
YMIN YMAX X 1MIN 1 X 1MAX 1
XX YY NPMESHGRIDNPARANGEXMIN XMAX PLOTSTEP
NPARANGEYMIN YMAX PLOTSTEP
PLTTIGHTLAYOUTHPAD05 WPAD05 PAD25
Z CLFPREDICTNPCXXRAVEL YYRAVEL
Z ZRESHAPEXXSHAPE
CS PLTCONTOURFXX YY Z CMAPPLTCMRDYLBU
PLTXLABELIRISFEATURENAMESPAIR0
PLTYLABELIRISFEATURENAMESPAIR1
PLOT THE TRAINING POINTS
FORI COLOR INZIPRANGENCLASSES PLOTCOLORS
IDX NPWHEREY I
PLTSCATTERXIDX 0 XIDX 1 CCOLOR LABELIRISTARGETNAMESI
CMAPPLTCMRDYLBU EDGECOLORBLACK S15
PLTSUPTITLEDECISION SURFACE OF A DECISION TREE USING PAIRED FEATURES
PLTLEGENDLOCLOWER RIGHT BORDERPAD0 HANDLETEXTPAD0
PLTAXISTIGHT
PLTFigure
CLF DECISIONTREECLASSIFIERFITIRISDATA IRISTARGET
PLOTTREECLF FILLEDTRUE
PLTSHOW
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0934 SECONDS
NOTE CLICK HERE TO DOWNLOAD THE FULL EXAMPLE CODE
5294 UNDERSTANDING THE DECISION TREE STRUCTURE
THE DECISION TREE STRUCTURE CAN BE ANALYSED TO GAIN FURTHER INSIGHT ON THE RELATION BETWEEN THE FEATURES AND THE TARGET
TO PREDICT IN THIS EXAMPLE WE SHOW HOW TO RETRIEVE
• THE BINARY TREE STRUCTURE
• THE DEPTH OF EACH NODE AND WHETHER OR NOT IT'S A LEAF
• THE NODES THAT WERE REACHED BY A SAMPLE USING THE DECISIONPATH METHOD
• THE LEAF THAT WAS REACHED BY A SAMPLE USING THE APPLY METHOD
• THE RULES THAT WERE USED TO PREDICT A SAMPLE
• THE DECISION PATH SHARED BY A GROUP OF SAMPLES
OUT
THE BINARY TREE STRUCTURE HAS 5 NODES AND HAS THE FOLLOWING TREE STRUCTURE
NODE0 TEST NODE GO TO NODE 1 IF X 3 0800000011920929 ELSE TO NODE 2
NODE1 LEAF NODE
NODE2 TEST NODE GO TO NODE 3 IF X 2 4950000047683716 ELSE TO NODE 4
NODE3 LEAF NODE
1454 CHAPTER 5 EXAMPLES
```

SCIKITLEARN USER GUIDE RELEASE 0213

NODE4 LEAF NODE

RULES USED TO PREDICT SAMPLE 0

DECISION ID NODE 0 XTEST0 3 24 0800000011920929

DECISION ID NODE 2 XTEST0 2 51 4950000047683716

THE FOLLOWING SAMPLES 0 1 SHARE THE NODE 0 2 IN THE TREE

IT IS 400 OF ALL NODES

IMPORT NUMPY AS NP

FROM SKLEARNMODELSELECTION IMPORT TRAINTESTSPLIT

FROM SKLEARNDATASETS IMPORT LOADIRIS

FROM SKLEARNTREE IMPORT DECISIONTREECLASSIFIER

IRIS LOADIRIS

X IRISDATA

Y IRISTARGET

XTRAIN XTEST YTRAIN YTEST TRAINTESTSPLITX Y RANDOMSTATE0

ESTIMATOR DECISIONTREECLASSIFIERMAXLEAFNODES3 RANDOMSTATE0

ESTIMATORFITXTRAIN YTRAIN

THE DECISION ESTIMATOR HAS AN ATTRIBUTE CALLED TREE WHICH STORES THE ENTIRE TREE STRUCTURE AND ALLOWS ACCESS TO LOW LEVEL ATTRIBUTES THE BINARY TREE TREE IS REPRESENTED AS A NUMBER OF PARALLEL ARRAYS THE ITH ELEMENT OF EACH ARRAY HOLDS INFORMATION ABOUT THE NODE I NODE 0 IS THE TREES ROOT NOTE SOME OF THE ARRAYS ONLY APPLY TO EITHER LEAVES OR SPLIT NODES RESP IN THIS CASE THE VALUES OF NODES OF THE OTHER TYPE ARE ARBITRARY

AMONG THOSE ARRAYS WE HAVE

LEFTCHILD ID OF THE LEFT CHILD OF THE NODE

RIGHTCHILD ID OF THE RIGHT CHILD OF THE NODE

FEATURE FEATURE USED FOR SPLITTING THE NODE

THRESHOLD THRESHOLD VALUE AT THE NODE

USING THOSE ARRAYS WE CAN PARSE THE TREE STRUCTURE

NNODES ESTIMATORTREENODECOUNT

CHILDRENLEFT ESTIMATORTREECHILDRENLEFT

CHILDRENRIGHT ESTIMATORTREECHILDRENRIGHT

FEATURE ESTIMATORTREEFEATURE

THRESHOLD ESTIMATORTREETHRESHOLD

THE TREE STRUCTURE CAN BE TRAVERSED TO COMPUTE VARIOUS PROPERTIES SUCH AS THE DEPTH OF EACH NODE AND WHETHER OR NOT IT IS A LEAF

NODEDEPTH NPZEROSSHAPENNODES DTYPEINT64

ISLEAVES NPZEROSSHAPENNODES DTYPEBOOL

STACK 0 1 SEED IS THE ROOT NODE ID AND ITS PARENT DEPTH

WHILELENSTACK 0

NODEID PARENTDEPTH STACKPOP

529 DECISION TREES 1455

```
SCIKITLEARN USER GUIDE RELEASE 0213
NODEDEPTHNODEID PARENTDEPTH 1
IF WE HAVE A TEST NODE
IFCHILDRENLEFTNODEID CHILDRENRIGHTNODEID
STACKAPPENDCHILDRENLEFTNODEID PARENTDEPTH 1
STACKAPPENDCHILDRENRIGHTNODEID PARENTDEPTH 1
ELSE
ISLEAVESNODEID TRUE
PRINTTHE BINARY TREE STRUCTURE HAS SNODES AND HAS
THE FOLLOWING TREE STRUCTURE
NNODES
FORIINRANGENNODES
IFISLEAVESI
PRINTSNODESLEAF NODE NODEDEPTHI T I
ELSE
PRINTSNODESTEST NODE GO TO NODE SIF X S SELSE TO
NODES
NODEDEPTHI T
I
CHILDRENLEFTI
FEATUREI
THRESHOLDI
CHILDRENRIGHTI

PRINT
FIRST LETS RETRIEVE THE DECISION PATH OF EACH SAMPLE THE DECISIONPATH
METHOD ALLOWS TO RETRIEVE THE NODE INDICATOR FUNCTIONS A NON ZERO ELEMENT OF
INDICATOR MATRIX AT THE POSITION I J INDICATES THAT THE SAMPLE I GOES
THROUGH THE NODE J
NODEINDICATOR ESTIMATORDECISIONPATHXTEST
SIMILARLY WE CAN ALSO HAVE THE LEAVES IDS REACHED BY EACH SAMPLE
LEAVEID ESTIMATORAPPLYXTEST
NOW ITS POSSIBLE TO GET THE TESTS THAT WERE USED TO PREDICT A SAMPLE OR
A GROUP OF SAMPLES FIRST LETS MAKE IT FOR THE SAMPLE
SAMPLEID 0
NODEINDEX NODEINDICATORINDICESNODEINDICATORINDPTRSAMPLEID
NODEINDICATORINDPTRSAMPLEID 1
PRINTRULES USED TO PREDICT SAMPLE S SAMPLEID
FORNODEID INNODEINDEX
IFLEAVEIDSAMPLEID NODEID
CONTINUE
IFXTESTSAMPLEID FEATURENODEID THRESHOLDNODEID
THRESHOLDSIGN
ELSE
THRESHOLDSIGN
PRINTDECISION ID NODE S XTEST SS SS S
NODEID
SAMPLEID
1456 CHAPTER 5 EXAMPLES
```

SCIKITLEARN USER GUIDE RELEASE 0213  
FEATURENODEID  
XTESTSAMPLEID FEATURENODEID  
THRESHOLDSIGN  
THRESHOLDNODEID  
FOR A GROUP OF SAMPLES WE HAVE THE FOLLOWING COMMON NODE  
SAMPLEIDS 0 1  
COMMONNODES NODEINDICATOROTOARRAYSAMPLEIDSSUMAXISO  
LENSAMPLEIDS  
COMMONNODEID NPARANGENNODESCOMMONNODES  
PRINTNTHE FOLLOWING SAMPLES SSHARE THE NODE SIN THE TREE  
SAMPLEIDS COMMONNODEID  
PRINTIT IS S OF ALL NODES 100 LENCOMMONNODEID NNODES  
TOTAL RUNNING TIME OF THE SCRIPT 0 MINUTES 0003 SECONDS  
529 DECISION TREES 1457





CHAPTER  
SIX  
API REFERENCE  
THIS IS THE CLASS AND FUNCTION REFERENCE OF SCIKITLEARN PLEASE REFER TO THE FULL USER GUIDE FOR FURTHER DETAILS AS THE CLASS  
AND FUNCTION RAW SPECIFICATIONS MAY NOT BE ENOUGH TO GIVE FULL GUIDELINES ON THEIR USES FOR REFERENCE ON CONCEPTS  
REPEATED ACROSS THE API SEE GLOSSARY OF COMMON TERMS AND API ELEMENTS  
61SKLEARNBASE BASE CLASSES AND UTILITY FUNCTIONS  
BASE CLASSES FOR ALL ESTIMATORS  
611 BASE CLASSES  
BASEBASEESTIMATOR BASE CLASS FOR ALL ESTIMATORS IN SCIKITLEARN  
BASEBICCLUSTERMIXIN MIXIN CLASS FOR ALL BICLUSTER ESTIMATORS IN SCIKITLEARN  
BASECLASSIFIERMIXIN MIXIN CLASS FOR ALL CLASSIFIERS IN SCIKITLEARN  
BASECLUSTERMIXIN MIXIN CLASS FOR ALL CLUSTER ESTIMATORS IN SCIKITLEARN  
BASEDENSITYMIXIN MIXIN CLASS FOR ALL DENSITY ESTIMATORS IN SCIKITLEARN  
BASEREGRESSORMIXIN MIXIN CLASS FOR ALL REGRESSION ESTIMATORS IN SCIKITLEARN  
BASETRANSFORMERMIXIN MIXIN CLASS FOR ALL TRANSFORMERS IN SCIKITLEARN  
SKLEARNBASE BASEESTIMATOR  
CLASSSSKLEARNBASE BASEESTIMATOR  
BASE CLASS FOR ALL ESTIMATORS IN SCIKITLEARN  
NOTES  
ALL ESTIMATORS SHOULD SPECIFY ALL THE PARAMETERS THAT CAN BE SET AT THE CLASS LEVEL IN THEIR INIT AS EXPLICIT  
KEYWORD ARGUMENTS NO ARGS ORKWARGS  
METHODS  
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR  
SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR  
1459

SCIKITLEARN USER GUIDE RELEASE 0213

INIT SELFARGS KWARGS

INITIALIZE SELF SEE HELPTYPESELF FOR ACCURATE SIGNATURE

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE

PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT

OF A NESTED OBJECT

RETURNS

SELF

EXAMPLES USING SKLEARNBASEBASEESTIMATOR

- INDUCTIVE CLUSTERING
- COLUMN TRANSFORMER WITH HETEROGENEOUS DATA SOURCES

SKLEARNBASE BICLUSTERMIXIN

CLASSSKLEARNBASE BICLUSTERMIXIN

MIXIN CLASS FOR ALL BICLUSTER ESTIMATORS IN SCIKITLEARN

ATTRIBUTES

BICLUSTERS CONVENIENT WAY TO GET ROW AND COLUMN INDICATORS TOGETHER

METHODS

GETINDICES SELF I ROW AND COLUMN INDICES OF THE I'TH BICLUSTER

GETSHAPE SELF I SHAPE OF THE I'TH BICLUSTER

GETSUBMATRIX SELF I DATA RETURNS THE SUBMATRIX CORRESPONDING TO BICLUSTER I

INIT SELFARGS KWARGS

INITIALIZE SELF SEE HELPTYPESELF FOR ACCURATE SIGNATURE

BICLUSTERS

CONVENIENT WAY TO GET ROW AND COLUMN INDICATORS TOGETHER

RETURNS THE ROWS ANDCOLUMNS MEMBERS

GETINDICES SELF I

ROW AND COLUMN INDICES OF THE I'TH BICLUSTER

1460 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

ONLY WORKS IF ROWS ANDCOLUMNS ATTRIBUTES EXIST

PARAMETERS

IINT THE INDEX OF THE CLUSTER

RETURNS

ROWIND NPARRAY DTYPENPINTP INDICES OF ROWS IN THE DATASET THAT BELONG TO THE BICLUSTER

COLIND NPARRAY DTYPENPINTP INDICES OF COLUMNS IN THE DATASET THAT BELONG TO THE BICLUSTER

GETSHAPE SELF I

SHAPE OF THE I'TH BICLUSTER

PARAMETERS

IINT THE INDEX OF THE CLUSTER

RETURNS

SHAPE INT INT NUMBER OF ROWS AND COLUMNS RESP IN THE BICLUSTER

GETSUBMATRIX SELFIDATA

RETURNS THE SUBMATRIX CORRESPONDING TO BICLUSTER I

PARAMETERS

IINT THE INDEX OF THE CLUSTER

DATA ARRAY THE DATA

RETURNS

SUBMATRIX ARRAY THE SUBMATRIX CORRESPONDING TO BICLUSTER I

NOTES

WORKS WITH SPARSE MATRICES ONLY WORKS IF ROWS ANDCOLUMNS ATTRIBUTES EXIST

SKLEARNBASE CLASSIFIERMIXIN

CLASSSSKLEARNBASE CLASSIFIERMIXIN

MIXIN CLASS FOR ALL CLASSIFIERS IN SCIKITLEARN

METHODS

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS

INIT SELFARGS KWARGS

INITIALIZE SELF SEE HELPTYPESELF FOR ACCURATE SIGNATURE

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS

IN MULTILABEL CLASSIFICATION THIS IS THE SUBSET ACCURACY WHICH IS A HARSH METRIC SINCE YOU REQUIRE FOR EACH SAMPLE THAT EACH LABEL SET BE CORRECTLY PREDICTED

61SKLEARNBASE BASE CLASSES AND UTILITY FUNCTIONS 1461

SCIKITLEARN USER GUIDE RELEASE 0213

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE LABELS FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT MEAN ACCURACY OF SELF PREDICT X WRT Y

SKLEARNBASE CLUSTER MIXIN

CLASS SKLEARNBASE CLUSTER MIXIN

MIXIN CLASS FOR ALL CLUSTER ESTIMATORS IN SCIKITLEARN

METHODS

FIT PREDICT SELF X Y PERFORMS CLUSTERING ON X AND RETURNS CLUSTER LABELS

INIT SELF ARGS KWARGS

INITIALIZE SELF SEE HELPTYPE SELF FOR ACCURATE SIGNATURE

FIT PREDICT SELF X Y NONE

PERFORMS CLUSTERING ON X AND RETURNS CLUSTER LABELS

PARAMETERS

X NDARRAY SHAPE NSAMPLES NFEATURES INPUT DATA

Y IGNORED NOT USED PRESENT FOR API CONSISTENCY BY CONVENTION

RETURNS

LABELS NDARRAY SHAPE NSAMPLES CLUSTER LABELS

SKLEARNBASE DENSITY MIXIN

CLASS SKLEARNBASE DENSITY MIXIN

MIXIN CLASS FOR ALL DENSITY ESTIMATORS IN SCIKITLEARN

METHODS

SCORE SELF X Y RETURNS THE SCORE OF THE MODEL ON THE DATA X

INIT SELF ARGS KWARGS

INITIALIZE SELF SEE HELPTYPE SELF FOR ACCURATE SIGNATURE

SCORE SELF X Y NONE

RETURNS THE SCORE OF THE MODEL ON THE DATA X

PARAMETERS

X ARRAYLIKE SHAPE NSAMPLES NFEATURES

1462 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

RETURNS

SCORE FLOAT

SKLEARNBASE REGRESSORMIXIN

CLASSSSKLEARNBASE REGRESSORMIXIN

MIXIN CLASS FOR ALL REGRESSION ESTIMATORS IN SCIKITLEARN

METHODS

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

INIT SELFARGS KWARGS

INITIALIZE SELF SEE HELPTYPESELF FOR ACCURATE SIGNATURE

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{U}{V}$  WHERE U IS THE RESIDUAL SUM OF SQUARES  $Y_{TRUE} - Y_{PRED}$   $2SUM$  AND V IS THE TOTAL SUM OF SQUARES  $Y_{TRUE} - Y_{TRUEMEAN}$   $2SUM$  THE BEST POSSIBLE SCORE IS 10 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 00

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES FOR SOME ESTIMATORS THIS MAY BE A PRECOMPUTED KERNEL MATRIX INSTEAD SHAPE NSAMPLES NSAMPLESFITTED WHERE NSAMPLESFITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE VALUES FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT R2 OF SELF PREDICTX WRT Y

NOTES

THE R2 SCORE USED WHEN CALLING SCORE ON A REGRESSOR WILL USE MULTIOUTPUTUNIFORMAVERAGE FROM VERSION 023 TO KEEP CONSISTENT WITH METRICSR2SCORE THIS WILL INFLUENCE THE SCORE METHOD OF ALL THE MULTIOUTPUT REGRESSORS EXCEPT FOR MULTIOUTPUTMULTIOUTPUTREGRESSOR TO SPECIFY THE DEFAULT VALUE MANUALLY AND AVOID THE WARNING PLEASE EITHER CALL METRICSR2SCORE DIRECTLY OR MAKE A CUSTOM SCORER WITH METRICSMAKESCORER THE BUILTIN SCORER R2 USES MULTIOUTPUTUNIFORMAVERAGE

SKLEARNBASE TRANSFORMERMIXIN

CLASSSSKLEARNBASE TRANSFORMERMIXIN

MIXIN CLASS FOR ALL TRANSFORMERS IN SCIKITLEARN

61SKLEARNBASE BASE CLASSES AND UTILITY FUNCTIONS 1463

SCIKITLEARN USER GUIDE RELEASE 0213

METHODS

FITTRANSFORM SELF X Y FIT TO DATA THEN TRANSFORM IT

INIT SELFARGS KWARGS

INITIALIZE SELF SEE HELPTYPESELF FOR ACCURATE SIGNATURE

FITTRANSFORM SELFXYNONE FITPARAMS

FIT TO DATA THEN TRANSFORM IT

FITS TRANSFORMER TO X AND Y WITH OPTIONAL PARAMETERS FITPARAMS AND RETURNS A TRANSFORMED VERSION OF X

PARAMETERS

XNUMPY ARRAY OF SHAPE NSAMPLES NFEATURES TRAINING SET

YNUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES

RETURNS

XNEW NUMPY ARRAY OF SHAPE NSAMPLES NFEATURESNEW TRANSFORMED ARRAY

EXAMPLES USING SKLEARNBASETRANSFORMERMIXIN

- COLUMN TRANSFORMER WITH HETEROGENEOUS DATA SOURCES

612 FUNCTIONS

BASECLONE ESTIMATOR SAFE CONSTRUCTS A NEW ESTIMATOR WITH THE SAME PARAMETERS

BASEISCLASSIFIER ESTIMATOR RETURNS TRUE IF THE GIVEN ESTIMATOR IS PROBABLY A CLASSIFIER

BASEISREGRESSOR ESTIMATOR RETURNS TRUE IF THE GIVEN ESTIMATOR IS PROBABLY A REGRESSOR

CONFIGCONTEXT NEWCONFIG CONTEXT MANAGER FOR GLOBAL SCIKITLEARN CONFIGURATION

GETCONFIG RETRIEVE CURRENT VALUES FOR CONFIGURATION SET BY

SETCONFIG

SETCONFIG ASSUMEFINITE WORKINGMEMORY SET GLOBAL SCIKITLEARN CONFIGURATION

SHOWVERSIONS PRINT USEFUL DEBUGGING INFORMATION

SKLEARNBASE CLONE

SKLEARNBASE CLONEESTIMATOR SAFETRUE

CONSTRUCTS A NEW ESTIMATOR WITH THE SAME PARAMETERS

CLONE DOES A DEEP COPY OF THE MODEL IN AN ESTIMATOR WITHOUT ACTUALLY COPYING ATTACHED DATA IT YIELDS A NEW ESTIMATOR WITH THE SAME PARAMETERS THAT HAS NOT BEEN FIT ON ANY DATA

PARAMETERS

ESTIMATOR ESTIMATOR OBJECT OR LIST TUPLE OR SET OF OBJECTS THE ESTIMATOR OR GROUP OF ESTIMATORS TO BE CLONED

SAFE BOOLEAN OPTIONAL IF SAFE IS FALSE CLONE WILL FALL BACK TO A DEEP COPY ON OBJECTS THAT ARE NOT ESTIMATORS

1464 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

SKLEARNBASE ISCLASSIFIER

SKLEARNBASE ISCLASSIFIER ESTIMATOR

RETURNS TRUE IF THE GIVEN ESTIMATOR IS PROBABLY A CLASSIFIER

PARAMETERS

ESTIMATOR OBJECT ESTIMATOR OBJECT TO TEST

RETURNS

OUT BOOL TRUE IF ESTIMATOR IS A CLASSIFIER AND FALSE OTHERWISE

SKLEARNBASE ISREGRESSOR

SKLEARNBASE ISREGRESSOR ESTIMATOR

RETURNS TRUE IF THE GIVEN ESTIMATOR IS PROBABLY A REGRESSOR

PARAMETERS

ESTIMATOR OBJECT ESTIMATOR OBJECT TO TEST

RETURNS

OUT BOOL TRUE IF ESTIMATOR IS A REGRESSOR AND FALSE OTHERWISE

SKLEARN CONFIGCONTEXT

SKLEARN CONFIGCONTEXT NEWCONFIG

CONTEXT MANAGER FOR GLOBAL SCIKITLEARN CONFIGURATION

PARAMETERS

ASSUMEFINITE BOOL OPTIONAL IF TRUE VALIDATION FOR FINITENESS WILL BE SKIPPED SAVING TIME BUT LEADING TO POTENTIAL CRASHES IF FALSE VALIDATION FOR FINITENESS WILL BE PERFORMED AVOIDING ERROR GLOBAL DEFAULT FALSE

WORKINGMEMORY INT OPTIONAL IF SET SCIKITLEARN WILL ATTEMPT TO LIMIT THE SIZE OF TEMPORARY ARRAYS TO THIS NUMBER OF MIB PER JOB WHEN PARALLELISED OFTEN SAVING BOTH COMPUTATION TIME AND MEMORY ON EXPENSIVE OPERATIONS THAT CAN BE PERFORMED IN CHUNKS GLOBAL DEFAULT 1024

SEE ALSO

SETCONFIG SET GLOBAL SCIKITLEARN CONFIGURATION

GETCONFIG RETRIEVE CURRENT VALUES OF THE GLOBAL CONFIGURATION

NOTES

ALL SETTINGS NOT JUST THOSE PRESENTLY MODIFIED WILL BE RETURNED TO THEIR PREVIOUS VALUES WHEN THE CONTEXT MANAGER IS EXITED THIS IS NOT THREADSAFE

EXAMPLES

61SKLEARNBASE BASE CLASSES AND UTILITY FUNCTIONS 1465

SCIKITLEARN USER GUIDE RELEASE 0213

```
import sklearn
from sklearn.utils.validation import assert_all_finite
with sklearn.config_context(assume_finite=True):
    assert_all_finite(float_nan)
with sklearn.config_context(assume_finite=True):
    with sklearn.config_context(assume_finite=False):
        assert_all_finite(float_nan)
```

Traceback most recent call last:

ValueError: Input contains NaN

sklearn.get\_config

sklearn.get\_config

Retrieve current values for configuration set by set\_config

Returns

config\_dict: Keys are parameter names that can be passed to set\_config

See also

ConfigContext: Context manager for global scikitlearn configuration

set\_config: Set global scikitlearn configuration

sklearn.set\_config

sklearn.set\_config(assume\_finite=None, working\_memory=None, print\_changed\_only=None)

Set global scikitlearn configuration

New in version 0.19

Parameters

assume\_finite: bool, optional. If True validation for finiteness will be skipped saving time but leading to potential crashes. If False validation for finiteness will be performed avoiding error. Global default: False

New in version 0.19

working\_memory: int, optional. If set scikitlearn will attempt to limit the size of temporary arrays to this number of MiB per job when parallelised often saving both computation time and memory on expensive operations that can be performed in chunks. Global default: 1024

New in version 0.20

print\_changed\_only: bool, optional. If True only the parameters that were set to nondefault values will be printed when printing an estimator. For example print\_svc while True will only print 'svc' while the default behaviour would be to print 'svcc10'.

Cachesize: 200 ' ' with all the nonchanged parameters

New in version 0.21

See also

ConfigContext: Context manager for global scikitlearn configuration

1466 Chapter 6: API Reference



SCIKITLEARN USER GUIDE RELEASE 0213  
GETCONFIG RETRIEVE CURRENT VALUES OF THE GLOBAL CONFIGURATION  
EXAMPLES USING SKLEARNSETCONFIG  
•COMPACT ESTIMATOR REPRESENTATIONS  
SKLEARN SHOWVERSIONS  
SKLEARN SHOWVERSIONS  
PRINT USEFUL DEBUGGING INFORMATION  
62SKLEARNCALIBRATION PROBABILITY CALIBRATION  
CALIBRATION OF PREDICTED PROBABILITIES  
USER GUIDE SEE THE PROBABILITY CALIBRATION SECTION FOR FURTHER DETAILS  
CALIBRATIONCALIBRATEDCLASSIFIERCV PROBABILITY CALIBRATION WITH ISOTONIC REGRESSION OR SIGMOID  
621SKLEARNCALIBRATION CALIBRATEDCLASSIFIERCV  
CLASSSSKLEARNCALIBRATION CALIBRATEDCLASSIFIERCV BASEESTIMATORNONE  
METHOD'SIGMOID' CV'WARN'  
PROBABILITY CALIBRATION WITH ISOTONIC REGRESSION OR SIGMOID  
SEE GLOSSARY ENTRY FOR CROSSVALIDATION ESTIMATOR  
WITH THIS CLASS THE BASEESTIMATOR IS FIT ON THE TRAIN SET OF THE CROSSVALIDATION GENERATOR AND THE TEST SET IS USED  
FOR CALIBRATION THE PROBABILITIES FOR EACH OF THE FOLDS ARE THEN AVERAGED FOR PREDICTION IN CASE THAT CV"PREFIT"  
IS PASSED TO INIT IT IS ASSUMED THAT BASEESTIMATOR HAS BEEN FITTED ALREADY AND ALL DATA IS USED FOR CALIBRATION  
NOTE THAT DATA FOR FITTING THE CLASSIFIER AND FOR CALIBRATING IT MUST BE DISJOINT  
READ MORE IN THE USER GUIDE  
PARAMETERS  
BASEESTIMATOR INSTANCE BASEESTIMATOR THE CLASSIFIER WHOSE OUTPUT DECISION FUNCTION NEEDS  
TO BE CALIBRATED TO OFFER MORE ACCURATE PREDICTPROBA OUTPUTS IF CVPREFIT THE CLASSIFIER MUST  
HAVE BEEN FIT ALREADY ON DATA  
METHOD 'SIGMOID' OR 'ISOTONIC' THE METHOD TO USE FOR CALIBRATION CAN BE 'SIGMOID' WHICH  
CORRESPONDS TO PLATT'S METHOD OR 'ISOTONIC' WHICH IS A NONPARAMETRIC APPROACH IT IS NOT  
ADVISED TO USE ISOTONIC CALIBRATION WITH TOO FEW CALIBRATION SAMPLES 1000 SINCE IT  
TENDS TO OVERFIT USE SIGMOIDS PLATT'S CALIBRATION IN THIS CASE  
CVINTEGER CROSSVALIDATION GENERATOR ITERABLE OR "PREFIT" OPTIONAL DETERMINES THE CROSS  
VALIDATION SPLITTING STRATEGY POSSIBLE INPUTS FOR CV ARE  
• NONE TO USE THE DEFAULT 3FOLD CROSSVALIDATION  
• INTEGER TO SPECIFY THE NUMBER OF FOLDS  
•CV SPLITTER  
• AN ITERABLE YIELDING TRAIN TEST SPLITS AS ARRAYS OF INDICES  
62SKLEARNCALIBRATION PROBABILITY CALIBRATION 1467

SCIKITLEARN USER GUIDE RELEASE 0213  
FOR INTEGERNONE INPUTS IF YIS BINARY OR MULTICLASS SKLEARNMODELSELECTION  
STRATIFIEDKFOLD IS USED IF YIS NEITHER BINARY NOR MULTICLASS SKLEARN  
MODELSELECTIONKFOLD IS USED  
REFER USER GUIDE FOR THE VARIOUS CROSSVALIDATION STRATEGIES THAT CAN BE USED HERE  
IF “PREFIT” IS PASSED IT IS ASSUMED THAT BASEESTIMATOR HAS BEEN FITTED ALREADY AND ALL DATA IS  
USED FOR CALIBRATION  
CHANGED IN VERSION 020 CVDEFAULT VALUE IF NONE WILL CHANGE FROM 3FOLD TO 5FOLD IN  
V022  
ATTRIBUTES  
CLASSES ARRAY SHAPE NCLASSES THE CLASS LABELS  
CALIBRATEDCLASSIFIERS LIST LEN EQUAL TO CV OR 1 IF CV “PREFIT” THE LIST OF CALIBRATED  
CLASSIFIERS ONE FOR EACH CROSSVALIDATION FOLD WHICH HAS BEEN FITTED ON ALL BUT THE VALIDATION  
FOLD AND CALIBRATED ON THE VALIDATION FOLD  
REFERENCES  
R57CF438D70601 R57CF438D70602 R57CF438D70603 R57CF438D70604  
METHODS  
FITSELF X Y SAMPLEWEIGHT FIT THE CALIBRATED MODEL  
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR  
PREDICT SELF X PREDICT THE TARGET OF NEW SAMPLES  
PREDICTPROBA SELF X POSTERIOR PROBABILITIES OF CLASSIFICATION  
SCORE SELF X Y SAMPLEWEIGHT RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND  
LABELS  
SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR  
INIT SELFBASEESTIMATORNONE METHOD‘SIGMOID’ CV‘WARN’  
FITSELFXYSAMPLEWEIGHTNONE  
FIT THE CALIBRATED MODEL  
PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA  
YARRAYLIKE SHAPE NSAMPLES TARGET VALUES  
SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OR NONE SAMPLE WEIGHTS IF NONE THEN  
SAMPLES ARE EQUALLY WEIGHTED  
RETURNS  
SELF OBJECT RETURNS AN INSTANCE OF SELF  
GETPARAMS SELFDEEPTURE  
GET PARAMETERS FOR THIS ESTIMATOR  
PARAMETERS  
DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBOBJECTS THAT ARE ESTIMATORS  
1468 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELF

PREDICT THE TARGET OF NEW SAMPLES CAN BE DIFFERENT FROM THE PREDICTION OF THE UNCALIBRATED CLASSIFIER

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES THE SAMPLES

RETURNS

CARRAY SHAPE NSAMPLES THE PREDICTED CLASS

PREDICTPROBA SELF

POSTERIOR PROBABILITIES OF CLASSIFICATION

THIS FUNCTION RETURNS POSTERIOR PROBABILITIES OF CLASSIFICATION ACCORDING TO EACH CLASS ON AN ARRAY OF TEST VECTORS X

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES THE SAMPLES

RETURNS

CARRAY SHAPE NSAMPLES NCLASSES THE PREDICTED PROBAS

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS

IN MULTILABEL CLASSIFICATION THIS IS THE SUBSET ACCURACY WHICH IS A HARSH METRIC SINCE YOU REQUIRE FOR EACH SAMPLE THAT EACH LABEL SET BE CORRECTLY PREDICTED

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE LABELS FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT MEAN ACCURACY OF SELF PREDICTX WRT Y

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

EXAMPLES USING SKLEARN CALIBRATION CALIBRATED CLASSIFIER CV

- PROBABILITY CALIBRATION CURVES
- PROBABILITY CALIBRATION OF CLASSIFIERS
- PROBABILITY CALIBRATION FOR 3 CLASS CLASSIFICATION

62 SKLEARN CALIBRATION PROBABILITY CALIBRATION 1469

SCIKITLEARN USER GUIDE RELEASE 0213  
CALIBRATIONCALIBRATIONCURVE YTRUE  
YPROB  
COMPUTE TRUE AND PREDICTED PROBABILITIES FOR A CALIBRATION  
CURVE  
622SKLEARNCALIBRATION CALIBRATIONCURVE  
SKLEARNCALIBRATION CALIBRATIONCURVE YTRUE YPROB NORMALIZEFALSE NBINS5 STRAT  
EGY'UNIFORM'  
COMPUTE TRUE AND PREDICTED PROBABILITIES FOR A CALIBRATION CURVE  
THE METHOD ASSUMES THE INPUTS COME FROM A BINARY CLASSIFIER  
CALIBRATION CURVES MAY ALSO BE REFERRED TO AS RELIABILITY DIAGRAMS  
READ MORE IN THE USER GUIDE  
PARAMETERS  
YTRUE ARRAY SHAPE NSAMPLES TRUE TARGETS  
YPROB ARRAY SHAPE NSAMPLES PROBABILITIES OF THE POSITIVE CLASS  
NORMALIZE BOOL OPTIONAL DEFAULTFALSE WHETHER YPROB NEEDS TO BE NORMALIZED INTO THE BIN  
0 1 IE IS NOT A PROPER PROBABILITY IF TRUE THE SMALLEST VALUE IN YPROB IS MAPPED ONTO  
0 AND THE LARGEST ONE ONTO 1  
NBINS INT NUMBER OF BINS A BIGGER NUMBER REQUIRES MORE DATA BINS WITH NO DATA POINTS  
IE WITHOUT CORRESPONDING VALUES IN YPROB WILL NOT BE RETURNED THUS THERE MAY BE FEWER  
THAN NBINS IN THE RETURN VALUE  
STRATEGY 'UNIFORM' 'QUANTILE' DEFAULT'UNIFORM' STRATEGY USED TO DEFINE THE WIDTHS OF THE  
BINS  
UNIFORM ALL BINS HAVE IDENTICAL WIDTHS  
QUANTILE ALL BINS HAVE THE SAME NUMBER OF POINTS  
RETURNS  
PROBTRUE ARRAY SHAPE NBINS OR SMALLER THE TRUE PROBABILITY IN EACH BIN FRACTION OF POSI  
TIVES  
PROBPRED ARRAY SHAPE NBINS OR SMALLER THE MEAN PREDICTED PROBABILITY IN EACH BIN  
REFERENCES  
ALEXANDRU NICULESCUMIZIL AND RICH CARUANA 2005 PREDICTING GOOD PROBABILITIES WITH SUPERVISED LEARNING  
IN PROCEEDINGS OF THE 22ND INTERNATIONAL CONFERENCE ON MACHINE LEARNING ICML SEE SECTION 4 QUALITATIVE  
ANALYSIS OF PREDICTIONS  
EXAMPLES USING SKLEARNCALIBRATIONCALIBRATIONCURVE  
•COMPARISON OF CALIBRATION OF CLASSIFIERS  
•PROBABILITY CALIBRATION CURVES  
1470 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213  
63SKLEARNCLUSTER CLUSTERING  
THESKLEARNCLUSTER MODULE GATHERS POPULAR UNSUPERVISED CLUSTERING ALGORITHMS  
USER GUIDE SEE THE CLUSTERING SECTION FOR FURTHER DETAILS  
631 CLASSES  
CLUSTERAFFINITYPROPAGATION DAMPING PERFORM AFFINITY PROPAGATION CLUSTERING OF DATA  
CLUSTERAGGLOMERATIVECLUSTERING AGGLOMERATIVE CLUSTERING  
CLUSTERBIRCH THRESHOLD BRANCHINGFACTOR IMPLEMENTS THE BIRCH CLUSTERING ALGORITHM  
CLUSTERDBSCAN EPS MINSAMPLES METRIC PERFORM DBSCAN CLUSTERING FROM VECTOR ARRAY OR DISTANCE  
MATRIX  
CLUSTEROPTICS MINSAMPLES MAXEPS ESTIMATE CLUSTERING STRUCTURE FROM VECTOR ARRAY  
CLUSTERFEATUREAGGLOMERATION NCLUSTERS  
AGGLOMERATE FEATURES  
CLUSTERKMEANS NCLUSTERS INIT NINIT KMEANS CLUSTERING  
CLUSTERMINIBATCHKMEANS NCLUSTERS INIT MINIBATCH KMEANS CLUSTERING  
CLUSTERMEANSHIFT BANDWIDTH SEEDS MEAN SHIFT CLUSTERING USING A FLAT KERNEL  
CLUSTERSPECTRALCLUSTERING NCLUSTERS APPLY CLUSTERING TO A PROJECTION OF THE NORMALIZED LAPLA  
CIAN  
SKLEARNCLUSTER AFFINITYPROPAGATION  
CLASSSSKLEARNCLUSTER AFFINITYPROPAGATION DAMPING05 MAXITER200 CONVER  
GENCEITER15 COPYTRUE PREFERENCECONE  
AFFINITY'EUCLIDEAN' VERBOSEFALSE  
PERFORM AFFINITY PROPAGATION CLUSTERING OF DATA  
READ MORE IN THE USER GUIDE  
PARAMETERS  
DAMPING FLOAT OPTIONAL DEFAULT 05 DAMPING FACTOR BETWEEN 05 AND 1 IS THE EXTENT TO  
WHICH THE CURRENT VALUE IS MAINTAINED RELATIVE TO INCOMING VALUES WEIGHTED 1 DAMPING  
THIS IN ORDER TO AVOID NUMERICAL OSCILLATIONS WHEN UPDATING THESE VALUES MESSAGES  
MAXITER INT OPTIONAL DEFAULT 200 MAXIMUM NUMBER OF ITERATIONS  
CONVERGENCEITER INT OPTIONAL DEFAULT 15 NUMBER OF ITERATIONS WITH NO CHANGE IN THE NUMBER  
OF ESTIMATED CLUSTERS THAT STOPS THE CONVERGENCE  
COPY BOOLEAN OPTIONAL DEFAULT TRUE MAKE A COPY OF INPUT DATA  
PREFERENCE ARRAYLIKE SHAPE NSAMPLES OR FLOAT OPTIONAL PREFERENCES FOR EACH POINT POINTS  
WITH LARGER VALUES OF PREFERENCES ARE MORE LIKELY TO BE CHOSEN AS EXEMPLARS THE NUMBER OF  
EXEMPLARS IE OF CLUSTERS IS INFLUENCED BY THE INPUT PREFERENCES VALUE IF THE PREFERENCES ARE  
NOT PASSED AS ARGUMENTS THEY WILL BE SET TO THE MEDIAN OF THE INPUT SIMILARITIES  
AFFINITY STRING OPTIONAL DEFAULT"EUCLIDEAN" WHICH AFFINITY TO USE AT THE MOMENT  
PRECOMPUTED ANDEUCLIDEAN ARE SUPPORTED EUCLIDEAN USES THE NEGATIVE SQUARED  
EUCLIDEAN DISTANCE BETWEEN POINTS  
VERBOSE BOOLEAN OPTIONAL DEFAULT FALSE WHETHER TO BE VERBOSE  
ATTRIBUTES  
63SKLEARNCLUSTER CLUSTERING 1471

SCIKITLEARN USER GUIDE RELEASE 0213

CLUSTERCENTERSINDICES ARRAY SHAPE NCLUSTERS INDICES OF CLUSTER CENTERS  
CLUSTERCENTERS ARRAY SHAPE NCLUSTERS NFEATURES CLUSTER CENTERS IF AFFINITY  
PRECOMPUTED

LABELS ARRAY SHAPE NSAMPLES LABELS OF EACH POINT  
AFFINITYMATRIX ARRAY SHAPE NSAMPLES NSAMPLES STORES THE AFFINITY MATRIX USED IN FIT  
NITER INT NUMBER OF ITERATIONS TAKEN TO CONVERGE  
NOTES

FOR AN EXAMPLE SEE EXAMPLESCUSTERPLOTAFFINITYPROPAGATIONPY  
THE ALGORITHMIC COMPLEXITY OF AFFINITY PROPAGATION IS QUADRATIC IN THE NUMBER OF POINTS  
WHENFIT DOES NOT CONVERGE CLUSTERCENTERS BECOMES AN EMPTY ARRAY AND ALL TRAINING SAMPLES WILL BE  
LABELLED AS1 IN ADDITION PREDICT WILL THEN LABEL EVERY SAMPLE AS 1  
WHEN ALL TRAINING SAMPLES HAVE EQUAL SIMILARITIES AND EQUAL PREFERENCES THE ASSIGNMENT OF CLUSTER CENTERS AND  
LABELS DEPENDS ON THE PREFERENCE IF THE PREFERENCE IS SMALLER THAN THE SIMILARITIES FIT WILL RESULT IN A SINGLE  
CLUSTER CENTER AND LABEL 0FOR EVERY SAMPLE OTHERWISE EVERY TRAINING SAMPLE BECOMES ITS OWN CLUSTER CENTER  
AND IS ASSIGNED A UNIQUE LABEL

REFERENCES  
BRENDAN J FREY AND DELBERT DUECK “CLUSTERING BY PASSING MESSAGES BETWEEN DATA POINTS” SCIENCE FEB 2007  
EXAMPLES

```
FROM SKLEARNCLUSTER IMPORT AFFINITYPROPAGATION
IMPORT NUMPY AS NP
X NPARRAY1 2 1 4 1 0
  4 2 4 4 4 0
CLUSTERING AFFINITYPROPAGATIONFITX
CLUSTERING
AFFINITYPROPAGATIONAFFINITYEUCLIDEAN CONVERGENCEITER15 COPYTRUE
DAMPING05 MAXITER200 PREFERENCENONE VERBOSEFALSE
CLUSTERINGLABELS
ARRAY0 0 0 1 1 1
CLUSTERINGPREDICT0 0 4 4
ARRAY0 1
CLUSTERINGCLUSTERCENTERS
ARRAY1 2
  4 2
```

METHODS  
FITSELF X Y CREATE AFFINITY MATRIX FROM NEGATIVE EUCLIDEAN DIS  
TANCES THEN APPLY AFFINITY PROPAGATION CLUSTERING  
FITPREDICT SELF X Y PERFORMS CLUSTERING ON X AND RETURNS CLUSTER LABELS  
CONTINUED ON NEXT PAGE  
1472 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

TABLE 614 – CONTINUED FROM PREVIOUS PAGE

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PREDICT SELF X PREDICT THE CLOSEST CLUSTER EACH SAMPLE IN X BELONGS TO

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELF DAMPING 0.5 MAX ITER 200 CONVERGENCE ITER 15 COPY TRUE PREFERENCE NONE

AFFINITY ‘EUCLIDEAN’ VERBOSE FALSE

FIT SELF FX NONE

CREATE AFFINITY MATRIX FROM NEGATIVE EUCLIDEAN DISTANCES THEN APPLY AFFINITY PROPAGATION CLUSTERING

PARAMETERS

X ARRAY LIKE SHAPE NSAMPLES NFEATURES OR NSAMPLES NSAMPLES DATA MATRIX OR IF

AFFINITY IS PRECOMPUTED MATRIX OF SIMILARITIES AFFINITIES

Y IGNORED

FIT PREDICT SELF FX NONE

PERFORMS CLUSTERING ON X AND RETURNS CLUSTER LABELS

PARAMETERS

X NDARRAY SHAPE NSAMPLES NFEATURES INPUT DATA

Y IGNORED NOT USED PRESENT FOR API CONSISTENCY BY CONVENTION

RETURNS

LABELS NDARRAY SHAPE NSAMPLES CLUSTER LABELS

GETPARAMS SELF DEEP TRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICT SELF X

PREDICT THE CLOSEST CLUSTER EACH SAMPLE IN X BELONGS TO

PARAMETERS

X ARRAY LIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES NEW DATA TO PREDICT

RETURNS

LABELS ARRAY SHAPE NSAMPLES INDEX OF THE CLUSTER EACH SAMPLE BELONGS TO

SETPARAMS SELF PARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE

PARAMETERS OF THE FORM COMPONENT PARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT

OF A NESTED OBJECT

RETURNS

SELF

63 SKLEARN CLUSTER CLUSTERING 1473

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES USING SKLEARNCLUSTERAFFINITYPROPAGATION

- DEMO OF AFFINITY PROPAGATION CLUSTERING ALGORITHM
- COMPARING DIFFERENT CLUSTERING ALGORITHMS ON TOY DATASETS

SKLEARNCLUSTER AGGLOMERATIVECLUSTERING

CLASSSSKLEARNCLUSTER AGGLOMERATIVECLUSTERING NCLUSTERS2 AFFINITY‘EUCLIDEAN’ MEM

ORYNONE CONNECTIVITYNONE COM

PUTEFULLTREE‘AUTO’ LINKAGE‘WARD’

POOLINGFUNC‘DEPRECATED’ DIS

TANCETHRESHOLDNONE

AGGLOMERATIVE CLUSTERING

RECURSIVELY MERGES THE PAIR OF CLUSTERS THAT MINIMALLY INCREASES A GIVEN LINKAGE DISTANCE

READ MORE IN THE USER GUIDE

PARAMETERS

NCLUSTERS INT OR NONE OPTIONAL DEFAULT2 THE NUMBER OF CLUSTERS TO FIND IT MUST BE NONE

IFDISTANCETHRESHOLD IS NOTNONE

AFFINITY STRING OR CALLABLE DEFAULT “EUCLIDEAN” METRIC USED TO COMPUTE THE LINKAGE CAN BE “EUCLIDEAN” “L1” “L2” “MANHATTAN” “COSINE” OR “PRECOMPUTED” IF LINKAGE IS “WARD” ONLY “EUCLIDEAN” IS ACCEPTED IF “PRECOMPUTED” A DISTANCE MATRIX INSTEAD OF A SIMILARITY MATRIX IS NEEDED AS INPUT FOR THE FIT METHOD

MEMORY NONE STR OR OBJECT WITH THE JOBLIBMEMORY INTERFACE OPTIONAL USED TO CACHE THE OUTPUT OF THE COMPUTATION OF THE TREE BY DEFAULT NO CACHING IS DONE IF A STRING IS GIVEN IT IS THE PATH TO THE CACHING DIRECTORY

CONNECTIVITY ARRAYLIKE OR CALLABLE OPTIONAL CONNECTIVITY MATRIX DEFINES FOR EACH SAMPLE THE NEIGHBORING SAMPLES FOLLOWING A GIVEN STRUCTURE OF THE DATA THIS CAN BE A CONNECTIVITY MATRIX ITSELF OR A CALLABLE THAT TRANSFORMS THE DATA INTO A CONNECTIVITY MATRIX SUCH AS DERIVED FROM KNEIGHBORSGRAPH DEFAULT IS NONE IE THE HIERARCHICAL CLUSTERING ALGORITHM IS UNSTRUCTURED

COMPUTEFULLTREE BOOL OR ‘AUTO’ OPTIONAL STOP EARLY THE CONSTRUCTION OF THE TREE AT NCLUSTERS THIS IS USEFUL TO DECREASE COMPUTATION TIME IF THE NUMBER OF CLUSTERS IS NOT SMALL COMPARED TO THE NUMBER OF SAMPLES THIS OPTION IS USEFUL ONLY WHEN SPECIFYING A CONNECTIVITY MATRIX NOTE ALSO THAT WHEN VARYING THE NUMBER OF CLUSTERS AND USING CACHING IT MAY BE ADVANTAGEOUS TO COMPUTE THE FULL TREE IT MUST BE TRUE IFDISTANCETHRESHOLD IS NOTNONE

LINKAGE “WARD” “COMPLETE” “AVERAGE” “SINGLE” OPTIONAL DEFAULT“WARD” WHICH LINKAGE CRITERION TO USE THE LINKAGE CRITERION DETERMINES WHICH DISTANCE TO USE BETWEEN SETS OF OBSERVATION THE ALGORITHM WILL MERGE THE PAIRS OF CLUSTER THAT MINIMIZE THIS CRITERION

- WARD MINIMIZES THE VARIANCE OF THE CLUSTERS BEING MERGED
- AVERAGE USES THE AVERAGE OF THE DISTANCES OF EACH OBSERVATION OF THE TWO SETS
- COMPLETE OR MAXIMUM LINKAGE USES THE MAXIMUM DISTANCES BETWEEN ALL OBSERVATIONS OF THE TWO SETS
- SINGLE USES THE MINIMUM OF THE DISTANCES BETWEEN ALL OBSERVATIONS OF THE TWO SETS

1474 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213  
POOLINGFUNC CALLABLE DEFAULT'DEPRECATED' IGNORED  
DEPRECATED SINCE VERSION 020 POOLINGFUNC HAS BEEN DEPRECATED IN 020 AND WILL BE  
REMOVED IN 022  
DISTANCETHRESHOLD FLOAT OPTIONAL DEFAULTNONE THE LINKAGE DISTANCE THRESHOLD ABOVE  
WHICH CLUSTERS WILL NOT BE MERGED IF NOT NONE NCLUSTERS MUST BENONE AND  
COMPUTEFULLTREE MUST BETRUE  
NEW IN VERSION 021  
ATTRIBUTES  
NCLUSTERS INT THE NUMBER OF CLUSTERS FOUND BY THE ALGORITHM IF  
DISTANCETHRESHOLDNONE IT WILL BE EQUAL TO THE GIVEN NCLUSTERS  
LABELS ARRAY NSAMPLES CLUSTER LABELS FOR EACH POINT  
NLEAVES INT NUMBER OF LEAVES IN THE HIERARCHICAL TREE  
NCONNECTEDCOMPONENTS INT THE ESTIMATED NUMBER OF CONNECTED COMPONENTS IN THE GRAPH  
CHILDREN ARRAYLIKE SHAPE NSAMPLES1 2 THE CHILDREN OF EACH NONLEAF NODE VALUES LESS  
THAN NSAMPLES CORRESPOND TO LEAVES OF THE TREE WHICH ARE THE ORIGINAL SAMPLES A NODE  
IGREATER THAN OR EQUAL TO NSAMPLES IS A NONLEAF NODE AND HAS CHILDREN CHILDRENI  
NSAMPLES ALTERNATIVELY AT THE ITH ITERATION CHILDRENI0 AND CHILDRENI1 ARE  
MERGED TO FORM NODE NSAMPLES I  
EXAMPLES  
FROM SKLEARNCLUSTER IMPORT AGGLOMERATIVECLUSTERING  
IMPORT NUMPY AS NP  
X NPARRAY1 2 1 4 1 0  
4 2 4 4 4 0  
CLUSTERING AGGLOMERATIVECLUSTERINGFITX  
CLUSTERING  
AGGLOMERATIVECLUSTERINGAFFINITYEUCLIDEAN COMPUTEFULLTREEAUTO  
CONNECTIVITYNONE DISTANCETHRESHOLDNONE  
LINKAGEWARD MEMORYNONE NCLUSTERS2  
POOLINGFUNCDEPRECATED  
CLUSTERINGLABELS  
ARRAY1 1 1 0 0 0  
METHODS  
FITSELF X Y FIT THE HIERARCHICAL CLUSTERING ON THE DATA  
FITPREDICT SELF X Y PERFORMS CLUSTERING ON X AND RETURNS CLUSTER LABELS  
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR  
SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR  
INIT SELF NCLUSTERS2 AFFINITY'EUCLIDEAN' MEMORYNONE CONNECTIVITYNONE  
COMPUTEFULLTREE'AUTO' LINKAGE'WARD' POOLINGFUNC'DEPRECATED' DIS  
TANCETHRESHOLDNONE  
FITSELFXYNONE  
FIT THE HIERARCHICAL CLUSTERING ON THE DATA  
63SKLEARNCLUSTER CLUSTERING 1475

SCIKITLEARN USER GUIDE RELEASE 0213

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA SHAPE NSAMPLES  
NFEATURES OR NSAMPLES NSAMPLES IF AFFINITY'PRECOMPUTED'

YIGNORED

RETURNS

SELF

FITPREDICT SELFXYNONE

PERFORMS CLUSTERING ON X AND RETURNS CLUSTER LABELS

PARAMETERS

XNDARRAY SHAPE NSAMPLES NFEATURES INPUT DATA

YIGNORED NOT USED PRESENT FOR API CONSISTENCY BY CONVENTION

RETURNS

LABELS NDARRAY SHAPE NSAMPLES CLUSTER LABELS

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE  
PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT

OF A NESTED OBJECT

RETURNS

SELF

EXAMPLES USING SKLEARNCLUSTERAGGLOMERATIVECLUSTERING

- AGGLOMERATIVE CLUSTERING WITH AND WITHOUT STRUCTURE
  - VARIOUS AGGLOMERATIVE CLUSTERING ON A 2D EMBEDDING OF DIGITS
  - A DEMO OF STRUCTURED WARD HIERARCHICAL CLUSTERING ON AN IMAGE OF COINS
  - HIERARCHICAL CLUSTERING STRUCTURED VS UNSTRUCTURED WARD
  - AGGLOMERATIVE CLUSTERING WITH DIFFERENT METRICS
  - INDUCTIVE CLUSTERING
  - COMPARING DIFFERENT HIERARCHICAL LINKAGE METHODS ON TOY DATASETS
  - COMPARING DIFFERENT CLUSTERING ALGORITHMS ON TOY DATASETS
- 1476 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

SKLEARNCLUSTER BIRCH

CLASSSSKLEARNCLUSTER BIRCHTHRESHOLD05 BRANCHINGFACTOR50 NCLUSTERS3 COM

PUTELABELSTRUE COPYTRUE

IMPLEMENTS THE BIRCH CLUSTERING ALGORITHM

IT IS A MEMORYEFFICIENT ONLINELEARNING ALGORITHM PROVIDED AS AN ALTERNATIVE TO MINIBATCHKMEANS IT CON

STRUCTS A TREE DATA STRUCTURE WITH THE CLUSTER CENTROIDS BEING READ OFF THE LEAF THESE CAN BE EITHER THE FINAL CLUSTER

CENTROIDS OR CAN BE PROVIDED AS INPUT TO ANOTHER CLUSTERING ALGORITHM SUCH AS AGGLOMERATIVECLUSTERING

READ MORE IN THE USER GUIDE

PARAMETERS

THRESHOLD FLOAT DEFAULT 05 THE RADIUS OF THE SUBCLUSTER OBTAINED BY MERGING A NEW SAMPLE

AND THE CLOSEST SUBCLUSTER SHOULD BE LESSER THAN THE THRESHOLD OTHERWISE A NEW SUBCLUSTER IS

STARTED SETTING THIS VALUE TO BE VERY LOW PROMOTES SPLITTING AND VICEVERSA

BRANCHINGFACTOR INT DEFAULT 50 MAXIMUM NUMBER OF CF SUBCLUSTERS IN EACH NODE IF A NEW

SAMPLES ENTERS SUCH THAT THE NUMBER OF SUBCLUSTERS EXCEED THE BRANCHINGFACTOR THEN THAT

NODE IS SPLIT INTO TWO NODES WITH THE SUBCLUSTERS REDISTRIBUTED IN EACH THE PARENT SUBCLUSTER

OF THAT NODE IS REMOVED AND TWO NEW SUBCLUSTERS ARE ADDED AS PARENTS OF THE 2 SPLIT NODES

NCLUSTERS INT INSTANCE OF SKLEARNCLUSTER MODEL DEFAULT 3 NUMBER OF CLUSTERS AFTER THE FINAL

CLUSTERING STEP WHICH TREATS THE SUBCLUSTERS FROM THE LEAVES AS NEW SAMPLES

- NONE THE FINAL CLUSTERING STEP IS NOT PERFORMED AND THE SUBCLUSTERS ARE RETURNED AS THEY ARE
- SKLEARNCLUSTER ESTIMATOR IF A MODEL IS PROVIDED THE MODEL IS FIT TREATING THE SUBCLUSTERS AS NEW SAMPLES AND THE INITIAL DATA IS MAPPED TO THE LABEL OF THE CLOSEST SUBCLUSTER
- INT THE MODEL FIT IS AGGLOMERATIVECLUSTERING WITHNCLUSTERS SET TO BE EQUAL TO THE INT

COMPUTELABELS BOOL DEFAULT TRUE WHETHER OR NOT TO COMPUTE LABELS FOR EACH FIT

COPY BOOL DEFAULT TRUE WHETHER OR NOT TO MAKE A COPY OF THE GIVEN DATA IF SET TO FALSE THE

INITIAL DATA WILL BE OVERWRITTEN

ATTRIBUTES

ROOT CFNODE ROOT OF THE CFTREE

DUMMYLEAF CFNODE START POINTER TO ALL THE LEAVES

SUBCLUSTERCENTERS NDARRAY CENTROIDS OF ALL SUBCLUSTERS READ DIRECTLY FROM THE LEAVES

SUBCLUSTERLABELS NDARRAY LABELS ASSIGNED TO THE CENTROIDS OF THE SUBCLUSTERS AFTER THEY ARE

CLUSTERED GLOBALLY

LABELS NDARRAY SHAPE NSAMPLES ARRAY OF LABELS ASSIGNED TO THE INPUT DATA IF PARTIALFIT IS

USED INSTEAD OF FIT THEY ARE ASSIGNED TO THE LAST BATCH OF DATA

NOTES

THE TREE DATA STRUCTURE CONSISTS OF NODES WITH EACH NODE CONSISTING OF A NUMBER OF SUBCLUSTERS THE MAXIMUM

NUMBER OF SUBCLUSTERS IN A NODE IS DETERMINED BY THE BRANCHING FACTOR EACH SUBCLUSTER MAINTAINS A LINEAR SUM

SQUARED SUM AND THE NUMBER OF SAMPLES IN THAT SUBCLUSTER IN ADDITION EACH SUBCLUSTER CAN ALSO HAVE A NODE AS

ITS CHILD IF THE SUBCLUSTER IS NOT A MEMBER OF A LEAF NODE

63SKLEARNCLUSTER CLUSTERING 1477

SCIKITLEARN USER GUIDE RELEASE 0213

FOR A NEW POINT ENTERING THE ROOT IT IS MERGED WITH THE SUBCLUSTER CLOSEST TO IT AND THE LINEAR SUM SQUARED SUM AND THE NUMBER OF SAMPLES OF THAT SUBCLUSTER ARE UPDATED THIS IS DONE RECURSIVELY TILL THE PROPERTIES OF THE LEAF NODE ARE UPDATED

REFERENCES

- TIAN ZHANG RAGHU RAMAKRISHNAN MARON LIVNY BIRCH AN EFFICIENT DATA CLUSTERING METHOD FOR LARGE DATABASES [HTTPSWWWCSFUCACOURSECENTRAL459HANPAPERSZHANG96PDF](https://www.cse.cmu.edu/~raghu/papers/birch.pdf)
- ROBERTO PERDISCI JBIRCH JAVA IMPLEMENTATION OF BIRCH CLUSTERING ALGORITHM [HTTPSCODEGOOGLECOMARCHIVEPJBIRCH](http://code.google.com/archive/p/jbirch/)

EXAMPLES

```
FROM SKLEARNCLUSTER IMPORT BIRCH
X = 0 1 03 1 03 1 0 1 03 1 03 1
BRC = BIRCH(BRANCHINGFACTOR=50, NCLUSTERS=None, THRESHOLD=0.5,
COMPUTE_LABELS=True,
BRC_FIT_X=None)
BIRCH(BRANCHINGFACTOR=50, COMPUTE_LABELS=True, COPY=True, NCLUSTERS=None,
THRESHOLD=0.5,
BRC_PREDICT_X=None)
ARRAY([[0, 0, 1, 1, 1]])
METHODS
FIT(self, X, Y=None) BUILD A CF TREE FOR THE INPUT DATA
FITPREDICT(self, X, Y=None) PERFORMS CLUSTERING ON X AND RETURNS CLUSTER LABELS
FITTRANSFORM(self, X, Y=None) FIT TO DATA THEN TRANSFORM IT
GETPARAMS(self) DEEP GET PARAMETERS FOR THIS ESTIMATOR
PARTIALFIT(self, X, Y=None) ONLINE LEARNING
PREDICT(self, X) PREDICT DATA USING THE CENTROIDS OF SUBCLUSTERS
SETPARAMS(self, PARAMS) SET THE PARAMETERS OF THIS ESTIMATOR
TRANSFORM(self, X) TRANSFORM X INTO SUBCLUSTER CENTROIDS DIMENSION
INIT(self, THRESHOLD=0.5, BRANCHINGFACTOR=50, NCLUSTERS=3, COMPUTE_LABELS=True,
COPY=True,
FIT_SELF_X=None)
BUILD A CF TREE FOR THE INPUT DATA
PARAMETERS
X: ARRAY-LIKE SPARSE MATRIX SHAPE: (NSAMPLES, NFEATURES) INPUT DATA
Y: IGNORED
FITPREDICT(self, X, Y=None)
PERFORMS CLUSTERING ON X AND RETURNS CLUSTER LABELS
PARAMETERS
X: NDARRAY SHAPE: (NSAMPLES, NFEATURES) INPUT DATA
```

1478 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213  
YIGNORED NOT USED PRESENT FOR API CONSISTENCY BY CONVENTION  
RETURNS  
LABELS NDARRAY SHAPE NSAMPLES CLUSTER LABELS  
FITTRANSFORM SELFXYNONE FITPARAMS  
FIT TO DATA THEN TRANSFORM IT  
FITS TRANSFORMER TO X AND Y WITH OPTIONAL PARAMETERS FITPARAMS AND RETURNS A TRANSFORMED VERSION OF X  
PARAMETERS  
XNUMPY ARRAY OF SHAPE NSAMPLES NFEATURES TRAINING SET  
YNUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES  
RETURNS  
XNEW NUMPY ARRAY OF SHAPE NSAMPLES NFEATURESNEW TRANSFORMED ARRAY  
GETPARAMS SELFDEEPTTRUE  
GET PARAMETERS FOR THIS ESTIMATOR  
PARAMETERS  
DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBOBJECTS THAT ARE ESTIMATORS  
RETURNS  
PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES  
PARTIALFIT SELFXYNONE YNONE  
ONLINE LEARNING PREVENTS REBUILDING OF CFTREE FROM SCRATCH  
PARAMETERS  
XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES NONE INPUT DATA IF X IS NOT  
PROVIDED ONLY THE GLOBAL CLUSTERING STEP IS DONE  
YIGNORED  
PREDICTSELF  
PREDICT DATA USING THE CENTROIDS OF SUBCLUSTERS  
AVOID COMPUTATION OF THE ROW NORMS OF X  
PARAMETERS  
XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES INPUT DATA  
RETURNS  
LABELS NDARRAY SHAPENSAMPLES LABELLED DATA  
SETPARAMS SELFPARAMS  
SET THE PARAMETERS OF THIS ESTIMATOR  
THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE  
PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT  
OF A NESTED OBJECT  
RETURNS  
SELF  
63SKLEARNCLUSTER CLUSTERING 1479

SCIKITLEARN USER GUIDE RELEASE 0213

TRANSFORM SELF X

TRANSFORM X INTO SUBCLUSTER CENTROIDS DIMENSION

EACH DIMENSION REPRESENTS THE DISTANCE FROM THE SAMPLE POINT TO EACH CLUSTER CENTROID

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES INPUT DATA

RETURNS

XTRANS ARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NCLUSTERS TRANSFORMED DATA

EXAMPLES USING SKLEARNCLUSTERBIRCH

- COMPARE BIRCH AND MINIBATCHKMEANS
- COMPARING DIFFERENT CLUSTERING ALGORITHMS ON TOY DATASETS

SKLEARNCLUSTER DBSCAN

CLASSSSKLEARNCLUSTER DBSCANEPS05 MINSAMPLES5 METRIC'EUCLIDEAN' METRICPARAMSNONE

ALGORITHM'AUTO' LEAFSIZE30 PNONE NJOBSNONE

PERFORM DBSCAN CLUSTERING FROM VECTOR ARRAY OR DISTANCE MATRIX

DBSCAN DENSITYBASED SPATIAL CLUSTERING OF APPLICATIONS WITH NOISE FINDS CORE SAMPLES OF HIGH DENSITY

AND EXPANDS CLUSTERS FROM THEM GOOD FOR DATA WHICH CONTAINS CLUSTERS OF SIMILAR DENSITY

READ MORE IN THE USER GUIDE

PARAMETERS

EPS FLOAT OPTIONAL THE MAXIMUM DISTANCE BETWEEN TWO SAMPLES FOR ONE TO BE CONSIDERED AS

IN THE NEIGHBORHOOD OF THE OTHER THIS IS NOT A MAXIMUM BOUND ON THE DISTANCES OF POINTS

WITHIN A CLUSTER THIS IS THE MOST IMPORTANT DBSCAN PARAMETER TO CHOOSE APPROPRIATELY FOR

YOUR DATA SET AND DISTANCE FUNCTION

MINSAMPLES INT OPTIONAL THE NUMBER OF SAMPLES OR TOTAL WEIGHT IN A NEIGHBORHOOD FOR A

POINT TO BE CONSIDERED AS A CORE POINT THIS INCLUDES THE POINT ITSELF

METRIC STRING OR CALLABLE THE METRIC TO USE WHEN CALCULATING DISTANCE BETWEEN INSTANCES IN

A FEATURE ARRAY IF METRIC IS A STRING OR CALLABLE IT MUST BE ONE OF THE OPTIONS ALLOWED

BYSKLEARNMETRICSPAIRWISEDISTANCES FOR ITS METRIC PARAMETER IF METRIC

IS "PRECOMPUTED" X IS ASSUMED TO BE A DISTANCE MATRIX AND MUST BE SQUARE X MAY BE

A SPARSE MATRIX IN WHICH CASE ONLY "NONZERO" ELEMENTS MAY BE CONSIDERED NEIGHBORS FOR

DBSCAN

NEW IN VERSION 017 METRIC PRECOMPUTED TO ACCEPT PRECOMPUTED SPARSE MATRIX

METRICPARAMS DICT OPTIONAL ADDITIONAL KEYWORD ARGUMENTS FOR THE METRIC FUNCTION

NEW IN VERSION 019

ALGORITHM 'AUTO' 'BALLTREE' 'KDTree' 'BRUTE' OPTIONAL THE ALGORITHM TO BE USED BY THE

NEARESTNEIGHBORS MODULE TO COMPUTE POINTWISE DISTANCES AND FIND NEAREST NEIGHBORS SEE

NEARESTNEIGHBORS MODULE DOCUMENTATION FOR DETAILS

LEAFSIZE INT OPTIONAL DEFAULT 30 LEAF SIZE PASSED TO BALLTREE OR CKDTree THIS CAN AFFECT

THE SPEED OF THE CONSTRUCTION AND QUERY AS WELL AS THE MEMORY REQUIRED TO STORE THE TREE

THE OPTIMAL VALUE DEPENDS ON THE NATURE OF THE PROBLEM

1480 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

PFLOAT OPTIONAL THE POWER OF THE MINKOWSKI METRIC TO BE USED TO CALCULATE DISTANCE BETWEEN POINTS

NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF PARALLEL JOBS TO RUN NONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

ATTRIBUTES

CORESAMPLEINDICES ARRAY SHAPE NCORESAMPLES INDICES OF CORE SAMPLES

COMPONENTS ARRAY SHAPE NCORESAMPLES NFEATURES COPY OF EACH CORE SAMPLE FOUND BY TRAINING

LABELS ARRAY SHAPE NSAMPLES CLUSTER LABELS FOR EACH POINT IN THE DATASET GIVEN TO FIT

NOISY SAMPLES ARE GIVEN THE LABEL 1

SEE ALSO

OPTICS A SIMILAR CLUSTERING AT MULTIPLE VALUES OF EPS OUR IMPLEMENTATION IS OPTIMIZED FOR MEMORY USAGE

NOTES

FOR AN EXAMPLE SEE EXAMPLESCUSTERPLOTDBSCANPY

THIS IMPLEMENTATION BULKCOMPUTES ALL NEIGHBORHOOD QUERIES WHICH INCREASES THE MEMORY COMPLEXITY TO  $O(n^2)$  WHERE D IS THE AVERAGE NUMBER OF NEIGHBORS WHILE ORIGINAL DBSCAN HAD MEMORY COMPLEXITY  $O(n)$  ON IT MAY ATTRACT A HIGHER MEMORY COMPLEXITY WHEN QUERYING THESE NEAREST NEIGHBORHOODS DEPENDING ON THE ALGORITHM

ONE WAY TO AVOID THE QUERY COMPLEXITY IS TO PRECOMPUTE SPARSE NEIGHBORHOODS IN CHUNKS USINGNEARESTNEIGHBORSRADIUSNEIGHBORSGRAPH WITHMODEDISTANCE THEN USING METRICPRECOMPUTED HERE

ANOTHER WAY TO REDUCE MEMORY AND COMPUTATION TIME IS TO REMOVE NEARDUPLICATE POINTS AND USE SAMPLEWEIGHT INSTEAD

CLUSTEROPTICS PROVIDES A SIMILAR CLUSTERING WITH LOWER MEMORY USAGE

REFERENCES

ESTER M H P KRIEGEL J SANDER AND X XU "A DENSITYBASED ALGORITHM FOR DISCOVERING CLUSTERS IN LARGE SPATIAL DATABASES WITH NOISE" IN PROCEEDINGS OF THE 2ND INTERNATIONAL CONFERENCE ON KNOWLEDGE DISCOVERY AND DATA MINING PORTLAND OR AAAI PRESS PP 226231 1996

SCHUBERT E SANDER J ESTER M KRIEGEL H P XU X 2017 DBSCAN REVISITED REVISITED WHY AND HOW YOU SHOULD STILL USE DBSCAN ACM TRANSACTIONS ON DATABASE SYSTEMS TODS 423 19

EXAMPLES

```
FROM SKLEARNCLUSTER IMPORT DBSCAN
IMPORT NUMPY AS NP
X = NPARRAY(1 2 2 2 2 3
            8 7 8 8 25 80)
CLUSTERING = DBSCAN(EPS=3, MIN_SAMPLES=2, FIT_X=True)
CLUSTERING.labels_
63SKLEARNCLUSTER CLUSTERING 1481
```

SCIKITLEARN USER GUIDE RELEASE 0213  
ARRAY 0 0 0 1 1 1  
CLUSTERING  
DBSCANALGORITHM AUTO EPS3 LEAF SIZE30 METRIC EUCLIDEAN  
METRIC PARAMS NONE MIN SAMPLES2 NJOBS NONE P NONE  
METHODS  
FITSELF X Y SAMPLEWEIGHT PERFORM DBSCAN CLUSTERING FROM FEATURES OR DISTANCE  
MATRIX  
FITPREDICT SELF X Y SAMPLEWEIGHT PERFORMS CLUSTERING ON X AND RETURNS CLUSTER LABELS  
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR  
SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR  
INIT SELF EPS05 MIN SAMPLES5 METRIC 'EUCLIDEAN' METRIC PARAMS NONE ALGO  
RITHM 'AUTO' LEAF SIZE30 P NONE NJOBS NONE  
FITSELF X Y NONE SAMPLEWEIGHT NONE  
PERFORM DBSCAN CLUSTERING FROM FEATURES OR DISTANCE MATRIX  
PARAMETERS  
X ARRAY OR SPARSE CSR MATRIX OF SHAPE NSAMPLES NFEATURES OR ARRAY OF SHAPE  
NSAMPLES NSAMPLES A FEATURE ARRAY OR ARRAY OF DISTANCES BETWEEN SAMPLES IF  
METRIC PRECOMPUTED  
SAMPLEWEIGHT ARRAY SHAPE NSAMPLES OPTIONAL WEIGHT OF EACH SAMPLE SUCH THAT A  
SAMPLE WITH A WEIGHT OF AT LEAST MIN SAMPLES IS BY ITSELF A CORE SAMPLE A SAMPLE  
WITH NEGATIVE WEIGHT MAY INHIBIT ITS EPS NEIGHBOR FROM BEING CORE NOTE THAT WEIGHTS ARE  
ABSOLUTE AND DEFAULT TO 1  
Y IGNORED  
FITPREDICT SELF X Y NONE SAMPLEWEIGHT NONE  
PERFORMS CLUSTERING ON X AND RETURNS CLUSTER LABELS  
PARAMETERS  
X ARRAY OR SPARSE CSR MATRIX OF SHAPE NSAMPLES NFEATURES OR ARRAY OF SHAPE  
NSAMPLES NSAMPLES A FEATURE ARRAY OR ARRAY OF DISTANCES BETWEEN SAMPLES IF  
METRIC PRECOMPUTED  
SAMPLEWEIGHT ARRAY SHAPE NSAMPLES OPTIONAL WEIGHT OF EACH SAMPLE SUCH THAT A  
SAMPLE WITH A WEIGHT OF AT LEAST MIN SAMPLES IS BY ITSELF A CORE SAMPLE A SAMPLE  
WITH NEGATIVE WEIGHT MAY INHIBIT ITS EPS NEIGHBOR FROM BEING CORE NOTE THAT WEIGHTS ARE  
ABSOLUTE AND DEFAULT TO 1  
Y IGNORED  
RETURNS  
Y NDARRAY SHAPE NSAMPLES CLUSTER LABELS  
GETPARAMS SELF DEEP TRUE  
GET PARAMETERS FOR THIS ESTIMATOR  
PARAMETERS  
DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBJECTS THAT ARE ESTIMATORS  
1482 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE

PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT

OF A NESTED OBJECT

RETURNS

SELF

EXAMPLES USING SKLEARNCLUSTERDBSCAN

•DEMO OF DBSCAN CLUSTERING ALGORITHM

•COMPARING DIFFERENT CLUSTERING ALGORITHMS ON TOY DATASETS

SKLEARNCLUSTER OPTICS

CLASSSKLEARNCLUSTER OPTICSMINSAMPLES5 MAXEPSINF METRIC’MINKOWSKI’ P2MET

RICPARAMSNONE CLUSTERMETHOD’XI’ EPSNONE XI005 PREDE

CESSORCORRECTIONTRUE MINCLUSTERSIZENONE ALGORITHM’AUTO’

LEAFSIZE30 NJOBSNONE

ESTIMATE CLUSTERING STRUCTURE FROM VECTOR ARRAY

OPTICS ORDERING POINTS TO IDENTIFY THE CLUSTERING STRUCTURE CLOSELY RELATED TO DBSCAN FINDS CORE SAMPLE

OF HIGH DENSITY AND EXPANDS CLUSTERS FROM THEM R2C55E37003FE1 UNLIKE DBSCAN KEEPS CLUSTER HIERARCHY

FOR A VARIABLE NEIGHBORHOOD RADIUS BETTER SUITED FOR USAGE ON LARGE DATASETS THAN THE CURRENT SKLEARN IMPLEMEN

TATION OF DBSCAN

CLUSTERS ARE THEN EXTRACTED USING A DBSCANLIKE METHOD CLUSTERMETHOD ‘DBSCAN’ OR AN AUTOMATIC TECHNIQUE

PROPOSED IN R2C55E37003FE1 CLUSTERMETHOD ‘XI’

THIS IMPLEMENTATION DEVIATES FROM THE ORIGINAL OPTICS BY FIRST PERFORMING KNEARESTNEIGHBORHOOD SEARCHES

ON ALL POINTS TO IDENTIFY CORE SIZES THEN COMPUTING ONLY THE DISTANCES TO UNPROCESSED POINTS WHEN CONSTRUCTING

THE CLUSTER ORDER NOTE THAT WE DO NOT EMPLOY A HEAP TO MANAGE THE EXPANSION CANDIDATES SO THE TIME COMPLEXITY

WILL BE ON2

READ MORE IN THE USER GUIDE

PARAMETERS

MINSAMPLES INT 1 OR FLOAT BETWEEN 0 AND 1 DEFAULT5 THE NUMBER OF SAMPLES IN A NEIGH

BORHOOD FOR A POINT TO BE CONSIDERED AS A CORE POINT ALSO UP AND DOWN STEEP REGIONS CAN’T

HAVE MORE THEN MINSAMPLES CONSECUTIVE NONSTEEP POINTS EXPRESSED AS AN ABSOLUTE

NUMBER OR A FRACTION OF THE NUMBER OF SAMPLES ROUNDED TO BE AT LEAST 2

MAXEPS FLOAT OPTIONAL DEFAULTNPINF THE MAXIMUM DISTANCE BETWEEN TWO SAMPLES FOR ONE

TO BE CONSIDERED AS IN THE NEIGHBORHOOD OF THE OTHER DEFAULT VALUE OF NPINF WILL IDENTIFY

CLUSTERS ACROSS ALL SCALES REDUCING MAXEPS WILL RESULT IN SHORTER RUN TIMES

METRIC STRING OR CALLABLE OPTIONAL DEFAULT’MINKOWSKI’ METRIC TO USE FOR DISTANCE COMPUTA

TION ANY METRIC FROM SCIKITLEARN OR SCIPYSPATIALDISTANCE CAN BE USED

IF METRIC IS A CALLABLE FUNCTION IT IS CALLED ON EACH PAIR OF INSTANCES ROWS AND THE RESULTING

VALUE RECORDED THE CALLABLE SHOULD TAKE TWO ARRAYS AS INPUT AND RETURN ONE VALUE INDICATING

63SKLEARNCLUSTER CLUSTERING 1483

SCIKITLEARN USER GUIDE RELEASE 0213

THE DISTANCE BETWEEN THEM THIS WORKS FOR SCIPY’S METRICS BUT IS LESS EFFICIENT THAN PASSING THE METRIC NAME AS A STRING IF METRIC IS “PRECOMPUTED” X IS ASSUMED TO BE A DISTANCE MATRIX AND MUST BE SQUARE

VALID VALUES FOR METRIC ARE

- FROM SCIKITLEARN ‘CITYBLOCK’ ‘COSINE’ ‘EUCLIDEAN’ ‘L1’ ‘L2’ ‘MANHATTAN’
- FROM SCIPYSPATIALDISTANCE ‘BRAYCURTIS’ ‘CANBERRA’ ‘CHEBYSHEV’ ‘CORRELATION’ ‘DICE’ ‘HAMMING’ ‘JACCARD’ ‘KULSINSKI’ ‘MAHALANOBIS’ ‘MINKOWSKI’ ‘ROGERSTANIMOTO’ ‘RUSSELLRAO’ ‘SEUCLIDEAN’ ‘SOKALMICHERNER’ ‘SOKALSNEATH’ ‘SQUEUCLIDEAN’ ‘YULE’

SEE THE DOCUMENTATION FOR SCIPYSPATIALDISTANCE FOR DETAILS ON THESE METRICS

PINTEGRER OPTIONAL DEFAULT2 PARAMETER FOR THE MINKOWSKI METRIC FROM SKLEARN

METRICSPAIRWISEDISTANCES WHEN P 1 THIS IS EQUIVALENT TO USING MANHAT

TANDISTANCE L1 AND EUCLIDEANDISTANCE L2 FOR P 2 FOR ARBITRARY P MINKOWSKIDISTANCE

LP IS USED

METRICPARAMS DICT OPTIONAL DEFAULTNONE ADDITIONAL KEYWORD ARGUMENTS FOR THE METRIC

FUNCTION

CLUSTERMETHOD STRING OPTIONAL DEFAULT’XI’ THE EXTRACTION METHOD USED TO EXTRACT CLUSTERS

USING THE CALCULATED REACHABILITY AND ORDERING POSSIBLE VALUES ARE “XI” AND “DBSCAN”

EPS FLOAT OPTIONAL DEFAULTNONE THE MAXIMUM DISTANCE BETWEEN TWO SAMPLES FOR ONE TO BE

CONSIDERED AS IN THE NEIGHBORHOOD OF THE OTHER BY DEFAULT IT ASSUMES THE SAME VALUE AS

MAXEPS USED ONLY WHEN CLUSTERMETHODDBSCAN

XIFLOAT BETWEEN 0 AND 1 OPTIONAL DEFAULT005 DETERMINES THE MINIMUM STEEPNESS ON THE

REACHABILITY PLOT THAT CONSTITUTES A CLUSTER BOUNDARY FOR EXAMPLE AN UPWARDS POINT IN THE

REACHABILITY PLOT IS DEFINED BY THE RATIO FROM ONE POINT TO ITS SUCCESSOR BEING AT MOST 1XI

USED ONLY WHEN CLUSTERMETHODXI

PREDECESSORCORRECTION BOOL OPTIONAL DEFAULTTRUE CORRECT CLUSTERS ACCORDING TO THE PRE

DECESSORS CALCULATED BY OPTICS R2C55E37003FE2 THIS PARAMETER HAS MINIMAL EFFECT ON

MOST DATASETS USED ONLY WHEN CLUSTERMETHODXI

MINCLUSTERSIZE INT 1 OR FLOAT BETWEEN 0 AND 1 DEFAULTNONE MINIMUM NUMBER OF SAM

PLES IN AN OPTICS CLUSTER EXPRESSED AS AN ABSOLUTE NUMBER OR A FRACTION OF THE NUMBER OF

SAMPLES ROUNDED TO BE AT LEAST 2 IF NONE THE VALUE OF MINSAMPLES IS USED INSTEAD

USED ONLY WHEN CLUSTERMETHODXI

ALGORITHM ‘AUTO’ ‘BALLTREE’ ‘KDTREE’ ‘BRUTE’ OPTIONAL ALGORITHM USED TO COMPUTE THE

NEAREST NEIGHBORS

- ‘BALLTREE’ WILL USE BALLTREE
- ‘KDTREE’ WILL USE KDTREE
- ‘BRUTE’ WILL USE A BRUTEFORCE SEARCH
- ‘AUTO’ WILL ATTEMPT TO DECIDE THE MOST APPROPRIATE ALGORITHM BASED ON THE VALUES PASSED

TOFIT METHOD DEFAULT

NOTE FITTING ON SPARSE INPUT WILL OVERRIDE THE SETTING OF THIS PARAMETER USING BRUTE FORCE

LEAFSIZE INT OPTIONAL DEFAULT30 LEAF SIZE PASSED TO BALLTREE ORKDTREE THIS CAN

AFFECT THE SPEED OF THE CONSTRUCTION AND QUERY AS WELL AS THE MEMORY REQUIRED TO STORE THE

TREE THE OPTIMAL VALUE DEPENDS ON THE NATURE OF THE PROBLEM

1484 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF PARALLEL JOBS TO RUN FOR NEIGHBORS

SEARCHNONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS

USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

ATTRIBUTES

LABELS ARRAY SHAPE NSAMPLES CLUSTER LABELS FOR EACH POINT IN THE DATASET GIVEN

TO FIT NOISY SAMPLES AND POINTS WHICH ARE NOT INCLUDED IN A LEAF CLUSTER OF

CLUSTERHIERARCHY ARE LABELED AS 1

REACHABILITY ARRAY SHAPE NSAMPLES REACHABILITY DISTANCES PER SAMPLE INDEXED BY OBJECT

ORDER USECLUSTREACHABILITYCLUSTORDERING TO ACCESS IN CLUSTER ORDER

ORDERING ARRAY SHAPE NSAMPLES THE CLUSTER ORDERED LIST OF SAMPLE INDICES

COREDISTANCES ARRAY SHAPE NSAMPLES DISTANCE AT WHICH EACH SAMPLE BECOMES A CORE

POINT INDEXED BY OBJECT ORDER POINTS WHICH WILL NEVER BE CORE HAVE A DISTANCE OF INF USE

CLUSTCOREDISTANCESCLUSTORDERING TO ACCESS IN CLUSTER ORDER

PREDECESSOR ARRAY SHAPE NSAMPLES POINT THAT A SAMPLE WAS REACHED FROM INDEXED BY

OBJECT ORDER SEED POINTS HAVE A PREDECESSOR OF 1

CLUSTERHIERARCHY ARRAY SHAPE NCLUSTERS 2 THE LIST OF CLUSTERS IN THE FORM OF START

END IN EACH ROW WITH ALL INDICES INCLUSIVE THE CLUSTERS ARE ORDERED ACCORDING

TOEND START ASCENDING SO THAT LARGER CLUSTERS ENCOMPASSING SMALLER CLUSTERS

COME AFTER THOSE SMALLER ONES SINCE LABELS DOES NOT REFLECT THE HIERARCHY USUALLY

LENCLUSTERHIERARCHY NPUNIQUEOPTICSLABELS PLEASE ALSO

NOTE THAT THESE INDICES ARE OF THE ORDERING IEXORDERINGSTARTEND

1FORM A CLUSTER ONLY AVAILABLE WHEN CLUSTERMETHODXI

SEE ALSO

DBSCAN A SIMILAR CLUSTERING FOR A SPECIFIED NEIGHBORHOOD RADIUS EPS OUR IMPLEMENTATION IS OPTIMIZED FOR

RUNTIME

REFERENCES

R2C55E37003FE1 R2C55E37003FE2

METHODS

FITSELF X Y PERFORM OPTICS CLUSTERING

FITPREDICT SELF X Y PERFORMS CLUSTERING ON X AND RETURNS CLUSTER LABELS

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELFMINSAMPLES5 MAXEPSINF METRIC'MINKOWSKI' P2METRICPARAMSNONE

CLUSTERMETHOD'XI' EPSNONE XI005 PREDECESSORCORRECTIONTRUE

MINCLUSTERSIZENONE ALGORITHM'AUTO' LEAFSIZE30 NJOBSNONE

FITSELFXYNONE

PERFORM OPTICS CLUSTERING

EXTRACTS AN ORDERED LIST OF POINTS AND REACHABILITY DISTANCES AND PERFORMS INITIAL CLUSTERING USING MAXEPS

DISTANCE SPECIFIED AT OPTICS OBJECT INSTANTIATION

PARAMETERS

63SKLEARNCLUSTER CLUSTERING 1485

SCIKITLEARN USER GUIDE RELEASE 0213

XARRAY SHAPE NSAMPLES NFEATURES OR NSAMPLES NSAMPLES IF MET

RIC'PRECOMPUTED' A FEATURE ARRAY OR ARRAY OF DISTANCES BETWEEN SAMPLES IF MET

RIC'PRECOMPUTED'

YIGNORED

RETURNS

SELF INSTANCE OF OPTICS THE INSTANCE

FITPREDICT SELFXYNONE

PERFORMS CLUSTERING ON X AND RETURNS CLUSTER LABELS

PARAMETERS

XNDARRAY SHAPE NSAMPLES NFEATURES INPUT DATA

YIGNORED NOT USED PRESENT FOR API CONSISTENCY BY CONVENTION

RETURNS

LABELS NDARRAY SHAPE NSAMPLES CLUSTER LABELS

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE

PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT

OF A NESTED OBJECT

RETURNS

SELF

EXAMPLES USING SKLEARNCLUSTEROPTICS

- DEMO OF OPTICS CLUSTERING ALGORITHM
- COMPARING DIFFERENT CLUSTERING ALGORITHMS ON TOY DATASETS

SKLEARNCLUSTER FEATUREAGGLOMERATION

CLASSSSKLEARNCLUSTER FEATUREAGGLOMERATION NCLUSTERS2 AFFINITY'EUCLIDEAN' MEM

ORYNONE CONNECTIVITYNONE COM

PUTEFULLTREE'AUTO' LINKAGE'WARD'

POOLINGFUNCFUNCTION MEAN DIS

TANCETHRESHOLDNONE

AGGLOMERATE FEATURES

SIMILAR TO AGGLOMERATIVECLUSTERING BUT RECURSIVELY MERGES FEATURES INSTEAD OF SAMPLES

1486 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

READ MORE IN THE USER GUIDE

PARAMETERS

NCLUSTERS INT OR NONE OPTIONAL DEFAULT2 THE NUMBER OF CLUSTERS TO FIND IT MUST BE NONE

IFDISTANCETHRESHOLD IS NOTNONE

AFFINITY STRING OR CALLABLE DEFAULT “EUCLIDEAN” METRIC USED TO COMPUTE THE LINKAGE CAN BE

“EUCLIDEAN” “L1” “L2” “MANHATTAN” “COSINE” OR ‘PRECOMPUTED’ IF LINKAGE IS “WARD” ONLY

“EUCLIDEAN” IS ACCEPTED

MEMORY NONE STR OR OBJECT WITH THE JOBLIBMEMORY INTERFACE OPTIONAL USED TO CACHE THE OUTPUT OF THE COMPUTATION OF THE TREE BY DEFAULT NO CACHING IS DONE IF A STRING IS GIVEN IT IS THE PATH TO THE CACHING DIRECTORY

CONNECTIVITY ARRAYLIKE OR CALLABLE OPTIONAL CONNECTIVITY MATRIX DEFINES FOR EACH FEATURE THE NEIGHBORING FEATURES FOLLOWING A GIVEN STRUCTURE OF THE DATA THIS CAN BE A CONNECTIVITY MATRIX ITSELF OR A CALLABLE THAT TRANSFORMS THE DATA INTO A CONNECTIVITY MATRIX SUCH AS DERIVED FROM KNEIGHBORSGRAPH DEFAULT IS NONE IE THE HIERARCHICAL CLUSTERING ALGORITHM IS UNSTRUCTURED

COMPUTEFULLTREE BOOL OR ‘AUTO’ OPTIONAL DEFAULT “AUTO” STOP EARLY THE CONSTRUCTION OF THE TREE AT NCLUSTERS THIS IS USEFUL TO DECREASE COMPUTATION TIME IF THE NUMBER OF CLUSTERS IS NOT SMALL COMPARED TO THE NUMBER OF FEATURES THIS OPTION IS USEFUL ONLY WHEN SPECIFYING A CONNECTIVITY MATRIX NOTE ALSO THAT WHEN VARYING THE NUMBER OF CLUSTERS AND USING CACHING IT MAY BE ADVANTAGEOUS TO COMPUTE THE FULL TREE IT MUST BE TRUE IF

DISTANCETHRESHOLD IS NOTNONE

LINKAGE “WARD” “COMPLETE” “AVERAGE” “SINGLE” OPTIONAL DEFAULT”WARD” WHICH LINKAGE CRITERION TO USE THE LINKAGE CRITERION DETERMINES WHICH DISTANCE TO USE BETWEEN SETS OF FEATURES THE ALGORITHM WILL MERGE THE PAIRS OF CLUSTER THAT MINIMIZE THIS CRITERION

- WARD MINIMIZES THE VARIANCE OF THE CLUSTERS BEING MERGED
  - AVERAGE USES THE AVERAGE OF THE DISTANCES OF EACH FEATURE OF THE TWO SETS
  - COMPLETE OR MAXIMUM LINKAGE USES THE MAXIMUM DISTANCES BETWEEN ALL FEATURES OF THE TWO SETS
  - SINGLE USES THE MINIMUM OF THE DISTANCES BETWEEN ALL OBSERVATIONS OF THE TWO SETS
- POOLINGFUNC CALLABLE DEFAULT NPMEAN THIS COMBINES THE VALUES OF AGGLOMERATED FEATURES INTO A SINGLE VALUE AND SHOULD ACCEPT AN ARRAY OF SHAPE M N AND THE KEYWORD ARGUMENT AXIS1 AND REDUCE IT TO AN ARRAY OF SIZE M

DISTANCETHRESHOLD FLOAT OPTIONAL DEFAULTNONE THE LINKAGE DISTANCE THRESHOLD ABOVE WHICH CLUSTERS WILL NOT BE MERGED IF NOT NONE NCLUSTERS MUST BENONE AND

COMPUTEFULLTREE MUST BETRUE

NEW IN VERSION 021

ATTRIBUTES

NCLUSTERS INT THE NUMBER OF CLUSTERS FOUND BY THE ALGORITHM IF DISTANCETHRESHOLDNONE IT WILL BE EQUAL TO THE GIVEN NCLUSTERS

LABELS ARRAYLIKE NFEATURES CLUSTER LABELS FOR EACH FEATURE

NLEAVES INT NUMBER OF LEAVES IN THE HIERARCHICAL TREE

NCONNECTEDCOMPONENTS INT THE ESTIMATED NUMBER OF CONNECTED COMPONENTS IN THE GRAPH CHILDREN ARRAYLIKE SHAPE NNODES1 2 THE CHILDREN OF EACH NONLEAF NODE VALUES LESS THAN NFEATURES CORRESPOND TO LEAVES OF THE TREE WHICH ARE THE ORIGINAL SAMPLES A NODE

63SKLEARNCLUSTER CLUSTERING 1487

SCIKITLEARN USER GUIDE RELEASE 0213

IGREATER THAN OR EQUAL TO NFEATURES IS A NONLEAF NODE AND HAS CHILDREN CHILDREN1  
NFEATURES ALTERNATIVELY AT THE ITH ITERATION CHILDREN10 AND CHILDREN11 ARE  
MERGED TO FORM NODE NFEATURES I

EXAMPLES

```
import numpy as np
from sklearn import datasets
cluster = ClusterFeatureAgglomeration(n_clusters=32,
agglo='fit',
feature_agglomeration='euclidean',
compute_full_tree='auto',
connectivity=None,
distance_threshold=None,
linkage='ward',
memory=None,
n_clusters=32,
pooling_func=None,
x_reduced=agglo.transform(X_reduced.shape[0:1797, 32])
```

METHODS

`fit(X, Y)` Fit the hierarchical clustering on the data

`fit_transform(X, Y)` Fit to data then transform it

`get_params(self, deep=True)` Get parameters for this estimator

`inverse_transform(self, X_reduced)` Inverse the transformation

`pooling_func(X, axis=0, dtype=None, out=None, keepdims=False)` Compute the arithmetic mean along the specified axis

`set_params(self, **kwargs)` Set the parameters of this estimator

`transform(self, X)` Transform a new matrix using the built clustering

`init(self, n_clusters=2, affinity='euclidean', memory=None, connectivity=None, compute_full_tree='auto', linkage='ward', pooling_func=None, function_mean_at=0x7f3c23df3400, distance_threshold=None)`

`fit(self, X, Y=None, params=None)`

Fit the hierarchical clustering on the data

PARAMETERS

`X` array-like shape (n\_samples, n\_features) the data

`Y` ignored

RETURNS

`self`

`fit_transform(self, X, Y=None, fit_params=None)`

Fit to data then transform it

Fits transformer to `X` and `Y` with optional parameters `fit_params` and returns a transformed version of `X`

PARAMETERS

`X` numpy array of shape (n\_samples, n\_features) training set

1488 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

YNUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES

RETURNS

XNEW NUMPY ARRAY OF SHAPE NSAMPLES NFEATURESNEW TRANSFORMED ARRAY

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

INVERSETRANSFORM SELFRED

INVERSE THE TRANSFORMATION RETURN A VECTOR OF SIZE NBFEATURES WITH THE VALUES OF RED ASSIGNED TO EACH

GROUP OF FEATURES

PARAMETERS

RED ARRAYLIKE SHAPENSAMPLES NCLUSTERS OR NCLUSTERS THE VALUES TO BE ASSIGNED

TO EACH CLUSTER OF SAMPLES

RETURNS

XARRAY SHAPENSAMPLES NFEATURES OR NFEATURES A VECTOR OF SIZE NSAMPLES WITH

THE VALUES OF RED ASSIGNED TO EACH OF THE CLUSTER OF SAMPLES

POOLINGFUNC AAXISNONE DYPENONE OUTNONE KEEPDIRMSNO VALUE

COMPUTE THE ARITHMETIC MEAN ALONG THE SPECIFIED AXIS

RETURNS THE AVERAGE OF THE ARRAY ELEMENTS THE AVERAGE IS TAKEN OVER THE FLATTENED ARRAY BY DEFAULT OTHER

WISE OVER THE SPECIFIED AXIS FLOAT64 INTERMEDIATE AND RETURN VALUES ARE USED FOR INTEGER INPUTS

PARAMETERS

AARRAYLIKE ARRAY CONTAINING NUMBERS WHOSE MEAN IS DESIRED IF AIS NOT AN ARRAY A

CONVERSION IS ATTEMPTED

AXIS NONE OR INT OR TUPLE OF INTS OPTIONAL AXIS OR AXES ALONG WHICH THE MEANS ARE COM

PUTED THE DEFAULT IS TO COMPUTE THE MEAN OF THE FLATTENED ARRAY

NEW IN VERSION 170

IF THIS IS A TUPLE OF INTS A MEAN IS PERFORMED OVER MULTIPLE AXES INSTEAD OF A SINGLE AXIS OR

ALL THE AXES AS BEFORE

DTYPE DATATYPE OPTIONAL DTYPE TO USE IN COMPUTING THE MEAN FOR INTEGER INPUTS THE

DEFAULT ISFLOAT64 FOR FLOATING POINT INPUTS IT IS THE SAME AS THE INPUT DTYPE

OUT NDARRAY OPTIONAL ALTERNATE OUTPUT ARRAY IN WHICH TO PLACE THE RESULT THE DEFAULT IS

NONE IF PROVIDED IT MUST HAVE THE SAME SHAPE AS THE EXPECTED OUTPUT BUT THE TYPE WILL

BE CAST IF NECESSARY SEE DOCUFUNCS FOR DETAILS

KEEPDIRMS BOOL OPTIONAL IF THIS IS SET TO TRUE THE AXES WHICH ARE REDUCED ARE LEFT IN THE

RESULT AS DIMENSIONS WITH SIZE ONE WITH THIS OPTION THE RESULT WILL BROADCAST CORRECTLY

AGAINST THE INPUT ARRAY

IF THE DEFAULT VALUE IS PASSED THEN KEEPDIRMS WILL NOT BE PASSED THROUGH TO THE MEAN

METHOD OF SUBCLASSES OF NDARRAY HOWEVER ANY NONDEFAULT VALUE WILL BE IF THE SUB

CLASS' METHOD DOES NOT IMPLEMENT KEEPDIRMS ANY EXCEPTIONS WILL BE RAISED

63SKLEARNCLUSTER CLUSTERING 1489

SCIKITLEARN USER GUIDE RELEASE 0213

RETURNS

MNDARRAY SEE DTYPE PARAMETER ABOVE IF OUTNONE RETURNS A NEW ARRAY CONTAINING THE MEAN VALUES OTHERWISE A REFERENCE TO THE OUTPUT ARRAY IS RETURNED

SEE ALSO

AVERAGE WEIGHTED AVERAGE

STDVARNANMEAN NANSTD NANVAR

NOTES

THE ARITHMETIC MEAN IS THE SUM OF THE ELEMENTS ALONG THE AXIS DIVIDED BY THE NUMBER OF ELEMENTS

NOTE THAT FOR FLOATINGPOINT INPUT THE MEAN IS COMPUTED USING THE SAME PRECISION THE INPUT HAS DEPENDING ON THE INPUT DATA THIS CAN CAUSE THE RESULTS TO BE INACCURATE ESPECIALLY FOR FLOAT32 SEE EXAMPLE BELOW

SPECIFYING A HIGHERPRECISION ACCUMULATOR USING THE DTYPE KEYWORD CAN ALLEVIATE THIS ISSUE

BY DEFAULT FLOAT16 RESULTS ARE COMPUTED USING FLOAT32 INTERMEDIATES FOR EXTRA PRECISION

EXAMPLES

A NPARRAY1 2 3 4

NPMEANA

25

NPMEANA AXIS0

ARRAY 2 3

NPMEANA AXIS1

ARRAY 15 35

IN SINGLE PRECISION MEAN CAN BE INACCURATE

A NPZEROS2 512 512 DYPENPFLOAT32

A0 10

A1 01

NPMEANA

054999924

COMPUTING THE MEAN IN FLOAT64 IS MORE ACCURATE

NPMEANA DYPENPFLOAT64

055000000074505806

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELF

TRANSFORM A NEW MATRIX USING THE BUILT CLUSTERING

1490 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES OR NFEATURES A M BY N ARRAY OF M OBSERVATIONS IN N DIMENSIONS OR A LENGTH M ARRAY OF M ONEDIMENSIONAL OBSERVATIONS

RETURNS

YARRAY SHAPE NSAMPLES NCLUSTERS OR NCLUSTERS THE POOLED VALUES FOR EACH FEATURE

CLUSTER

EXAMPLES USING SKLEARNCLUSTERFEATUREAGGLOMERATION

- FEATURE AGGLOMERATION
- FEATURE AGGLOMERATION VS UNIVARIATE SELECTION

SKLEARNCLUSTER KMEANS

CLASSSSKLEARNCLUSTER KMEANSNCLUSTERS8 INIT'KMEANS' NINIT10 MAXITER300

TOL00001 PRECOMPUTEDISTANCES'AUTO' VERBOSE0 RAN

DOMSTATENONE COPYXTRUE NJOBSNONE ALGORITHM'AUTO'

KMEANS CLUSTERING

READ MORE IN THE USER GUIDE

PARAMETERS

NCLUSTERS INT OPTIONAL DEFAULT 8 THE NUMBER OF CLUSTERS TO FORM AS WELL AS THE NUMBER OF CENTROIDS TO GENERATE

INIT 'KMEANS' 'RANDOM' OR AN NDARRAY METHOD FOR INITIALIZATION DEFAULTS TO 'K

MEANS'

'KMEANS' SELECTS INITIAL CLUSTER CENTERS FOR KMEAN CLUSTERING IN A SMART WAY TO SPEED

UP CONVERGENCE SEE SECTION NOTES IN KINIT FOR MORE DETAILS

'RANDOM' CHOOSE K OBSERVATIONS ROWS AT RANDOM FROM DATA FOR THE INITIAL CENTROIDS

IF AN NDARRAY IS PASSED IT SHOULD BE OF SHAPE NCLUSTERS NFEATURES AND GIVES THE INITIAL

CENTERS

NINIT INT DEFAULT 10 NUMBER OF TIME THE KMEANS ALGORITHM WILL BE RUN WITH DIFFERENT CEN TROID SEEDS THE FINAL RESULTS WILL BE THE BEST OUTPUT OF NINIT CONSECUTIVE RUNS IN TERMS OF

INERTIA

MAXITER INT DEFAULT 300 MAXIMUM NUMBER OF ITERATIONS OF THE KMEANS ALGORITHM FOR A SINGLE RUN

TOLFLOAT DEFAULT 1E4 RELATIVE TOLERANCE WITH REGARDS TO INERTIA TO DECLARE CONVERGENCE

PRECOMPUTEDISTANCES 'AUTO' TRUE FALSE PRECOMPUTE DISTANCES FASTER BUT TAKES MORE

MEMORY

'AUTO' DO NOT PRECOMPUTE DISTANCES IF NSAMPLES NCLUSTERS 12 MILLION THIS CORRE

SPONDS TO ABOUT 100MB OVERHEAD PER JOB USING DOUBLE PRECISION

TRUE ALWAYS PRECOMPUTE DISTANCES

FALSE NEVER PRECOMPUTE DISTANCES

VERBOSE INT DEFAULT 0 VERBOSITY MODE

63SKLEARNCLUSTER CLUSTERING 1491

SCIKITLEARN USER GUIDE RELEASE 0213

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE DEFAULT DETERMINES RANDOM NUMBER GENERATION FOR CENTROID INITIALIZATION USE AN INT TO MAKE THE RANDOMNESS DETERMINISTIC SEE GLOSSARY

COPYX BOOLEAN OPTIONAL WHEN PRECOMPUTING DISTANCES IT IS MORE NUMERICALLY ACCURATE TO CENTER THE DATA FIRST IF COPYX IS TRUE DEFAULT THEN THE ORIGINAL DATA IS NOT MODIFIED ENSURING X IS CCONTIGUOUS IF FALSE THE ORIGINAL DATA IS MODIFIED AND PUT BACK BEFORE THE FUNCTION RETURNS BUT SMALL NUMERICAL DIFFERENCES MAY BE INTRODUCED BY SUBTRACTING AND THEN ADDING THE DATA MEAN IN THIS CASE IT WILL ALSO NOT ENSURE THAT DATA IS CCONTIGUOUS WHICH MAY CAUSE A SIGNIFICANT SLOWDOWN

NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF JOBS TO USE FOR THE COMPUTATION THIS WORKS BY COMPUTING EACH OF THE NINIT RUNS IN PARALLEL NONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

ALGORITHM “AUTO” “FULL” OR “ELKAN” DEFAULT”AUTO” KMEANS ALGORITHM TO USE THE CLASSICAL EMSTYLE ALGORITHM IS “FULL” THE “ELKAN” VARIATION IS MORE EFFICIENT BY USING THE TRIANGLE INEQUALITY BUT CURRENTLY DOESN’T SUPPORT SPARSE DATA “AUTO” CHOOSES “ELKAN” FOR DENSE DATA AND “FULL” FOR SPARSE DATA

ATTRIBUTES

CLUSTERCENTERS ARRAY NCLUSTERS NFEATURES COORDINATES OF CLUSTER CENTERS IF THE ALGORITHM STOPS BEFORE FULLY CONVERGING SEE TOL ANDMAXITER THESE WILL NOT BE CONSISTENT WITHLABELS

LABELS LABELS OF EACH POINT

INERTIA FLOAT SUM OF SQUARED DISTANCES OF SAMPLES TO THEIR CLOSEST CLUSTER CENTER

NITER INT NUMBER OF ITERATIONS RUN

SEE ALSO

MINIBATCHKMEANS ALTERNATIVE ONLINE IMPLEMENTATION THAT DOES INCREMENTAL UPDATES OF THE CENTERS POSITIONS USING MINIBATCHES FOR LARGE SCALE LEARNING SAY NSAMPLES 10K MINIBATCHKMEANS IS PROBABLY MUCH FASTER THAN THE DEFAULT BATCH IMPLEMENTATION

NOTES

THE KMEANS PROBLEM IS SOLVED USING EITHER LLOYD’S OR ELKAN’S ALGORITHM THE AVERAGE COMPLEXITY IS GIVEN BY  $O(K N T)$  WERE N IS THE NUMBER OF SAMPLES AND T IS THE NUMBER OF ITERATION THE WORST CASE COMPLEXITY IS GIVEN BY  $O(NK^2P)$  WITH N NSAMPLES P NFEATURES D ARTHUR AND S VASSILVITSKII ‘HOW SLOW IS THE KMEANS METHOD’ SOCG2006

IN PRACTICE THE KMEANS ALGORITHM IS VERY FAST ONE OF THE FASTEST CLUSTERING ALGORITHMS AVAILABLE BUT IT FALLS IN LOCAL MINIMA THAT’S WHY IT CAN BE USEFUL TO RESTART IT SEVERAL TIMES

IF THE ALGORITHM STOPS BEFORE FULLY CONVERGING BECAUSE OF TOL ORMAXITER LABELS AND CLUSTERCENTERS WILL NOT BE CONSISTENT IE THE CLUSTERCENTERS WILL NOT BE THE MEANS OF THE POINTS IN EACH CLUSTER ALSO THE ESTIMATOR WILL REASSIGN LABELS AFTER THE LAST ITERATION TO MAKE LABELS CONSISTENT WITH PREDICT ON THE TRAINING SET

```
SCIKITLEARN USER GUIDE RELEASE 0213
EXAMPLES
FROM SKLEARNCLUSTER IMPORT KMEANS
IMPORT NUMPY AS NP
X NPARRAY1 2 1 4 1 0
10 2 10 4 10 0
KMEANS KMEANSNCLUSTERS2 RANDOMSTATE0FITX
KMEANSLABELS
ARRAY1 1 1 0 0 0 DTYPEINT32
KMEANSPREDICT0 0 12 3
ARRAY1 0 DTYPEINT32
KMEANSCLUSTERCENTERS
ARRAY10 2
1 2
METHODS
FITSELF X Y SAMPLEWEIGHT COMPUTE KMEANS CLUSTERING
FITPREDICT SELF X Y SAMPLEWEIGHT COMPUTE CLUSTER CENTERS AND PREDICT CLUSTER INDEX FOR
EACH SAMPLE
FITTRANSFORM SELF X Y SAMPLEWEIGHT COMPUTE CLUSTERING AND TRANSFORM X TO CLUSTERDISTANCE
SPACE
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR
PREDICT SELF X SAMPLEWEIGHT PREDICT THE CLOSEST CLUSTER EACH SAMPLE IN X BELONGS TO
SCORE SELF X Y SAMPLEWEIGHT OPPOSITE OF THE VALUE OF X ON THE KMEANS OBJECTIVE
SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR
TRANSFORM SELF X TRANSFORM X TO A CLUSTERDISTANCE SPACE
INIT SELFNCLUSTERS8 INIT'KMEANS' NINIT10 MAXITER300 TOL00001 PRECOM
PUTEDISTANCES'AUTO' VERBOSE0 RANDOMSTATENONE COPYXTRUE NJOBSNONE AL
GORITHM'AUTO'
FITSELFXYNONE SAMPLEWEIGHTNONE
COMPUTE KMEANS CLUSTERING
PARAMETERS
XARRAYLIKE OR SPARSE MATRIX SHAPENSAMPLES NFEATURES TRAINING INSTANCES TO CLUSTER
IT MUST BE NOTED THAT THE DATA WILL BE CONVERTED TO C ORDERING WHICH WILL CAUSE A MEMORY
COPY IF THE GIVEN DATA IS NOT CCONTIGUOUS
YIGNORED NOT USED PRESENT HERE FOR API CONSISTENCY BY CONVENTION
SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL THE WEIGHTS FOR EACH OBSERVATION
IN X IF NONE ALL OBSERVATIONS ARE ASSIGNED EQUAL WEIGHT DEFAULT NONE
FITPREDICT SELFXYNONE SAMPLEWEIGHTNONE
COMPUTE CLUSTER CENTERS AND PREDICT CLUSTER INDEX FOR EACH SAMPLE
CONVENIENCE METHOD EQUIVALENT TO CALLING FITX FOLLOWED BY PREDICTX
PARAMETERS
XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES NEW DATA TO TRANSFORM
YIGNORED NOT USED PRESENT HERE FOR API CONSISTENCY BY CONVENTION
63SKLEARNCLUSTER CLUSTERING 1493
```

SCIKITLEARN USER GUIDE RELEASE 0213

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL THE WEIGHTS FOR EACH OBSERVATION  
IN X IF NONE ALL OBSERVATIONS ARE ASSIGNED EQUAL WEIGHT DEFAULT NONE

RETURNS

LABELS ARRAY SHAPE NSAMPLES INDEX OF THE CLUSTER EACH SAMPLE BELONGS TO

FITTRANSFORM SELFXYNONE SAMPLEWEIGHTNONE

COMPUTE CLUSTERING AND TRANSFORM X TO CLUSTERDISTANCE SPACE  
EQUIVALENT TO FITXTRANSFORMX BUT MORE EFFICIENTLY IMPLEMENTED

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES NEW DATA TO TRANSFORM  
YIGNORED NOT USED PRESENT HERE FOR API CONSISTENCY BY CONVENTION

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL THE WEIGHTS FOR EACH OBSERVATION  
IN X IF NONE ALL OBSERVATIONS ARE ASSIGNED EQUAL WEIGHT DEFAULT NONE

RETURNS

XNEW ARRAY SHAPE NSAMPLES K X TRANSFORMED IN THE NEW SPACE

GETPARAMS SELFDEEPTTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELFXSAMPLEWEIGHTNONE

PREDICT THE CLOSEST CLUSTER EACH SAMPLE IN X BELONGS TO  
IN THE VECTOR QUANTIZATION LITERATURE CLUSTERCENTERS IS CALLED THE CODE BOOK AND EACH VALUE RE  
TURNED BYPREDICT IS THE INDEX OF THE CLOSEST CODE IN THE CODE BOOK

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES NEW DATA TO PREDICT

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL THE WEIGHTS FOR EACH OBSERVATION  
IN X IF NONE ALL OBSERVATIONS ARE ASSIGNED EQUAL WEIGHT DEFAULT NONE

RETURNS

LABELS ARRAY SHAPE NSAMPLES INDEX OF THE CLUSTER EACH SAMPLE BELONGS TO

SCORESELFXYNONE SAMPLEWEIGHTNONE

OPPOSITE OF THE VALUE OF X ON THE KMEANS OBJECTIVE

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES NEW DATA  
YIGNORED NOT USED PRESENT HERE FOR API CONSISTENCY BY CONVENTION

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL THE WEIGHTS FOR EACH OBSERVATION  
IN X IF NONE ALL OBSERVATIONS ARE ASSIGNED EQUAL WEIGHT DEFAULT NONE

RETURNS

1494 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

SCORE FLOAT OPPOSITE OF THE VALUE OF X ON THE KMEANS OBJECTIVE

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELF X

TRANSFORM X TO A CLUSTERDISTANCE SPACE

IN THE NEW SPACE EACH DIMENSION IS THE DISTANCE TO THE CLUSTER CENTERS NOTE THAT EVEN IF X IS SPARSE THE ARRAY RETURNED BY TRANSFORM WILL TYPICALLY BE DENSE

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES NEW DATA TO TRANSFORM

RETURNS

XNEW ARRAY SHAPE NSAMPLES K X TRANSFORMED IN THE NEW SPACE

EXAMPLES USING SKLEARNCLUSTERKMEANS

- DEMONSTRATION OF KMEANS ASSUMPTIONS
- VECTOR QUANTIZATION EXAMPLE
- KMEANS CLUSTERING
- COLOR QUANTIZATION USING KMEANS
- EMPIRICAL EVALUATION OF THE IMPACT OF KMEANS INITIALIZATION
- COMPARISON OF THE KMEANS AND MINIBATCHKMEANS CLUSTERING ALGORITHMS
- A DEMO OF KMEANS CLUSTERING ON THE HANDWRITTEN DIGITS DATA
- SELECTING THE NUMBER OF CLUSTERS WITH SILHOUETTE ANALYSIS ON KMEANS CLUSTERING
- CLUSTERING TEXT DOCUMENTS USING KMEANS

SKLEARNCLUSTER MINIBATCHKMEANS

CLASSSSKLEARNCLUSTER MINIBATCHKMEANS NCLUSTERS8 INIT'KMEANS' MAXITER100

BATCHSIZE100 VERBOSE0 COMPUTELABELSTRUE

RANDOMSTATENONE TOL00 MAXNOIMPROVEMENT10

INITSIZENONE NINIT3 REASSIGNMENTRATIO001

MINIBATCH KMEANS CLUSTERING

READ MORE IN THE USER GUIDE

PARAMETERS

NCLUSTERS INT OPTIONAL DEFAULT 8 THE NUMBER OF CLUSTERS TO FORM AS WELL AS THE NUMBER OF CENTROIDS TO GENERATE

63SKLEARNCLUSTER CLUSTERING 1495

SCIKITLEARN USER GUIDE RELEASE 0213

INIT ‘KMEANS’ ‘RANDOM’ OR AN NDARRAY DEFAULT ‘KMEANS’ METHOD FOR INITIALIZATION  
DEFAULTS TO ‘KMEANS’

‘KMEANS’ SELECTS INITIAL CLUSTER CENTERS FOR KMEAN CLUSTERING IN A SMART WAY TO SPEED  
UP CONVERGENCE SEE SECTION NOTES IN KINIT FOR MORE DETAILS

‘RANDOM’ CHOOSE K OBSERVATIONS ROWS AT RANDOM FROM DATA FOR THE INITIAL CENTROIDS  
IF AN NDARRAY IS PASSED IT SHOULD BE OF SHAPE NCLUSTERS NFEATURES AND GIVES THE INITIAL  
CENTERS

MAXITER INT OPTIONAL MAXIMUM NUMBER OF ITERATIONS OVER THE COMPLETE DATASET BEFORE STOP  
PING INDEPENDENTLY OF ANY EARLY STOPPING CRITERION HEURISTICS

BATCHSIZE INT OPTIONAL DEFAULT 100 SIZE OF THE MINI BATCHES

VERBOSE BOOLEAN OPTIONAL VERBOSITY MODE

COMPUTELABELS BOOLEAN DEFAULTTRUE COMPUTE LABEL ASSIGNMENT AND INERTIA FOR THE COM  
PLETE DATASET ONCE THE MINIBATCH OPTIMIZATION HAS CONVERGED IN FIT

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE DEFAULT DETERMINES RANDOM NUMBER GEN  
ERATION FOR CENTROID INITIALIZATION AND RANDOM REASSIGNMENT USE AN INT TO MAKE THE RANDOM  
NESS DETERMINISTIC SEE GLOSSARY

TOLFLOAT DEFAULT 00 CONTROL EARLY STOPPING BASED ON THE RELATIVE CENTER CHANGES AS MEASURED  
BY A SMOOTHED VARIANCENORMALIZED OF THE MEAN CENTER SQUARED POSITION CHANGES THIS  
EARLY STOPPING HEURISTICS IS CLOSER TO THE ONE USED FOR THE BATCH VARIANT OF THE ALGORITHMS BUT  
INDUCES A SLIGHT COMPUTATIONAL AND MEMORY OVERHEAD OVER THE INERTIA HEURISTIC

TO DISABLE CONVERGENCE DETECTION BASED ON NORMALIZED CENTER CHANGE SET TOL TO 00 DEFAULT  
MAXNOIMPROVEMENT INT DEFAULT 10 CONTROL EARLY STOPPING BASED ON THE CONSECUTIVE NUM  
BER OF MINI BATCHES THAT DOES NOT YIELD AN IMPROVEMENT ON THE SMOOTHED INERTIA

TO DISABLE CONVERGENCE DETECTION BASED ON INERTIA SET MAXNOIMPROVEMENT TO NONE  
INITSIZE INT OPTIONAL DEFAULT 3 BATCHSIZE NUMBER OF SAMPLES TO RANDOMLY SAMPLE FOR  
SPEEDING UP THE INITIALIZATION SOMETIMES AT THE EXPENSE OF ACCURACY THE ONLY ALGORITHM  
IS INITIALIZED BY RUNNING A BATCH KMEANS ON A RANDOM SUBSET OF THE DATA THIS NEEDS TO BE  
LARGER THAN NCLUSTERS

NINIT INT DEFAULT3 NUMBER OF RANDOM INITIALIZATIONS THAT ARE TRIED IN CONTRAST TO KMEANS  
THE ALGORITHM IS ONLY RUN ONCE USING THE BEST OF THE NINIT INITIALIZATIONS AS MEASURED BY  
INERTIA

REASSIGNMENTRATIO FLOAT DEFAULT 001 CONTROL THE FRACTION OF THE MAXIMUM NUMBER OF  
COUNTS FOR A CENTER TO BE REASSIGNED A HIGHER VALUE MEANS THAT LOW COUNT CENTERS ARE MORE  
EASILY REASSIGNED WHICH MEANS THAT THE MODEL WILL TAKE LONGER TO CONVERGE BUT SHOULD CON  
VERGE IN A BETTER CLUSTERING

ATTRIBUTES

CLUSTERCENTERS ARRAY NCLUSTERS NFEATURES COORDINATES OF CLUSTER CENTERS

LABELS LABELS OF EACH POINT IF COMPUTELABELS IS SET TO TRUE

INERTIA FLOAT THE VALUE OF THE INERTIA CRITERION ASSOCIATED WITH THE CHOSEN PARTITION IF COM  
PUTELABELS IS SET TO TRUE THE INERTIA IS DEFINED AS THE SUM OF SQUARE DISTANCES OF SAMPLES  
TO THEIR NEAREST NEIGHBOR

SEE ALSO

1496 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

KMEANS THE CLASSIC IMPLEMENTATION OF THE CLUSTERING METHOD BASED ON THE LLOYD’S ALGORITHM IT CONSUMES THE WHOLE SET OF INPUT DATA AT EACH ITERATION

NOTES

SEE [HTTPSWWWEECSTUFTSEDUDSCULLEYPAPERSFASTKMEANS.PDF](https://www.eecstufsedudsculleypapers/fastkmeans.pdf)

EXAMPLES

```
FROM SKLEARNCLUSTER IMPORT MINIBATCHKMEANS
IMPORT NUMPY AS NP
X = NPARRAY(1 2 1 4 1 0
4 2 4 0 4 4
4 5 0 1 2 2
3 2 5 5 1 1)
MANUALLY FIT ON BATCHES
KMEANS = MINIBATCHKMEANS(NCLUSTERS=2
RANDOMSTATE=0
BATCHSIZE=6)
KMEANS.fit(X, KMEANS.PARTIALFIT(X, 0.6))
KMEANS = KMEANS.PARTIALFIT(X, 0.6)
KMEANS.CLUSTERCENTERS
ARRAY([1
3 4])
KMEANS.PREDICT(0 0 4 4)
ARRAY([1 DTYPE=INT32])
FIT ON THE WHOLE DATA
KMEANS = MINIBATCHKMEANS(NCLUSTERS=2
RANDOMSTATE=0
BATCHSIZE=6
MAXITER=10)
KMEANS.CLUSTERCENTERS
ARRAY([395918367 240816327
112195122 13902439])
KMEANS.PREDICT(0 0 4 4)
ARRAY([1 DTYPE=INT32])
```

METHODS

`fit(self, X, Y, sample_weight)` COMPUTE THE CENTROIDS ON X BY CHUNKING IT INTO MINI BATCHES

`fit_predict(self, X, Y, sample_weight)` COMPUTE CLUSTER CENTERS AND PREDICT CLUSTER INDEX FOR EACH SAMPLE

`fit_transform(self, X, Y, sample_weight)` COMPUTE CLUSTERING AND TRANSFORM X TO CLUSTERDISTANCE SPACE

`get_params(self, deep=True)` GET PARAMETERS FOR THIS ESTIMATOR

`partial_fit(self, X, Y, sample_weight)` UPDATE K MEANS ESTIMATE ON A SINGLE MINIBATCH X

`predict(self, X, sample_weight)` PREDICT THE CLOSEST CLUSTER EACH SAMPLE IN X BELONGS TO

`score(self, X, Y, sample_weight)` OPPOSITE OF THE VALUE OF X ON THE KMEANS OBJECTIVE

`set_params(self, **params)` SET THE PARAMETERS OF THIS ESTIMATOR

`transform(self, X)` TRANSFORM X TO A CLUSTERDISTANCE SPACE

63SKLEARNCLUSTER CLUSTERING 1497

SCIKITLEARN USER GUIDE RELEASE 0213

INIT SELFNCCLUSTERS8 INIT'KMEANS' MAXITER100 BATCHSIZE100 VERBOSE0  
COMPUTELABELSTRUE RANDOMSTATENONE TOL00 MAXNOIMPROVEMENT10  
INITSIZEONE NINIT3 REASSIGNMENTRATIO001  
FITSELFXYNONE SAMPLEWEIGHTNONE  
COMPUTE THE CENTROIDS ON X BY CHUNKING IT INTO MINIBATCHES  
PARAMETERS  
XARRAYLIKE OR SPARSE MATRIX SHAPENSAMPLES NFEATURES TRAINING INSTANCES TO CLUSTER  
IT MUST BE NOTED THAT THE DATA WILL BE CONVERTED TO C ORDERING WHICH WILL CAUSE A MEMORY  
COPY IF THE GIVEN DATA IS NOT CCONTIGUOUS  
YIGNORED NOT USED PRESENT HERE FOR API CONSISTENCY BY CONVENTION  
SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL THE WEIGHTS FOR EACH OBSERVATION  
IN X IF NONE ALL OBSERVATIONS ARE ASSIGNED EQUAL WEIGHT DEFAULT NONE  
FITPREDICT SELFXYNONE SAMPLEWEIGHTNONE  
COMPUTE CLUSTER CENTERS AND PREDICT CLUSTER INDEX FOR EACH SAMPLE  
CONVENIENCE METHOD EQUIVALENT TO CALLING FITX FOLLOWED BY PREDICTX  
PARAMETERS  
XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES NEW DATA TO TRANSFORM  
YIGNORED NOT USED PRESENT HERE FOR API CONSISTENCY BY CONVENTION  
SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL THE WEIGHTS FOR EACH OBSERVATION  
IN X IF NONE ALL OBSERVATIONS ARE ASSIGNED EQUAL WEIGHT DEFAULT NONE  
RETURNS  
LABELS ARRAY SHAPE NSAMPLES INDEX OF THE CLUSTER EACH SAMPLE BELONGS TO  
FITTRANSFORM SELFXYNONE SAMPLEWEIGHTNONE  
COMPUTE CLUSTERING AND TRANSFORM X TO CLUSTERDISTANCE SPACE  
EQUIVALENT TO FITXTRANSFORMX BUT MORE EFFICIENTLY IMPLEMENTED  
PARAMETERS  
XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES NEW DATA TO TRANSFORM  
YIGNORED NOT USED PRESENT HERE FOR API CONSISTENCY BY CONVENTION  
SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL THE WEIGHTS FOR EACH OBSERVATION  
IN X IF NONE ALL OBSERVATIONS ARE ASSIGNED EQUAL WEIGHT DEFAULT NONE  
RETURNS  
XNEW ARRAY SHAPE NSAMPLES K X TRANSFORMED IN THE NEW SPACE  
GETPARAMS SELFDEEPTTRUE  
GET PARAMETERS FOR THIS ESTIMATOR  
PARAMETERS  
DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBOBJECTS THAT ARE ESTIMATORS  
RETURNS  
PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES  
PARTIALFIT SELFXYNONE SAMPLEWEIGHTNONE  
UPDATE K MEANS ESTIMATE ON A SINGLE MINIBATCH X

1498 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES COORDINATES OF THE DATA POINTS TO CLUSTER IT  
MUST BE NOTED THAT X WILL BE COPIED IF IT IS NOT CCONTIGUOUS

YIGNORED NOT USED PRESENT HERE FOR API CONSISTENCY BY CONVENTION

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL THE WEIGHTS FOR EACH OBSERVATION

IN X IF NONE ALL OBSERVATIONS ARE ASSIGNED EQUAL WEIGHT DEFAULT NONE

PREDICTSELFXSAMPLEWEIGHTNONE

PREDICT THE CLOSEST CLUSTER EACH SAMPLE IN X BELONGS TO

IN THE VECTOR QUANTIZATION LITERATURE CLUSTERCENTERS IS CALLED THE CODE BOOK AND EACH VALUE RE  
TURNED BYPREDICT IS THE INDEX OF THE CLOSEST CODE IN THE CODE BOOK

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES NEW DATA TO PREDICT

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL THE WEIGHTS FOR EACH OBSERVATION

IN X IF NONE ALL OBSERVATIONS ARE ASSIGNED EQUAL WEIGHT DEFAULT NONE

RETURNS

LABELS ARRAY SHAPE NSAMPLES INDEX OF THE CLUSTER EACH SAMPLE BELONGS TO

SCORESELFXYNONE SAMPLEWEIGHTNONE

OPPOSITE OF THE VALUE OF X ON THE KMEANS OBJECTIVE

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES NEW DATA

YIGNORED NOT USED PRESENT HERE FOR API CONSISTENCY BY CONVENTION

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL THE WEIGHTS FOR EACH OBSERVATION

IN X IF NONE ALL OBSERVATIONS ARE ASSIGNED EQUAL WEIGHT DEFAULT NONE

RETURNS

SCORE FLOAT OPPOSITE OF THE VALUE OF X ON THE KMEANS OBJECTIVE

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE  
PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT

OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELF

TRANSFORM X TO A CLUSTERDISTANCE SPACE

IN THE NEW SPACE EACH DIMENSION IS THE DISTANCE TO THE CLUSTER CENTERS NOTE THAT EVEN IF X IS SPARSE THE  
ARRAY RETURNED BY TRANSFORM WILL TYPICALLY BE DENSE

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES NEW DATA TO TRANSFORM

RETURNS

XNEW ARRAY SHAPE NSAMPLES K X TRANSFORMED IN THE NEW SPACE

63SKLEARNCLUSTER CLUSTERING 1499

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES USING SKLEARNCLUSTERMINIBATCHKMEANS

- BICCLUSTERING DOCUMENTS WITH THE SPECTRAL COCLUSTERING ALGORITHM
- ONLINE LEARNING OF A DICTIONARY OF PARTS OF FACES
- COMPARE BIRCH AND MINIBATCHKMEANS
- EMPIRICAL EVALUATION OF THE IMPACT OF KMEANS INITIALIZATION
- COMPARISON OF THE KMEANS AND MINIBATCHKMEANS CLUSTERING ALGORITHMS
- COMPARING DIFFERENT CLUSTERING ALGORITHMS ON TOY DATASETS
- FACES DATASET DECOMPOSITIONS
- CLUSTERING TEXT DOCUMENTS USING KMEANS

SKLEARNCLUSTER MEANSHIFT

CLASSSSKLEARNCLUSTER MEANSHIFT BANDWIDTHNONE SEEDSNONE BINSEEDINGFALSE

MINBINFREQ1 CLUSTERALLTRUE NJOBSNONE

MEAN SHIFT CLUSTERING USING A FLAT KERNEL

MEAN SHIFT CLUSTERING AIMS TO DISCOVER “BLOBS” IN A SMOOTH DENSITY OF SAMPLES IT IS A CENTROIDBASED ALGORITHM WHICH WORKS BY UPDATING CANDIDATES FOR CENTROIDS TO BE THE MEAN OF THE POINTS WITHIN A GIVEN REGION THESE CANDIDATES ARE THEN FILTERED IN A POSTPROCESSING STAGE TO ELIMINATE NEARDUPPLICATES TO FORM THE FINAL SET OF CENTROIDS

SEEDING IS PERFORMED USING A BINNING TECHNIQUE FOR SCALABILITY

READ MORE IN THE USER GUIDE

PARAMETERS

BANDWIDTH FLOAT OPTIONAL BANDWIDTH USED IN THE RBF KERNEL

IF NOT GIVEN THE BANDWIDTH IS ESTIMATED USING SKLEARNCLUSTERESTIMATEBANDWIDTH SEE THE DOCUMENTATION FOR THAT FUNCTION FOR HINTS ON SCALABILITY SEE ALSO THE NOTES BELOW

SEEDS ARRAY SHAPENSAMPLES NFEATURES OPTIONAL SEEDS USED TO INITIALIZE KERNELS IF NOT SET THE SEEDS ARE CALCULATED BY CLUSTERINGGETBINSEEDS WITH BANDWIDTH AS THE GRID SIZE AND DEFAULT VALUES FOR OTHER PARAMETERS

BINSEEDING BOOLEAN OPTIONAL IF TRUE INITIAL KERNEL LOCATIONS ARE NOT LOCATIONS OF ALL POINTS BUT RATHER THE LOCATION OF THE DISCRETIZED VERSION OF POINTS WHERE POINTS ARE BINNED ONTO A GRID WHOSE COARSENESS CORRESPONDS TO THE BANDWIDTH SETTING THIS OPTION TO TRUE WILL SPEED UP THE ALGORITHM BECAUSE FEWER SEEDS WILL BE INITIALIZED DEFAULT VALUE FALSE IGNORED IF SEEDS ARGUMENT IS NOT NONE

MINBINFREQ INT OPTIONAL TO SPEED UP THE ALGORITHM ACCEPT ONLY THOSE BINS WITH AT LEAST MINBINFREQ POINTS AS SEEDS IF NOT DEFINED SET TO 1

CLUSTERALL BOOLEAN DEFAULT TRUE IF TRUE THEN ALL POINTS ARE CLUSTERED EVEN THOSE ORPHANS THAT ARE NOT WITHIN ANY KERNEL ORPHANS ARE ASSIGNED TO THE NEAREST KERNEL IF FALSE THEN ORPHANS ARE GIVEN CLUSTER LABEL 1

NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF JOBS TO USE FOR THE COMPUTATION THIS WORKS BY COMPUTING EACH OF THE NINIT RUNS IN PARALLEL

NONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

1500 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

ATTRIBUTES

CLUSTERCENTERS ARRAY NCLUSTERS NFEATURES COORDINATES OF CLUSTER CENTERS

LABELS LABELS OF EACH POINT

NOTES

SCALABILITY

BECAUSE THIS IMPLEMENTATION USES A FLAT KERNEL AND A BALL TREE TO LOOK UP MEMBERS OF EACH KERNEL THE COMPLEXITY WILL TEND TOWARDS OTNLOGN IN LOWER DIMENSIONS WITH N THE NUMBER OF SAMPLES AND T THE NUMBER OF POINTS IN HIGHER DIMENSIONS THE COMPLEXITY WILL TEND TOWARDS OTN2

SCALABILITY CAN BE BOOSTED BY USING FEWER SEEDS FOR EXAMPLE BY USING A HIGHER VALUE OF MINBINFREQ IN THE GETBINSEEDS FUNCTION

NOTE THAT THE ESTIMATEBANDWIDTH FUNCTION IS MUCH LESS SCALABLE THAN THE MEAN SHIFT ALGORITHM AND WILL BE THE BOTTLENECK IF IT IS USED

REFERENCES

DORIN COMANICIU AND PETER MEER “MEAN SHIFT A ROBUST APPROACH TOWARD FEATURE SPACE ANALYSIS” IEEE TRANS ACTIONS ON PATTERN ANALYSIS AND MACHINE INTELLIGENCE 2002 PP 603619

EXAMPLES

```
FROM SKLEARNCLUSTER IMPORT MEANSHIFT
IMPORT NUMPY AS NP
X NPARRAY1 1 2 1 1 0
  4 7 3 5 3 6
CLUSTERING MEANSHIFTBANDWIDTH2FITX
CLUSTERINGLABELS
ARRAY1 1 1 0 0 0
CLUSTERINGPREDICT0 0 5 5
ARRAY1 0
CLUSTERING
MEANSHIFTBANDWIDTH2 BINSEEDINGFALSE CLUSTERALLTRUE MINBINFREQ1
NJOBSNONE SEEDSNONE
METHODS
FITSELF X Y PERFORM CLUSTERING
FITPREDICT SELF X Y PERFORMS CLUSTERING ON X AND RETURNS CLUSTER LABELS
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR
PREDICT SELF X PREDICT THE CLOSEST CLUSTER EACH SAMPLE IN X BELONGS TO
SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR
INIT SELF BANDWIDTHNONE SEEDSNONE BINSEEDINGFALSE MINBINFREQ1 CLUSTERALLTRUE
NJOBSNONE
FITSELFXYNONE
63SKLEARNCLUSTER CLUSTERING 1501
```

SCIKITLEARN USER GUIDE RELEASE 0213

PERFORM CLUSTERING

PARAMETERS

XARRAYLIKE SHAPENSAMPLES NFEATURES SAMPLES TO CLUSTER

YIGNORED

FITPREDICT SELFXYNONE

PERFORMS CLUSTERING ON X AND RETURNS CLUSTER LABELS

PARAMETERS

XNDARRAY SHAPE NSAMPLES NFEATURES INPUT DATA

YIGNORED NOT USED PRESENT FOR API CONSISTENCY BY CONVENTION

RETURNS

LABELS NDARRAY SHAPE NSAMPLES CLUSTER LABELS

GETPARAMS SELFDEEPTTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELF

PREDICT THE CLOSEST CLUSTER EACH SAMPLE IN X BELONGS TO

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPENSAMPLES NFEATURES NEW DATA TO PREDICT

RETURNS

LABELS ARRAY SHAPE NSAMPLES INDEX OF THE CLUSTER EACH SAMPLE BELONGS TO

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE

PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT

OF A NESTED OBJECT

RETURNS

SELF

EXAMPLES USING SKLEARNCLUSTERMEANSHIFT

- A DEMO OF THE MEANSHIFT CLUSTERING ALGORITHM
- COMPARING DIFFERENT CLUSTERING ALGORITHMS ON TOY DATASETS

1502 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

SKLEARNCLUSTER SPECTRALCLUSTERING

CLASSSSKLEARNCLUSTER SPECTRALCLUSTERING NCLUSTERS8 EIGENSOLVERNONE RAN

DOMSTATENONE NINIT10 GAMMA10 AFFIN

ITY'RBF' NNEIGHBORS10 EIGENTOL00 AS

SIGNLABELS'KMEANS' DEGREE3 COEF01 KER

NELPARAMSNONE NJOBSNONE

APPLY CLUSTERING TO A PROJECTION OF THE NORMALIZED LAPLACIAN

IN PRACTICE SPECTRAL CLUSTERING IS VERY USEFUL WHEN THE STRUCTURE OF THE INDIVIDUAL CLUSTERS IS HIGHLY NONCONVEX

OR MORE GENERALLY WHEN A MEASURE OF THE CENTER AND SPREAD OF THE CLUSTER IS NOT A SUITABLE DESCRIPTION OF THE

COMPLETE CLUSTER FOR INSTANCE WHEN CLUSTERS ARE NESTED CIRCLES ON THE 2D PLANE

IF AFFINITY IS THE ADJACENCY MATRIX OF A GRAPH THIS METHOD CAN BE USED TO FIND NORMALIZED GRAPH CUTS

WHEN CALLING FIT AN AFFINITY MATRIX IS CONSTRUCTED USING EITHER KERNEL FUNCTION SUCH THE GAUSSIAN AKA RBF

KERNEL OF THE EUCLIDEAN DISTANCED DX X

NPEXP GAMMA DXX2

OR A KNEAREST NEIGHBORS CONNECTIVITY MATRIX

ALTERNATIVELY USING PRECOMPUTED A USERPROVIDED AFFINITY MATRIX CAN BE USED

READ MORE IN THE USER GUIDE

PARAMETERS

NCLUSTERS INTEGER OPTIONAL THE DIMENSION OF THE PROJECTION SUBSPACE

EIGENSOLVER NONE 'ARPACK' 'LOBPCG' OR 'AMG' THE EIGENVALUE DECOMPOSITION STRATEGY TO

USE AMG REQUIRES PYAMG TO BE INSTALLED IT CAN BE FASTER ON VERY LARGE SPARSE PROBLEMS

BUT MAY ALSO LEAD TO INSTABILITIES

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE DEFAULT A PSEUDO RANDOM NUMBER

GENERATOR USED FOR THE INITIALIZATION OF THE LOBPCG EIGEN VECTORS DECOMPOSITION WHEN

EIGENSOLVERAMG AND BY THE KMEANS INITIALIZATION USE AN INT TO MAKE THE RAN

DOMNESS DETERMINISTIC SEE GLOSSARY

NINIT INT OPTIONAL DEFAULT 10 NUMBER OF TIME THE KMEANS ALGORITHM WILL BE RUN WITH DIF

FERENT CENTROID SEEDS THE FINAL RESULTS WILL BE THE BEST OUTPUT OF NINIT CONSECUTIVE RUNS IN

TERMS OF INERTIA

GAMMA FLOAT DEFAULT10 KERNEL COEFFICIENT FOR RBF POLY SIGMOID LAPLACIAN AND CHI2 KERNELS

IGNORED FOR AFFINITYNEARESTNEIGHBORS

AFFINITY STRING ARRAYLIKE OR CALLABLE DEFAULT 'RBF' IF A STRING THIS MAY BE ONE OF 'NEAR

ESTNEIGHBORS' 'PRECOMPUTED' 'RBF' OR ONE OF THE KERNELS SUPPORTED BY SKLEARN

METRICSPAIRWISEKERNELS

ONLY KERNELS THAT PRODUCE SIMILARITY SCORES NONNEGATIVE VALUES THAT INCREASE WITH SIMILAR

ITY SHOULD BE USED THIS PROPERTY IS NOT CHECKED BY THE CLUSTERING ALGORITHM

NNEIGHBORS INTEGER NUMBER OF NEIGHBORS TO USE WHEN CONSTRUCTING THE AFFINITY MATRIX USING

THE NEAREST NEIGHBORS METHOD IGNORED FOR AFFINITYRBF

EIGENTOL FLOAT OPTIONAL DEFAULT 00 STOPPING CRITERION FOR EIGENDECOMPOSITION OF THE LAPLA

CIAN MATRIX WHEN EIGENSOLVERARPACK

ASSIGNLABELS 'KMEANS' 'DISCRETIZE' DEFAULT 'KMEANS' THE STRATEGY TO USE TO ASSIGN LABELS

IN THE EMBEDDING SPACE THERE ARE TWO WAYS TO ASSIGN LABELS AFTER THE LAPLACIAN EMBEDDING

63SKLEARNCLUSTER CLUSTERING 1503

SCIKITLEARN USER GUIDE RELEASE 0213

KMEANS CAN BE APPLIED AND IS A POPULAR CHOICE BUT IT CAN ALSO BE SENSITIVE TO INITIALIZATION  
DISCRETIZATION IS ANOTHER APPROACH WHICH IS LESS SENSITIVE TO RANDOM INITIALIZATION  
DEGREE FLOAT DEFAULT3 DEGREE OF THE POLYNOMIAL KERNEL IGNORED BY OTHER KERNELS  
COEF0 FLOAT DEFAULT1 ZERO COEFFICIENT FOR POLYNOMIAL AND SIGMOID KERNELS IGNORED BY OTHER  
KERNELS

KERNELPARAMS DICTIONARY OF STRING TO ANY OPTIONAL PARAMETERS KEYWORD ARGUMENTS AND  
VALUES FOR KERNEL PASSED AS CALLABLE OBJECT IGNORED BY OTHER KERNELS

NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF PARALLEL JOBS TO RUN NONE MEANS

1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL PROCESSORS SEE

GLOSSARY FOR MORE DETAILS

ATTRIBUTES

AFFINITYMATRIX ARRAYLIKE SHAPE NSAMPLES NSAMPLES AFFINITY MATRIX USED FOR CLUSTER

ING AVAILABLE ONLY IF AFTER CALLING FIT

LABELS LABELS OF EACH POINT

NOTES

IF YOU HAVE AN AFFINITY MATRIX SUCH AS A DISTANCE MATRIX FOR WHICH 0 MEANS IDENTICAL ELEMENTS AND HIGH VALUES  
MEANS VERY DISSIMILAR ELEMENTS IT CAN BE TRANSFORMED IN A SIMILARITY MATRIX THAT IS WELL SUITED FOR THE ALGORITHM  
BY APPLYING THE GAUSSIAN RBF HEAT KERNEL

NPEXP DISTMATRIX 2 2 DELTA2

WHEREDELTA IS A FREE PARAMETER REPRESENTING THE WIDTH OF THE GAUSSIAN KERNEL

ANOTHER ALTERNATIVE IS TO TAKE A SYMMETRIC VERSION OF THE K NEAREST NEIGHBORS CONNECTIVITY MATRIX OF THE POINTS

IF THE PYAMG PACKAGE IS INSTALLED IT IS USED THIS GREATLY SPEEDS UP COMPUTATION

REFERENCES

• NORMALIZED CUTS AND IMAGE SEGMENTATION 2000 JIANBO SHI JITENDRA MALIK HTTPCITESEERISTPSUEDU  
VIEWDOCSUMMARYDOI10111602324

• A TUTORIAL ON SPECTRAL CLUSTERING 2007 ULRIKE VON LUXBURG HTTPCITESEERXISTPSUEDUVIEWDOC  
SUMMARYDOI10111659323

• MULTICLASS SPECTRAL CLUSTERING 2003 STELLA X YU JIANBO SHI HTTPSWWW1ICSIBERKELEYEDUSTELLAYU  
PUBLICATIONDOC2003KWAYICCV PDF

EXAMPLES

FROM SKLEARNCLUSTER IMPORT SPECTRALCLUSTERING

IMPORT NUMPY AS NP

X NPARRAY1 1 2 1 1 0

4 7 3 5 3 6

CLUSTERING SPECTRALCLUSTERINGNCLUSTERS2

ASSIGNLABELSDISCRETIZE

RANDOMSTATE0FITX

CLUSTERINGLABELS

1504 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

ARRAY1 1 1 0 0 0

CLUSTERING

SPECTRALCLUSTERINGAFFINITYRBF ASSIGNLABELSDISCRETIZE COEF01

DEGREE3 EIGENSOLVERNONE EIGENTOL00 GAMMA10

KERNELPARAMSNONE NCLUSTERS2 NINIT10 NJOBSNONE

NNEIGHBORS10 RANDOMSTATE0

METHODS

FITSELF X Y CREATES AN AFFINITY MATRIX FOR X USING THE SELECTED AFFINITY THEN APPLIES SPECTRAL CLUSTERING TO THIS AFFINITY MATRIX

FITPREDICT SELF X Y PERFORMS CLUSTERING ON X AND RETURNS CLUSTER LABELS

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELFNCLUSTERS8 EIGENSOLVERNONE RANDOMSTATENONE NINIT10 GAMMA10 AFFINITY'RBF' NNEIGHBORS10 EIGENTOL00 ASSIGNLABELS'KMEANS' DEGREE3 COEF01

KERNELPARAMSNONE NJOBSNONE

FITSELFXYNONE

CREATES AN AFFINITY MATRIX FOR X USING THE SELECTED AFFINITY THEN APPLIES SPECTRAL CLUSTERING TO THIS AFFINITY MATRIX

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES OR IF AFFINITY'PRECOMPUTED' A PRECOMPUTED AFFINITY MATRIX OF SHAPE NSAMPLES NSAMPLES

YIGNORED

FITPREDICT SELFXYNONE

PERFORMS CLUSTERING ON X AND RETURNS CLUSTER LABELS

PARAMETERS

XNDARRAY SHAPE NSAMPLES NFEATURES INPUT DATA

YIGNORED NOT USED PRESENT FOR API CONSISTENCY BY CONVENTION

RETURNS

LABELS NDARRAY SHAPE NSAMPLES CLUSTER LABELS

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

63SKLEARNCLUSTER CLUSTERING 1505

SCIKITLEARN USER GUIDE RELEASE 0213

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS  
SELF

EXAMPLES USING SKLEARNCLUSTERSPECTRALCLUSTERING

- COMPARING DIFFERENT CLUSTERING ALGORITHMS ON TOY DATASETS

632 FUNCTIONS

CLUSTERAFFINITYPROPAGATION S    PERFORM AFFINITY PROPAGATION CLUSTERING OF DATA

CLUSTERCLUSTEROPTICSDBCAN REACHABILITY  
    PERFORMS DBSCAN EXTRACTION FOR AN ARBITRARY EPSILON

CLUSTERCLUSTEROPTICSXI REACHABILITY    AUTOMATICALLY EXTRACT CLUSTERS ACCORDING TO THE XISTEEP METHOD

CLUSTERCOMPUTEOPTICSGRAPH X  
MINSAMPLES    COMPUTES THE OPTICS REACHABILITY GRAPH

CLUSTERDBSCAN X EPS MINSAMPLES    PERFORM DBSCAN CLUSTERING FROM VECTOR ARRAY OR DISTANCE MATRIX

CLUSTERESTIMATEBANDWIDTH X QUANTILE    ESTIMATE THE BANDWIDTH TO USE WITH THE MEANSHIFT ALGORITHM

CLUSTERKMEANS X NCLUSTERS    KMEANS CLUSTERING ALGORITHM

CLUSTERMEANSHIFT X BANDWIDTH SEEDS    PERFORM MEAN SHIFT CLUSTERING OF DATA USING A FLAT KERNEL

CLUSTERSPECTRALCLUSTERING AFFINITY    APPLY CLUSTERING TO A PROJECTION OF THE NORMALIZED LAPLACIAN

CLUSTERWARDTREE X CONNECTIVITY    WARD CLUSTERING BASED ON A FEATURE MATRIX

SKLEARNCLUSTER AFFINITYPROPAGATION  
SKLEARNCLUSTER AFFINITYPROPAGATION S PREFERENCENONE    CONVERGENCEITER15  
MAXITER200 DAMPING05 COPYTRUE VER  
BOSEFALSE RETURNNITERFALSE  
PERFORM AFFINITY PROPAGATION CLUSTERING OF DATA  
READ MORE IN THE USER GUIDE  
PARAMETERS

SARRAYLIKE SHAPE NSAMPLES NSAMPLES MATRIX OF SIMILARITIES BETWEEN POINTS

PREFERENCE ARRAYLIKE SHAPE NSAMPLES OR FLOAT OPTIONAL PREFERENCES FOR EACH POINT    POINTS WITH LARGER VALUES OF PREFERENCES ARE MORE LIKELY TO BE CHOSEN AS EXEMPLARS THE NUMBER OF EXEMPLARS IE OF CLUSTERS IS INFLUENCED BY THE INPUT PREFERENCES VALUE IF THE PREFERENCES ARE NOT PASSED AS ARGUMENTS THEY WILL BE SET TO THE MEDIAN OF THE INPUT SIMILARITIES RESULTING IN A MODERATE NUMBER OF CLUSTERS FOR A SMALLER AMOUNT OF CLUSTERS THIS CAN BE SET TO THE MINIMUM VALUE OF THE SIMILARITIES

CONVERGENCEITER INT OPTIONAL DEFAULT 15 NUMBER OF ITERATIONS WITH NO CHANGE IN THE NUMBER OF ESTIMATED CLUSTERS THAT STOPS THE CONVERGENCE

MAXITER INT OPTIONAL DEFAULT 200 MAXIMUM NUMBER OF ITERATIONS

1506 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

DAMPING FLOAT OPTIONAL DEFAULT 0.5 DAMPING FACTOR BETWEEN 0.5 AND 1  
COPY BOOLEAN OPTIONAL DEFAULT TRUE IF COPY IS FALSE THE AFFINITY MATRIX IS MODIFIED INPLACE  
BY THE ALGORITHM FOR MEMORY EFFICIENCY  
VERBOSE BOOLEAN OPTIONAL DEFAULT FALSE THE VERBOSITY LEVEL  
RETURN\_NITER BOOL DEFAULT FALSE WHETHER OR NOT TO RETURN THE NUMBER OF ITERATIONS

RETURNS  
CLUSTER\_CENTERS INDICES ARRAY SHAPE NCLUSTERS INDEX OF CLUSTERS CENTERS  
LABELS ARRAY SHAPE NSAMPLES CLUSTER LABELS FOR EACH POINT  
NITER INT NUMBER OF ITERATIONS RUN RETURNED ONLY IF RETURN\_NITER IS SET TO TRUE

NOTES  
FOR AN EXAMPLE SEE EXAMPLES/CLUSTER\_PLOT\_AFFINITY\_PROPAGATION.PY  
WHEN THE ALGORITHM DOES NOT CONVERGE IT RETURNS AN EMPTY ARRAY AS CLUSTER\_CENTERS INDICES AND 1 AS  
LABEL FOR EACH TRAINING SAMPLE  
WHEN ALL TRAINING SAMPLES HAVE EQUAL SIMILARITIES AND EQUAL PREFERENCES THE ASSIGNMENT OF CLUSTER CENTERS AND  
LABELS DEPENDS ON THE PREFERENCE IF THE PREFERENCE IS SMALLER THAN THE SIMILARITIES A SINGLE CLUSTER CENTER AND  
LABEL FOR EVERY SAMPLE WILL BE RETURNED OTHERWISE EVERY TRAINING SAMPLE BECOMES ITS OWN CLUSTER CENTER AND  
IS ASSIGNED A UNIQUE LABEL

REFERENCES  
BRENDAN J FREY AND DELBERT DUECK “CLUSTERING BY PASSING MESSAGES BETWEEN DATA POINTS” SCIENCE FEB 2007  
EXAMPLES USING SKLEARN/CLUSTER/AFFINITY\_PROPAGATION

•VISUALIZING THE STOCK MARKET STRUCTURE  
SKLEARN/CLUSTER/CLUSTER\_OPTICS\_DBSCAN  
SKLEARN/CLUSTER/CLUSTER\_OPTICS\_DBSCAN/REACHABILITY\_CORE\_DISTANCES\_ORDERING\_EPS  
PERFORMS DBSCAN EXTRACTION FOR AN ARBITRARY EPSILON  
EXTRACTING THE CLUSTERS RUNS IN LINEAR TIME NOTE THAT THIS RESULTS IN LABELS WHICH ARE CLOSE TO A DBSCAN  
WITH SIMILAR SETTINGS AND EPS ONLY IF EPS IS CLOSE TO MAX\_EPS

PARAMETERS  
REACHABILITY ARRAY SHAPE NSAMPLES REACHABILITY DISTANCES CALCULATED BY OPTICS  
REACHABILITY  
CORE\_DISTANCES ARRAY SHAPE NSAMPLES DISTANCES AT WHICH POINTS BECOME CORE  
CORE\_DISTANCES  
ORDERING ARRAY SHAPE NSAMPLES OPTICS ORDERED POINT INDICES ORDERING  
EPS FLOAT DBSCAN EPS PARAMETER MUST BE SET TO MAX\_EPS RESULTS WILL BE CLOSE TO  
DBSCAN ALGORITHM IF EPS AND MAX\_EPS ARE CLOSE TO ONE ANOTHER  
63 SKLEARN/CLUSTER CLUSTERING 1507

SCIKITLEARN USER GUIDE RELEASE 0213

RETURNS

LABELS ARRAY SHAPE NSAMPLES THE ESTIMATED LABELS

EXAMPLES USING SKLEARNCLUSTERCLUSTEROPTICSDBSCAN

•DEMO OF OPTICS CLUSTERING ALGORITHM

SKLEARNCLUSTER CLUSTEROPTICSXI

SKLEARNCLUSTER CLUSTEROPTICSXI REACHABILITY PREDECESSOR ORDERING MINSAMPLES

MINCLUSTERSIZENONE XI005 PREDECES

SORCORRECTIONTRUE

AUTOMATICALLY EXTRACT CLUSTERS ACCORDING TO THE XISTEEP METHOD

PARAMETERS

REACHABILITY ARRAY SHAPE NSAMPLES REACHABILITY DISTANCES CALCULATED BY OPTICS

REACHABILITY

PREDECESSOR ARRAY SHAPE NSAMPLES PREDECESSORS CALCULATED BY OPTICS

ORDERING ARRAY SHAPE NSAMPLES OPTICS ORDERED POINT INDICES ORDERING

MINSAMPLES INT 1 OR FLOAT BETWEEN 0 AND 1 THE SAME AS THE MINSAMPLES GIVEN TO OPTICS

UP AND DOWN STEEP REGIONS CAN'T HAVE MORE THEN MINSAMPLES CONSECUTIVE NONSTEEP

POINTS EXPRESSED AS AN ABSOLUTE NUMBER OR A FRACTION OF THE NUMBER OF SAMPLES ROUNDED TO

BE AT LEAST 2

MINCLUSTERSIZE INT 1 OR FLOAT BETWEEN 0 AND 1 DEFAULTNONE MINIMUM NUMBER OF SAM

PLES IN AN OPTICS CLUSTER EXPRESSED AS AN ABSOLUTE NUMBER OR A FRACTION OF THE NUMBER OF

SAMPLES ROUNDED TO BE AT LEAST 2 IF NONE THE VALUE OF MINSAMPLES IS USED INSTEAD

XIFLOAT BETWEEN 0 AND 1 OPTIONAL DEFAULT005 DETERMINES THE MINIMUM STEEPNESS ON THE

REACHABILITY PLOT THAT CONSTITUTES A CLUSTER BOUNDARY FOR EXAMPLE AN UPWARDS POINT IN THE

REACHABILITY PLOT IS DEFINED BY THE RATIO FROM ONE POINT TO ITS SUCCESSOR BEING AT MOST 1XI

PREDECESSORCORRECTION BOOL OPTIONAL DEFAULTTRUE CORRECT CLUSTERS BASED ON THE CALCU

LATED PREDECESSORS

RETURNS

LABELS ARRAY SHAPE NSAMPLES THE LABELS ASSIGNED TO SAMPLES POINTS WHICH ARE NOT INCLUDED

IN ANY CLUSTER ARE LABELED AS 1

CLUSTERS ARRAY SHAPE NCLUSTERS 2 THE LIST OF CLUSTERS IN THE FORM OF START END

IN EACH ROW WITH ALL INDICES INCLUSIVE THE CLUSTERS ARE ORDERED ACCORDING TO END

START ASCENDING SO THAT LARGER CLUSTERS ENCOMPASSING SMALLER CLUSTERS COME AF

TER SUCH NESTED SMALLER CLUSTERS SINCE LABELS DOES NOT REFLECT THE HIERARCHY USUALLY

LENCLUSTERS NPUNIQUELABELS

SKLEARNCLUSTER COMPUTEOPTICSGRAPH

SKLEARNCLUSTER COMPUTEOPTICSGRAPH XMINSAMPLES MAXEPS METRIC PMETRICPARAMS AL

GORITHM LEAFSIZE NJOBS

COMPUTES THE OPTICS REACHABILITY GRAPH

READ MORE IN THE USER GUIDE

1508 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

PARAMETERS

XARRAY SHAPE NSAMPLES NFEATURES OR NSAMPLES NSAMPLES IF METRIC'PRECOMPUTED'  
A FEATURE ARRAY OR ARRAY OF DISTANCES BETWEEN SAMPLES IF METRIC'PRECOMPUTED'  
MINSAMPLES INT 1 OR FLOAT BETWEEN 0 AND 1 THE NUMBER OF SAMPLES IN A NEIGHBORHOOD FOR A  
POINT TO BE CONSIDERED AS A CORE POINT EXPRESSED AS AN ABSOLUTE NUMBER OR A FRACTION OF THE  
NUMBER OF SAMPLES ROUNDED TO BE AT LEAST 2  
MAXEPS FLOAT OPTIONAL DEFAULTNPINF THE MAXIMUM DISTANCE BETWEEN TWO SAMPLES FOR ONE  
TO BE CONSIDERED AS IN THE NEIGHBORHOOD OF THE OTHER DEFAULT VALUE OF NPINF WILL IDENTIFY  
CLUSTERS ACROSS ALL SCALES REDUCING MAXEPS WILL RESULT IN SHORTER RUN TIMES  
METRIC STRING OR CALLABLE OPTIONAL DEFAULT'MINKOWSKI' METRIC TO USE FOR DISTANCE COMPUTA  
TION ANY METRIC FROM SCIKITLEARN OR SCIPYSPATIALDISTANCE CAN BE USED  
IF METRIC IS A CALLABLE FUNCTION IT IS CALLED ON EACH PAIR OF INSTANCES ROWS AND THE RESULTING  
VALUE RECORDED THE CALLABLE SHOULD TAKE TWO ARRAYS AS INPUT AND RETURN ONE VALUE INDICATING  
THE DISTANCE BETWEEN THEM THIS WORKS FOR SCIPY'S METRICS BUT IS LESS EFFICIENT THAN PASSING  
THE METRIC NAME AS A STRING IF METRIC IS "PRECOMPUTED" X IS ASSUMED TO BE A DISTANCE MATRIX  
AND MUST BE SQUARE  
VALID VALUES FOR METRIC ARE  
• FROM SCIKITLEARN 'CITYBLOCK' 'COSINE' 'EUCLIDEAN' 'L1' 'L2' 'MANHATTAN'  
• FROM SCIPYSPATIALDISTANCE 'BRAYCURTIS' 'CANBERRA' 'CHEBYSHEV' 'CORRELATION' 'DICE'  
'HAMMING' 'JACCARD' 'KULSINSKI' 'MAHALANOBIS' 'MINKOWSKI' 'ROGERSTANIMOTO' 'RUS  
SELLRAO' 'SEUCLIDEAN' 'SOKALMICHENER' 'SOKALSNEATH' 'SQEUCLIDEAN' 'YULE'  
SEE THE DOCUMENTATION FOR SCIPYSPATIALDISTANCE FOR DETAILS ON THESE METRICS  
PINTEGER OPTIONAL DEFAULT2 PARAMETER FOR THE MINKOWSKI METRIC FROM SKLEARN  
METRICSPAIRWISEDISTANCES WHEN P 1 THIS IS EQUIVALENT TO USING MANHAT  
TANDISTANCE L1 AND EUCLIDEANDISTANCE L2 FOR P 2 FOR ARBITRARY P MINKOWSKIDISTANCE  
LP IS USED  
METRICPARAMS DICT OPTIONAL DEFAULTNONE ADDITIONAL KEYWORD ARGUMENTS FOR THE METRIC  
FUNCTION  
ALGORITHM 'AUTO' 'BALLTREE' 'KDTree' 'BRUTE' OPTIONAL ALGORITHM USED TO COMPUTE THE  
NEAREST NEIGHBORS  
• 'BALLTREE' WILL USE BALLTREE  
• 'KDTree' WILL USE KDTree  
• 'BRUTE' WILL USE A BRUTEFORCE SEARCH  
• 'AUTO' WILL ATTEMPT TO DECIDE THE MOST APPROPRIATE ALGORITHM BASED ON THE VALUES PASSED  
TOFIT METHOD DEFAULT  
NOTE FITTING ON SPARSE INPUT WILL OVERRIDE THE SETTING OF THIS PARAMETER USING BRUTE FORCE  
LEAF SIZE INT OPTIONAL DEFAULT30 LEAF SIZE PASSED TO BALLTREE ORKDTree THIS CAN  
AFFECT THE SPEED OF THE CONSTRUCTION AND QUERY AS WELL AS THE MEMORY REQUIRED TO STORE THE  
TREE THE OPTIMAL VALUE DEPENDS ON THE NATURE OF THE PROBLEM  
NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF PARALLEL JOBS TO RUN FOR NEIGHBORS  
SEARCHNONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS  
USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS  
RETURNS  
63SKLEARNCLUSTER CLUSTERING 1509

SCIKITLEARN USER GUIDE RELEASE 0213

ORDERING ARRAY SHAPE NSAMPLES THE CLUSTER ORDERED LIST OF SAMPLE INDICES

COREDISTANCES ARRAY SHAPE NSAMPLES DISTANCE AT WHICH EACH SAMPLE BECOMES A CORE POINT INDEXED BY OBJECT ORDER POINTS WHICH WILL NEVER BE CORE HAVE A DISTANCE OF INF USE CLUSTCOREDISTANCESCLUSTORDERING TO ACCESS IN CLUSTER ORDER

REACHABILITY ARRAY SHAPE NSAMPLES REACHABILITY DISTANCES PER SAMPLE INDEXED BY OBJECT ORDER USECLUSTREACHABILITYCLUSTORDERING TO ACCESS IN CLUSTER ORDER

PREDECESSOR ARRAY SHAPE NSAMPLES POINT THAT A SAMPLE WAS REACHED FROM INDEXED BY OBJECT ORDER SEED POINTS HAVE A PREDECESSOR OF 1

REFERENCES

1

SKLEARNCLUSTER DBSCAN

SKLEARNCLUSTER DBSCANXEPS05 MINSAMPLES5 METRIC'MINKOWSKI' METRICPARAMSNONE ALGORITHM'AUTO' LEAFSIZE30 P2SAMPLEWEIGHTNONE NJOBSNONE

PERFORM DBSCAN CLUSTERING FROM VECTOR ARRAY OR DISTANCE MATRIX

READ MORE IN THE USER GUIDE

PARAMETERS

XARRAY OR SPARSE CSR MATRIX OF SHAPE NSAMPLES NFEATURES OR ARRAY OF SHAPE NSAMPLES NSAMPLES A FEATURE ARRAY OR ARRAY OF DISTANCES BETWEEN SAMPLES IF METRICPRECOMPUTED

EPS FLOAT OPTIONAL THE MAXIMUM DISTANCE BETWEEN TWO SAMPLES FOR ONE TO BE CONSIDERED AS IN THE NEIGHBORHOOD OF THE OTHER THIS IS NOT A MAXIMUM BOUND ON THE DISTANCES OF POINTS WITHIN A CLUSTER THIS IS THE MOST IMPORTANT DBSCAN PARAMETER TO CHOOSE APPROPRIATELY FOR YOUR DATA SET AND DISTANCE FUNCTION

MINSAMPLES INT OPTIONAL THE NUMBER OF SAMPLES OR TOTAL WEIGHT IN A NEIGHBORHOOD FOR A POINT TO BE CONSIDERED AS A CORE POINT THIS INCLUDES THE POINT ITSELF

METRIC STRING OR CALLABLE THE METRIC TO USE WHEN CALCULATING DISTANCE BETWEEN INSTANCES IN A FEATURE ARRAY IF METRIC IS A STRING OR CALLABLE IT MUST BE ONE OF THE OPTIONS ALLOWED BYSKLEARNMETRICSPAIRWISEDISTANCES FOR ITS METRIC PARAMETER IF METRIC IS "PRECOMPUTED" X IS ASSUMED TO BE A DISTANCE MATRIX AND MUST BE SQUARE X MAY BE A SPARSE MATRIX IN WHICH CASE ONLY "NONZERO" ELEMENTS MAY BE CONSIDERED NEIGHBORS FOR DBSCAN

METRICPARAMS DICT OPTIONAL ADDITIONAL KEYWORD ARGUMENTS FOR THE METRIC FUNCTION

NEW IN VERSION 019

ALGORITHM 'AUTO' 'BALLTREE' 'KDTree' 'BRUTE' OPTIONAL THE ALGORITHM TO BE USED BY THE NEARESTNEIGHBORS MODULE TO COMPUTE POINTWISE DISTANCES AND FIND NEAREST NEIGHBORS SEE NEARESTNEIGHBORS MODULE DOCUMENTATION FOR DETAILS

LEAFSIZE INT OPTIONAL DEFAULT 30 LEAF SIZE PASSED TO BALLTREE OR CKDTree THIS CAN AFFECT THE SPEED OF THE CONSTRUCTION AND QUERY AS WELL AS THE MEMORY REQUIRED TO STORE THE TREE THE OPTIMAL VALUE DEPENDS ON THE NATURE OF THE PROBLEM

PFLOAT OPTIONAL THE POWER OF THE MINKOWSKI METRIC TO BE USED TO CALCULATE DISTANCE BETWEEN POINTS

1510 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

SAMPLEWEIGHT ARRAY SHAPE NSAMPLES OPTIONAL WEIGHT OF EACH SAMPLE SUCH THAT A SAMPLE WITH A WEIGHT OF AT LEAST MINSAMPLES IS BY ITSELF A CORE SAMPLE A SAMPLE WITH NEGATIVE WEIGHT MAY INHIBIT ITS EPSNEIGHBOR FROM BEING CORE NOTE THAT WEIGHTS ARE ABSOLUTE AND DEFAULT TO 1

NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF PARALLEL JOBS TO RUN FOR NEIGHBORS

SEARCHNONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

RETURNS

CORESAMPLES ARRAY NCORESAMPLES INDICES OF CORE SAMPLES

LABELS ARRAY NSAMPLES CLUSTER LABELS FOR EACH POINT NOISY SAMPLES ARE GIVEN THE LABEL 1

SEE ALSO

DBSCAN AN ESTIMATOR INTERFACE FOR THIS CLUSTERING ALGORITHM

OPTICS A SIMILAR ESTIMATOR INTERFACE CLUSTERING AT MULTIPLE VALUES OF EPS OUR IMPLEMENTATION IS OPTIMIZED FOR MEMORY USAGE

NOTES

FOR AN EXAMPLE SEE EXAMPLESCUSTERPLOTDBSCANPY

THIS IMPLEMENTATION BULKCOMPUTES ALL NEIGHBORHOOD QUERIES WHICH INCREASES THE MEMORY COMPLEXITY TO  $O(DN)$  WHERE D IS THE AVERAGE NUMBER OF NEIGHBORS WHILE ORIGINAL DBSCAN HAD MEMORY COMPLEXITY  $O(N)$  IT MAY ATTRACT A HIGHER MEMORY COMPLEXITY WHEN QUERYING THESE NEAREST NEIGHBORHOODS DEPENDING ON THE ALGORITHM

ONE WAY TO AVOID THE QUERY COMPLEXITY IS TO PRECOMPUTE SPARSE NEIGHBORHOODS IN CHUNKS USINGNEARESTNEIGHBORSRADIUSNEIGHBORSGRAPH WITHMODEDISTANCE THEN USING METRICPRECOMPUTED HERE

ANOTHER WAY TO REDUCE MEMORY AND COMPUTATION TIME IS TO REMOVE NEARDUPLICATE POINTS AND USE SAMPLEWEIGHT INSTEAD

CLUSTEROPTICS PROVIDES A SIMILAR CLUSTERING WITH LOWER MEMORY USAGE

REFERENCES

ESTER M H P KRIEGL J SANDER AND X XU “A DENSITYBASED ALGORITHM FOR DISCOVERING CLUSTERS IN LARGE SPATIAL DATABASES WITH NOISE” IN PROCEEDINGS OF THE 2ND INTERNATIONAL CONFERENCE ON KNOWLEDGE DISCOVERY AND DATA MINING PORTLAND OR AAAI PRESS PP 226231 1996

SCHUBERT E SANDER J ESTER M KRIEGL H P XU X 2017 DBSCAN REVISITED REVISITED WHY AND HOW YOU SHOULD STILL USE DBSCAN ACM TRANSACTIONS ON DATABASE SYSTEMS TODS 423 19

SKLEARNCLUSTER ESTIMATEBANDWIDTH

SKLEARNCLUSTER ESTIMATEBANDWIDTH XQUANTILE03 NSAMPLESNONE RANDOMSTATE0

NJOBSNONE

ESTIMATE THE BANDWIDTH TO USE WITH THE MEANSHIFT ALGORITHM

THAT THIS FUNCTION TAKES TIME AT LEAST QUADRATIC IN NSAMPLES FOR LARGE DATASETS IT’S WISE TO SET THAT PARAMETER TO A SMALL VALUE

63SKLEARNCLUSTER CLUSTERING 1511

SCIKITLEARN USER GUIDE RELEASE 0213

PARAMETERS

XARRAYLIKE SHAPENSAMPLES NFEATURES INPUT POINTS

QUANTILE FLOAT DEFAULT 0.3 SHOULD BE BETWEEN 0 1 0.5 MEANS THAT THE MEDIAN OF ALL PAIRWISE DISTANCES IS USED

NSAMPLES INT OPTIONAL THE NUMBER OF SAMPLES TO USE IF NOT GIVEN ALL SAMPLES ARE USED

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE DEFAULT THE GENERATOR USED TO RANDOMLY SELECT THE SAMPLES FROM INPUT POINTS FOR BANDWIDTH ESTIMATION USE AN INT TO MAKE THE

RANDOMNESS DETERMINISTIC SEE GLOSSARY

NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF PARALLEL JOBS TO RUN FOR NEIGHBORS

SEARCHNONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS

USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

RETURNS

BANDWIDTH FLOAT THE BANDWIDTH PARAMETER

EXAMPLES USING SKLEARNCLUSTERESTIMATEBANDWIDTH

•A DEMO OF THE MEANSHIFT CLUSTERING ALGORITHM

•COMPARING DIFFERENT CLUSTERING ALGORITHMS ON TOY DATASETS

SKLEARNCLUSTER KMEANS

SKLEARNCLUSTER KMEANSXNCLUSTERS SAMPLEWEIGHTNONE INIT'KMEANS' PRECOM

PUTEDISTANCES'AUTO' NINIT10 MAXITER300 VERBOSEFALSE

TOL0.0001 RANDOMSTATENONE COPYXTRUE NJOBSNONE ALGO

RITHM'AUTO' RETURNNNITERFALSE

KMEANS CLUSTERING ALGORITHM

READ MORE IN THE USER GUIDE

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE OBSERVATIONS TO CLUSTER IT

MUST BE NOTED THAT THE DATA WILL BE CONVERTED TO C ORDERING WHICH WILL CAUSE A MEMORY

COPY IF THE GIVEN DATA IS NOT CCONTIGUOUS

NCLUSTERS INT THE NUMBER OF CLUSTERS TO FORM AS WELL AS THE NUMBER OF CENTROIDS TO GENERATE

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL THE WEIGHTS FOR EACH OBSERVATION IN

X IF NONE ALL OBSERVATIONS ARE ASSIGNED EQUAL WEIGHT DEFAULT NONE

INIT 'KMEANS' 'RANDOM' OR NDARRAY OR A CALLABLE OPTIONAL METHOD FOR INITIALIZATION

DEFAULT TO 'KMEANS'

'KMEANS' SELECTS INITIAL CLUSTER CENTERS FOR KMEAN CLUSTERING IN A SMART WAY TO SPEED

UP CONVERGENCE SEE SECTION NOTES IN KINIT FOR MORE DETAILS

'RANDOM' CHOOSE K OBSERVATIONS ROWS AT RANDOM FROM DATA FOR THE INITIAL CENTROIDS

IF AN NDARRAY IS PASSED IT SHOULD BE OF SHAPE NCLUSTERS NFEATURES AND GIVES THE INITIAL

CENTERS  
IF A CALLABLE IS PASSED IT SHOULD TAKE ARGUMENTS X K AND AND A RANDOM STATE AND RETURN AN  
INITIALIZATION

1512 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

PRECOMPUTEDISTANCES ‘AUTO’ TRUE FALSE PRECOMPUTE DISTANCES FASTER BUT TAKES MORE MEMORY

‘AUTO’ DO NOT PRECOMPUTE DISTANCES IF NSAMPLES NCLUSTERS 12 MILLION THIS CORRESPONDS TO ABOUT 100MB OVERHEAD PER JOB USING DOUBLE PRECISION

TRUE ALWAYS PRECOMPUTE DISTANCES

FALSE NEVER PRECOMPUTE DISTANCES

NINIT INT OPTIONAL DEFAULT 10 NUMBER OF TIME THE KMEANS ALGORITHM WILL BE RUN WITH DIFFERENT CENTROID SEEDS THE FINAL RESULTS WILL BE THE BEST OUTPUT OF NINIT CONSECUTIVE RUNS IN TERMS OF INERTIA

MAXITER INT OPTIONAL DEFAULT 300 MAXIMUM NUMBER OF ITERATIONS OF THE KMEANS ALGORITHM TO RUN

VERBOSE BOOLEAN OPTIONAL VERBOSITY MODE

TOLFLOAT OPTIONAL THE RELATIVE INCREMENT IN THE RESULTS BEFORE DECLARING CONVERGENCE

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE DEFAULT DETERMINES RANDOM NUMBER GENERATION FOR CENTROID INITIALIZATION USE AN INT TO MAKE THE RANDOMNESS DETERMINISTIC SEE GLOSSARY

COPYX BOOLEAN OPTIONAL WHEN PRECOMPUTING DISTANCES IT IS MORE NUMERICALLY ACCURATE TO CENTER THE DATA FIRST IF COPYX IS TRUE DEFAULT THEN THE ORIGINAL DATA IS NOT MODIFIED ENSURING X IS CCONTIGUOUS IF FALSE THE ORIGINAL DATA IS MODIFIED AND PUT BACK BEFORE THE FUNCTION RETURNS BUT SMALL NUMERICAL DIFFERENCES MAY BE INTRODUCED BY SUBTRACTING AND THEN ADDING THE DATA MEAN IN THIS CASE IT WILL ALSO NOT ENSURE THAT DATA IS CCONTIGUOUS WHICH MAY CAUSE A SIGNIFICANT SLOWDOWN

NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF JOBS TO USE FOR THE COMPUTATION THIS WORKS BY COMPUTING EACH OF THE NINIT RUNS IN PARALLEL

NONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

ALGORITHM “AUTO” “FULL” OR “ELKAN” DEFAULT“AUTO” KMEANS ALGORITHM TO USE THE CLASSICAL EMSTYLE ALGORITHM IS “FULL” THE “ELKAN” VARIATION IS MORE EFFICIENT BY USING THE TRIANGLE INEQUALITY BUT CURRENTLY DOESN’T SUPPORT SPARSE DATA “AUTO” CHOOSES “ELKAN” FOR DENSE DATA AND “FULL” FOR SPARSE DATA

RETURNNITER BOOL OPTIONAL WHETHER OR NOT TO RETURN THE NUMBER OF ITERATIONS

RETURNS

CENTROID FLOAT NDARRAY WITH SHAPE K NFEATURES CENTROIDS FOUND AT THE LAST ITERATION OF K MEANS

LABEL INTEGER NDARRAY WITH SHAPE NSAMPLES LABELI IS THE CODE OR INDEX OF THE CENTROID THE I’TH OBSERVATION IS CLOSEST TO

INERTIA FLOAT THE FINAL VALUE OF THE INERTIA CRITERION SUM OF SQUARED DISTANCES TO THE CLOSEST CENTROID FOR ALL OBSERVATIONS IN THE TRAINING SET

BESTNITER INT NUMBER OF ITERATIONS CORRESPONDING TO THE BEST RESULTS RETURNED ONLY IF RETURNNITER IS SET TO TRUE

63SKLEARNCLUSTER CLUSTERING 1513

SCIKITLEARN USER GUIDE RELEASE 0213

SKLEARNCLUSTER MEANSHIFT

SKLEARNCLUSTER MEANSHIFT X BANDWIDTHNONE SEEDSNONE BINSEEDINGFALSE

MINBINFREQ1 CLUSTERALLTRUE MAXITER300 NJOBSNONE

PERFORM MEAN SHIFT CLUSTERING OF DATA USING A FLAT KERNEL

READ MORE IN THE USER GUIDE

PARAMETERS

XARRAYLIKE SHAPENSAMPLES NFEATURES INPUT DATA

BANDWIDTH FLOAT OPTIONAL KERNEL BANDWIDTH

IF BANDWIDTH IS NOT GIVEN IT IS DETERMINED USING A HEURISTIC BASED ON THE MEDIAN OF ALL PAIRWISE DISTANCES THIS WILL TAKE QUADRATIC TIME IN THE NUMBER OF SAMPLES THE SKLEARNCLUSTERESTIMATEBANDWIDTH FUNCTION CAN BE USED TO DO THIS MORE EFFICIENTLY

SEEDS ARRAYLIKE SHAPENSEEDS NFEATURES OR NONE POINT USED AS INITIAL KERNEL LOCATIONS IF NONE AND BINSEEDINGFALSE EACH DATA POINT IS USED AS A SEED IF NONE AND BINSEEDINGTRUE SEE BINSEEDING

BINSEEDING BOOLEAN DEFAULTFALSE IF TRUE INITIAL KERNEL LOCATIONS ARE NOT LOCATIONS OF ALL POINTS BUT RATHER THE LOCATION OF THE DISCRETIZED VERSION OF POINTS WHERE POINTS ARE BINNED ONTO A GRID WHOSE COARSENESS CORRESPONDS TO THE BANDWIDTH SETTING THIS OPTION TO TRUE WILL SPEED UP THE ALGORITHM BECAUSE FEWER SEEDS WILL BE INITIALIZED IGNORED IF SEEDS ARGUMENT IS NOT NONE

MINBINFREQ INT DEFAULT1 TO SPEED UP THE ALGORITHM ACCEPT ONLY THOSE BINS WITH AT LEAST MINBINFREQ POINTS AS SEEDS

CLUSTERALL BOOLEAN DEFAULT TRUE IF TRUE THEN ALL POINTS ARE CLUSTERED EVEN THOSE ORPHANS THAT ARE NOT WITHIN ANY KERNEL ORPHANS ARE ASSIGNED TO THE NEAREST KERNEL IF FALSE THEN ORPHANS ARE GIVEN CLUSTER LABEL 1

MAXITER INT DEFAULT 300 MAXIMUM NUMBER OF ITERATIONS PER SEED POINT BEFORE THE CLUSTERING OPERATION TERMINATES FOR THAT SEED POINT IF HAS NOT CONVERGED YET

NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF JOBS TO USE FOR THE COMPUTATION THIS WORKS BY COMPUTING EACH OF THE NINIT RUNS IN PARALLEL

NONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

NEW IN VERSION 017 PARALLEL EXECUTION USING NJOBS

RETURNS

CLUSTERCENTERS ARRAY SHAPENCLUSTERS NFEATURES COORDINATES OF CLUSTER CENTERS

LABELS ARRAY SHAPENSAMPLES CLUSTER LABELS FOR EACH POINT

NOTES

FOR AN EXAMPLE SEE EXAMPLESCUSTERPLOTMEANSHIFTPY

1514 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

SKLEARNCLUSTER SPECTRALCLUSTERING

SKLEARNCLUSTER SPECTRALCLUSTERING AFFINITY NCLUSTERS8 NCOMPONENTSNONE

EIGENSOLVERNONE RANDOMSTATENONE NINIT10

EIGENTOL00 ASSIGNLABELS'KMEANS'

APPLY CLUSTERING TO A PROJECTION OF THE NORMALIZED LAPLACIAN

IN PRACTICE SPECTRAL CLUSTERING IS VERY USEFUL WHEN THE STRUCTURE OF THE INDIVIDUAL CLUSTERS IS HIGHLY NONCONVEX OR MORE GENERALLY WHEN A MEASURE OF THE CENTER AND SPREAD OF THE CLUSTER IS NOT A SUITABLE DESCRIPTION OF THE COMPLETE CLUSTER FOR INSTANCE WHEN CLUSTERS ARE NESTED CIRCLES ON THE 2D PLANE

IF AFFINITY IS THE ADJACENCY MATRIX OF A GRAPH THIS METHOD CAN BE USED TO FIND NORMALIZED GRAPH CUTS

READ MORE IN THE USER GUIDE

PARAMETERS

AFFINITY ARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NSAMPLES THE AFFINITY MATRIX DESCRIBING THE RELATIONSHIP OF THE SAMPLES TO EMBED MUST BE SYMMETRIC

POSSIBLE EXAMPLES

- ADJACENCY MATRIX OF A GRAPH
- HEAT KERNEL OF THE PAIRWISE DISTANCE MATRIX OF THE SAMPLES
- SYMMETRIC KNEAREST NEIGHBOURS CONNECTIVITY MATRIX OF THE SAMPLES

NCLUSTERS INTEGER OPTIONAL NUMBER OF CLUSTERS TO EXTRACT

NCOMPONENTS INTEGER OPTIONAL DEFAULT IS NCLUSTERS NUMBER OF EIGEN VECTORS TO USE FOR THE SPECTRAL EMBEDDING

EIGENSOLVER NONE 'ARPACK' 'LOBPCG' OR 'AMG' THE EIGENVALUE DECOMPOSITION STRATEGY TO USE AMG REQUIRES PYAMG TO BE INSTALLED IT CAN BE FASTER ON VERY LARGE SPARSE PROBLEMS BUT MAY ALSO LEAD TO INSTABILITIES

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE DEFAULT A PSEUDO RANDOM NUMBER GENERATOR USED FOR THE INITIALIZATION OF THE LOBPCG EIGEN VECTORS DECOMPOSITION WHEN EIGENSOLVER 'AMG' AND BY THE KMEANS INITIALIZATION USE AN INT TO MAKE THE RANDOMNESS DETERMINISTIC SEE GLOSSARY

NINIT INT OPTIONAL DEFAULT 10 NUMBER OF TIME THE KMEANS ALGORITHM WILL BE RUN WITH DIFFERENT CENTROID SEEDS THE FINAL RESULTS WILL BE THE BEST OUTPUT OF NINIT CONSECUTIVE RUNS IN TERMS OF INERTIA

EIGENTOL FLOAT OPTIONAL DEFAULT 00 STOPPING CRITERION FOR EIGENDECOMPOSITION OF THE LAPLACIAN MATRIX WHEN USING ARPACK EIGENSOLVER

ASSIGNLABELS 'KMEANS' 'DISCRETIZE' DEFAULT 'KMEANS' THE STRATEGY TO USE TO ASSIGN LABELS IN THE EMBEDDING SPACE THERE ARE TWO WAYS TO ASSIGN LABELS AFTER THE LAPLACIAN EMBEDDING KMEANS CAN BE APPLIED AND IS A POPULAR CHOICE BUT IT CAN ALSO BE SENSITIVE TO INITIALIZATION DISCRETIZATION IS ANOTHER APPROACH WHICH IS LESS SENSITIVE TO RANDOM INITIALIZATION SEE THE 'MULTICLASS SPECTRAL CLUSTERING' PAPER REFERENCED BELOW FOR MORE DETAILS ON THE DISCRETIZATION APPROACH

RETURNS

LABELS ARRAY OF INTEGERS SHAPE NSAMPLES THE LABELS OF THE CLUSTERS

63SKLEARNCLUSTER CLUSTERING 1515

SCIKITLEARN USER GUIDE RELEASE 0213

NOTES

THE GRAPH SHOULD CONTAIN ONLY ONE CONNECT COMPONENT ELSEWHERE THE RESULTS MAKE LITTLE SENSE  
THIS ALGORITHM SOLVES THE NORMALIZED CUT FOR K2 IT IS A NORMALIZED SPECTRAL CLUSTERING

REFERENCES

- NORMALIZED CUTS AND IMAGE SEGMENTATION 2000 JIANBO SHI JITENDRA MALIK HTTPCITESEERISTPSUEDUVIEWDOC SUMMARYDOI10111602324
- A TUTORIAL ON SPECTRAL CLUSTERING 2007 ULRIKE VON LUXBURG HTTPCITESEERXISTPSUEDUVIEWDOC SUMMARYDOI10111659323
- MULTICLASS SPECTRAL CLUSTERING 2003 STELLA X YU JIANBO SHI HTTPSWWW1ICSIBERKELEYEDUSTELLAYUPUBLICATIONDOC2003KWAYICCV PDF

EXAMPLES USING SKLEARNCLUSTERSPECTRALCLUSTERING

•SEGMENTING THE PICTURE OF GREEK COINS IN REGIONS

•SPECTRAL CLUSTERING FOR IMAGE SEGMENTATION

SKLEARNCLUSTER WARDTREE

SKLEARNCLUSTER WARDTREE XCONNECTIVITYNONE NCLUSTERSNONE RETURNDISTANCEFALSE

WARD CLUSTERING BASED ON A FEATURE MATRIX

RECURSIVELY MERGES THE PAIR OF CLUSTERS THAT MINIMALLY INCREASES WITHINCLUSTER VARIANCE

THE INERTIA MATRIX USES A HEAPQBASED REPRESENTATION

THIS IS THE STRUCTURED VERSION THAT TAKES INTO ACCOUNT SOME TOPOLOGICAL STRUCTURE BETWEEN SAMPLES

READ MORE IN THE USER GUIDE

PARAMETERS

XARRAY SHAPE NSAMPLES NFEATURES FEATURE MATRIX REPRESENTING NSAMPLES SAMPLES TO BE CLUSTERED

CONNECTIVITY SPARSE MATRIX OPTIONAL CONNECTIVITY MATRIX DEFINES FOR EACH SAMPLE THE NEIGHBORING SAMPLES FOLLOWING A GIVEN STRUCTURE OF THE DATA THE MATRIX IS ASSUMED TO BE SYMMETRIC AND ONLY THE UPPER TRIANGULAR HALF IS USED DEFAULT IS NONE IE THE WARD ALGORITHM IS UNSTRUCTURED

NCLUSTERS INT OPTIONAL STOP EARLY THE CONSTRUCTION OF THE TREE AT NCLUSTERS THIS IS USEFUL TO DECREASE COMPUTATION TIME IF THE NUMBER OF CLUSTERS IS NOT SMALL COMPARED TO THE NUMBER OF SAMPLES IN THIS CASE THE COMPLETE TREE IS NOT COMPUTED THUS THE 'CHILDREN' OUTPUT IS OF LIMITED USE AND THE 'PARENTS' OUTPUT SHOULD RATHER BE USED THIS OPTION IS VALID ONLY WHEN SPECIFYING A CONNECTIVITY MATRIX

RETURNDISTANCE BOOL OPTIONAL IF TRUE RETURN THE DISTANCE BETWEEN THE CLUSTERS

RETURNS

1516 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

CHILDREN 2D ARRAY SHAPE NNODES1 2 THE CHILDREN OF EACH NONLEAF NODE VALUES LESS THAN NSAMPLES CORRESPOND TO LEAVES OF THE TREE WHICH ARE THE ORIGINAL SAMPLES A NODE I GREATER THAN OR EQUAL TO NSAMPLES IS A NONLEAF NODE AND HAS CHILDREN CHILDRENI NSAMPLES ALTERNATIVELY AT THE ITH ITERATION CHILDRENI0 AND CHILDRENI1 ARE MERGED TO FORM NODE NSAMPLES I

NCONNECTEDCOMPONENTS INT THE NUMBER OF CONNECTED COMPONENTS IN THE GRAPH NLEAVES INT THE NUMBER OF LEAVES IN THE TREE

PARENTS 1D ARRAY SHAPE NNODES OR NONE THE PARENT OF EACH NODE ONLY RETURNED WHEN A CONNECTIVITY MATRIX IS SPECIFIED ELSEWHERE 'NONE' IS RETURNED

DISTANCES 1D ARRAY SHAPE NNODES1 ONLY RETURNED IF RETURNDISTANCE IS SET TO TRUE FOR COMPATIBILITY THE DISTANCES BETWEEN THE CENTERS OF THE NODES DISTANCESI COR RESPONDS TO A WEIGHTED EUCLIDEAN DISTANCE BETWEEN THE NODES CHILDRENI 1 AND CHILDRENI 2 IF THE NODES REFER TO LEAVES OF THE TREE THEN DISTANCESI IS THEIR UNWEIGHTED EUCLIDEAN DISTANCE DISTANCES ARE UPDATED IN THE FOLLOWING WAY FROM SCIPYHIERARCHYLINKAGE

THE NEW ENTRY  $\sqrt{c_1^2 + c_2^2}$  IS COMPUTED AS FOLLOWS

WHERE  $c_1$  IS THE NEWLY JOINED CLUSTER CONSISTING OF CLUSTERS  $c_1$  AND  $c_2$  IS AN UNUSED CLUSTER IN THE FOREST  $n$  AND  $n$  IS THE CARDINALITY OF ITS ARGUMENT THIS IS ALSO KNOWN AS THE INCREMENTAL ALGORITHM

64SKLEARNCLUSTERBICLUSTER BICLUSTERING

SPECTRAL BICLUSTERING ALGORITHMS

AUTHORS KEMAL EREN LICENSE BSD 3 CLAUSE

USER GUIDE SEE THE BICLUSTERING SECTION FOR FURTHER DETAILS

641 CLASSES

SPECTRALBICLUSTERING NCLUSTERS METHOD SPECTRAL BICLUSTERING KLUGER 2003

SPECTRALCOCLUSTERING NCLUSTERS SPECTRAL COCLUSTERING ALGORITHM DHILLON 2001

SKLEARNCLUSTERBICLUSTER SPECTRALBICLUSTERING

CLASSSSKLEARNCLUSTERBICLUSTER SPECTRALBICLUSTERING NCLUSTERS3

METHOD'BISTOCHASTIC'

NCOMPONENTS6 NBEST3

SVDMETHOD'RANDOMIZED'

NSVDVECSNONE

MINIBATCHFALSE INIT'K

MEANS' NINIT10

NJOBSNONE RAN

DOMSTATENONE

SPECTRAL BICLUSTERING KLUGER 2003

64SKLEARNCLUSTERBICLUSTER BICLUSTERING 1517

SCIKITLEARN USER GUIDE RELEASE 0213

PARTITIONS ROWS AND COLUMNS UNDER THE ASSUMPTION THAT THE DATA HAS AN UNDERLYING CHECKERBOARD STRUCTURE FOR INSTANCE IF THERE ARE TWO ROW PARTITIONS AND THREE COLUMN PARTITIONS EACH ROW WILL BELONG TO THREE BICLUSTERS AND EACH COLUMN WILL BELONG TO TWO BICLUSTERS THE OUTER PRODUCT OF THE CORRESPONDING ROW AND COLUMN LABEL VECTORS GIVES THIS CHECKERBOARD STRUCTURE  
READ MORE IN THE USER GUIDE

PARAMETERS

NCLUSTERS INTEGER OR TUPLE NROWCLUSTERS NCOLUMNCLUSTERS THE NUMBER OF ROW AND COLUMN CLUSTERS IN THE CHECKERBOARD STRUCTURE

METHOD STRING OPTIONAL DEFAULT 'BISTOCHASTIC' METHOD OF NORMALIZING AND CONVERTING SINGULAR VECTORS INTO BICLUSTERS MAY BE ONE OF 'SCALE' 'BISTOCHASTIC' OR 'LOG' THE AUTHORS RECOMMEND USING 'LOG' IF THE DATA IS SPARSE HOWEVER LOG NORMALIZATION WILL NOT WORK WHICH IS WHY THE DEFAULT IS 'BISTOCHASTIC' CAUTION IF METHODLOG THE DATA MUST NOT BE SPARSE

NCOMPONENTS INTEGER OPTIONAL DEFAULT 6 NUMBER OF SINGULAR VECTORS TO CHECK

NBEST INTEGER OPTIONAL DEFAULT 3 NUMBER OF BEST SINGULAR VECTORS TO WHICH TO PROJECT THE DATA FOR CLUSTERING

SVDMETHOD STRING OPTIONAL DEFAULT 'RANDOMIZED' SELECTS THE ALGORITHM FOR FINDING SINGULAR VECTORS MAY BE 'RANDOMIZED' OR 'ARPACK' IF 'RANDOMIZED' USES SKLEARNUTILS EXTMATHRANDOMIZEDSVD WHICH MAY BE FASTER FOR LARGE MATRICES IF 'ARPACK' USES SCIPYSPARSELINALGSVDS WHICH IS MORE ACCURATE BUT POSSIBLY SLOWER IN SOME CASES

NSVDVECS INT OPTIONAL DEFAULT NONE NUMBER OF VECTORS TO USE IN CALCULATING THE SVD CORRESPONDS TO NCV WHENSVDMETHODARPACK ANDNOVERSAMPLES WHEN SVDMETHOD IS 'RANDOMIZED'

MINIBATCH BOOL OPTIONAL DEFAULT FALSE WHETHER TO USE MINIBATCH KMEANS WHICH IS FASTER BUT MAY GET DIFFERENT RESULTS

INIT 'KMEANS' 'RANDOM' OR AN NDARRAY METHOD FOR INITIALIZATION OF KMEANS ALGORITHM DEFAULTS TO 'KMEANS'

NINIT INT OPTIONAL DEFAULT 10 NUMBER OF RANDOM INITIALIZATIONS THAT ARE TRIED WITH THE KMEANS ALGORITHM

IF MINIBATCH KMEANS IS USED THE BEST INITIALIZATION IS CHOSEN AND THE ALGORITHM RUNS ONCE OTHERWISE THE ALGORITHM IS RUN FOR EACH INITIALIZATION AND THE BEST SOLUTION CHOSEN

NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF JOBS TO USE FOR THE COMPUTATION THIS WORKS BY BREAKING DOWN THE PAIRWISE MATRIX INTO NJOBS EVEN SLICES AND COMPUTING THEM IN PARALLEL

NONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE DEFAULT USED FOR RANDOMIZING THE SINGULAR VALUE DECOMPOSITION AND THE KMEANS INITIALIZATION USE AN INT TO MAKE THE RANDOMNESS DETERMINISTIC SEE GLOSSARY

ATTRIBUTES

ROWS ARRAYLIKE SHAPE NROWCLUSTERS NROWS RESULTS OF THE CLUSTERING ROWS I R IS TRUE IF CLUSTER ICONTAINS ROW R AVAILABLE ONLY AFTER CALLING FIT

COLUMNS ARRAYLIKE SHAPE NCOLUMNCLUSTERS NCOLUMNS RESULTS OF THE CLUSTERING LIKE ROWS

SCIKITLEARN USER GUIDE RELEASE 0213

ROWLABELS ARRAYLIKE SHAPE NROWS ROW PARTITION LABELS

COLUMNLABELS ARRAYLIKE SHAPE NCOLS COLUMN PARTITION LABELS

REFERENCES

- KLUGER YUVAL ET AL 2003 SPECTRAL BICLUSTERING OF MICROARRAY DATA COCLUSTERING GENES AND CONDITIONS

EXAMPLES

```
FROM SKLEARNCLUSTER IMPORT SPECTRALBICLUSTERING
IMPORT NUMPY AS NP
X NPARRAY1 1 2 1 1 0
4 7 3 5 3 6
CLUSTERING SPECTRALBICLUSTERINGNCLUSTERS2 RANDOMSTATE0FITX
CLUSTERINGROWLABELS
ARRAY1 1 1 0 0 0 DTYPEINT32
CLUSTERINGCOLUMNLABELS
ARRAY0 1 DTYPEINT32
CLUSTERING
SPECTRALBICLUSTERINGINITKMEANS METHODBISTOCHASTIC
MINIBATCHFALSE NBEST3 NCLUSTERS2 NCOMPONENTS6
NINIT10 NJOBSNONE NSVDVECSNONE RANDOMSTATE0
SVDMETHODRANDOMIZED
METHODS
FITSELF X Y CREATES A BICLUSTERING FOR X
GETINDICES SELF I ROW AND COLUMN INDICES OF THE I'TH BICLUSTER
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR
GETSHAPE SELF I SHAPE OF THE I'TH BICLUSTER
GETSUBMATRIX SELF I DATA RETURNS THE SUBMATRIX CORRESPONDING TO BICLUSTER I
SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR
INIT SELF NCLUSTERS3 METHOD'BISTOCHASTIC' NCOMPONENTS6 NBEST3
SVDMETHOD'RANDOMIZED' NSVDVECSNONE MINIBATCHFALSE INIT'KMEANS'
NINIT10 NJOBSNONE RANDOMSTATENONE
BICLUSTERS
CONVENIENT WAY TO GET ROW AND COLUMN INDICATORS TOGETHER
RETURNS THE ROWS ANDCOLUMNS MEMBERS
FITSELFXYNONE
CREATES A BICLUSTERING FOR X
PARAMETERS
XARRAYLIKE SHAPE NSAMPLES NFEATURES
YIGNORED
GETINDICES SELF I
ROW AND COLUMN INDICES OF THE I'TH BICLUSTER
64SKLEARNCLUSTERBICLUSTER BICLUSTERING 1519
```

SCIKITLEARN USER GUIDE RELEASE 0213

ONLY WORKS IF ROWS ANDCOLUMNS ATTRIBUTES EXIST

PARAMETERS

IINT THE INDEX OF THE CLUSTER

RETURNS

ROWIND NPARRAY DTYPENPINTP INDICES OF ROWS IN THE DATASET THAT BELONG TO THE BICLUSTER

COLIND NPARRAY DTYPENPINTP INDICES OF COLUMNS IN THE DATASET THAT BELONG TO THE BICLUSTER

TER

GETPARAMS SELFDEEPTTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

GETSHAPE SELF

SHAPE OF THE I'TH BICLUSTER

PARAMETERS

IINT THE INDEX OF THE CLUSTER

RETURNS

SHAPE INT INT NUMBER OF ROWS AND COLUMNS RESP IN THE BICLUSTER

GETSUBMATRIX SELFIDATA

RETURNS THE SUBMATRIX CORRESPONDING TO BICLUSTER I

PARAMETERS

IINT THE INDEX OF THE CLUSTER

DATA ARRAY THE DATA

RETURNS

SUBMATRIX ARRAY THE SUBMATRIX CORRESPONDING TO BICLUSTER I

NOTES

WORKS WITH SPARSE MATRICES ONLY WORKS IF ROWS ANDCOLUMNS ATTRIBUTES EXIST

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

1520 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

SKLEARNCLUSTERBICLUSTER SPECTRALCOCLUSTERING

CLASSSSKLEARNCLUSTERBICLUSTER SPECTRALCOCLUSTERING NCLUSTERS3

SVDMETHOD'RANDOMIZED'

NSVDVECSNONE

MINIBATCHFALSE INIT'K

MEANS' NINIT10

NJOBSNONE RAN

DOMSTATENONE

SPECTRAL COCLUSTERING ALGORITHM DHILLON 2001

CLUSTERS ROWS AND COLUMNS OF AN ARRAY XTO SOLVE THE RELAXED NORMALIZED CUT OF THE BIPARTITE GRAPH CREATED FROM

XAS FOLLOWS THE EDGE BETWEEN ROW VERTEX IAND COLUMN VERTEX JHAS WEIGHT XI J

THE RESULTING BICLUSTER STRUCTURE IS BLOCKDIAGONAL SINCE EACH ROW AND EACH COLUMN BELONGS TO EXACTLY ONE

BICLUSTER

SUPPORTS SPARSE MATRICES AS LONG AS THEY ARE NONNEGATIVE

READ MORE IN THE USER GUIDE

PARAMETERS

NCLUSTERS INTEGER OPTIONAL DEFAULT 3 THE NUMBER OF BICLUSTERS TO FIND

SVDMETHOD STRING OPTIONAL DEFAULT 'RANDOMIZED' SELECTS THE ALGORITHM FOR FINDING SINGU

LAR VECTORS MAY BE 'RANDOMIZED' OR 'ARPACK' IF 'RANDOMIZED' USE SKLEARNUTILS

EXTMATHRANDOMIZEDSVD WHICH MAY BE FASTER FOR LARGE MATRICES IF 'ARPACK' USE

SCIPYSPARSELINALGSVDS WHICH IS MORE ACCURATE BUT POSSIBLY SLOWER IN SOME

CASES

NSVDVECS INT OPTIONAL DEFAULT NONE NUMBER OF VECTORS TO USE IN CALCULATING THE

SVD CORRESPONDS TO NCV WHENSVDMETHODARPACK ANDNOVERSAMPLES WHEN

SVDMETHOD IS 'RANDOMIZED'

MINIBATCH BOOL OPTIONAL DEFAULT FALSE WHETHER TO USE MINIBATCH KMEANS WHICH IS FASTER

BUT MAY GET DIFFERENT RESULTS

INIT 'KMEANS' 'RANDOM' OR AN NDARRAY METHOD FOR INITIALIZATION OF KMEANS ALGORITHM

DEFAULTS TO 'KMEANS'

NINIT INT OPTIONAL DEFAULT 10 NUMBER OF RANDOM INITIALIZATIONS THAT ARE TRIED WITH THE K

MEANS ALGORITHM

IF MINIBATCH KMEANS IS USED THE BEST INITIALIZATION IS CHOSEN AND THE ALGORITHM RUNS ONCE

OTHERWISE THE ALGORITHM IS RUN FOR EACH INITIALIZATION AND THE BEST SOLUTION CHOSEN

NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF JOBS TO USE FOR THE COMPUTATION

THIS WORKS BY BREAKING DOWN THE PAIRWISE MATRIX INTO NJOBS EVEN SLICES AND COMPUTING

THEM IN PARALLEL

NONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL

PROCESSORS SEE GLOSSARY FOR MORE DETAILS

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE DEFAULT USED FOR RANDOMIZING THE SINGULAR

VALUE DECOMPOSITION AND THE KMEANS INITIALIZATION USE AN INT TO MAKE THE RANDOMNESS

DETERMINISTIC SEE GLOSSARY

ATTRIBUTES

ROWS ARRAYLIKE SHAPE NROWCLUSTERS NROWS RESULTS OF THE CLUSTERING ROWSI R IS

TRUE IF CLUSTER ICONTAINS ROW R AVAILABLE ONLY AFTER CALLING FIT

64SKLEARNCLUSTERBICLUSTER BICLUSTERING 1521

SCIKITLEARN USER GUIDE RELEASE 0213

COLUMNS ARRAYLIKE SHAPE NCOLUMNCLUSTERS NCOLUMNS RESULTS OF THE CLUSTERING LIKE ROWS

ROWLABELS ARRAYLIKE SHAPE NROWS THE BICLUSTER LABEL OF EACH ROW

COLUMNLABELS ARRAYLIKE SHAPE NCOLS THE BICLUSTER LABEL OF EACH COLUMN

REFERENCES

- DHILLON Inderjit S 2001 COCLUSTERING DOCUMENTS AND WORDS USING BIPARTITE SPECTRAL GRAPH PARTITIONING

EXAMPLES

```
FROM SKLEARNCLUSTER IMPORT SPECTRALCOCLUSTERING
IMPORT NUMPY AS NP
X = NPARRAY1 1 2 1 1 0
4 7 3 5 3 6
CLUSTERING = SPECTRALCOCLUSTERING(NCLUSTERS2, RANDOMSTATE=FITX,
CLUSTERINGROWLABELS,
ARRAY0 1 1 0 0 0, DTYPE=INT32,
CLUSTERINGCOLUMNLABELS,
ARRAY0 0, DTYPE=INT32)
CLUSTERING =
SPECTRALCOCLUSTERING(INIT='KMEANS', MINIBATCH=False, NCLUSTERS2,
NINIT=10, NJOBS=None, NSVDVECS=None, RANDOMSTATE=0,
SVDMETHOD='RANDOMIZED')
METHODS
FITSELF(X, Y) CREATES A BICLUSTERING FOR X
GETINDICES(self, i) ROW AND COLUMN INDICES OF THE I'TH BICLUSTER
GETPARAMS(self) DEEP GET PARAMETERS FOR THIS ESTIMATOR
GETSHAPE(self, i) SHAPE OF THE I'TH BICLUSTER
GETSUBMATRIX(self, i) DATA RETURNS THE SUBMATRIX CORRESPONDING TO BICLUSTER I
SETPARAMS(self, params) SET THE PARAMETERS OF THIS ESTIMATOR
INIT(self, nclusters=3, svdmethod='RANDOMIZED', nsvdvecs=None, minibatch=False,
init='KMEANS', ninit=10, njobs=None, randomstate=None)
BICLUSTERS
CONVENIENT WAY TO GET ROW AND COLUMN INDICATORS TOGETHER
RETURNS THE ROWS AND COLUMNS MEMBERS
FITSELF(X, Y, None)
CREATES A BICLUSTERING FOR X
PARAMETERS
X ARRAYLIKE SHAPE NSAMPLES NFEATURES
Y IGNORED
GETINDICES(self, i)
ROW AND COLUMN INDICES OF THE I'TH BICLUSTER
```

1522 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

ONLY WORKS IF ROWS ANDCOLUMNS ATTRIBUTES EXIST

PARAMETERS

IINT THE INDEX OF THE CLUSTER

RETURNS

ROWIND NPARRAY DTYPENPINTP INDICES OF ROWS IN THE DATASET THAT BELONG TO THE BICLUSTER

COLIND NPARRAY DTYPENPINTP INDICES OF COLUMNS IN THE DATASET THAT BELONG TO THE BICLUSTER

TER

GETPARAMS SELFDEEPTTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

GETSHAPE SELF

SHAPE OF THE I'TH BICLUSTER

PARAMETERS

IINT THE INDEX OF THE CLUSTER

RETURNS

SHAPE INT INT NUMBER OF ROWS AND COLUMNS RESP IN THE BICLUSTER

GETSUBMATRIX SELFIDATA

RETURNS THE SUBMATRIX CORRESPONDING TO BICLUSTER I

PARAMETERS

IINT THE INDEX OF THE CLUSTER

DATA ARRAY THE DATA

RETURNS

SUBMATRIX ARRAY THE SUBMATRIX CORRESPONDING TO BICLUSTER I

NOTES

WORKS WITH SPARSE MATRICES ONLY WORKS IF ROWS ANDCOLUMNS ATTRIBUTES EXIST

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

645KLEARNCLUSTERBICLUSTER BICLUSTERING 1523

SCIKITLEARN USER GUIDE RELEASE 0213

65SKLEARNCOMPOSE COMPOSITE ESTIMATORS

METAESTIMATORS FOR BUILDING COMPOSITE MODELS WITH TRANSFORMERS

IN ADDITION TO ITS CURRENT CONTENTS THIS MODULE WILL EVENTUALLY BE HOME TO REFURBISHED VERSIONS OF PIPELINE AND FEATUREUNION

USER GUIDE SEE THE PIPELINES AND COMPOSITE ESTIMATORS SECTION FOR FURTHER DETAILS

COMPOSECOLUMNTRANSFORMER TRANSFORMERS APPLIES TRANSFORMERS TO COLUMNS OF AN ARRAY OR PANDAS DATAFRAME

COMPOSETRANSFORMEDTARGETREGRESSOR METAESTIMATOR TO REGRESS ON A TRANSFORMED TARGET

651SKLEARNCOMPOSE COLUMNTRANSFORMER

CLASSSSKLEARNCOMPOSE COLUMNTRANSFORMER TRANSFORMERS REMAINDER'DROP'

SPARSETHRESHOLD03 NJOBSNONE TRANS

FORMERWEIGHTSNONE VERBOSEFALSE

APPLIES TRANSFORMERS TO COLUMNS OF AN ARRAY OR PANDAS DATAFRAME

THIS ESTIMATOR ALLOWS DIFFERENT COLUMNS OR COLUMN SUBSETS OF THE INPUT TO BE TRANSFORMED SEPARATELY AND THE FEATURES GENERATED BY EACH TRANSFORMER WILL BE CONCATENATED TO FORM A SINGLE FEATURE SPACE THIS IS USEFUL FOR HETEROGENEOUS OR COLUMNAR DATA TO COMBINE SEVERAL FEATURE EXTRACTION MECHANISMS OR TRANSFORMATIONS INTO A SINGLE TRANSFORMER

READ MORE IN THE USER GUIDE

NEW IN VERSION 020

PARAMETERS

TRANSFORMERS LIST OF TUPLES LIST OF NAME TRANSFORMER COLUMNS TUPLES SPECIFYING THE TRANSFORMER OBJECTS TO BE APPLIED TO SUBSETS OF THE DATA

NAME STRING LIKE IN PIPELINE AND FEATUREUNION THIS ALLOWS THE TRANSFORMER AND ITS PARAMETERS TO BE SET USING SETPARAMS AND SEARCHED IN GRID SEARCH

TRANSFORMER ESTIMATOR OR 'PASSTHROUGH' 'DROP' ESTIMATOR MUST SUPPORT FIT AND TRANSFORM SPECIALCASED STRINGS 'DROP' AND 'PASSTHROUGH' ARE ACCEPTED AS WELL TO INDICATE TO DROP THE COLUMNS OR TO PASS THEM THROUGH UNTRANSFORMED RESPECTIVELY

COLUMNS STRING OR INT ARRAYLIKE OF STRING OR INT SLICE BOOLEAN MASK ARRAY OR CALLABLE INDEXES THE DATA ON ITS SECOND AXIS INTEGERS ARE INTERPRETED AS POSITIONAL COLUMNS WHILE STRINGS CAN REFERENCE DATAFRAME COLUMNS BY NAME A SCALAR STRING OR INT SHOULD BE USED WHERETRANSFORMER EXPECTS X TO BE A 1D ARRAYLIKE VECTOR OTHERWISE A 2D ARRAY WILL BE PASSED TO THE TRANSFORMER A CALLABLE IS PASSED THE INPUT DATA XAND CAN RETURN ANY OF THE ABOVE

REMAINDER 'DROP' 'PASSTHROUGH' OR ESTIMATOR DEFAULT 'DROP' BY DEFAULT ONLY THE SPECIFIED COLUMNS IN TRANSFORMERS ARE TRANSFORMED AND COMBINED IN THE OUTPUT AND THE NONSPECIFIED COLUMNS ARE DROPPED DEFAULT OF DROP BY SPECIFYING

REMAINDERPASSTHROUGH ALL REMAINING COLUMNS THAT WERE NOT SPECIFIED IN TRANSFORMERS WILL BE AUTOMATICALLY PASSED THROUGH THIS SUBSET OF COLUMNS IS CONCATENATED WITH THE OUTPUT OF THE TRANSFORMERS BY SETTING REMAINDER TO BE AN ESTIMATOR THE REMAINING NONSPECIFIED COLUMNS WILL USE THE REMAINDER ESTIMATOR THE ESTIMATOR MUST SUPPORT FIT ANDTRANSFORM NOTE THAT USING THIS FEATURE REQUIRES THAT THE DATAFRAME COLUMNS INPUT AT FIT ANDTRANSFORM HAVE IDENTICAL ORDER

1524 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

SPARSETHRESHOLD FLOAT DEFAULT 0.3 IF THE OUTPUT OF THE DIFFERENT TRANSFORMERS CONTAINS SPARSE MATRICES THESE WILL BE STACKED AS A SPARSE MATRIX IF THE OVERALL DENSITY IS LOWER THAN THIS VALUE USESPARSETHRESHOLD0 TO ALWAYS RETURN DENSE WHEN THE TRANSFORMED OUTPUT CONSISTS OF ALL DENSE DATA THE STACKED RESULT WILL BE DENSE AND THIS KEYWORD WILL BE IGNORED NJOBS INT OR NONE OPTIONAL DEFAULTNONE NUMBER OF JOBS TO RUN IN PARALLEL NONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

TRANSFORMERWEIGHTS DICT OPTIONAL MULTIPLICATIVE WEIGHTS FOR FEATURES PER TRANSFORMER THE OUTPUT OF THE TRANSFORMER IS MULTIPLIED BY THESE WEIGHTS KEYS ARE TRANSFORMER NAMES VALUES THE WEIGHTS

VERBOSE BOOLEAN OPTIONALDEFAULTFALSE IF TRUE THE TIME ELAPSED WHILE FITTING EACH TRANSFORMER WILL BE PRINTED AS IT IS COMPLETED

ATTRIBUTES

TRANSFORMERS LIST THE COLLECTION OF FITTED TRANSFORMERS AS TUPLES OF NAME FITTEDTRANSFORMER COLUMNFITTEDTRANSFORMER CAN BE AN ESTIMATOR 'DROP' OR 'PASSTHROUGH' IN

CASE THERE WERE NO COLUMNS SELECTED THIS WILL BE THE UNFITTED TRANSFORMER IF THERE ARE REMAINING COLUMNS THE FINAL ELEMENT IS A TUPLE OF THE FORM 'REMAINDER' TRANSFORMER REMAININGCOLUMNS CORRESPONDING TO THE REMAINDER PARAMETER IF THERE ARE REMAINING COLUMNS THEN LENTRANSFORMERSLENTTRANSFORMERS1 OTHER

WISELENTTRANSFORMERSLENTTRANSFORMERS

NAMEDTRANSFORMERS BUNCH OBJECT A DICTIONARY WITH ATTRIBUTE ACCESS ACCESS THE FITTED

TRANSFORMER BY NAME

SPARSEOUTPUT BOOLEAN BOOLEAN FLAG INDICATING WETHER THE OUTPUT OF TRANSFORM IS A SPARSE MATRIX OR A DENSE NUMPY ARRAY WHICH DEPENDS ON THE OUTPUT OF THE INDIVIDUAL TRANSFORMERS AND THESPARSETHRESHOLD KEYWORD

SEE ALSO

SKLEARNCOMPOSEMAKECOLUMNTRANSFORMER CONVENIENCE FUNCTION FOR COMBINING THE OUTPUTS OF MULTIPLE TRANSFORMER OBJECTS APPLIED TO COLUMN SUBSETS OF THE ORIGINAL FEATURE SPACE

NOTES

THE ORDER OF THE COLUMNS IN THE TRANSFORMED FEATURE MATRIX FOLLOWS THE ORDER OF HOW THE COLUMNS ARE SPECIFIED IN THE TRANSFORMERS LIST COLUMNS OF THE ORIGINAL FEATURE MATRIX THAT ARE NOT SPECIFIED ARE DROPPED FROM THE RESULTING TRANSFORMED FEATURE MATRIX UNLESS SPECIFIED IN THE PASSTHROUGH KEYWORD THOSE COLUMNS SPECIFIED WITHPASSTHROUGH ARE ADDED AT THE RIGHT TO THE OUTPUT OF THE TRANSFORMERS

EXAMPLES

```
import numpy as np
from sklearn.compose import ColumnTransformer
from sklearn.preprocessing import Normalizer
ct = ColumnTransformer(
    [ ('norm1', Normalizer(), 0),
      ('norm2', Normalizer(), 1),
      ('x', np.array([1, 2, 2],
                    [1, 1, 0],
                    [1, 1, 0])
        ),
    ], remainder='passthrough')
normalizer_scales = ct.get_feature_names_out()
normalizer_scales
```

SCIKITLEARN USER GUIDE RELEASE 0213

IS APPLIED FOR THE TWO FIRST AND TWO LAST ELEMENTS OF EACH  
ROW INDEPENDENTLY

CTFITTRANSFORMX

ARRAY0 1 05 05

05 05 0 1

METHODS

FITSELF X Y FIT ALL TRANSFORMERS USING X

FITTRANSFORM SELF X Y FIT ALL TRANSFORMERS TRANSFORM THE DATA AND CONCATENATE  
RESULTS

GETFEATURENAMES SELF GET FEATURE NAMES FROM ALL TRANSFORMERS

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

SETPARAMS SELF KWARGS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF X TRANSFORM X SEPARATELY BY EACH TRANSFORMER CONCATENATE  
RESULTS

INIT SELFTRANSFORMERS REMAINDER'DROP' SPARSETHRESHOLD03 NJOBSNONE TRANSFORMERWEIGHTSNONE VERBOSEFALSE

FITSELFXYNONE

FIT ALL TRANSFORMERS USING X

PARAMETERS

XARRAYLIKE OR DATAFRAME OF SHAPE NSAMPLES NFEATURES INPUT DATA OF WHICH SPECIFIED  
SUBSETS ARE USED TO FIT THE TRANSFORMERS

YARRAYLIKE SHAPE NSAMPLES OPTIONAL TARGETS FOR SUPERVISED LEARNING

RETURNS

SELF COLUMNTRANSFORMER THIS ESTIMATOR

FITTRANSFORM SELFXYNONE

FIT ALL TRANSFORMERS TRANSFORM THE DATA AND CONCATENATE RESULTS

PARAMETERS

XARRAYLIKE OR DATAFRAME OF SHAPE NSAMPLES NFEATURES INPUT DATA OF WHICH SPECIFIED  
SUBSETS ARE USED TO FIT THE TRANSFORMERS

YARRAYLIKE SHAPE NSAMPLES OPTIONAL TARGETS FOR SUPERVISED LEARNING

RETURNS

XT ARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES SUMNCOMPONENTS HSTACK OF RESULTS  
OF TRANSFORMERS SUMNCOMPONENTS IS THE SUM OF NCOMPONENTS OUTPUT DIMENSION OVER  
TRANSFORMERS IF ANY RESULT IS A SPARSE MATRIX EVERYTHING WILL BE CONVERTED TO SPARSE MATRICES

GETFEATURENAMES SELF

GET FEATURE NAMES FROM ALL TRANSFORMERS

RETURNS

FEATURENAMES LIST OF STRINGS NAMES OF THE FEATURES PRODUCED BY TRANSFORM

1526 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

GETPARAMS SELFDEEPTREE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

NAMEDTRANSFORMERS

ACCESS THE FITTED TRANSFORMER BY NAME

READONLY ATTRIBUTE TO ACCESS ANY TRANSFORMER BY GIVEN NAME KEYS ARE TRANSFORMER NAMES AND VALUES ARE THE FITTED TRANSFORMER OBJECTS

SETPARAMS SELFKWARGS

SET THE PARAMETERS OF THIS ESTIMATOR

VALID PARAMETER KEYS CAN BE LISTED WITH GETPARAMS

RETURNS

SELF

TRANSFORM SELF

TRANSFORM X SEPARATELY BY EACH TRANSFORMER CONCATENATE RESULTS

PARAMETERS

XARRAYLIKE OR DATAFRAME OF SHAPE NSAMPLES NFEATURES THE DATA TO BE TRANSFORMED

BY SUBSET

RETURNS

XT ARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES SUMNCOMPONENTS HSTACK OF RESULTS

OF TRANSFORMERS SUMNCOMPONENTS IS THE SUM OF NCOMPONENTS OUTPUT DIMENSION OVER

TRANSFORMERS IF ANY RESULT IS A SPARSE MATRIX EVERYTHING WILL BE CONVERTED TO SPARSE MATRICES

EXAMPLES USING SKLEARNCOMPOSECOLUMNTRANSFORMER

- COLUMN TRANSFORMER WITH MIXED TYPES
- COLUMN TRANSFORMER WITH HETEROGENEOUS DATA SOURCES

652SKLEARNCOMPOSE TRANSFORMEDTARGETREGRESSOR

CLASSSKLEARNCOMPOSE TRANSFORMEDTARGETREGRESSOR REGRESSORNONE TRANSFORMERNONE

FUNCNONE INVERSEFUNCNONE

CHECKINVERSETRUE

METAESTIMATOR TO REGRESS ON A TRANSFORMED TARGET

USEFUL FOR APPLYING A NONLINEAR TRANSFORMATION IN REGRESSION PROBLEMS THIS TRANSFORMATION CAN BE GIVEN AS A TRANSFORMER SUCH AS THE QUANTILETRANSFORMER OR AS A FUNCTION AND ITS INVERSE SUCH AS LOG AND EXP

THE COMPUTATION DURING FIT IS

REGRESSORFITX FUNCY

655SKLEARNCOMPOSE COMPOSITE ESTIMATORS 1527

SCIKITLEARN USER GUIDE RELEASE 0213

OR

REGRESSORFITX TRANSFORMERTRANSFORMY

THE COMPUTATION DURING PREDICT IS

INVERSEFUNCREGRESSORPREDICTX

OR

TRANSFORMERINVERSETRANSFORMREGRESSORPREDICTX

READ MORE IN THE USER GUIDE

PARAMETERS

REGRESSOR OBJECT DEFAULTLINEARREGRESSION REGRESSOR OBJECT SUCH AS DERIVED FROM

REGRESSORMIXIN THIS REGRESSOR WILL AUTOMATICALLY BE CLONED EACH TIME PRIOR TO FITTING

TRANSFORMER OBJECT DEFAULTNONE ESTIMATOR OBJECT SUCH AS DERIVED FROM

TRANSFORMERMIXIN CANNOT BE SET AT THE SAME TIME AS FUNC ANDINVERSEFUNC IF

TRANSFORMER ISNONE AS WELL ASFUNC ANDINVERSEFUNC THE TRANSFORMER WILL BE

AN IDENTITY TRANSFORMER NOTE THAT THE TRANSFORMER WILL BE CLONED DURING FITTING ALSO THE

TRANSFORMER IS RESTRICTING YTO BE A NUMPY ARRAY

FUNC FUNCTION OPTIONAL FUNCTION TO APPLY TO YBEFORE PASSING TO FIT CANNOT BE SET AT THE

SAME TIME AS TRANSFORMER THE FUNCTION NEEDS TO RETURN A 2DIMENSIONAL ARRAY IF FUNC

ISNONE THE FUNCTION USED WILL BE THE IDENTITY FUNCTION

INVERSEFUNC FUNCTION OPTIONAL FUNCTION TO APPLY TO THE PREDICTION OF THE REGRESSOR CAN

NOT BE SET AT THE SAME TIME AS TRANSFORMER AS WELL THE FUNCTION NEEDS TO RETURN A

2DIMENSIONAL ARRAY THE INVERSE FUNCTION IS USED TO RETURN PREDICTIONS TO THE SAME SPACE OF

THE ORIGINAL TRAINING LABELS

CHECKINVERSE BOOL DEFAULTTRUE WHETHER TO CHECK THAT TRANSFORM FOLLOWED BY

INVERSETRANSFORM ORFUNC FOLLOWED BY INVERSEFUNC LEADS TO THE ORIGINAL TAR

GETS

ATTRIBUTES

REGRESSOR OBJECT FITTED REGRESSOR

TRANSFORMER OBJECT TRANSFORMER USED IN FIT ANDPREDICT

NOTES

INTERNALLY THE TARGET YIS ALWAYS CONVERTED INTO A 2DIMENSIONAL ARRAY TO BE USED BY SCIKITLEARN TRANSFORMERS

AT THE TIME OF PREDICTION THE OUTPUT WILL BE RESHAPED TO A HAVE THE SAME NUMBER OF DIMENSIONS AS Y

SEEEXAMPLESCOMPOSEPLOTTRANSFORMEDTARGETPY

EXAMPLES

```
import numpy as np
from sklearn.linear_model import LinearRegression
from sklearn.compose import TransformTargetRegressor
tt = TransformTargetRegressor(LinearRegression(),
                              func=np.log,
                              inverse_func=np.exp)
```

1528 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

X NPARANGE4RESHAPE1 1

Y NPEXP2 XRAVEL

TTFITX Y

TRANSFORMEDTARGETREGRESSOR

TTSCOREX Y

10

TTREGRESSORCOEF

ARRAY2

METHODS

FITSELF X Y SAMPLEWEIGHT FIT THE MODEL ACCORDING TO THE GIVEN TRAINING DATA

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PREDICT SELF X PREDICT USING THE BASE REGRESSOR APPLYING INVERSE

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PRE

DICTION

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELF REGRESSORNONE TRANSFORMERNONE FUNCNONE INVERSEFUNCNONE

CHECKINVERSETRUE

FITSELFXYSAMPLEWEIGHTNONE

FIT THE MODEL ACCORDING TO THE GIVEN TRAINING DATA

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES TRAINING VECTOR WHERE

NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

YARRAYLIKE SHAPE NSAMPLES TARGET VALUES

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL ARRAY OF WEIGHTS THAT ARE ASSIGNED

TO INDIVIDUAL SAMPLES IF NOT PROVIDED THEN EACH SAMPLE IS GIVEN UNIT WEIGHT

RETURNS

SELF OBJECT

GETPARAMS SELFDEEPTTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELF

PREDICT USING THE BASE REGRESSOR APPLYING INVERSE

THE REGRESSOR IS USED TO PREDICT AND THE INVERSEFUNC ORINVERSESTRANSFORM IS APPLIED BEFORE

RETURNING THE PREDICTION

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES SAMPLES

65SKLEARNCOMPOSE COMPOSITE ESTIMATORS 1529

SCIKITLEARN USER GUIDE RELEASE 0213

RETURNS

YHAT ARRAY SHAPE NSAMPLES PREDICTED VALUES

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{U}{V}$  WHERE U IS THE RESIDUAL SUM OF SQUARES  $Y_{TRUE} - Y_{PRED}$

2SUM AND V IS THE TOTAL SUM OF SQUARES  $Y_{TRUE} - Y_{TRUEMEAN}$  2SUM THE BEST POSSIBLE SCORE

IS 10 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS

PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 00

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES FOR SOME ESTIMATORS THIS MAY

BE A PRECOMPUTED KERNEL MATRIX INSTEAD SHAPE NSAMPLES NSAMPLESFITTED WHERE

NSAMPLESFITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE VALUES FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT R2 OF SELFpredictX wrt Y

NOTES

THE R2 SCORE USED WHEN CALLING SCORE ON A REGRESSOR WILL USE MULTIOUTPUTUNIFORMAVERAGE

FROM VERSION 023 TO KEEP CONSISTENT WITH METRICSR2SCORE THIS WILL INFLUENCE THE SCORE

METHOD OF ALL THE MULTIOUTPUT REGRESSORS EXCEPT FOR MULTIOUTPUTMULTIOUTPUTREGRESSOR

TO SPECIFY THE DEFAULT VALUE MANUALLY AND AVOID THE WARNING PLEASE EITHER CALL METRICSR2SCORE

DIRECTLY OR MAKE A CUSTOM SCORER WITH METRICSMakescorer THE BUILTIN SCORER R2 USES

MULTIOUTPUTUNIFORMAVERAGE

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE

PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT

OF A NESTED OBJECT

RETURNS

SELF

EXAMPLES USING SKLEARNCOMPOSETRANSFORMEDTARGETREGRESSOR

•EFFECT OF TRANSFORMING THE TARGETS IN REGRESSION MODEL

COMPOSEMAKECOLUMNTRANSFORMER CONSTRUCT A COLUMNTRANSFORMER FROM THE GIVEN TRANSFORM

ERS

653SKLEARNCOMPOSE MAKECOLUMNTRANSFORMER

SKLEARNCOMPOSE MAKECOLUMNTRANSFORMER TRANSFORMERS KWARGS

CONSTRUCT A COLUMNTRANSFORMER FROM THE GIVEN TRANSFORMERS

1530 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

THIS IS A SHORTHAND FOR THE COLUMNTRANSFORMER CONSTRUCTOR IT DOES NOT REQUIRE AND DOES NOT PERMIT NAMING THE TRANSFORMERS INSTEAD THEY WILL BE GIVEN NAMES AUTOMATICALLY BASED ON THEIR TYPES IT ALSO DOES NOT ALLOW WEIGHTING WITH TRANSFORMERWEIGHTS

PARAMETERS

TRANSFORMERS TUPLES OF TRANSFORMERS AND COLUMN SELECTIONS

REMAINDER ‘DROP’ ‘PASSTHROUGH’ OR ESTIMATOR DEFAULT ‘DROP’ BY DEFAULT ONLY THE SPECIFIED COLUMNS IN TRANSFORMERS ARE TRANSFORMED AND COMBINED IN THE OUTPUT AND THE NONSPECIFIED COLUMNS ARE DROPPED DEFAULT OF DROP BY SPECIFYING

REMAINDERPASSTHROUGH ALL REMAINING COLUMNS THAT WERE NOT SPECIFIED IN TRANSFORMERS WILL BE AUTOMATICALLY PASSED THROUGH THIS SUBSET OF COLUMNS IS CONCATENATED WITH THE OUTPUT OF THE TRANSFORMERS BY SETTING REMAINDER TO BE AN ESTIMATOR THE REMAINING NONSPECIFIED COLUMNS WILL USE THE REMAINDER ESTIMATOR THE ESTIMATOR MUST SUPPORT FIT ANDTRANSFORM

SPARSETHRESHOLD FLOAT DEFAULT 03 IF THE TRANSFORMED OUTPUT CONSISTS OF A MIX OF SPARSE AND DENSE DATA IT WILL BE STACKED AS A SPARSE MATRIX IF THE DENSITY IS LOWER THAN THIS VALUE USE SPARSETHRESHOLD0 TO ALWAYS RETURN DENSE WHEN THE TRANSFORMED OUTPUT CONSISTS OF ALL SPARSE OR ALL DENSE DATA THE STACKED RESULT WILL BE SPARSE OR DENSE RESPECTIVELY AND THIS KEYWORD WILL BE IGNORED

NJOBS INT OR NONE OPTIONAL DEFAULTNONE NUMBER OF JOBS TO RUN IN PARALLEL NONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL PROCESSORS SEE

GLOSSARY FOR MORE DETAILS

VERBOSE BOOLEAN OPTIONALDEFAULTFALSE IF TRUE THE TIME ELAPSED WHILE FITTING EACH TRANSFORMER WILL BE PRINTED AS IT IS COMPLETED

RETURNS

CTCOLUMNTRANSFORMER

SEE ALSO

SKLEARNCOMPOSECOLUMNTRANSFORMER CLASS THAT ALLOWS COMBINING THE OUTPUTS OF MULTIPLE TRANSFORMER OBJECTS USED ON COLUMN SUBSETS OF THE DATA INTO A SINGLE FEATURE SPACE

EXAMPLES

FROM SKLEARNPREPROCESSING IMPORT STANDARDSCALER ONEHOTENCODER

FROM SKLEARNCOMPOSE IMPORT MAKECOLUMNTRANSFORMER

MAKECOLUMNTRANSFORMER

STANDARDSCALER NUMERICALCOLUMN

ONEHOTENCODER CATEGORICALCOLUMN

COLUMNTRANSFORMERNJOBSNONE REMAINDERDROP SPARSETHRESHOLD03

TRANSFORMERWEIGHTSNONE

TRANSFORMERSSTANDARDSCALER

STANDARDSCALER

NUMERICALCOLUMN

ONEHOTENCODER

ONEHOTENCODER

CATEGORICALCOLUMN VERBOSEFALSE

65SKLEARNCOMPOSE COMPOSITE ESTIMATORS 1531

SCIKITLEARN USER GUIDE RELEASE 0213

66SKLEARNCOVARIANCE COVARIANCE ESTIMATORS

THE SKLEARNCOVARIANCE MODULE INCLUDES METHODS AND ALGORITHMS TO ROBUSTLY ESTIMATE THE COVARIANCE OF FEATURES GIVEN A SET OF POINTS THE PRECISION MATRIX DEFINED AS THE INVERSE OF THE COVARIANCE IS ALSO ESTIMATED COVARIANCE ESTIMATION IS CLOSELY RELATED TO THE THEORY OF GAUSSIAN GRAPHICAL MODELS

USER GUIDE SEE THE COVARIANCE ESTIMATION SECTION FOR FURTHER DETAILS

COVARIANCEEMPIRICALCOVARIANCE MAXIMUM LIKELIHOOD COVARIANCE ESTIMATOR

COVARIANCEELLIPTICENVELOPE AN OBJECT FOR DETECTING OUTLIERS IN A GAUSSIAN DISTRIBUTED DATASET

COVARIANCEGRAPHICALASSO ALPHA MODE SPARSE INVERSE COVARIANCE ESTIMATION WITH AN L1PENALIZED ESTIMATOR

COVARIANCEGRAPHICALASSOCV ALPHAS SPARSE INVERSE COVARIANCE W CROSSVALIDATED CHOICE OF THE L1 PENALTY

COVARIANCELEDOITWOLF STOREPRECISION LEDOITWOLF ESTIMATOR

COVARIANCEMINCOVDDET STOREPRECISION MINIMUM COVARIANCE DETERMINANT MCD ROBUST ESTIMATOR OF COVARIANCE

COVARIANCEOAS STOREPRECISION ORACLE APPROXIMATING SHRINKAGE ESTIMATOR

COVARIANCESHRUNKCOVARIANCE COVARIANCE ESTIMATOR WITH SHRINKAGE

661SKLEARNCOVARIANCE EMPIRICALCOVARIANCE

CLASSSKLEARNCOVARIANCE EMPIRICALCOVARIANCE STOREPRECISIONTRUE AS

SUMECENTEREDFALSE

MAXIMUM LIKELIHOOD COVARIANCE ESTIMATOR

READ MORE IN THE USER GUIDE

PARAMETERS

STOREPRECISION BOOL SPECIFIES IF THE ESTIMATED PRECISION IS STORED

ASSUMECENTERED BOOL IF TRUE DATA ARE NOT CENTERED BEFORE COMPUTATION USEFUL WHEN WORKING WITH DATA WHOSE MEAN IS ALMOST BUT NOT EXACTLY ZERO IF FALSE DEFAULT DATA ARE CENTERED BEFORE COMPUTATION

ATTRIBUTES

LOCATION ARRAYLIKE SHAPE NFEATURES ESTIMATED LOCATION IE THE ESTIMATED MEAN

COVARIANCE 2D NDARRAY SHAPE NFEATURES NFEATURES ESTIMATED COVARIANCE MATRIX

PRECISION 2D NDARRAY SHAPE NFEATURES NFEATURES ESTIMATED PSEUDOINVERSE MATRIX

STORED ONLY IF STOREPRECISION IS TRUE

EXAMPLES

```
import numpy as np
from sklearn.covariance import EmpiricalCovariance
from sklearn.datasets import make_gaussian_quantiles
realcov = np.array(
    [
        [3.4, 0.0, 0.0],
        [0.0, 3.4, 0.0],
        [0.0, 0.0, 3.4]
    ])
rng = np.random.RandomState(0)
X = rng.multivariate_normal(mean0, realcov, 1000)
```

1532 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

COVREALCOV  
SIZE500

COV EMPIRICALCOVARIANCEFITX  
COVCOVARIANCE  
ARRAY07569 02818  
02818 03928  
COVLOCATION  
ARRAY00622 00193

METHODS

ERRORNORM SELF COMPCOV NORM SCALING COMPUTES THE MEAN SQUARED ERROR BETWEEN TWO COVARIANCE ESTIMATORS

FITSELF X Y FITS THE MAXIMUM LIKELIHOOD ESTIMATOR COVARIANCE MODEL ACCORDING TO THE GIVEN TRAINING DATA AND PARAMETERS

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

GETPRECISION SELF GETTER FOR THE PRECISION MATRIX

MAHALANOBIS SELF X COMPUTES THE SQUARED MAHALANOBIS DISTANCES OF GIVEN OBSERVATIONS

SCORE SELF XTEST Y COMPUTES THE LOGLIKELIHOOD OF A GAUSSIAN DATA SET WITH SELFCOVARIANCE AS AN ESTIMATOR OF ITS COVARIANCE MATRIX

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELFSTOREPRECISIONTRUE ASSUMECENTEREDFALSE

ERRORNORM SELFCOMPCOV NORM'FROBENIUS' SCALINGTRUE SQUAREDTRUE  
COMPUTES THE MEAN SQUARED ERROR BETWEEN TWO COVARIANCE ESTIMATORS IN THE SENSE OF THE FROBENIUS NORM

PARAMETERS

COMPCOV ARRAYLIKE SHAPE NFEATURES NFEATURES THE COVARIANCE TO COMPARE WITH

NORM STR THE TYPE OF NORM USED TO COMPUTE THE ERROR AVAILABLE ERROR TYPES 'FROBENIUS' DEFAULT SQRTTRATA 'SPECTRAL' SQRTMAXEIGENVALUESATA WHERE A IS THE ERROR

COMPCOV SELFCOVARIANCE

SCALING BOOL IF TRUE DEFAULT THE SQUARED ERROR NORM IS DIVIDED BY NFEATURES IF FALSE THE SQUARED ERROR NORM IS NOT RESCALED

SQUARED BOOL WHETHER TO COMPUTE THE SQUARED ERROR NORM OR THE ERROR NORM IF TRUE DEFAULT THE SQUARED ERROR NORM IS RETURNED IF FALSE THE ERROR NORM IS RETURNED

RETURNS

THE MEAN SQUARED ERROR IN THE SENSE OF THE FROBENIUS NORM BETWEEN SELF ANDCOMPCOV COVARIANCE ESTIMATORS

FITSELFXYNONE

FITS THE MAXIMUM LIKELIHOOD ESTIMATOR COVARIANCE MODEL ACCORDING TO THE GIVEN TRAINING DATA AND PARAMETERS

PARAMETERS

66SKLEARNCOVARIANCE COVARIANCE ESTIMATORS 1533

SCIKITLEARN USER GUIDE RELEASE 0213

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

YNOT USED PRESENT FOR API CONSISTENCE PURPOSE

RETURNS

SELF OBJECT

GETPARAMS SELFDEEPTREE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

GETPRECISION SELF

GETTER FOR THE PRECISION MATRIX

RETURNS

PRECISION ARRAYLIKE THE PRECISION MATRIX ASSOCIATED TO THE CURRENT COVARIANCE OBJECT

MAHALANOBIS SELF

COMPUTES THE SQUARED MAHALANOBIS DISTANCES OF GIVEN OBSERVATIONS

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES THE OBSERVATIONS THE MAHALANOBIS DISTANCES OF THE WHICH WE COMPUTE OBSERVATIONS ARE ASSUMED TO BE DRAWN FROM THE SAME DISTRIBUTION THAN THE DATA USED IN FIT

RETURNS

DIST ARRAY SHAPE NSAMPLES SQUARED MAHALANOBIS DISTANCES OF THE OBSERVATIONS

SCORESELFXTTEST YNONE

COMPUTES THE LOGLIKELIHOOD OF A GAUSSIAN DATA SET WITH SELFCOVARIANCE AS AN ESTIMATOR OF ITS COVARIANCE MATRIX

PARAMETERS

XTTEST ARRAYLIKE SHAPE NSAMPLES NFEATURES TEST DATA OF WHICH WE COMPUTE THE LIKELIHOOD WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES XTTEST IS ASSUMED TO BE DRAWN FROM THE SAME DISTRIBUTION THAN THE DATA USED IN FIT INCLUDING CENTERING

YNOT USED PRESENT FOR API CONSISTENCE PURPOSE

RETURNS

RESFLOAT THE LIKELIHOOD OF THE DATA SET WITH SELFCOVARIANCE AS AN ESTIMATOR OF ITS COVARIANCE MATRIX

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

1534 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

RETURNS

SELF

EXAMPLES USING SKLEARNCOVARIANCEEMPIRICALCOVARIANCE

•ROBUST COVARIANCE ESTIMATION AND MAHALANOBIS DISTANCES RELEVANCE

•ROBUST VS EMPIRICAL COVARIANCE ESTIMATE

662SKLEARNCOVARIANCE ELLIPTICENVELOPE

CLASSSSKLEARNCOVARIANCE ELLIPTICENVELOPE STOREPRECISIONTRUE ASSUMECENTEREDFALSE

SUPPORTFRACTIONNONE CONTAMINATION01

RANDOMSTATENONE

AN OBJECT FOR DETECTING OUTLIERS IN A GAUSSIAN DISTRIBUTED DATASET

READ MORE IN THE USER GUIDE

PARAMETERS

STOREPRECISION BOOLEAN OPTIONAL DEFAULTTRUE SPECIFY IF THE ESTIMATED PRECISION IS STORED

ASSUMECENTERED BOOLEAN OPTIONAL DEFAULTFALSE IF TRUE THE SUPPORT OF ROBUST LOCATION AND

COVARIANCE ESTIMATES IS COMPUTED AND A COVARIANCE ESTIMATE IS RECOMPUTED FROM IT WITHOUT

CENTERING THE DATA USEFUL TO WORK WITH DATA WHOSE MEAN IS SIGNIFICANTLY EQUAL TO ZERO BUT IS

NOT EXACTLY ZERO IF FALSE THE ROBUST LOCATION AND COVARIANCE ARE DIRECTLY COMPUTED WITH THE

FASTMCD ALGORITHM WITHOUT ADDITIONAL TREATMENT

SUPPORTFRACTION FLOAT IN 0 1 OPTIONAL DEFAULTNONE THE PROPORTION OF POINTS TO BE

INCLUDED IN THE SUPPORT OF THE RAW MCD ESTIMATE IF NONE THE MINIMUM VALUE OF SUP

PORTFRACTION WILL BE USED WITHIN THE ALGORITHM NSAMPLE NFEATURES 1

2

CONTAMINATION FLOAT IN 0 05 OPTIONAL DEFAULT01 THE AMOUNT OF CONTAMINATION OF THE

DATA SET IE THE PROPORTION OF OUTLIERS IN THE DATA SET

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE THE SEED OF THE

PSEUDO RANDOM NUMBER GENERATOR TO USE WHEN SHUFFLING THE DATA IF INT RANDOMSTATE IS

THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS

THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE

INSTANCE USED BY NPRANDOM

ATTRIBUTES

LOCATION ARRAYLIKE SHAPE NFEATURES ESTIMATED ROBUST LOCATION

COVARIANCE ARRAYLIKE SHAPE NFEATURES NFEATURES ESTIMATED ROBUST COVARIANCE MATRIX

PRECISION ARRAYLIKE SHAPE NFEATURES NFEATURES ESTIMATED PSEUDO INVERSE MATRIX STORED

ONLY IF STOREPRECISION IS TRUE

SUPPORT ARRAYLIKE SHAPE NSAMPLES A MASK OF THE OBSERVATIONS THAT HAVE BEEN USED TO

COMPUTE THE ROBUST ESTIMATES OF LOCATION AND SHAPE

OFFSET FLOAT OFFSET USED TO DEFINE THE DECISION FUNCTION FROM THE RAW SCORES WE HAVE THE RELA

TIONDECISIONFUNCTION SCORESAMPLES OFFSET THE OFFSET DEPENDS

ON THE CONTAMINATION PARAMETER AND IS DEFINED IN SUCH A WAY WE OBTAIN THE EXPECTED NUMBER

OF OUTLIERS SAMPLES WITH DECISION FUNCTION 0 IN TRAINING

SEE ALSO

665SKLEARNCOVARIANCE COVARIANCE ESTIMATORS 1535

SCIKITLEARN USER GUIDE RELEASE 0213  
EMPIRICALCOVARIANCE MINCOVDET  
NOTES  
OUTLIER DETECTION FROM COVARIANCE ESTIMATION MAY BREAK OR NOT PERFORM WELL IN HIGHDIMENSIONAL SETTINGS IN PARTICULAR ONE WILL ALWAYS TAKE CARE TO WORK WITH NSAMPLES NFEATURES 2  
REFERENCES  
R68AE096DA0E41  
EXAMPLES  
IMPORT NUMPY AS NP  
FROM SKLEARNCOVARIANCE IMPORT ELLIPTICENVELOPE  
TRUECOV NPARRAY8 3  
3 4  
X NPRANDOMRANDOMSTATE0MULTIVARIATENORMALMEAN0 0  
COVTRUECOV  
SIZE500  
COV ELLIPTICENVELOPERANDOMSTATE0FITX  
PREDICT RETURNS 1 FOR AN INLIER AND 1 FOR AN OUTLIER  
COVPREDICT0 0  
3 3  
ARRAY 1 1  
COVCOVARIANCE  
ARRAY07411 02535  
02535 03053  
COVLOCATION  
ARRAY00813 00427  
METHODS  
CORRECTCOVARIANCE SELF DATA APPLY A CORRECTION TO RAW MINIMUM COVARIANCE DETERMINANT ESTIMATES  
DECISIONFUNCTION SELF X RAWVALUES COMPUTE THE DECISION FUNCTION OF THE GIVEN OBSERVATIONS  
ERRORNORM SELF COMPCOV NORM SCALING COMPUTES THE MEAN SQUARED ERROR BETWEEN TWO COVARIANCE ESTIMATORS  
FITSELF X Y FIT THE ELLIPTICENVELOPE MODEL  
FITPREDICT SELF X Y PERFORMS FIT ON X AND RETURNS LABELS FOR X  
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR  
GETPRECISION SELF GETTER FOR THE PRECISION MATRIX  
MAHALANOBIS SELF X COMPUTES THE SQUARED MAHALANOBIS DISTANCES OF GIVEN OBSERVATIONS  
PREDICT SELF X PREDICT THE LABELS 1 INLIER 1 OUTLIER OF X ACCORDING TO THE FITTED MODEL  
REWEIGHTCOVARIANCE SELF DATA REWEIGHT RAW MINIMUM COVARIANCE DETERMINANT ESTIMATES  
CONTINUED ON NEXT PAGE  
1536 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

TABLE 634 – CONTINUED FROM PREVIOUS PAGE

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS

SCORESAMPLES SELF X COMPUTE THE NEGATIVE MAHALANOBIS DISTANCES

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELFSTOREPRECISIONTRUE ASSUMECENTEREDFALSE SUPPORTFRACTIONNONE CONTAMINATION01 RANDOMSTATENONE

CORRECTCOVARIANCE SELFDATA

APPLY A CORRECTION TO RAW MINIMUM COVARIANCE DETERMINANT ESTIMATES

CORRECTION USING THE EMPIRICAL CORRECTION FACTOR SUGGESTED BY ROUSSEEUW AND VAN DRIESSEN IN RVD PARAMETERS

DATA ARRAYLIKE SHAPE NSAMPLES NFEATURES THE DATA MATRIX WITH P FEATURES AND N SAMPLES THE DATA SET MUST BE THE ONE WHICH WAS USED TO COMPUTE THE RAW ESTIMATES

RETURNS

COVARIANCECORRECTED ARRAYLIKE SHAPE NFEATURES NFEATURES CORRECTED ROBUST COVARIANCE ESTIMATE

REFERENCES

RVD

DECISIONFUNCTION SELFXXRAWVALUESNONE

COMPUTE THE DECISION FUNCTION OF THE GIVEN OBSERVATIONS

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES

RAWVALUES BOOL OPTIONAL WHETHER OR NOT TO CONSIDER RAW MAHALANOBIS DISTANCES AS THE DECISION FUNCTION MUST BE FALSE DEFAULT FOR COMPATIBILITY WITH THE OTHERS OUTLIER DETECTION TOOLS

DEPRECATED SINCE VERSION 020 RAWVALUES HAS BEEN DEPRECATED IN 020 AND WILL BE REMOVED IN 022

RETURNS

DECISION ARRAYLIKE SHAPE NSAMPLES DECISION FUNCTION OF THE SAMPLES IT IS EQUAL TO THE SHIFTED MAHALANOBIS DISTANCES THE THRESHOLD FOR BEING AN OUTLIER IS 0 WHICH ENSURES A COMPATIBILITY WITH OTHER OUTLIER DETECTION ALGORITHMS

ERRORNORM SELFCOMPCOV NORM'FROBENIUS' SCALINGTRUE SQUAREDTRUE

COMPUTES THE MEAN SQUARED ERROR BETWEEN TWO COVARIANCE ESTIMATORS IN THE SENSE OF THE FROBENIUS NORM

PARAMETERS

COMPCOV ARRAYLIKE SHAPE NFEATURES NFEATURES THE COVARIANCE TO COMPARE WITH

NORM STR THE TYPE OF NORM USED TO COMPUTE THE ERROR AVAILABLE ERROR TYPES 'FROBENIUS' DEFAULT SQRTTRATA 'SPECTRAL' SQRTMAXEIGENVALUESATA WHERE A IS THE ERROR

COMPCOV SELFCOVARIANCE

SCALING BOOL IF TRUE DEFAULT THE SQUARED ERROR NORM IS DIVIDED BY NFEATURES IF FALSE THE SQUARED ERROR NORM IS NOT RESCALED

66SKLEARNCOVARIANCE COVARIANCE ESTIMATORS 1537

SCIKITLEARN USER GUIDE RELEASE 0213

SQUARED BOOL WHETHER TO COMPUTE THE SQUARED ERROR NORM OR THE ERROR NORM IF TRUE  
DEFAULT THE SQUARED ERROR NORM IS RETURNED IF FALSE THE ERROR NORM IS RETURNED

RETURNS  
THE MEAN SQUARED ERROR IN THE SENSE OF THE FROBENIUS NORM BETWEEN  
SELF ANDCOMP COV COVARIANCE ESTIMATORS

FITSELFXYNONE

FIT THE ELLIPTICENVELOPE MODEL

PARAMETERS  
XNUMPY ARRAY OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES TRAINING DATA  
YIGNORED NOT USED PRESENT FOR API CONSISTENCY BY CONVENTION

FITPREDICT SELFXYNONE

PERFORMS FIT ON X AND RETURNS LABELS FOR X

RETURNS 1 FOR OUTLIERS AND 1 FOR INLIERS

PARAMETERS  
XNDARRAY SHAPE NSAMPLES NFEATURES INPUT DATA  
YIGNORED NOT USED PRESENT FOR API CONSISTENCY BY CONVENTION

RETURNS  
YNDARRAY SHAPE NSAMPLES 1 FOR INLIERS 1 FOR OUTLIERS

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS  
DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBOBJECTS THAT ARE ESTIMATORS

RETURNS  
PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

GETPRECISION SELF

GETTER FOR THE PRECISION MATRIX

RETURNS  
PRECISION ARRAYLIKE THE PRECISION MATRIX ASSOCIATED TO THE CURRENT COVARIANCE OBJECT

MAHALANOBIS SELF

COMPUTES THE SQUARED MAHALANOBIS DISTANCES OF GIVEN OBSERVATIONS

PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES THE OBSERVATIONS THE MAHALANOBIS DISTANCES  
OF THE WHICH WE COMPUTE OBSERVATIONS ARE ASSUMED TO BE DRAWN FROM THE SAME DISTRIBU  
TION THAN THE DATA USED IN FIT

RETURNS  
DIST ARRAY SHAPE NSAMPLES SQUARED MAHALANOBIS DISTANCES OF THE OBSERVATIONS

PREDICTSELF

PREDICT THE LABELS 1 INLIER 1 OUTLIER OF X ACCORDING TO THE FITTED MODEL

1538 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES

RETURNS

ISINLIER ARRAY SHAPE NSAMPLES RETURNS 1 FOR ANOMALIESOUTLIERS AND 1 FOR INLIERS

REWEIGHTCOVARIANCE SELFDATA

REWEIGHT RAW MINIMUM COVARIANCE DETERMINANT ESTIMATES

REWEIGHT OBSERVATIONS USING ROUSSEEUW'S METHOD EQUIVALENT TO DELETING OUTLYING OBSERVATIONS FROM THE DATA SET BEFORE COMPUTING LOCATION AND COVARIANCE ESTIMATES DESCRIBED IN RVDRIESSEN

PARAMETERS

DATA ARRAYLIKE SHAPE NSAMPLES NFEATURES THE DATA MATRIX WITH P FEATURES AND N SAMPLES THE DATA SET MUST BE THE ONE WHICH WAS USED TO COMPUTE THE RAW ESTIMATES

RETURNS

LOCATIONREWEIGHTED ARRAYLIKE SHAPE NFEATURES REWEIGHTED ROBUST LOCATION ESTIMATE

COVARIANCEREWEIGHTED ARRAYLIKE SHAPE NFEATURES NFEATURES REWEIGHTED ROBUST COVARIANCE ESTIMATE

SUPPORTREWEIGHTED ARRAYLIKE TYPE BOOLEAN SHAPE NSAMPLES A MASK OF THE OBSERVATIONS THAT HAVE BEEN USED TO COMPUTE THE REWEIGHTED ROBUST LOCATION AND COVARIANCE ESTIMATES

REFERENCES

RVDRIESSEN

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS

IN MULTILABEL CLASSIFICATION THIS IS THE SUBSET ACCURACY WHICH IS A HARSH METRIC SINCE YOU REQUIRE FOR EACH SAMPLE THAT EACH LABEL SET BE CORRECTLY PREDICTED

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE LABELS FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT MEAN ACCURACY OF SELF PREDICT X WRT Y

SCORES SAMPLES SELF X

COMPUTE THE NEGATIVE MAHALANOBIS DISTANCES

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES

RETURNS

NEGATIVE MAHALANOBIS DISTANCES ARRAYLIKE SHAPE NSAMPLES OPPOSITE OF THE MAHALANOBIS DISTANCES

66SKLEARN COVARIANCE COVARIANCE ESTIMATORS 1539

SCIKITLEARN USER GUIDE RELEASE 0213

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

EXAMPLES USING SKLEARNCOVARIANCEELLIPTICENVELOPE

- COMPARING ANOMALY DETECTION ALGORITHMS FOR OUTLIER DETECTION ON TOY DATASETS
- OUTLIER DETECTION ON A REAL DATA SET

663SKLEARNCOVARIANCE GRAPHICALASSO

CLASSSSKLEARNCOVARIANCE GRAPHICALASSO ALPHA001 MODE'CD' TOL00001

ENETTOL00001 MAXITER100 VERBOSEFALSE

ASSUMECENTEREDFALSE

SPARSE INVERSE COVARIANCE ESTIMATION WITH AN L1PENALIZED ESTIMATOR

READ MORE IN THE USER GUIDE

PARAMETERS

ALPHA POSITIVE FLOAT DEFAULT 001 THE REGULARIZATION PARAMETER THE HIGHER ALPHA THE MORE REGULARIZATION THE SPARSER THE INVERSE COVARIANCE

MODE 'CD' 'LARS' DEFAULT 'CD' THE LASSO SOLVER TO USE COORDINATE DESCENT OR LARS USE LARS FOR VERY SPARSE UNDERLYING GRAPHS WHERE  $p \gg n$  ELSEWHERE PREFER CD WHICH IS MORE NUMERICALLY STABLE

TOLPOSITIVE FLOAT DEFAULT 1E4 THE TOLERANCE TO DECLARE CONVERGENCE IF THE DUAL GAP GOES BELOW THIS VALUE ITERATIONS ARE STOPPED

ENETTOL POSITIVE FLOAT OPTIONAL THE TOLERANCE FOR THE ELASTIC NET SOLVER USED TO CALCULATE THE DESCENT DIRECTION THIS PARAMETER CONTROLS THE ACCURACY OF THE SEARCH DIRECTION FOR A GIVEN COLUMN UPDATE NOT OF THE OVERALL PARAMETER ESTIMATE ONLY USED FOR MODE'CD'

MAXITER INTEGER DEFAULT 100 THE MAXIMUM NUMBER OF ITERATIONS

VERBOSE BOOLEAN DEFAULT FALSE IF VERBOSE IS TRUE THE OBJECTIVE FUNCTION AND DUAL GAP ARE PLOTTED AT EACH ITERATION

ASSUMECENTERED BOOLEAN DEFAULT FALSE IF TRUE DATA ARE NOT CENTERED BEFORE COMPUTATION USEFUL WHEN WORKING WITH DATA WHOSE MEAN IS ALMOST BUT NOT EXACTLY ZERO IF FALSE DATA ARE CENTERED BEFORE COMPUTATION

ATTRIBUTES

LOCATION ARRAYLIKE SHAPE NFEATURES ESTIMATED LOCATION IE THE ESTIMATED MEAN

COVARIANCE ARRAYLIKE SHAPE NFEATURES NFEATURES ESTIMATED COVARIANCE MATRIX

PRECISION ARRAYLIKE SHAPE NFEATURES NFEATURES ESTIMATED PSEUDO INVERSE MATRIX

NITER INT NUMBER OF ITERATIONS RUN

SEE ALSO

1540 CHAPTER 6 API REFERENCE

```
SCIKITLEARN USER GUIDE RELEASE 0213
GRAPHICALASSO GRAPHICALASSOCV
EXAMPLES
IMPORT NUMPY AS NP
FROM SKLEARNCOVARIANCE IMPORT GRAPHICALASSO
TRUECOV NPARRAY08 00 02 00
00 04 00 00
02 00 03 01
00 00 01 07
NPRANDOMSEED0
X NPRANDOMMULTIVARIATENORMALMEAN0 0 0 0
COVTRUECOV
SIZE200
COV GRAPHICALASSOFITX
NPAROUNDCOVCOVARIANCE DECIMALS3
ARRAY0816 0049 0218 0019
0049 0364 0017 0034
0218 0017 0322 0093
0019 0034 0093 069
NPAROUNDCOVLOCATION DECIMALS3
ARRAY0073 004 0038 0143
METHODS
ERRORNORM SELF COMPCOV NORM SCALING COMPUTES THE MEAN SQUARED ERROR BETWEEN TWO COVARIANCE ESTIMATORS
FITSELF X Y FITS THE GRAPHICALASSO MODEL TO X
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR
GETPRECISION SELF GETTER FOR THE PRECISION MATRIX
MAHALANOBIS SELF X COMPUTES THE SQUARED MAHALANOBIS DISTANCES OF GIVEN OBSERVATIONS
SCORE SELF XTEST Y COMPUTES THE LOGLIKELIHOOD OF A GAUSSIAN DATA SET WITH SELFCOVARIANCE AS AN ESTIMATOR OF ITS COVARIANCE MATRIX
SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR
INIT SELFALPHA001 MODE'CD' TOL00001 ENETTOL00001 MAXITER100 VERBOSEFALSE
ASSUMECENTEREDFALSE
ERRORNORM SELFCOMPCOV NORM'FROBENIUS' SCALINGTRUE SQUAREDTRUE
COMPUTES THE MEAN SQUARED ERROR BETWEEN TWO COVARIANCE ESTIMATORS IN THE SENSE OF THE FROBENIUS NORM
PARAMETERS
COMPCOV ARRAYLIKE SHAPE NFEATURES NFEATURES THE COVARIANCE TO COMPARE WITH
NORM STR THE TYPE OF NORM USED TO COMPUTE THE ERROR AVAILABLE ERROR TYPES 'FROBENIUS'
DEFAULT SQRTTRATA 'SPECTRAL' SQRTMAXEIGENVALUESATA WHERE A IS THE ERROR
COMPCOV SELFCOVARIANCE
SCALING BOOL IF TRUE DEFAULT THE SQUARED ERROR NORM IS DIVIDED BY NFEATURES IF FALSE
THE SQUARED ERROR NORM IS NOT RESCALED
66SKLEARNCOVARIANCE COVARIANCE ESTIMATORS 1541
```

SCIKITLEARN USER GUIDE RELEASE 0213

SQUARED BOOL WHETHER TO COMPUTE THE SQUARED ERROR NORM OR THE ERROR NORM IF TRUE  
DEFAULT THE SQUARED ERROR NORM IS RETURNED IF FALSE THE ERROR NORM IS RETURNED

RETURNS

THE MEAN SQUARED ERROR IN THE SENSE OF THE FROBENIUS NORM BETWEEN  
SELF ANDCOMP COV COVARIANCE ESTIMATORS

FITSELFXYNONE

FITS THE GRAPHICALASSO MODEL TO X

PARAMETERS

XNDARRAY SHAPE NSAMPLES NFEATURES DATA FROM WHICH TO COMPUTE THE COVARIANCE ES  
TIMATE

YIGNORED

GETPARAMS SELFDEEPTTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

GETPRECISION SELF

GETTER FOR THE PRECISION MATRIX

RETURNS

PRECISION ARRAYLIKE THE PRECISION MATRIX ASSOCIATED TO THE CURRENT COVARIANCE OBJECT

MAHALANOBIS SELF

COMPUTES THE SQUARED MAHALANOBIS DISTANCES OF GIVEN OBSERVATIONS

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES THE OBSERVATIONS THE MAHALANOBIS DISTANCES  
OF THE WHICH WE COMPUTE OBSERVATIONS ARE ASSUMED TO BE DRAWN FROM THE SAME DISTRIBU  
TION THAN THE DATA USED IN FIT

RETURNS

DIST ARRAY SHAPE NSAMPLES SQUARED MAHALANOBIS DISTANCES OF THE OBSERVATIONS

SCORESELFXTTEST YNONE

COMPUTES THE LOGLIKELIHOOD OF A GAUSSIAN DATA SET WITH SELF COVARIANCE AS AN ESTIMATOR OF ITS  
COVARIANCE MATRIX

PARAMETERS

XTTEST ARRAYLIKE SHAPE NSAMPLES NFEATURES TEST DATA OF WHICH WE COMPUTE THE  
LIKELIHOOD WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF  
FEATURES XTTEST IS ASSUMED TO BE DRAWN FROM THE SAME DISTRIBUTION THAN THE DATA USED IN  
FIT INCLUDING CENTERING

YNOT USED PRESENT FOR API CONSISTENCE PURPOSE

RETURNS

1542 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213  
RESFLOAT THE LIKELIHOOD OF THE DATA SET WITH SELF COVARIANCE AS AN ESTIMATOR OF ITS  
COVARIANCE MATRIX  
SETPARAMS SELFPARAMS  
SET THE PARAMETERS OF THIS ESTIMATOR  
THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE  
PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT  
OF A NESTED OBJECT  
RETURNS  
SELF  
664SKLEARNCOVARIANCE GRAPHICALLASSOCV  
CLASSSSKLEARNCOVARIANCE GRAPHICALLASSOCV ALPHAS4 NREFINEMENTS4 CV'WARN'  
TOL00001 ENETTOL00001 MAXITER100  
MODE'CD' NJOBSNONE VERBOSEFALSE AS  
SUMECENTEREDFALSE  
SPARSE INVERSE COVARIANCE W CROSSVALIDATED CHOICE OF THE L1 PENALTY  
SEE GLOSSARY ENTRY FOR CROSSVALIDATION ESTIMATOR  
READ MORE IN THE USER GUIDE  
PARAMETERS  
ALPHAS INTEGER OR LIST POSITIVE FLOAT OPTIONAL IF AN INTEGER IS GIVEN IT FIXES THE NUMBER OF POINTS  
ON THE GRIDS OF ALPHA TO BE USED IF A LIST IS GIVEN IT GIVES THE GRID TO BE USED SEE THE NOTES  
IN THE CLASS DOCSTRING FOR MORE DETAILS  
NREFINEMENTS STRICTLY POSITIVE INTEGER THE NUMBER OF TIMES THE GRID IS REFINED NOT USED IF  
EXPLICIT VALUES OF ALPHAS ARE PASSED  
CVINT CROSSVALIDATION GENERATOR OR AN ITERABLE OPTIONAL DETERMINES THE CROSSVALIDATION SPLIT  
TING STRATEGY POSSIBLE INPUTS FOR CV ARE  
• NONE TO USE THE DEFAULT 3FOLD CROSSVALIDATION  
• INTEGER TO SPECIFY THE NUMBER OF FOLDS  
• CV SPLITTER  
• AN ITERABLE YIELDING TRAIN TEST SPLITS AS ARRAYS OF INDICES  
FOR INTEGERNONE INPUTS KFOLD IS USED  
REFER USER GUIDE FOR THE VARIOUS CROSSVALIDATION STRATEGIES THAT CAN BE USED HERE  
CHANGED IN VERSION 020 CVDEFAULT VALUE IF NONE WILL CHANGE FROM 3FOLD TO 5FOLD IN  
V022  
TOLPOSITIVE FLOAT OPTIONAL THE TOLERANCE TO DECLARE CONVERGENCE IF THE DUAL GAP GOES BELOW  
THIS VALUE ITERATIONS ARE STOPPED  
ENETTOL POSITIVE FLOAT OPTIONAL THE TOLERANCE FOR THE ELASTIC NET SOLVER USED TO CALCULATE THE  
DESCENT DIRECTION THIS PARAMETER CONTROLS THE ACCURACY OF THE SEARCH DIRECTION FOR A GIVEN  
COLUMN UPDATE NOT OF THE OVERALL PARAMETER ESTIMATE ONLY USED FOR MODE'CD'  
MAXITER INTEGER OPTIONAL MAXIMUM NUMBER OF ITERATIONS  
66SKLEARNCOVARIANCE COVARIANCE ESTIMATORS 1543

SCIKITLEARN USER GUIDE RELEASE 0213

MODE 'CD' 'LARS' THE LASSO SOLVER TO USE COORDINATE DESCENT OR LARS USE LARS FOR  
VERY SPARSE UNDERLYING GRAPHS WHERE NUMBER OF FEATURES IS GREATER THAN NUMBER OF SAMPLES  
ELSEWHERE PREFER CD WHICH IS MORE NUMERICALLY STABLE

NJOBS INT OR NONE OPTIONAL DEFAULTNONE NUMBER OF JOBS TO RUN IN PARALLEL NONE MEANS 1  
UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL PROCESSORS SEE

GLOSSARY FOR MORE DETAILS

VERBOSE BOOLEAN OPTIONAL IF VERBOSE IS TRUE THE OBJECTIVE FUNCTION AND DUALITY GAP ARE PRINTED  
AT EACH ITERATION

ASSUMECENTERED BOOLEAN IF TRUE DATA ARE NOT CENTERED BEFORE COMPUTATION USEFUL WHEN  
WORKING WITH DATA WHOSE MEAN IS ALMOST BUT NOT EXACTLY ZERO IF FALSE DATA ARE CENTERED  
BEFORE COMPUTATION

ATTRIBUTES

LOCATION ARRAYLIKE SHAPE NFEATURES ESTIMATED LOCATION IE THE ESTIMATED MEAN

COVARIANCE NUMPYNDARRAY SHAPE NFEATURES NFEATURES ESTIMATED COVARIANCE MATRIX

PRECISION NUMPYNDARRAY SHAPE NFEATURES NFEATURES ESTIMATED PRECISION MATRIX INVERSE  
COVARIANCE

ALPHA FLOAT PENALIZATION PARAMETER SELECTED

CVALPHAS LIST OF FLOAT ALL PENALIZATION PARAMETERS EXPLORED

GRIDSCORES 2D NUMPYNDARRAY NALPHAS NFOLDS LOGLIKELIHOOD SCORE ON LEFTOUT DATA  
ACROSS FOLDS

NITER INT NUMBER OF ITERATIONS RUN FOR THE OPTIMAL ALPHA

SEE ALSO

GRAPHICALASSO GRAPHICALASSO

NOTES

THE SEARCH FOR THE OPTIMAL PENALIZATION PARAMETER ALPHA IS DONE ON AN ITERATIVELY REFINED GRID FIRST THE CROSS  
VALIDATED SCORES ON A GRID ARE COMPUTED THEN A NEW REFINED GRID IS CENTERED AROUND THE MAXIMUM AND SO ON  
ONE OF THE CHALLENGES WHICH IS FACED HERE IS THAT THE SOLVERS CAN FAIL TO CONVERGE TO A WELLCONDITIONED ESTIMATE  
THE CORRESPONDING VALUES OF ALPHA THEN COME OUT AS MISSING VALUES BUT THE OPTIMUM MAY BE CLOSE TO THESE  
MISSING VALUES

EXAMPLES

```
import numpy as np
from sklearn.covariance import GraphicalLassoCV
truecov = np.array([0.8, 0.0, 0.2, 0.0,
                    0.0, 0.4, 0.0, 0.0,
                    0.2, 0.0, 0.3, 0.1,
                    0.0, 0.0, 0.1, 0.7])
np.random.seed(0)
X = np.random.multivariate_normal(mean=[0, 0, 0, 0],
                                  cov=truecov,
                                  size=200)
cov = GraphicalLassoCV(cv=5).fit(X)
```

1544 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213  
NPAROUNDCOVCOVARIANCE DECIMALS3  
ARRAY0816 0051 022 0017  
0051 0364 0018 0036  
022 0018 0322 0094  
0017 0036 0094 069  
NPAROUNDCOVLOCATION DECIMALS3  
ARRAY0073 004 0038 0143  
METHODS  
ERRORNORM SELF COMPCOV NORM SCALING COMPUTES THE MEAN SQUARED ERROR BETWEEN TWO COVARIANCE ESTIMATORS  
FITSELF X Y FITS THE GRAPHICALASSO COVARIANCE MODEL TO X  
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR  
GETPRECISION SELF GETTER FOR THE PRECISION MATRIX  
MAHALANOBIS SELF X COMPUTES THE SQUARED MAHALANOBIS DISTANCES OF GIVEN OBSERVATIONS  
SCORE SELF XTEST Y COMPUTES THE LOGLIKELIHOOD OF A GAUSSIAN DATA SET WITH SELFCOVARIANCE AS AN ESTIMATOR OF ITS COVARIANCE MATRIX  
SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR  
INIT SELFALPHAS4 NREFINEMENTS4 CV'WARN' TOL00001 ENETTOL00001 MAXITER100  
MODE'CD' NJOBSNONE VERBOSEFALSE ASSUMECENTEREDFALSE  
ERRORNORM SELFCOMPCOV NORM'FROBENIUS' SCALINGTRUE SQUAREDTRUE  
COMPUTES THE MEAN SQUARED ERROR BETWEEN TWO COVARIANCE ESTIMATORS IN THE SENSE OF THE FROBENIUS NORM  
PARAMETERS  
COMPCOV ARRAYLIKE SHAPE NFEATURES NFEATURES THE COVARIANCE TO COMPARE WITH  
NORM STR THE TYPE OF NORM USED TO COMPUTE THE ERROR AVAILABLE ERROR TYPES 'FROBENIUS'  
DEFAULT SQRTTRATA 'SPECTRAL' SQRTMAXEIGENVALUESATA WHERE A IS THE ERROR  
COMPCOV SELFCOVARIANCE  
SCALING BOOL IF TRUE DEFAULT THE SQUARED ERROR NORM IS DIVIDED BY NFEATURES IF FALSE THE SQUARED ERROR NORM IS NOT RESCALED  
SQUARED BOOL WHETHER TO COMPUTE THE SQUARED ERROR NORM OR THE ERROR NORM IF TRUE  
DEFAULT THE SQUARED ERROR NORM IS RETURNED IF FALSE THE ERROR NORM IS RETURNED  
RETURNS  
THE MEAN SQUARED ERROR IN THE SENSE OF THE FROBENIUS NORM BETWEEN  
SELF ANDCOMPCOV COVARIANCE ESTIMATORS  
FITSELFXYNONE  
FITS THE GRAPHICALASSO COVARIANCE MODEL TO X  
PARAMETERS  
XNDARRAY SHAPE NSAMPLES NFEATURES DATA FROM WHICH TO COMPUTE THE COVARIANCE ESTIMATE  
YIGNORED  
66SKLEARNCOVARIANCE COVARIANCE ESTIMATORS 1545

SCIKITLEARN USER GUIDE RELEASE 0213

GETPARAMS SELFDEEPTURE  
GET PARAMETERS FOR THIS ESTIMATOR  
PARAMETERS  
DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBOBJECTS THAT ARE ESTIMATORS  
RETURNS  
PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES  
GETPRECISION SELF  
GETTER FOR THE PRECISION MATRIX  
RETURNS  
PRECISION ARRAYLIKE THE PRECISION MATRIX ASSOCIATED TO THE CURRENT COVARIANCE OBJECT  
MAHALANOBIS SELF  
COMPUTES THE SQUARED MAHALANOBIS DISTANCES OF GIVEN OBSERVATIONS  
PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES THE OBSERVATIONS THE MAHALANOBIS DISTANCES  
OF THE WHICH WE COMPUTE OBSERVATIONS ARE ASSUMED TO BE DRAWN FROM THE SAME DISTRIBUTION  
THAN THE DATA USED IN FIT  
RETURNS  
DIST ARRAY SHAPE NSAMPLES SQUARED MAHALANOBIS DISTANCES OF THE OBSERVATIONS  
SCORESELFXTTEST YNONE  
COMPUTES THE LOGLIKELIHOOD OF A GAUSSIAN DATA SET WITH SELFCOVARIANCE AS AN ESTIMATOR OF ITS  
COVARIANCE MATRIX  
PARAMETERS  
XTTEST ARRAYLIKE SHAPE NSAMPLES NFEATURES TEST DATA OF WHICH WE COMPUTE THE  
LIKELIHOOD WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF  
FEATURES XTTEST IS ASSUMED TO BE DRAWN FROM THE SAME DISTRIBUTION THAN THE DATA USED IN  
FIT INCLUDING CENTERING  
YNOT USED PRESENT FOR API CONSISTENCE PURPOSE  
RETURNS  
RESFLOAT THE LIKELIHOOD OF THE DATA SET WITH SELFCOVARIANCE AS AN ESTIMATOR OF ITS  
COVARIANCE MATRIX  
SETPARAMS SELFPARAMS  
SET THE PARAMETERS OF THIS ESTIMATOR  
THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE  
PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT  
OF A NESTED OBJECT  
RETURNS  
SELF  
EXAMPLES USING SKLEARNCOVARIANCEGRAPHICALASSOCV  
•VISUALIZING THE STOCK MARKET STRUCTURE  
1546 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

•SPARSE INVERSE COVARIANCE ESTIMATION

665SKLEARNCOVARIANCE LEDOITWOLF

CLASSSSKLEARNCOVARIANCE LEDOITWOLF STOREPRECISIONTRUE ASSUMECENTEREDFALSE

BLOCKSIZE1000

LEDOITWOLF ESTIMATOR

LEDOITWOLF IS A PARTICULAR FORM OF SHRINKAGE WHERE THE SHRINKAGE COEFFICIENT IS COMPUTED USING O LEDOIT AND M WOLF’S FORMULA AS DESCRIBED IN “A WELLCONDITIONED ESTIMATOR FOR LARGEDIMENSIONAL COVARIANCE MATRICES”

LEDOIT AND WOLF JOURNAL OF MULTIVARIATE ANALYSIS V OLUME 88 ISSUE 2 FEBRUARY 2004 PAGES 365411

READ MORE IN THE USER GUIDE

PARAMETERS

STOREPRECISION BOOL DEFAULTTRUE SPECIFY IF THE ESTIMATED PRECISION IS STORED

ASSUMECENTERED BOOL DEFAULTFALSE IF TRUE DATA WILL NOT BE CENTERED BEFORE COMPUTATION

USEFUL WHEN WORKING WITH DATA WHOSE MEAN IS ALMOST BUT NOT EXACTLY ZERO IF FALSE DE

FAULT DATA WILL BE CENTERED BEFORE COMPUTATION

BLOCKSIZE INT DEFAULT1000 SIZE OF THE BLOCKS INTO WHICH THE COVARIANCE MATRIX WILL BE SPLIT

DURING ITS LEDOITWOLF ESTIMATION THIS IS PURELY A MEMORY OPTIMIZATION AND DOES NOT AFFECT

RESULTS

ATTRIBUTES

LOCATION ARRAYLIKE SHAPE NFEATURES ESTIMATED LOCATION IE THE ESTIMATED MEAN

COVARIANCE ARRAYLIKE SHAPE NFEATURES NFEATURES ESTIMATED COVARIANCE MATRIX

PRECISION ARRAYLIKE SHAPE NFEATURES NFEATURES ESTIMATED PSEUDO INVERSE MATRIX STORED

ONLY IF STOREPRECISION IS TRUE

SHRINKAGE FLOAT 0 SHRINKAGE 1 COEFFICIENT IN THE CONVEX COMBINATION USED FOR THE

COMPUTATION OF THE SHRUNK ESTIMATE

NOTES

THE REGULARISED COVARIANCE IS

$$\frac{1}{n} \text{SHRINKAGE} \text{ COV} \text{SHRINKAGE} \mu \text{ NPIDENTITY} \text{NFEATURES}$$

WHERE  $\mu$   $\text{TRACECOV} \text{NFEATURES}$  AND SHRINKAGE IS GIVEN BY THE LEDOIT AND WOLF FORMULA SEE REFERENCES

REFERENCES

“A WELLCONDITIONED ESTIMATOR FOR LARGEDIMENSIONAL COVARIANCE MATRICES” LEDOIT AND WOLF JOURNAL OF MUL

TIVARIATE ANALYSIS V OLUME 88 ISSUE 2 FEBRUARY 2004 PAGES 365411

EXAMPLES

```
import numpy as np
from sklearn.covariance import LedoitWolf
real_cov = np.array([[2, 0.5], [0.5, 1]])
lw = LedoitWolf()
lw.fit(real_cov)
```

2 8

665SKLEARNCOVARIANCE COVARIANCE ESTIMATORS 1547

SCIKITLEARN USER GUIDE RELEASE 0213

NPRANDOMSEED0

X NPRANDOMMULTIVARIATENORMALMEAN0 0

COVREALCOV

SIZE50

COV LEDOITWOLFFITX

COVCOVARIANCE

ARRAY04406 01616

01616 08022

COVLOCATION

ARRAY 00595 00075

METHODS

ERRORNORM SELF COMPCOV NORM SCALING COMPUTES THE MEAN SQUARED ERROR BETWEEN TWO COVARIANCE ESTIMATORS

FITSELF X Y FITS THE LEDOITWOLF SHRUNK COVARIANCE MODEL ACCORDING TO THE GIVEN TRAINING DATA AND PARAMETERS

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

GETPRECISION SELF GETTER FOR THE PRECISION MATRIX

MAHALANOBIS SELF X COMPUTES THE SQUARED MAHALANOBIS DISTANCES OF GIVEN OBSERVATIONS

SCORE SELF XTEST Y COMPUTES THE LOGLIKELIHOOD OF A GAUSSIAN DATA SET WITH SELFCOVARIANCE AS AN ESTIMATOR OF ITS COVARIANCE MATRIX

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELFSTOREPRECISIONTRUE ASSUMECENTEREDFALSE BLOCKSIZE1000

ERRORNORM SELFCOMPCOV NORM'FROBENIUS' SCALINGTRUE SQUAREDTRUE

COMPUTES THE MEAN SQUARED ERROR BETWEEN TWO COVARIANCE ESTIMATORS IN THE SENSE OF THE FROBENIUS NORM

PARAMETERS

COMPCOV ARRAYLIKE SHAPE NFEATURES NFEATURES THE COVARIANCE TO COMPARE WITH

NORM STR THE TYPE OF NORM USED TO COMPUTE THE ERROR AVAILABLE ERROR TYPES 'FROBENIUS'

DEFAULT SQRTTRATA 'SPECTRAL' SQRTMAXEIGENVALUESATA WHERE A IS THE ERROR

COMPCOV SELFCOVARIANCE

SCALING BOOL IF TRUE DEFAULT THE SQUARED ERROR NORM IS DIVIDED BY NFEATURES IF FALSE

THE SQUARED ERROR NORM IS NOT RESCALED

SQUARED BOOL WHETHER TO COMPUTE THE SQUARED ERROR NORM OR THE ERROR NORM IF TRUE

DEFAULT THE SQUARED ERROR NORM IS RETURNED IF FALSE THE ERROR NORM IS RETURNED

RETURNS

THE MEAN SQUARED ERROR IN THE SENSE OF THE FROBENIUS NORM BETWEEN

SELF ANDCOMPCOV COVARIANCE ESTIMATORS

FITSELFXYNONE

FITS THE LEDOITWOLF SHRUNK COVARIANCE MODEL ACCORDING TO THE GIVEN TRAINING DATA AND PARAMETERS

PARAMETERS

1548 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

YNOT USED PRESENT FOR API CONSISTENCE PURPOSE

RETURNS

SELF OBJECT

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

GETPRECISION SELF

GETTER FOR THE PRECISION MATRIX

RETURNS

PRECISION ARRAYLIKE THE PRECISION MATRIX ASSOCIATED TO THE CURRENT COVARIANCE OBJECT

MAHALANOBIS SELF

COMPUTES THE SQUARED MAHALANOBIS DISTANCES OF GIVEN OBSERVATIONS

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES THE OBSERVATIONS THE MAHALANOBIS DISTANCES OF THE WHICH WE COMPUTE OBSERVATIONS ARE ASSUMED TO BE DRAWN FROM THE SAME DISTRIBUTION THAN THE DATA USED IN FIT

RETURNS

DIST ARRAY SHAPE NSAMPLES SQUARED MAHALANOBIS DISTANCES OF THE OBSERVATIONS

SCORESELFXTTEST YNONE

COMPUTES THE LOGLIKELIHOOD OF A GAUSSIAN DATA SET WITH SELFCOVARIANCE AS AN ESTIMATOR OF ITS COVARIANCE MATRIX

PARAMETERS

XTTEST ARRAYLIKE SHAPE NSAMPLES NFEATURES TEST DATA OF WHICH WE COMPUTE THE LIKELIHOOD WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES XTTEST IS ASSUMED TO BE DRAWN FROM THE SAME DISTRIBUTION THAN THE DATA USED IN FIT INCLUDING CENTERING

YNOT USED PRESENT FOR API CONSISTENCE PURPOSE

RETURNS

RESFLOAT THE LIKELIHOOD OF THE DATA SET WITH SELFCOVARIANCE AS AN ESTIMATOR OF ITS COVARIANCE MATRIX

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

66SKLEARNCOVARIANCE COVARIANCE ESTIMATORS 1549

SCIKITLEARN USER GUIDE RELEASE 0213

RETURNS

SELF

EXAMPLES USING SKLEARNCOVARIANCELEDOITWOLF

- LEDOITWOLF VS OAS ESTIMATION
- SHRINKAGE COVARIANCE ESTIMATION LEDOITWOLF VS OAS AND MAXLIKELIHOOD
- MODEL SELECTION WITH PROBABILISTIC PCA AND FACTOR ANALYSIS FA

666SKLEARNCOVARIANCE MINCOVDET

CLASSSSKLEARNCOVARIANCE MINCOVDET STOREPRECISIONTRUE ASSUMECENTEREDFALSE SUP

PORTFRACTIONNONE RANDOMSTATENONE

MINIMUM COVARIANCE DETERMINANT MCD ROBUST ESTIMATOR OF COVARIANCE

THE MINIMUM COVARIANCE DETERMINANT COVARIANCE ESTIMATOR IS TO BE APPLIED ON GAUSSIANDISTRIBUTED DATA BUT COULD STILL BE RELEVANT ON DATA DRAWN FROM A UNIMODAL SYMMETRIC DISTRIBUTION IT IS NOT MEANT TO BE USED WITH MULTIMODAL DATA THE ALGORITHM USED TO FIT A MINCOVDET OBJECT IS LIKELY TO FAIL IN SUCH A CASE ONE SHOULD CONSIDER PROJECTION PURSUIT METHODS TO DEAL WITH MULTIMODAL DATASETS

READ MORE IN THE USER GUIDE

PARAMETERS

STOREPRECISION BOOL SPECIFY IF THE ESTIMATED PRECISION IS STORED

ASSUMECENTERED BOOL IF TRUE THE SUPPORT OF THE ROBUST LOCATION AND THE COVARIANCE ESTIMATES IS COMPUTED AND A COVARIANCE ESTIMATE IS RECOMPUTED FROM IT WITHOUT CENTERING THE DATA

USEFUL TO WORK WITH DATA WHOSE MEAN IS SIGNIFICANTLY EQUAL TO ZERO BUT IS NOT EXACTLY ZERO IF FALSE THE ROBUST LOCATION AND COVARIANCE ARE DIRECTLY COMPUTED WITH THE FASTMCD ALGORITHM

WITHOUT ADDITIONAL TREATMENT

SUPPORTFRACTION FLOAT 0 SUPPORTFRACTION 1 THE PROPORTION OF POINTS TO BE INCLUDED IN

THE SUPPORT OF THE RAW MCD ESTIMATE DEFAULT IS NONE WHICH IMPLIES THAT THE MINIMUM

VALUE OF SUPPORTFRACTION WILL BE USED WITHIN THE ALGORITHM NSAMPLE NFEATURES 1 2

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RAN

DOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE

RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS

THE RANDOMSTATE INSTANCE USED BY NPRANDOM

ATTRIBUTES

RAWLOCATION ARRAYLIKE SHAPE NFEATURES THE RAW ROBUST ESTIMATED LOCATION BEFORE COR

RECTION AND REWEIGHTING

RAWCOVARIANCE ARRAYLIKE SHAPE NFEATURES NFEATURES THE RAW ROBUST ESTIMATED COVARI

ANCE BEFORE CORRECTION AND REWEIGHTING

RAWSUPPORT ARRAYLIKE SHAPE NSAMPLES A MASK OF THE OBSERVATIONS THAT HAVE BEEN

USED TO COMPUTE THE RAW ROBUST ESTIMATES OF LOCATION AND SHAPE BEFORE CORRECTION AND REWEIGHTING

LOCATION ARRAYLIKE SHAPE NFEATURES ESTIMATED ROBUST LOCATION

COVARIANCE ARRAYLIKE SHAPE NFEATURES NFEATURES ESTIMATED ROBUST COVARIANCE MATRIX

1550 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

PRECISION ARRAYLIKE SHAPE NFEATURES NFEATURES ESTIMATED PSEUDO INVERSE MATRIX STORED ONLY IF STOREPRECISION IS TRUE

SUPPORT ARRAYLIKE SHAPE NSAMPLES A MASK OF THE OBSERVATIONS THAT HAVE BEEN USED TO COMPUTE THE ROBUST ESTIMATES OF LOCATION AND SHAPE

DIST ARRAYLIKE SHAPE NSAMPLES MAHALANOBIS DISTANCES OF THE TRAINING SET ON WHICH FIT IS CALLED OBSERVATIONS

REFERENCES

R9F63E655F7BDROUSEEUW1984 R9F63E655F7BDROUSSEEUW R9F63E655F7BDBUTLERDAVIES

EXAMPLES

```
import numpy as np
from sklearn.covariance import MinCovDet
from sklearn.datasets import make_gaussian_quantiles
real_cov = np.array([
    [3, 4],
    [4, 3]
])
rng = np.random.RandomState(0)
X = rng.multivariate_normal(mean_0, cov=real_cov, size=500)
cov = MinCovDet(random_state=0).fit(X)
cov.covariance_
array([[0.7411, 0.2535],
       [0.2535, 0.3053]])
cov.location_
array([[0.0813, 0.0427]])
```

METHODS

CORRECTCOVARIANCE SELF DATA APPLY A CORRECTION TO RAW MINIMUM COVARIANCE DETERMINANT ESTIMATES

ERRORNORM SELF COMPCOV NORM SCALING COMPUTES THE MEAN SQUARED ERROR BETWEEN TWO COVARIANCE ESTIMATORS

FITSELF X Y FITS A MINIMUM COVARIANCE DETERMINANT WITH THE FASTMCD ALGORITHM

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

GETPRECISION SELF GETTER FOR THE PRECISION MATRIX

MAHALANOBIS SELF X COMPUTES THE SQUARED MAHALANOBIS DISTANCES OF GIVEN OBSERVATIONS

REWEIGHTCOVARIANCE SELF DATA REWEIGHT RAW MINIMUM COVARIANCE DETERMINANT ESTIMATES

SCORE SELF XTEST Y COMPUTES THE LOGLIKELIHOOD OF A GAUSSIAN DATA SET WITH SELF COVARIANCE AS AN ESTIMATOR OF ITS COVARIANCE MATRIX

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

66SKLEARNCOVARIANCE COVARIANCE ESTIMATORS 1551

SCIKITLEARN USER GUIDE RELEASE 0213

INIT SELFSTOREPRECISIONTRUE ASSUMECENTEREDFALSE SUPPORTFRACTIONNONE RAN  
DOMSTATENONE

CORRECTCOVARIANCE SELFDATA

APPLY A CORRECTION TO RAW MINIMUM COVARIANCE DETERMINANT ESTIMATES

CORRECTION USING THE EMPIRICAL CORRECTION FACTOR SUGGESTED BY ROUSSEEUW AND VAN DRIESSEN IN RVD  
PARAMETERS

DATA ARRAYLIKE SHAPE NSAMPLES NFEATURES THE DATA MATRIX WITH P FEATURES AND N SAM  
PLES THE DATA SET MUST BE THE ONE WHICH WAS USED TO COMPUTE THE RAW ESTIMATES

RETURNS

COVARIANCECORRECTED ARRAYLIKE SHAPE NFEATURES NFEATURES CORRECTED ROBUST COVARI  
ANCE ESTIMATE

REFERENCES

RVD

ERRORNORM SELFCOMPCOV NORM'FROBENIUS' SCALINGTRUE SQUAREDTRUE

COMPUTES THE MEAN SQUARED ERROR BETWEEN TWO COVARIANCE ESTIMATORS IN THE SENSE OF THE FROBENIUS  
NORM

PARAMETERS

COMPCOV ARRAYLIKE SHAPE NFEATURES NFEATURES THE COVARIANCE TO COMPARE WITH

NORM STR THE TYPE OF NORM USED TO COMPUTE THE ERROR AVAILABLE ERROR TYPES 'FROBENIUS'  
DEFAULT SQRTTRATA 'SPECTRAL' SQRTMAXEIGENVALUESATA WHERE A IS THE ERROR

COMPCOV SELFCOVARIANCE

SCALING BOOL IF TRUE DEFAULT THE SQUARED ERROR NORM IS DIVIDED BY NFEATURES IF FALSE  
THE SQUARED ERROR NORM IS NOT RESCALED

SQUARED BOOL WHETHER TO COMPUTE THE SQUARED ERROR NORM OR THE ERROR NORM IF TRUE  
DEFAULT THE SQUARED ERROR NORM IS RETURNED IF FALSE THE ERROR NORM IS RETURNED

RETURNS

THE MEAN SQUARED ERROR IN THE SENSE OF THE FROBENIUS NORM BETWEEN  
SELF ANDCOMPCOV COVARIANCE ESTIMATORS

FITSELFXYNONE

FITS A MINIMUM COVARIANCE DETERMINANT WITH THE FASTMCD ALGORITHM

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA WHERE NSAMPLES IS THE NUM  
BER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

YNOT USED PRESENT FOR API CONSISTENCE PURPOSE

RETURNS

SELF OBJECT

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

1552 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

GETPRECISION SELF

GETTER FOR THE PRECISION MATRIX

RETURNS

PRECISION ARRAYLIKE THE PRECISION MATRIX ASSOCIATED TO THE CURRENT COVARIANCE OBJECT

MAHALANOBIS SELF

COMPUTES THE SQUARED MAHALANOBIS DISTANCES OF GIVEN OBSERVATIONS

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES THE OBSERVATIONS THE MAHALANOBIS DISTANCES OF THE WHICH WE COMPUTE OBSERVATIONS ARE ASSUMED TO BE DRAWN FROM THE SAME DISTRIBUTION THAN THE DATA USED IN FIT

RETURNS

DIST ARRAY SHAPE NSAMPLES SQUARED MAHALANOBIS DISTANCES OF THE OBSERVATIONS

REWEIGHTCOVARIANCE SELF

REWEIGHT RAW MINIMUM COVARIANCE DETERMINANT ESTIMATES

REWEIGHT OBSERVATIONS USING ROUSSEEUW'S METHOD EQUIVALENT TO DELETING OUTLYING OBSERVATIONS FROM THE DATA SET BEFORE COMPUTING LOCATION AND COVARIANCE ESTIMATES DESCRIBED IN RVDRIESSEN

PARAMETERS

DATA ARRAYLIKE SHAPE NSAMPLES NFEATURES THE DATA MATRIX WITH P FEATURES AND N SAMPLES THE DATA SET MUST BE THE ONE WHICH WAS USED TO COMPUTE THE RAW ESTIMATES

RETURNS

LOCATIONREWEIGHTED ARRAYLIKE SHAPE NFEATURES REWEIGHTED ROBUST LOCATION ESTIMATE

COVARIANCEREWEIGHTED ARRAYLIKE SHAPE NFEATURES NFEATURES REWEIGHTED ROBUST COVARIANCE ESTIMATE

SUPPORTREWEIGHTED ARRAYLIKE TYPE BOOLEAN SHAPE NSAMPLES A MASK OF THE OBSERVATIONS THAT HAVE BEEN USED TO COMPUTE THE REWEIGHTED ROBUST LOCATION AND COVARIANCE ESTIMATES

REFERENCES

RVDRIESSEN

SCORESELFXTTEST YNONE

COMPUTES THE LOGLIKELIHOOD OF A GAUSSIAN DATA SET WITH SELFCOVARIANCE AS AN ESTIMATOR OF ITS COVARIANCE MATRIX

PARAMETERS

XTTEST ARRAYLIKE SHAPE NSAMPLES NFEATURES TEST DATA OF WHICH WE COMPUTE THE LIKELIHOOD WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF

66SKLEARNCOVARIANCE COVARIANCE ESTIMATORS 1553

SCIKITLEARN USER GUIDE RELEASE 0213

FEATURES XTEST IS ASSUMED TO BE DRAWN FROM THE SAME DISTRIBUTION THAN THE DATA USED IN  
FIT INCLUDING CENTERING

YNOT USED PRESENT FOR API CONSISTENCE PURPOSE

RETURNS

RESFLOAT THE LIKELIHOOD OF THE DATA SET WITH SELFCOVARIANCE AS AN ESTIMATOR OF ITS  
COVARIANCE MATRIX

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE  
PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT  
OF A NESTED OBJECT

RETURNS

SELF

EXAMPLES USING SKLEARNCOVARIANCEMINCOVDET

- ROBUST COVARIANCE ESTIMATION AND MAHALANOBIS DISTANCES RELEVANCE
- ROBUST VS EMPIRICAL COVARIANCE ESTIMATE

667SKLEARNCOVARIANCE OAS

CLASSSSKLEARNCOVARIANCE OASSTOREPRECISIONTRUE ASSUMECENTEREDFALSE

ORACLE APPROXIMATING SHRINKAGE ESTIMATOR

READ MORE IN THE USER GUIDE

OAS IS A PARTICULAR FORM OF SHRINKAGE DESCRIBED IN “SHRINKAGE ALGORITHMS FOR MMSE COVARIANCE ESTIMATION”  
CHEN ET AL IEEE TRANS ON SIGN PROC V OLUME 58 ISSUE 10 OCTOBER 2010

THE FORMULA USED HERE DOES NOT CORRESPOND TO THE ONE GIVEN IN THE ARTICLE IN THE ORIGINAL ARTICLE FORMULA 23  
STATES THAT 2P IS MULTIPLIED BY TRACECOVCOV IN BOTH THE NUMERATOR AND DENOMINATOR BUT THIS OPERATION IS  
OMITTED BECAUSE FOR A LARGE P THE VALUE OF 2P IS SO SMALL THAT IT DOESN'T AFFECT THE VALUE OF THE ESTIMATOR  
PARAMETERS

STOREPRECISION BOOL DEFAULTTRUE SPECIFY IF THE ESTIMATED PRECISION IS STORED

ASSUMECENTERED BOOL DEFAULTFALSE IF TRUE DATA WILL NOT BE CENTERED BEFORE COMPUTATION

USEFUL WHEN WORKING WITH DATA WHOSE MEAN IS ALMOST BUT NOT EXACTLY ZERO IF FALSE DE  
FAULT DATA WILL BE CENTERED BEFORE COMPUTATION

ATTRIBUTES

COVARIANCE ARRAYLIKE SHAPE NFEATURES NFEATURES ESTIMATED COVARIANCE MATRIX

PRECISION ARRAYLIKE SHAPE NFEATURES NFEATURES ESTIMATED PSEUDO INVERSE MATRIX STORED  
ONLY IF STOREPRECISION IS TRUE

SHRINKAGE FLOAT 0 SHRINKAGE 1 COEFFICIENT IN THE CONVEX COMBINATION USED FOR THE  
COMPUTATION OF THE SHRUNK ESTIMATE

1554 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

NOTES

THE REGULARISED COVARIANCE IS

1 SHRINKAGE COV SHRINKAGE MU NPIDENTITYNFEATURES

WHERE MU TRACECOV NFEATURES AND SHRINKAGE IS GIVEN BY THE OAS FORMULA SEE REFERENCES

REFERENCES

“SHRINKAGE ALGORITHMS FOR MMSE COVARIANCE ESTIMATION” CHEN ET AL IEEE TRANS ON SIGN PROC V OLUME 58  
ISSUE 10 OCTOBER 2010

METHODS

ERRORNORM SELF COMPCOV NORM SCALING COMPUTES THE MEAN SQUARED ERROR BETWEEN TWO COVARIANCE ESTIMATORS

FITSELF X Y FITS THE ORACLE APPROXIMATING SHRINKAGE COVARIANCE

MODEL ACCORDING TO THE GIVEN TRAINING DATA AND PARAMETERS

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

GETPRECISION SELF GETTER FOR THE PRECISION MATRIX

MAHALANOBIS SELF X COMPUTES THE SQUARED MAHALANOBIS DISTANCES OF GIVEN OBSERVATIONS

SCORE SELF XTEST Y COMPUTES THE LOGLIKELIHOOD OF A GAUSSIAN DATA SET WITH SELFCOVARIANCE AS AN ESTIMATOR OF ITS COVARIANCE MATRIX

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELFSTOREPRECISIONTRUE ASSUMECENTEREDFALSE

ERRORNORM SELFCOMPCOV NORM‘FROBENIUS’ SCALINGTRUE SQUAREDTRUE

COMPUTES THE MEAN SQUARED ERROR BETWEEN TWO COVARIANCE ESTIMATORS IN THE SENSE OF THE FROBENIUS NORM

PARAMETERS

COMPCOV ARRAYLIKE SHAPE NFEATURES NFEATURES THE COVARIANCE TO COMPARE WITH

NORM STR THE TYPE OF NORM USED TO COMPUTE THE ERROR AVAILABLE ERROR TYPES ‘FROBENIUS’

DEFAULT SQRTTRATA ‘SPECTRAL’ SQRTMAXEIGENVALUESATA WHERE A IS THE ERROR

COMPCOV SELFCOVARIANCE

SCALING BOOL IF TRUE DEFAULT THE SQUARED ERROR NORM IS DIVIDED BY NFEATURES IF FALSE THE SQUARED ERROR NORM IS NOT RESCALED

SQUARED BOOL WHETHER TO COMPUTE THE SQUARED ERROR NORM OR THE ERROR NORM IF TRUE DEFAULT THE SQUARED ERROR NORM IS RETURNED IF FALSE THE ERROR NORM IS RETURNED

RETURNS

THE MEAN SQUARED ERROR IN THE SENSE OF THE FROBENIUS NORM BETWEEN SELF ANDCOMPCOV COVARIANCE ESTIMATORS

66SKLEARNCOVARIANCE COVARIANCE ESTIMATORS 1555

SCIKITLEARN USER GUIDE RELEASE 0213

**FITSELFXYNONE**  
FITS THE ORACLE APPROXIMATING SHRINKAGE COVARIANCE MODEL ACCORDING TO THE GIVEN TRAINING DATA AND PA  
RAMETERS  
PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA WHERE NSAMPLES IS THE NUM  
BER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES  
YNOT USED PRESENT FOR API CONSISTENCE PURPOSE  
RETURNS  
SELF OBJECT  
GETPARAMS SELFDEEPTTRUE  
GET PARAMETERS FOR THIS ESTIMATOR  
PARAMETERS  
DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBOBJECTS THAT ARE ESTIMATORS  
RETURNS  
PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES  
GETPRECISION SELF  
GETTER FOR THE PRECISION MATRIX  
RETURNS  
PRECISION ARRAYLIKE THE PRECISION MATRIX ASSOCIATED TO THE CURRENT COVARIANCE OBJECT  
MAHALANOBIS SELF  
COMPUTES THE SQUARED MAHALANOBIS DISTANCES OF GIVEN OBSERVATIONS  
PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES THE OBSERVATIONS THE MAHALANOBIS DISTANCES  
OF THE WHICH WE COMPUTE OBSERVATIONS ARE ASSUMED TO BE DRAWN FROM THE SAME DISTRIBU  
TION THAN THE DATA USED IN FIT  
RETURNS  
DIST ARRAY SHAPE NSAMPLES SQUARED MAHALANOBIS DISTANCES OF THE OBSERVATIONS  
SCORESELFXTTEST YNONE  
COMPUTES THE LOGLIKELIHOOD OF A GAUSSIAN DATA SET WITH SELFCOVARIANCE AS AN ESTIMATOR OF ITS  
COVARIANCE MATRIX  
PARAMETERS  
XTTEST ARRAYLIKE SHAPE NSAMPLES NFEATURES TEST DATA OF WHICH WE COMPUTE THE  
LIKELIHOOD WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF  
FEATURES XTTEST IS ASSUMED TO BE DRAWN FROM THE SAME DISTRIBUTION THAN THE DATA USED IN  
FIT INCLUDING CENTERING  
YNOT USED PRESENT FOR API CONSISTENCE PURPOSE  
RETURNS  
RESFLOAT THE LIKELIHOOD OF THE DATA SET WITH SELFCOVARIANCE AS AN ESTIMATOR OF ITS  
COVARIANCE MATRIX

1556 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

EXAMPLES USING SKLEARNCOVARIANCEOAS

- LEDOITWOLF VS OAS ESTIMATION
- SHRINKAGE COVARIANCE ESTIMATION LEDOITWOLF VS OAS AND MAXLIKELIHOOD

668SKLEARNCOVARIANCE SHRUNKCOVARIANCE

CLASSSSKLEARNCOVARIANCE SHRUNKCOVARIANCE STOREPRECISIONTRUE ASSUMECENTEREDFALSE

SHRINKAGE01

COVARIANCE ESTIMATOR WITH SHRINKAGE

READ MORE IN THE USER GUIDE

PARAMETERS

STOREPRECISION BOOLEAN DEFAULT TRUE SPECIFY IF THE ESTIMATED PRECISION IS STORED

ASSUMECENTERED BOOLEAN DEFAULT FALSE IF TRUE DATA WILL NOT BE CENTERED BEFORE COMPUTATION USEFUL WHEN WORKING WITH DATA WHOSE MEAN IS ALMOST BUT NOT EXACTLY ZERO IF FALSE DATA WILL BE CENTERED BEFORE COMPUTATION

SHRINKAGE FLOAT 0 SHRINKAGE 1 DEFAULT 01 COEFFICIENT IN THE CONVEX COMBINATION USED FOR THE COMPUTATION OF THE SHRUNK ESTIMATE

ATTRIBUTES

LOCATION ARRAYLIKE SHAPE NFEATURES ESTIMATED LOCATION IE THE ESTIMATED MEAN

COVARIANCE ARRAYLIKE SHAPE NFEATURES NFEATURES ESTIMATED COVARIANCE MATRIX

PRECISION ARRAYLIKE SHAPE NFEATURES NFEATURES ESTIMATED PSEUDO INVERSE MATRIX STORED ONLY IF STOREPRECISION IS TRUE

SHRINKAGE FLOAT 0 SHRINKAGE 1 COEFFICIENT IN THE CONVEX COMBINATION USED FOR THE COMPUTATION OF THE SHRUNK ESTIMATE

NOTES

THE REGULARIZED COVARIANCE IS GIVEN BY

$$\frac{1}{N} \text{SHRINKAGE} \text{ COV} \text{SHRINKAGE} \text{ MU} \text{ NPIDENTITY} \text{NFEATURES}$$

WHERE MU TRACECOV NFEATURES

66SKLEARNCOVARIANCE COVARIANCE ESTIMATORS 1557

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES

```
import numpy as np
from sklearn.covariance import shrunk_covariance
from sklearn.datasets import make_gaussian_quantiles
real_cov = np.array([[3, 4],
                      [4, 3]])
rng = np.random.RandomState(0)
X = rng.multivariate_normal(mean=[0, 0],
                             cov=real_cov,
                             size=500)
cov = shrunk_covariance(X,
                        cov=covariance,
                        array=[0.7387, 0.2536],
                        [0.2536, 0.4110],
                        cov_location=[0.0622, 0.0193])

METHODS
```

**error\_norm** self, compcov, norm, scaling: COMPUTES THE MEAN SQUARED ERROR BETWEEN TWO COVARIANCE ESTIMATORS

**fit** self, X, y: FITS THE SHRUNK COVARIANCE MODEL ACCORDING TO THE GIVEN TRAINING DATA AND PARAMETERS

**get\_params** self: DEEP GET PARAMETERS FOR THIS ESTIMATOR

**get\_precision** self: GETTER FOR THE PRECISION MATRIX

**mahalanobis** self, X: COMPUTES THE SQUARED MAHALANOBIS DISTANCES OF GIVEN OBSERVATIONS

**score** self, X, test\_y: COMPUTES THE LOGLIKELIHOOD OF A GAUSSIAN DATA SET WITH SELF COVARIANCE AS AN ESTIMATOR OF ITS COVARIANCE MATRIX

**set\_params** self, params: SET THE PARAMETERS OF THIS ESTIMATOR

**init** self, store\_precision, true, assume\_centered, false, shrinkage=0.1

**error\_norm** self, compcov, norm='frobenius', scaling=True, squared=True: COMPUTES THE MEAN SQUARED ERROR BETWEEN TWO COVARIANCE ESTIMATORS IN THE SENSE OF THE FROBENIUS NORM

PARAMETERS

**compcov** array-like, shape: (n\_features, n\_features) THE COVARIANCE TO COMPARE WITH

**norm** str: THE TYPE OF NORM USED TO COMPUTE THE ERROR AVAILABLE ERROR TYPES 'FROBENIUS'

**default\_sqrt\_trata** 'SPECTRAL' SQRTMAXEIGENVALUESATA WHERE A IS THE ERROR

**compcov** selfcovariance

**scaling** bool: IF TRUE DEFAULT THE SQUARED ERROR NORM IS DIVIDED BY NFEATURES IF FALSE THE SQUARED ERROR NORM IS NOT RESCALED

**squared** bool: WHETHER TO COMPUTE THE SQUARED ERROR NORM OR THE ERROR NORM IF TRUE DEFAULT THE SQUARED ERROR NORM IS RETURNED IF FALSE THE ERROR NORM IS RETURNED

RETURNS

THE MEAN SQUARED ERROR IN THE SENSE OF THE FROBENIUS NORM BETWEEN

1558 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

SELF ANDCOMPCOV COVARIANCE ESTIMATORS

FITSELFXYNONE

FITS THE SHRUNK COVARIANCE MODEL ACCORDING TO THE GIVEN TRAINING DATA AND PARAMETERS

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

YNOT USED PRESENT FOR API CONSISTENCE PURPOSE

RETURNS

SELF OBJECT

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

GETPRECISION SELF

GETTER FOR THE PRECISION MATRIX

RETURNS

PRECISION ARRAYLIKE THE PRECISION MATRIX ASSOCIATED TO THE CURRENT COVARIANCE OBJECT

MAHALANOBIS SELF

COMPUTES THE SQUARED MAHALANOBIS DISTANCES OF GIVEN OBSERVATIONS

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES THE OBSERVATIONS THE MAHALANOBIS DISTANCES OF THE WHICH WE COMPUTE OBSERVATIONS ARE ASSUMED TO BE DRAWN FROM THE SAME DISTRIBUTION THAN THE DATA USED IN FIT

RETURNS

DIST ARRAY SHAPE NSAMPLES SQUARED MAHALANOBIS DISTANCES OF THE OBSERVATIONS

SCORESELFXTTEST YNONE

COMPUTES THE LOGLIKELIHOOD OF A GAUSSIAN DATA SET WITH SELFCOVARIANCE AS AN ESTIMATOR OF ITS COVARIANCE MATRIX

PARAMETERS

XTTEST ARRAYLIKE SHAPE NSAMPLES NFEATURES TEST DATA OF WHICH WE COMPUTE THE LIKELIHOOD WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES XTTEST IS ASSUMED TO BE DRAWN FROM THE SAME DISTRIBUTION THAN THE DATA USED IN FIT INCLUDING CENTERING

YNOT USED PRESENT FOR API CONSISTENCE PURPOSE

RETURNS

RESFLOAT THE LIKELIHOOD OF THE DATA SET WITH SELFCOVARIANCE AS AN ESTIMATOR OF ITS COVARIANCE MATRIX

66SKLEARNCOVARIANCE COVARIANCE ESTIMATORS 1559

SCIKITLEARN USER GUIDE RELEASE 0213

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

EXAMPLES USING SKLEARNCOVARIANCESHRUNKCOVARIANCE

- SHRINKAGE COVARIANCE ESTIMATION LEDOITWOLF VS OAS AND MAXLIKELIHOOD
- MODEL SELECTION WITH PROBABILISTIC PCA AND FACTOR ANALYSIS FA

COVARIANCEEMPIRICALCOVARIANCE X COMPUTES THE MAXIMUM LIKELIHOOD COVARIANCE ESTIMATOR

COVARIANCEGRAPHICALASSO EMPCOV ALPHA

L1PENALIZED COVARIANCE ESTIMATOR

COVARIANCELEDOITWOLF X ASSUMECENTERED

ESTIMATES THE SHRUNK LEDOITWOLF COVARIANCE MATRIX

COVARIANCEOAS X ASSUMECENTERED ESTIMATE COVARIANCE WITH THE ORACLE APPROXIMATING

SHRINKAGE ALGORITHM

COVARIANCESHRUNKCOVARIANCE EMPCOV CALCULATES A COVARIANCE MATRIX SHRUNK ON THE DIAGONAL

669SKLEARNCOVARIANCE EMPIRICALCOVARIANCE

SKLEARNCOVARIANCE EMPIRICALCOVARIANCE XASSUMECENTEREDFALSE

COMPUTES THE MAXIMUM LIKELIHOOD COVARIANCE ESTIMATOR

PARAMETERS

XNDARRAY SHAPE NSAMPLES NFEATURES DATA FROM WHICH TO COMPUTE THE COVARIANCE ESTIMATE

ASSUMECENTERED BOOLEAN IF TRUE DATA WILL NOT BE CENTERED BEFORE COMPUTATION USEFUL WHEN WORKING WITH DATA WHOSE MEAN IS ALMOST BUT NOT EXACTLY ZERO IF FALSE DATA WILL BE CENTERED BEFORE COMPUTATION

RETURNS

COVARIANCE 2D NDARRAY SHAPE NFEATURES NFEATURES EMPIRICAL COVARIANCE MAXIMUM LIKE

LIHOOD ESTIMATOR

EXAMPLES USING SKLEARNCOVARIANCEEMPIRICALCOVARIANCE

- SHRINKAGE COVARIANCE ESTIMATION LEDOITWOLF VS OAS AND MAXLIKELIHOOD

6610SKLEARNCOVARIANCE GRAPHICALASSO

SKLEARNCOVARIANCE GRAPHICALASSO EMPCOV ALPHA COVINITNONE MODE’CD’ TOL00001

ENETTOL00001 MAXITER100 VERBOSEFALSE RE

TURNCOSTSFALSE EPS2220446049250313E16 RE

TURNNITERFALSE

L1PENALIZED COVARIANCE ESTIMATOR

1560 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

READ MORE IN THE USER GUIDE

PARAMETERS

EMPCOV 2D NDARRAY SHAPE NFEATURES NFEATURES EMPIRICAL COVARIANCE FROM WHICH TO COMPUTE THE COVARIANCE ESTIMATE

ALPHA POSITIVE FLOAT THE REGULARIZATION PARAMETER THE HIGHER ALPHA THE MORE REGULARIZATION THE SPARSER THE INVERSE COVARIANCE

COVINIT 2D ARRAY NFEATURES NFEATURES OPTIONAL THE INITIAL GUESS FOR THE COVARIANCE

MODE 'CD' 'LARS' THE LASSO SOLVER TO USE COORDINATE DESCENT OR LARS USE LARS FOR

VERY SPARSE UNDERLYING GRAPHS WHERE  $P \gg N$  ELSEWHERE PREFER CD WHICH IS MORE NUMERICALLY STABLE

TOLPOSITIVE FLOAT OPTIONAL THE TOLERANCE TO DECLARE CONVERGENCE IF THE DUAL GAP GOES BELOW THIS VALUE ITERATIONS ARE STOPPED

ENETTOL POSITIVE FLOAT OPTIONAL THE TOLERANCE FOR THE ELASTIC NET SOLVER USED TO CALCULATE THE DESCENT DIRECTION THIS PARAMETER CONTROLS THE ACCURACY OF THE SEARCH DIRECTION FOR A GIVEN COLUMN UPDATE NOT OF THE OVERALL PARAMETER ESTIMATE ONLY USED FOR MODE 'CD'

MAXITER INTEGER OPTIONAL THE MAXIMUM NUMBER OF ITERATIONS

VERBOSE BOOLEAN OPTIONAL IF VERBOSE IS TRUE THE OBJECTIVE FUNCTION AND DUAL GAP ARE PRINTED AT EACH ITERATION

RETURNCOSTS BOOLEAN OPTIONAL IF RETURNCOSTS IS TRUE THE OBJECTIVE FUNCTION AND DUAL GAP AT EACH ITERATION ARE RETURNED

EPS FLOAT OPTIONAL THE MACHINEPRECISION REGULARIZATION IN THE COMPUTATION OF THE CHOLESKY DIAGONAL FACTORS INCREASE THIS FOR VERY ILLCONDITIONED SYSTEMS

RETURNNITER BOOL OPTIONAL WHETHER OR NOT TO RETURN THE NUMBER OF ITERATIONS

RETURNS

COVARIANCE 2D NDARRAY SHAPE NFEATURES NFEATURES THE ESTIMATED COVARIANCE MATRIX

PRECISION 2D NDARRAY SHAPE NFEATURES NFEATURES THE ESTIMATED SPARSE PRECISION MATRIX

COSTS LIST OF OBJECTIVE DUALGAP PAIRS THE LIST OF VALUES OF THE OBJECTIVE FUNCTION AND THE DUAL GAP AT EACH ITERATION RETURNED ONLY IF RETURNCOSTS IS TRUE

NITER INT NUMBER OF ITERATIONS RETURNED ONLY IF RETURNNITER IS SET TO TRUE

SEE ALSO

GRAPHICALASSO GRAPHICALASSOCV

NOTES

THE ALGORITHM EMPLOYED TO SOLVE THIS PROBLEM IS THE GLASSO ALGORITHM FROM THE FRIEDMAN 2008 BIOSTATISTICS PAPER IT IS THE SAME ALGORITHM AS IN THE R GLASSO PACKAGE

ONE POSSIBLE DIFFERENCE WITH THE GLASSO R PACKAGE IS THAT THE DIAGONAL COEFFICIENTS ARE NOT PENALIZED

66SKLEARNCOVARIANCE COVARIANCE ESTIMATORS 1561

SCIKITLEARN USER GUIDE RELEASE 0213

6611SKLEARNCOVARIANCE LEDOITWOLF

SKLEARNCOVARIANCE LEDOITWOLF XASSUMECENTEREDFALSE BLOCKSIZE1000

ESTIMATES THE SHRUNK LEDOITWOLF COVARIANCE MATRIX

READ MORE IN THE USER GUIDE

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES DATA FROM WHICH TO COMPUTE THE COVARIANCE ESTIMATE

ASSUMECENTERED BOOLEAN DEFAULTFALSE IF TRUE DATA WILL NOT BE CENTERED BEFORE COMPUTATION USEFUL TO WORK WITH DATA WHOSE MEAN IS SIGNIFICANTLY EQUAL TO ZERO BUT IS NOT EXACTLY ZERO IF FALSE DATA WILL BE CENTERED BEFORE COMPUTATION

BLOCKSIZE INT DEFAULT1000 SIZE OF THE BLOCKS INTO WHICH THE COVARIANCE MATRIX WILL BE SPLIT THIS IS PURELY A MEMORY OPTIMIZATION AND DOES NOT AFFECT RESULTS

RETURNS

SHRUNKCOV ARRAYLIKE SHAPE NFEATURES NFEATURES SHRUNK COVARIANCE

SHRINKAGE FLOAT COEFFICIENT IN THE CONVEX COMBINATION USED FOR THE COMPUTATION OF THE SHRUNK ESTIMATE

NOTES

THE REGULARIZED SHRUNK COVARIANCE IS

$1 \text{ SHRINKAGE COV SHRINKAGE MU NPIDENTITYNFEATURES}$

WHERE  $\text{MU TRACECOV NFEATURES}$

EXAMPLES USING SKLEARNCOVARIANCELEDOITWOLF

- SPARSE INVERSE COVARIANCE ESTIMATION

6612SKLEARNCOVARIANCE OAS

SKLEARNCOVARIANCE OASXASSUMECENTEREDFALSE

ESTIMATE COVARIANCE WITH THE ORACLE APPROXIMATING SHRINKAGE ALGORITHM

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES DATA FROM WHICH TO COMPUTE THE COVARIANCE ESTIMATE

ASSUMECENTERED BOOLEAN IF TRUE DATA WILL NOT BE CENTERED BEFORE COMPUTATION USEFUL TO WORK WITH DATA WHOSE MEAN IS SIGNIFICANTLY EQUAL TO ZERO BUT IS NOT EXACTLY ZERO IF FALSE DATA WILL BE CENTERED BEFORE COMPUTATION

RETURNS

SHRUNKCOV ARRAYLIKE SHAPE NFEATURES NFEATURES SHRUNK COVARIANCE

SHRINKAGE FLOAT COEFFICIENT IN THE CONVEX COMBINATION USED FOR THE COMPUTATION OF THE SHRUNK ESTIMATE

1562 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

NOTES

THE REGULARISED SHRUNK COVARIANCE IS

1 SHRINKAGE COV SHRINKAGE MU NPIDENTITYNFEATURES

WHERE MU TRACECOV NFEATURES

THE FORMULA WE USED TO IMPLEMENT THE OAS IS SLIGHTLY MODIFIED COMPARED TO THE ONE GIVEN IN THE ARTICLE SEE

OAS FOR MORE DETAILS

6613SKLEARNCOVARIANCE SHRUNKCOVARIANCE

SKLEARNCOVARIANCE SHRUNKCOVARIANCE EMPCOV SHRINKAGE01

CALCULATES A COVARIANCE MATRIX SHRUNK ON THE DIAGONAL

READ MORE IN THE USER GUIDE

PARAMETERS

EMPCOV ARRAYLIKE SHAPE NFEATURES NFEATURES COVARIANCE MATRIX TO BE SHRUNK

SHRINKAGE FLOAT 0 SHRINKAGE 1 COEFFICIENT IN THE CONVEX COMBINATION USED FOR THE

COMPUTATION OF THE SHRUNK ESTIMATE

RETURNS

SHRUNKCOV ARRAYLIKE SHRUNK COVARIANCE

NOTES

THE REGULARIZED SHRUNK COVARIANCE IS GIVEN BY

1 SHRINKAGE COV SHRINKAGE MU NPIDENTITYNFEATURES

WHERE MU TRACECOV NFEATURES

67SKLEARNCROSSDECOMPOSITION CROSS DECOMPOSITION

USER GUIDE SEE THE CROSS DECOMPOSITION SECTION FOR FURTHER DETAILS

CROSSDECOMPOSITIONCCA NCOMPONENTS CCA CANONICAL CORRELATION ANALYSIS

CROSSDECOMPOSITIONPLSCANONICAL PLSCANONICAL IMPLEMENTS THE 2 BLOCKS CANONICAL PLS OF

THE ORIGINAL WOLD ALGORITHM TENENHAUS 1998 P204 RE

FERRED AS PLSC2A IN WEGELIN 2000

CROSSDECOMPOSITIONPLSREGRESSION PLS REGRESSION

CROSSDECOMPOSITIONPLSSVD NCOMPONENTS

PARTIAL LEAST SQUARE SVD

671SKLEARNCROSSDECOMPOSITION CCA

CLASSSKLEARNCROSSDECOMPOSITION CCANCOMPONENTS2 SCALETRUE MAXITER500 TOL1E06

COPYTRUE

CCA CANONICAL CORRELATION ANALYSIS

67SKLEARNCROSSDECOMPOSITION CROSS DECOMPOSITION 1563

SCIKITLEARN USER GUIDE RELEASE 0213  
CCA INHERITS FROM PLS WITH MODE" B" AND DEFLATIONMODE"CANONICAL"  
READ MORE IN THE USER GUIDE  
PARAMETERS  
NCOMPONENTS INT DEFAULT 2 NUMBER OF COMPONENTS TO KEEP  
SCALE BOOLEAN DEFAULT TRUE WHETHER TO SCALE THE DATA  
MAXITER AN INTEGER DEFAULT 500 THE MAXIMUM NUMBER OF ITERATIONS OF THE NIPALS INNER  
LOOP  
TOLNONNEGATIVE REAL DEFAULT 1E06 THE TOLERANCE USED IN THE ITERATIVE ALGORITHM  
COPY BOOLEAN WHETHER THE DEFLATION BE DONE ON A COPY LET THE DEFAULT VALUE TO TRUE UNLESS  
YOU DON'T CARE ABOUT SIDE EFFECTS  
ATTRIBUTES  
XWEIGHTS ARRAY P NCOMPONENTS X BLOCK WEIGHTS VECTORS  
YWEIGHTS ARRAY Q NCOMPONENTS Y BLOCK WEIGHTS VECTORS  
XLOADINGS ARRAY P NCOMPONENTS X BLOCK LOADINGS VECTORS  
YLOADINGS ARRAY Q NCOMPONENTS Y BLOCK LOADINGS VECTORS  
XSCORES ARRAY NSAMPLES NCOMPONENTS X SCORES  
YSCORES ARRAY NSAMPLES NCOMPONENTS Y SCORES  
XROTATIONS ARRAY P NCOMPONENTS X BLOCK TO LATENTS ROTATIONS  
YROTATIONS ARRAY Q NCOMPONENTS Y BLOCK TO LATENTS ROTATIONS  
NITER ARRAYLIKE NUMBER OF ITERATIONS OF THE NIPALS INNER LOOP FOR EACH COMPONENT  
SEE ALSO  
PLSCANONICAL  
PLSSVD  
NOTES  
FOR EACH COMPONENT K FIND THE WEIGHTS U V THAT MAXIMIZES MAX CORR<sub>XK</sub> U Y<sub>K</sub> V SUCH THAT U V  
1  
NOTE THAT IT MAXIMIZES ONLY THE CORRELATIONS BETWEEN THE SCORES  
THE RESIDUAL MATRIX OF X X<sub>K1</sub> BLOCK IS OBTAINED BY THE DEFLATION ON THE CURRENT X SCORE XSCORE  
THE RESIDUAL MATRIX OF Y Y<sub>K1</sub> BLOCK IS OBTAINED BY DEFLATION ON THE CURRENT Y SCORE  
REFERENCES  
JACOB A WEGELIN A SURVEY OF PARTIAL LEAST SQUARES PLS METHODS WITH EMPHASIS ON THE TWOBLOCK CASE  
TECHNICAL REPORT 371 DEPARTMENT OF STATISTICS UNIVERSITY OF WASHINGTON SEATTLE 2000  
IN FRENCH BUT STILL A REFERENCE TENENHAUS M 1998 LA REGRESSION PLS THEORIE ET PRATIQUE PARIS EDITIONS  
TECHNIC  
1564 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES

```
FROM SKLEARNCROSSDECOMPOSITION IMPORT CCA
X 0 0 1 100 222 354
Y 01 02 09 11 62 59 119 123
CCA CCANCOMPONENTS1
CCAFITX Y
```

CCACOPYTRUE MAXITER500 NCOMPONENTS1 SCALETRUE TOL1E06

XC YC CCATTRANSFORMX Y

METHODS

FITSELF X Y FIT MODEL TO DATA

FITTRANSFORM SELF X Y LEARN AND APPLY THE DIMENSION REDUCTION ON THE TRAIN DATA

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PREDICT SELF X COPY APPLY THE DIMENSION REDUCTION LEARNED ON THE TRAIN DATA

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF X Y COPY APPLY THE DIMENSION REDUCTION LEARNED ON THE TRAIN DATA

INIT SELFNCOMPONENTS2 SCALETRUE MAXITER500 TOL1E06 COPYTRUE

FITSELFXY

FIT MODEL TO DATA

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING VECTORS WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF PREDICTORS

YARRAYLIKE SHAPE NSAMPLES NTARGETS TARGET VECTORS WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NTARGETS IS THE NUMBER OF RESPONSE VARIABLES

FITTRANSFORM SELFXYNONE

LEARN AND APPLY THE DIMENSION REDUCTION ON THE TRAIN DATA

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING VECTORS WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF PREDICTORS

YARRAYLIKE SHAPE NSAMPLES NTARGETS TARGET VECTORS WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NTARGETS IS THE NUMBER OF RESPONSE VARIABLES

RETURNS

XSCORES IF Y IS NOT GIVEN XSCORES YSCORES OTHERWISE

GETPARAMS SELFDEEPTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

67SKLEARNCROSSDECOMPOSITION CROSS DECOMPOSITION 1565

SCIKITLEARN USER GUIDE RELEASE 0213

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELFXCOPYTRUE

APPLY THE DIMENSION REDUCTION LEARNED ON THE TRAIN DATA

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING VECTORS WHERE NSAMPLES IS THE

NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF PREDICTORS

COPY BOOLEAN DEFAULT TRUE WHETHER TO COPY X AND Y OR PERFORM INPLACE NORMALIZATION

NOTES

THIS CALL REQUIRES THE ESTIMATION OF A P X Q MATRIX WHICH MAY BE AN ISSUE IN HIGH DIMENSIONAL SPACE

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{U}{V}$  WHERE U IS THE RESIDUAL SUM OF SQUARES  $Y_{TRUE} - Y_{PRED}$

$2SUM$  AND V IS THE TOTAL SUM OF SQUARES  $Y_{TRUE} - Y_{TRUEMEAN}$   $2SUM$  THE BEST POSSIBLE SCORE

IS 10 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS

PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 00

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES FOR SOME ESTIMATORS THIS MAY

BE A PRECOMPUTED KERNEL MATRIX INSTEAD SHAPE NSAMPLES NSAMPLESFITTED WHERE

NSAMPLESFITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE VALUES FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT R2 OF SELF-PREDICTX WRT Y

NOTES

THE R2 SCORE USED WHEN CALLING SCORE ON A REGRESSOR WILL USE MULTIOUTPUTUNIFORMAVERAGE

FROM VERSION 023 TO KEEP CONSISTENT WITH METRICSR2SCORE THIS WILL INFLUENCE THE SCORE

METHOD OF ALL THE MULTIOUTPUT REGRESSORS EXCEPT FOR MULTIOUTPUTMULTIOUTPUTREGRESSOR

TO SPECIFY THE DEFAULT VALUE MANUALLY AND AVOID THE WARNING PLEASE EITHER CALL METRICSR2SCORE

DIRECTLY OR MAKE A CUSTOM SCORER WITH METRICSMAKESCORER THE BUILTIN SCORER R2 USES

MULTIOUTPUTUNIFORMAVERAGE

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE

PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT

OF A NESTED OBJECT

RETURNS

SELF

1566 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

TRANSFORM SELFXYNONE COPYTRUE

APPLY THE DIMENSION REDUCTION LEARNED ON THE TRAIN DATA

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING VECTORS WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF PREDICTORS

YARRAYLIKE SHAPE NSAMPLES NTARGETS TARGET VECTORS WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NTARGETS IS THE NUMBER OF RESPONSE VARIABLES

COPY BOOLEAN DEFAULT TRUE WHETHER TO COPY X AND Y OR PERFORM INPLACE NORMALIZATION

RETURNS

XSCORES IF Y IS NOT GIVEN XSCORES YSCORES OTHERWISE

EXAMPLES USING SKLEARNCROSSDECOMPOSITIONCCA

- MULTILABEL CLASSIFICATION
- COMPARE CROSS DECOMPOSITION METHODS

672SKLEARNCROSSDECOMPOSITION PLSCANONICAL

CLASSSSKLEARNCROSSDECOMPOSITION PLSCANONICAL NCOMPONENTS2 SCALETRUE ALGO

RITHM'NIPALS' MAXITER500 TOL1E06

COPYTRUE

PLSCANONICAL IMPLEMENTS THE 2 BLOCKS CANONICAL PLS OF THE ORIGINAL WOLD ALGORITHM TENENHAUS 1998 P204 REFERRED AS PLSC2A IN WEGELIN 2000

THIS CLASS INHERITS FROM PLS WITH MODE"A" AND DEFLATIONMODE"CANONICAL" NORMYWEIGHTSTRUE AND ALGORITHM"NIPALS" BUT SVD SHOULD PROVIDE SIMILAR RESULTS UP TO NUMERICAL ERRORS

READ MORE IN THE USER GUIDE

PARAMETERS

NCOMPONENTS INT DEFAULT 2 NUMBER OF COMPONENTS TO KEEP

SCALE BOOLEAN DEFAULT TRUE OPTION TO SCALE DATA

ALGORITHM STRING "NIPALS" OR "SVD" THE ALGORITHM USED TO ESTIMATE THE WEIGHTS IT WILL BE CALLED NCOMPONENTS TIMES IE ONCE FOR EACH ITERATION OF THE OUTER LOOP

MAXITER AN INTEGER DEFAULT 500 THE MAXIMUM NUMBER OF ITERATIONS OF THE NIPALS INNER LOOP USED ONLY IF ALGORITHM"NIPALS"

TOLNONNEGATIVE REAL DEFAULT 1E06 THE TOLERANCE USED IN THE ITERATIVE ALGORITHM

COPY BOOLEAN DEFAULT TRUE WHETHER THE DEFLATION SHOULD BE DONE ON A COPY LET THE DEFAULT VALUE TO TRUE UNLESS YOU DON'T CARE ABOUT SIDE EFFECT

ATTRIBUTES

XWEIGHTS ARRAY SHAPE P NCOMPONENTS X BLOCK WEIGHTS VECTORS

YWEIGHTS ARRAY SHAPE Q NCOMPONENTS Y BLOCK WEIGHTS VECTORS

XLOADINGS ARRAY SHAPE P NCOMPONENTS X BLOCK LOADINGS VECTORS

YLOADINGS ARRAY SHAPE Q NCOMPONENTS Y BLOCK LOADINGS VECTORS

XSCORES ARRAY SHAPE NSAMPLES NCOMPONENTS X SCORES

67SKLEARNCROSSDECOMPOSITION CROSS DECOMPOSITION 1567

SCIKITLEARN USER GUIDE RELEASE 0213

YSCORES ARRAY SHAPE NSAMPLES NCOMPONENTS Y SCORES

XROTATIONS ARRAY SHAPE P NCOMPONENTS X BLOCK TO LATENTS ROTATIONS

YROTATIONS ARRAY SHAPE Q NCOMPONENTS Y BLOCK TO LATENTS ROTATIONS

NITER ARRAYLIKE NUMBER OF ITERATIONS OF THE NIPALS INNER LOOP FOR EACH COMPONENT NOT USEFUL IF THE ALGORITHM PROVIDED IS "SVD"

SEE ALSO

CCA

PLSSVD

NOTES

MATRICES

T XSCORES

U YSCORES

W XWEIGHTS

C YWEIGHTS

P XLOADINGS

Q YLOADINGS

ARE COMPUTED SUCH THAT

$X - T P^T$  ERR ANDY  $U - Q T^T$  ERR

T K XK W K FORKINRANGENCOMPONENTS

U K YK C K FORKINRANGENCOMPONENTS

XROTATIONS W PT W1

YROTATIONS C QT C1

WHERE XK AND YK ARE RESIDUAL MATRICES AT ITERATION K

SLIDES EXPLAINING PLS

FOR EACH COMPONENT K FIND WEIGHTS U V THAT OPTIMIZE

$\max \text{CORR}(X - U Y, V^T X - U^T Y)$  SUCH THAT  $U^T U = 1$

NOTE THAT IT MAXIMIZES BOTH THE CORRELATIONS BETWEEN THE SCORES AND THE INTRABLOCK VARIANCES

THE RESIDUAL MATRIX OF X XK1 BLOCK IS OBTAINED BY THE DEFLATION ON THE CURRENT X SCORE XSCORE

THE RESIDUAL MATRIX OF Y YK1 BLOCK IS OBTAINED BY DEFLATION ON THE CURRENT Y SCORE THIS PERFORMS A CANONICAL SYMMETRIC VERSION OF THE PLS REGRESSION BUT SLIGHTLY DIFFERENT THAN THE CCA THIS IS MOSTLY USED FOR MODELING

THIS IMPLEMENTATION PROVIDES THE SAME RESULTS THAT THE "PLSPM" PACKAGE PROVIDED IN THE R LANGUAGE R

PROJECT USING THE FUNCTION PLSCAX Y RESULTS ARE EQUAL OR COLLINEAR WITH THE FUNCTION PLS MODE

CANONICAL OF THE "MIXOMICS" PACKAGE THE DIFFERENCE RELIES IN THE FACT THAT MIXOMICS IMPLEMENTATION DOES NOT EXACTLY IMPLEMENT THE WOLD ALGORITHM SINCE IT DOES NOT NORMALIZE YWEIGHTS TO ONE

REFERENCES

JACOB A WEGELIN A SURVEY OF PARTIAL LEAST SQUARES PLS METHODS WITH EMPHASIS ON THE TWOBLOCK CASE

TECHNICAL REPORT 371 DEPARTMENT OF STATISTICS UNIVERSITY OF WASHINGTON SEATTLE 2000

1568 CHAPTER 6 API REFERENCE

```
SCIKITLEARN USER GUIDE RELEASE 0213
TENENHAUS M 1998 LA REGRESSION PLS THEORIE ET PRATIQUE PARIS EDITIONS TECHNIC
EXAMPLES
FROM SKLEARNCROSSDECOMPOSITION IMPORT PLSCANONICAL
X 0 0 1 100 222 254
Y 01 02 09 11 62 59 119 123
PLSCA PLSCANONICALNCOMPONENTS2
PLSCAFITX Y

PLSCANONICALALGORITHMNIPALS COPYTRUE MAXITER500 NCOMPONENTS2
SCALETRUE TOL1E06
XC YC PLSCATTRANSFORMX Y
METHODS
FITSELF X Y FIT MODEL TO DATA
FITTRANSFORM SELF X Y LEARN AND APPLY THE DIMENSION REDUCTION ON THE TRAIN
DATA
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR
PREDICT SELF X COPY APPLY THE DIMENSION REDUCTION LEARNED ON THE TRAIN DATA
SCORE SELF X Y SAMPLEWEIGHT RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PRE
DICTION
SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR
TRANSFORM SELF X Y COPY APPLY THE DIMENSION REDUCTION LEARNED ON THE TRAIN DATA
INIT SELFNCOMPONENTS2 SCALETRUE ALGORITHM'NIPALS' MAXITER500 TOL1E06
COPYTRUE
FITSELFXY
FIT MODEL TO DATA
PARAMETERS
XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING VECTORS WHERE NSAMPLES IS THE
NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF PREDICTORS
YARRAYLIKE SHAPE NSAMPLES NTARGETS TARGET VECTORS WHERE NSAMPLES IS THE NUM
BER OF SAMPLES AND NTARGETS IS THE NUMBER OF RESPONSE VARIABLES
FITTRANSFORM SELFXYNONE
LEARN AND APPLY THE DIMENSION REDUCTION ON THE TRAIN DATA
PARAMETERS
XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING VECTORS WHERE NSAMPLES IS THE
NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF PREDICTORS
YARRAYLIKE SHAPE NSAMPLES NTARGETS TARGET VECTORS WHERE NSAMPLES IS THE NUMBER
OF SAMPLES AND NTARGETS IS THE NUMBER OF RESPONSE VARIABLES
RETURNS
XSCORES IF Y IS NOT GIVEN XSCORES YSCORES OTHERWISE
67SKLEARNCROSSDECOMPOSITION CROSS DECOMPOSITION 1569
```

SCIKITLEARN USER GUIDE RELEASE 0213

GETPARAMS SELFDEEPTREE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELFXYCOPYTRUE

APPLY THE DIMENSION REDUCTION LEARNED ON THE TRAIN DATA

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING VECTORS WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF PREDICTORS

COPY BOOLEAN DEFAULT TRUE WHETHER TO COPY X AND Y OR PERFORM INPLACE NORMALIZATION

NOTES

THIS CALL REQUIRES THE ESTIMATION OF A P X Q MATRIX WHICH MAY BE AN ISSUE IN HIGH DIMENSIONAL SPACE

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{U}{V}$  WHERE U IS THE RESIDUAL SUM OF SQUARES  $\sum (y_{true} - y_{pred})^2$  AND V IS THE TOTAL SUM OF SQUARES  $\sum (y_{true} - y_{true\_mean})^2$  THE BEST POSSIBLE SCORE IS 10 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 00

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES FOR SOME ESTIMATORS THIS MAY BE A PRECOMPUTED KERNEL MATRIX INSTEAD SHAPE NSAMPLES NSAMPLESFITTED WHERE NSAMPLESFITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE VALUES FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT R2 OF SELF-PREDICTX WRT Y

NOTES

THE R2 SCORE USED WHEN CALLING SCORE ON A REGRESSOR WILL USE MULTIOUTPUTUNIFORMAVERAGE FROM VERSION 023 TO KEEP CONSISTENT WITH METRICSR2SCORE THIS WILL INFLUENCE THE SCORE METHOD OF ALL THE MULTIOUTPUT REGRESSORS EXCEPT FOR MULTIOUTPUTMULTIOUTPUTREGRESSOR TO SPECIFY THE DEFAULT VALUE MANUALLY AND AVOID THE WARNING PLEASE EITHER CALL METRICSR2SCORE DIRECTLY OR MAKE A CUSTOM SCORER WITH METRICSMAKESCORER THE BUILTIN SCORER R2 USES MULTIOUTPUTUNIFORMAVERAGE

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

1570 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS  
SELF

TRANSFORM SELFXYNONE COPYTRUE  
APPLY THE DIMENSION REDUCTION LEARNED ON THE TRAIN DATA

PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING VECTORS WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF PREDICTORS  
YARRAYLIKE SHAPE NSAMPLES NTARGETS TARGET VECTORS WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NTARGETS IS THE NUMBER OF RESPONSE VARIABLES  
COPY BOOLEAN DEFAULT TRUE WHETHER TO COPY X AND Y OR PERFORM INPLACE NORMALIZATION

RETURNS  
XSCORES IF Y IS NOT GIVEN XSCORES YSCORES OTHERWISE

EXAMPLES USING SKLEARNCROSSDECOMPOSITIONPLSCANONICAL

- COMPARE CROSS DECOMPOSITION METHODS

673SKLEARNCROSSDECOMPOSITION PLSREGRESSION

CLASSSSKLEARNCROSSDECOMPOSITION PLSREGRESSION NCOMPONENTS2 SCALETRUE

MAXITER500 TOL1E06 COPYTRUE

PLS REGRESSION

PLSREGRESSION IMPLEMENTS THE PLS 2 BLOCKS REGRESSION KNOWN AS PLS2 OR PLS1 IN CASE OF ONE DIMENSIONAL RESPONSE THIS CLASS INHERITS FROM PLS WITH MODE“A” DEFLATIONMODE“REGRESSION” NORMYWEIGHTSFALSE AND ALGORITHM“NIPALS”

READ MORE IN THE USER GUIDE

PARAMETERS  
NCOMPONENTS INT DEFAULT 2 NUMBER OF COMPONENTS TO KEEP  
SCALE BOOLEAN DEFAULT TRUE WHETHER TO SCALE THE DATA  
MAXITER AN INTEGER DEFAULT 500 THE MAXIMUM NUMBER OF ITERATIONS OF THE NIPALS INNER LOOP USED ONLY IF ALGORITHM“NIPALS”  
TOLNONNEGATIVE REAL TOLERANCE USED IN THE ITERATIVE ALGORITHM DEFAULT 1E06  
COPY BOOLEAN DEFAULT TRUE WHETHER THE DEFLATION SHOULD BE DONE ON A COPY LET THE DEFAULT VALUE TO TRUE UNLESS YOU DON’T CARE ABOUT SIDE EFFECT

ATTRIBUTES  
XWEIGHTS ARRAY P NCOMPONENTS X BLOCK WEIGHTS VECTORS  
YWEIGHTS ARRAY Q NCOMPONENTS Y BLOCK WEIGHTS VECTORS  
XLOADINGS ARRAY P NCOMPONENTS X BLOCK LOADINGS VECTORS  
YLOADINGS ARRAY Q NCOMPONENTS Y BLOCK LOADINGS VECTORS

675SKLEARNCROSSDECOMPOSITION CROSS DECOMPOSITION 1571

SCIKITLEARN USER GUIDE RELEASE 0213

XSCORES ARRAY NSAMPLES NCOMPONENTS X SCORES  
YSCORES ARRAY NSAMPLES NCOMPONENTS Y SCORES  
XROTATIONS ARRAY P NCOMPONENTS X BLOCK TO LATENTS ROTATIONS  
YROTATIONS ARRAY Q NCOMPONENTS Y BLOCK TO LATENTS ROTATIONS  
COEF ARRAY P Q THE COEFFICIENTS OF THE LINEAR MODEL  $Y = X \text{ COEF} + \text{ERR}$   
NITER ARRAYLIKE NUMBER OF ITERATIONS OF THE NIPALS INNER LOOP FOR EACH COMPONENT

NOTES

MATRICES

T XSCORES  
U YSCORES  
W XWEIGHTS  
C YWEIGHTS  
P XLOADINGS  
Q YLOADINGS  
ARE COMPUTED SUCH THAT  
 $X = T \text{ PT} + \text{ERR}$  AND  $Y = U \text{ QT} + \text{ERR}$   
T K XK W K FORKINRANGENCOMPONENTS  
U K YK C K FORKINRANGENCOMPONENTS

XROTATIONS W PT W1  
YROTATIONS C QT C1  
WHERE XK AND YK ARE RESIDUAL MATRICES AT ITERATION K  
SLIDES EXPLAINING PLS

FOR EACH COMPONENT K FIND WEIGHTS U V THAT OPTIMIZES  $\text{MAX CORR}_K(U, YK) \sqrt{\text{STD}_K(U) \text{STD}_K(Y)}$  SUCH THAT  $U^T U = 1$

NOTE THAT IT MAXIMIZES BOTH THE CORRELATIONS BETWEEN THE SCORES AND THE INTRABLOCK VARIANCES  
THE RESIDUAL MATRIX OF X XK1 BLOCK IS OBTAINED BY THE DEFLATION ON THE CURRENT X SCORE XSCORE  
THE RESIDUAL MATRIX OF Y YK1 BLOCK IS OBTAINED BY DEFLATION ON THE CURRENT X SCORE THIS PERFORMS THE PLS  
REGRESSION KNOWN AS PLS2 THIS MODE IS PREDICTION ORIENTED

THIS IMPLEMENTATION PROVIDES THE SAME RESULTS THAT 3 PLS PACKAGES PROVIDED IN THE R LANGUAGE RPROJECT

- “MIXOMICS” WITH FUNCTION PLSX Y MODE “REGRESSION”
- “PLSPM ” WITH FUNCTION PLSREG2X Y
- “PLS” WITH FUNCTION OSCORESPLSFITX Y

REFERENCES

JACOB A WEGELIN A SURVEY OF PARTIAL LEAST SQUARES PLS METHODS WITH EMPHASIS ON THE TWOBLOCK CASE  
TECHNICAL REPORT 371 DEPARTMENT OF STATISTICS UNIVERSITY OF WASHINGTON SEATTLE 2000  
IN FRENCH BUT STILL A REFERENCE TENENHAUS M 1998 LA REGRESSION PLS THEORIE ET PRATIQUE PARIS EDITIONS  
TECHNIC  
1572 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES

FROM SKLEARNCROSSDECOMPOSITION IMPORT PLSREGRESSION

X 0 0 1 100 222 254

Y 01 02 09 11 62 59 119 123

PLS2 PLSREGRESSIONNCOMPONENTS2

PLS2FITX Y

PLSREGRESSIONCOPYTRUE MAXITER500 NCOMPONENTS2 SCALETRUE

TOL1E06

YPRED PLS2PREDICTX

METHODS

FITSELF X Y FIT MODEL TO DATA

FITTRANSFORM SELF X Y LEARN AND APPLY THE DIMENSION REDUCTION ON THE TRAIN

DATA

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PREDICT SELF X COPY APPLY THE DIMENSION REDUCTION LEARNED ON THE TRAIN DATA

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PRE  
DICTION

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF X Y COPY APPLY THE DIMENSION REDUCTION LEARNED ON THE TRAIN DATA

INIT SELFNCOMPONENTS2 SCALETRUE MAXITER500 TOL1E06 COPYTRUE

FITSELFXY

FIT MODEL TO DATA

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING VECTORS WHERE NSAMPLES IS THE  
NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF PREDICTORS

YARRAYLIKE SHAPE NSAMPLES NTARGETS TARGET VECTORS WHERE NSAMPLES IS THE NUM  
BER OF SAMPLES AND NTARGETS IS THE NUMBER OF RESPONSE VARIABLES

FITTRANSFORM SELFXYNONE

LEARN AND APPLY THE DIMENSION REDUCTION ON THE TRAIN DATA

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING VECTORS WHERE NSAMPLES IS THE  
NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF PREDICTORS

YARRAYLIKE SHAPE NSAMPLES NTARGETS TARGET VECTORS WHERE NSAMPLES IS THE NUMBER  
OF SAMPLES AND NTARGETS IS THE NUMBER OF RESPONSE VARIABLES

RETURNS

XSCORES IF Y IS NOT GIVEN XSCORES YSCORES OTHERWISE

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

67SKLEARNCROSSDECOMPOSITION CROSS DECOMPOSITION 1573

SCIKITLEARN USER GUIDE RELEASE 0213

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELFXCOPYTRUE

APPLY THE DIMENSION REDUCTION LEARNED ON THE TRAIN DATA

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING VECTORS WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF PREDICTORS

COPY BOOLEAN DEFAULT TRUE WHETHER TO COPY X AND Y OR PERFORM INPLACE NORMALIZATION

NOTES

THIS CALL REQUIRES THE ESTIMATION OF A P X Q MATRIX WHICH MAY BE AN ISSUE IN HIGH DIMENSIONAL SPACE

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{U}{V}$  WHERE U IS THE RESIDUAL SUM OF SQUARES  $Y_{TRUE} - Y_{PRED}$  2SUM AND V IS THE TOTAL SUM OF SQUARES  $Y_{TRUE} - Y_{TRUEMEAN}$  2SUM THE BEST POSSIBLE SCORE IS 10 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 00

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES FOR SOME ESTIMATORS THIS MAY BE A PRECOMPUTED KERNEL MATRIX INSTEAD SHAPE NSAMPLES NSAMPLESFITTED WHERE NSAMPLESFITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE VALUES FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT R2 OF SELF PREDICT X WRT Y

NOTES

THE R2 SCORE USED WHEN CALLING SCORE ON A REGRESSOR WILL USE MULTIOUTPUTUNIFORMAVERAGE FROM VERSION 023 TO KEEP CONSISTENT WITH METRICSR2SCORE THIS WILL INFLUENCE THE SCORE METHOD OF ALL THE MULTIOUTPUT REGRESSORS EXCEPT FOR MULTIOUTPUTMULTIOUTPUTREGRESSOR TO SPECIFY THE DEFAULT VALUE MANUALLY AND AVOID THE WARNING PLEASE EITHER CALL METRICSR2SCORE DIRECTLY OR MAKE A CUSTOM SCORER WITH METRICSMAKESCORER THE BUILTIN SCORER R2 USES MULTIOUTPUTUNIFORMAVERAGE

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

1574 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

SELF

TRANSFORM SELFXYNONE COPYTRUE

APPLY THE DIMENSION REDUCTION LEARNED ON THE TRAIN DATA

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING VECTORS WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF PREDICTORS

YARRAYLIKE SHAPE NSAMPLES NTARGETS TARGET VECTORS WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NTARGETS IS THE NUMBER OF RESPONSE VARIABLES

COPY BOOLEAN DEFAULT TRUE WHETHER TO COPY X AND Y OR PERFORM INPLACE NORMALIZATION

RETURNS

XSCORES IF Y IS NOT GIVEN XSCORES YSCORES OTHERWISE

EXAMPLES USING SKLEARNCROSSDECOMPOSITIONPLSREGRESSION

- COMPARE CROSS DECOMPOSITION METHODS

674SKLEARNCROSSDECOMPOSITION PLSSVD

CLASSSSKLEARNCROSSDECOMPOSITION PLSSVDNCOMPONENTS2 SCALETRUE COPYTRUE

PARTIAL LEAST SQUARE SVD

SIMPLY PERFORM A SVD ON THE CROSSCOVARIANCE MATRIX X\*Y THERE ARE NO ITERATIVE DEFLATION HERE

READ MORE IN THE USER GUIDE

PARAMETERS

NCOMPONENTS INT DEFAULT 2 NUMBER OF COMPONENTS TO KEEP

SCALE BOOLEAN DEFAULT TRUE WHETHER TO SCALE X AND Y

COPY BOOLEAN DEFAULT TRUE WHETHER TO COPY X AND Y OR PERFORM INPLACE COMPUTATIONS

ATTRIBUTES

XWEIGHTS ARRAY P NCOMPONENTS X BLOCK WEIGHTS VECTORS

YWEIGHTS ARRAY Q NCOMPONENTS Y BLOCK WEIGHTS VECTORS

XSCORES ARRAY NSAMPLES NCOMPONENTS X SCORES

YSCORES ARRAY NSAMPLES NCOMPONENTS Y SCORES

SEE ALSO

PLSCANONICAL

CCA

EXAMPLES

675SKLEARNCROSSDECOMPOSITION CROSS DECOMPOSITION 1575

SCIKITLEARN USER GUIDE RELEASE 0213

```
import numpy as np
from sklearn.cross_decomposition import PLSSVD
X = np.array(0, 1)
100
222
254
Y = np.array(0, 1, 2)
09 11
62 59
119 123
PLSCA = PLSSVD(n_components=2)
PLSCA.fit(X, Y)
PLSSVD(copy=True, n_components=2, scale=True)
XC, YC = PLSCA.transform(X, Y)
XC.shape, YC.shape
(4, 2), (4, 2)

METHODS
fit(self, X, Y): fit model to data
fit_transform(self, X, Y): learn and apply the dimension reduction on the train data
get_params(self): get parameters for this estimator
set_params(self, **kwargs): set the parameters of this estimator
transform(self, X, Y): apply the dimension reduction learned on the train data
init(self, n_components=2, scale=True, copy=True)
fit(self, X, Y): fit model to data
parameters
X: array-like shape (n_samples, n_features) training vectors where n_samples is the number of samples and n_features is the number of predictors
Y: array-like shape (n_samples, n_targets) target vectors where n_samples is the number of samples and n_targets is the number of response variables
fit_transform(self, X, Y): learn and apply the dimension reduction on the train data
parameters
X: array-like shape (n_samples, n_features) training vectors where n_samples is the number of samples and n_features is the number of predictors
Y: array-like shape (n_samples, n_targets) target vectors where n_samples is the number of samples and n_targets is the number of response variables
returns
X_scores if Y is not given, Y_scores otherwise
get_params(self): get parameters for this estimator
parameters
1576 CHAPTER 6 API REFERENCE
```

SCIKITLEARN USER GUIDE RELEASE 0213

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELFXYNONE

APPLY THE DIMENSION REDUCTION LEARNED ON THE TRAIN DATA

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING VECTORS WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF PREDICTORS

YARRAYLIKE SHAPE NSAMPLES NTARGETS TARGET VECTORS WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NTARGETS IS THE NUMBER OF RESPONSE VARIABLES

68SKLEARNDATASETS DATASETS

THESKLEARNDATASETS MODULE INCLUDES UTILITIES TO LOAD DATASETS INCLUDING METHODS TO LOAD AND FETCH POPULAR REFERENCE DATASETS IT ALSO FEATURES SOME ARTIFICIAL DATA GENERATORS

USER GUIDE SEE THE DATASET LOADING UTILITIES SECTION FOR FURTHER DETAILS

681 LOADERS

DATASETSCLEARDATAHOME DATAHOME DELETE ALL THE CONTENT OF THE DATA HOME CACHE

DATASETSDUMPSVMLIGHTFILE X Y F DUMP THE DATASET IN SVMLIGHT LIBSVM FILE FORMAT

DATASETSFETCH20NEWSGROUPS DATAHOME

LOAD THE FILENAMES AND DATA FROM THE 20 NEWSGROUPS

DATASET CLASSIFICATION

DATASETSFETCH20NEWSGROUPSVECTORIZED LOAD THE 20 NEWSGROUPS DATASET AND VECTORIZE IT INTO TOKEN COUNTS CLASSIFICATION

DATASETSFETCHCALIFORNIAHOUSING LOAD THE CALIFORNIA HOUSING DATASET REGRESSION

DATASETSFETCHCOVTYPE DATAHOME LOAD THE COVERTYPE DATASET CLASSIFICATION

DATASETSFETCHKDDCUP99 SUBSET DATAHOME

LOAD THE KDDCUP99 DATASET CLASSIFICATION

DATASETSFETCHLFWPAIRS SUBSET LOAD THE LABELED FACES IN THE WILD LFW PAIRS DATASET CLASSIFICATION

DATASETSFETCHLFWPEOPLE DATAHOME LOAD THE LABELED FACES IN THE WILD LFW PEOPLE DATASET CLASSIFICATION

DATASETSFETCHOLIVETTIFACES DATAHOME

LOAD THE OLIVETTI FACES DATASET FROM ATT CLASSIFICATION

DATASETSFETCHOPENML NAME VERSION FETCH DATASET FROM OPENML BY NAME OR DATASET ID

CONTINUED ON NEXT PAGE

68SKLEARNDATASETS DATASETS 1577

SCIKITLEARN USER GUIDE RELEASE 0213

TABLE 647 – CONTINUED FROM PREVIOUS PAGE

DATASETSFETCHRCV1 DATAHOME SUBSET    LOAD THE RCV1 MULTILABEL DATASET CLASSIFICATION

DATASETSFETCHSPECIESDISTRIBUTIONS    LOADER FOR SPECIES DISTRIBUTION DATASET FROM PHILLIPS ET

DATASETSGETDATAHOME DATAHOME RETURN THE PATH OF THE SCIKITLEARN DATA DIR

DATASETSLOADBOSTON RETURNXY LOAD AND RETURN THE BOSTON HOUSEPRICES DATASET REGRES  
SION

DATASETSLOADBREASTCANCER RETURNXY LOAD AND RETURN THE BREAST CANCER WISCONSIN DATASET CLAS  
SIFICATION

DATASETSLOADDIABETES RETURNXY LOAD AND RETURN THE DIABETES DATASET REGRESSION

DATASETSLOADDIGITS NCLASS RETURNXY LOAD AND RETURN THE DIGITS DATASET CLASSIFICATION

DATASETSLOADFILES CONTAINERPATH    LOAD TEXT FILES WITH CATEGORIES AS SUBFOLDER NAMES

DATASETSLOADIRIS RETURNXY LOAD AND RETURN THE IRIS DATASET CLASSIFICATION

DATASETSLOADLINNERUD RETURNXY LOAD AND RETURN THE LINNERUD DATASET MULTIVARIATE REGRES  
SION

DATASETSLOADSAMPLEIMAGE IMAGENAME LOAD THE NUMPY ARRAY OF A SINGLE SAMPLE IMAGE

DATASETSLOADSAMPLEIMAGES    LOAD SAMPLE IMAGES FOR IMAGE MANIPULATION

DATASETSLOADSVMLIGHTFILE F NFEATURES  
  LOAD DATASETS IN THE SVMLIGHT LIBSVM FORMAT INTO SPARSE  
CSR MATRIX

DATASETSLOADSVMLIGHTFILES FILES    LOAD DATASET FROM MULTIPLE FILES IN SVMLIGHT FORMAT

DATASETSLOADWINE RETURNXY LOAD AND RETURN THE WINE DATASET CLASSIFICATION

SKLEARNDATASETS CLEARDATAHOME

SKLEARNDATASETS CLEARDATAHOME DATAHOMENONE

DELETE ALL THE CONTENT OF THE DATA HOME CACHE

PARAMETERS

DATAHOME STR NONE THE PATH TO SCIKITLEARN DATA DIR

SKLEARNDATASETS DUMPSVMLIGHTFILE

SKLEARNDATASETS DUMPSVMLIGHTFILE XY F ZEROBASEDTRUE COMMENTNONE

QUERYIDNONE MULTILABELFALSE

DUMP THE DATASET IN SVMLIGHT LIBSVM FILE FORMAT

THIS FORMAT IS A TEXTBASED FORMAT WITH ONE SAMPLE PER LINE IT DOES NOT STORE ZERO VALUED FEATURES HENCE IS  
SUITABLE FOR SPARSE DATASET

THE FIRST ELEMENT OF EACH LINE CAN BE USED TO STORE A TARGET VARIABLE TO PREDICT

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES TRAINING VECTORS WHERE  
NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

YARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NLABELS TARGET VALUES CLASS LABELS  
MUST BE AN INTEGER OR FLOAT OR ARRAYLIKE OBJECTS OF INTEGER OR FLOAT FOR MULTILABEL CLASSIFICA  
TIONS

FSTRING OR FILELIKE IN BINARY MODE IF STRING SPECIFIES THE PATH THAT WILL CONTAIN THE DATA IF  
FILELIKE DATA WILL BE WRITTEN TO F F SHOULD BE OPENED IN BINARY MODE

ZEROBASED BOOLEAN OPTIONAL WHETHER COLUMN INDICES SHOULD BE WRITTEN ZEROBASED TRUE OR  
ONEBASED FALSE

COMMENT STRING OPTIONAL COMMENT TO INSERT AT THE TOP OF THE FILE THIS SHOULD BE EITHER A  
UNICODE STRING WHICH WILL BE ENCODED AS UTF8 OR AN ASCII BYTE STRING IF A COMMENT

1578 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

IS GIVEN THEN IT WILL BE PRECEDED BY ONE THAT IDENTIFIES THE FILE AS HAVING BEEN DUMPED BY

SCIKITLEARN NOTE THAT NOT ALL TOOLS GROK COMMENTS IN SVMLIGHT FILES

QUERYID ARRAYLIKE SHAPE NSAMPLES ARRAY CONTAINING PAIRWISE PREFERENCE CONSTRAINTS

QID IN SVMLIGHT FORMAT

MULTILABEL BOOLEAN OPTIONAL SAMPLES MAY HAVE SEVERAL LABELS EACH SEE [HTTPS://WWW.CS.CMU.EDU/ELKAN/EDUTWCJLINLIBSVMTOOLS/DATASETS/MULTILABEL.HTML](https://www.cs.cmu.edu/~elkan/EDUTWCJLINLIBSVMTOOLS/DATASETS/MULTILABEL.HTML)

NEW IN VERSION 0.17 PARAMETER MULTILABEL TO SUPPORT MULTILABEL DATASETS

EXAMPLES USING SKLEARN.DATASETS.DUMPSVMLIGHTFILE

- LIBSVM GUI

SKLEARN.DATASETS.FETCH20NEWSGROUPS

SKLEARN.DATASETS.FETCH20NEWSGROUPS.DAT.HOME.NONE.SUBSET='TRAIN'.CATEGORIES.NONE

SHUFFLE.TRUE.RANDOM.STATE.42.REMOVE.DOWN

LOAD.IF.MISSING.TRUE

LOAD THE FILENAMES AND DATA FROM THE 20 NEWSGROUPS DATASET CLASSIFICATION

DOWNLOAD IT IF NECESSARY

CLASSES 20

SAMPLES TOTAL 18846

DIMENSIONALITY 1

FEATURES TEXT

READ MORE IN THE USER GUIDE

PARAMETERS

DAT.HOME.OPTIONAL.DEFAULT.NONE.SPECIFY A DOWNLOAD AND CACHE FOLDER FOR THE DATASETS IF

NONE ALL SCIKITLEARN DATA IS STORED IN 'SCIKITLEARN.DAT.' SUBFOLDERS

SUBSET 'TRAIN' OR 'TEST' 'ALL' OPTIONAL SELECT THE DATASET TO LOAD 'TRAIN' FOR THE TRAINING SET

'TEST' FOR THE TEST SET 'ALL' FOR BOTH WITH SHUFFLED ORDERING

CATEGORIES NONE OR COLLECTION OF STRING OR UNICODE IF NONE DEFAULT LOAD ALL THE CATEGORIES IF

NOT NONE LIST OF CATEGORY NAMES TO LOAD OTHER CATEGORIES IGNORED

SHUFFLE.BOOL.OPTIONAL.WHETHER OR NOT TO SHUFFLE THE DATA MIGHT BE IMPORTANT FOR MODELS THAT

MAKE THE ASSUMPTION THAT THE SAMPLES ARE INDEPENDENT AND IDENTICALLY DISTRIBUTED IID

SUCH AS STOCHASTIC GRADIENT DESCENT

RANDOM.STATE.INT.RANDOM.STATE.INSTANCE OR NONE DEFAULT DETERMINES RANDOM NUMBER GEN

ERATION FOR DATASET SHUFFLING PASS AN INT FOR REPRODUCIBLE OUTPUT ACROSS MULTIPLE FUNCTION

CALLS SEE GLOSSARY

REMOVE.TUPLE.MAY.CONTAIN ANY SUBSET OF 'HEADERS' 'FOOTERS' 'QUOTES' EACH OF THESE ARE

KINDS OF TEXT THAT WILL BE DETECTED AND REMOVED FROM THE NEWSGROUP POSTS PREVENTING CLAS

SIFIERS FROM OVERFITTING ON METADATA

'HEADERS' REMOVES NEWSGROUP HEADERS 'FOOTERS' REMOVES BLOCKS AT THE ENDS OF POSTS THAT

LOOK LIKE SIGNATURES AND 'QUOTES' REMOVES LINES THAT APPEAR TO BE QUOTING ANOTHER POST

'HEADERS' FOLLOWS AN EXACT STANDARD THE OTHER FILTERS ARE NOT ALWAYS CORRECT

68SKLEARN.DATASETS.DATASETS.1579

SCIKITLEARN USER GUIDE RELEASE 0213

DOWNLOADIFMISSING OPTIONAL TRUE BY DEFAULT IF FALSE RAISE AN IOERROR IF THE DATA IS NOT  
LOCALLY AVAILABLE INSTEAD OF TRYING TO DOWNLOAD THE DATA FROM THE SOURCE SITE

RETURNS

BUNCH BUNCH OBJECT WITH THE FOLLOWING ATTRIBUTE

- BUNCHDATA LIST LENGTH NSAMPLES
- BUNCHTARGET ARRAY SHAPE NSAMPLES
- BUNCHFILENAMES LIST LENGTH NSAMPLES
- BUNCHDESCR A DESCRIPTION OF THE DATASET
- BUNCHTARGETNAMES A LIST OF CATEGORIES OF THE RETURNED DATA LENGTH NCLASSES THIS  
DEPENDS ON THE CATEGORIES PARAMETER

EXAMPLES USING SKLEARNDATASETSFETCH20NEWSGROUPS

- TOPIC EXTRACTION WITH NONNEGATIVE MATRIX FACTORIZATION AND LATENT DIRICHLET ALLOCATION
- BICLUSTERING DOCUMENTS WITH THE SPECTRAL COCLUSTERING ALGORITHM
- COLUMN TRANSFORMER WITH HETEROGENEOUS DATA SOURCES
- SAMPLE PIPELINE FOR TEXT FEATURE EXTRACTION AND EVALUATION
- FEATUREHASHER AND DICTVECTORIZER COMPARISON
- CLUSTERING TEXT DOCUMENTS USING KMEANS
- CLASSIFICATION OF TEXT DOCUMENTS USING SPARSE FEATURES

SKLEARNDATASETS FETCH20NEWSGROUPSVECTORIZED

SKLEARNDATASETS FETCH20NEWSGROUPSVECTORIZED SUBSET‘TRAIN’ REMOVE

DATAHOMENONE DOWN

LOADIFMISSINGTRUE RE

TURNXYFALSE

LOAD THE 20 NEWSGROUPS DATASET AND VECTORIZE IT INTO TOKEN COUNTS CLASSIFICATION

DOWNLOAD IT IF NECESSARY

THIS IS A CONVENIENCE FUNCTION THE TRANSFORMATION IS DONE USING THE DEFAULT SETTINGS FOR SKLEARN

FEATUREEXTRACTIONTEXTCOUNTVECTORIZER FOR MORE ADVANCED USAGE STOPWORD

FILTERING NGRAM EXTRACTION ETC COMBINE FETCH20NEWSGROUPS WITH A CUSTOM SKLEARN

FEATUREEXTRACTIONTEXTCOUNTVECTORIZER SKLEARNFEATUREEXTRACTIONTEXT

HASHINGVECTORIZER SKLEARNFEATUREEXTRACTIONTEXTTFIDFTRANSFORMER OR

SKLEARNFEATUREEXTRACTIONTEXTTFIDFVECTORIZER

CLASSES 20

SAMPLES TOTAL 18846

DIMENSIONALITY 130107

FEATURES REAL

READ MORE IN THE USER GUIDE

PARAMETERS

1580 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

SUBSET 'TRAIN' OR 'TEST' 'ALL' OPTIONAL SELECT THE DATASET TO LOAD 'TRAIN' FOR THE TRAINING SET 'TEST' FOR THE TEST SET 'ALL' FOR BOTH WITH SHUFFLED ORDERING

REMOVE TUPLE MAY CONTAIN ANY SUBSET OF 'HEADERS' 'FOOTERS' 'QUOTES' EACH OF THESE ARE KINDS OF TEXT THAT WILL BE DETECTED AND REMOVED FROM THE NEWSGROUP POSTS PREVENTING CLASSIFIERS FROM OVERFITTING ON METADATA

'HEADERS' REMOVES NEWSGROUP HEADERS 'FOOTERS' REMOVES BLOCKS AT THE ENDS OF POSTS THAT LOOK LIKE SIGNATURES AND 'QUOTES' REMOVES LINES THAT APPEAR TO BE QUOTING ANOTHER POST

DATAHOME OPTIONAL DEFAULT NONE SPECIFY AN DOWNLOAD AND CACHE FOLDER FOR THE DATASETS IF NONE ALL SCIKITLEARN DATA IS STORED IN 'SCIKITLEARNDATA' SUBFOLDERS

DOWNLOADIFMISSING OPTIONAL TRUE BY DEFAULT IF FALSE RAISE AN IOERROR IF THE DATA IS NOT LOCALLY AVAILABLE INSTEAD OF TRYING TO DOWNLOAD THE DATA FROM THE SOURCE SITE

RETURNXY BOOLEAN DEFAULT FALSE IF TRUE RETURNS DATADATA DATATARGET INSTEAD OF A BUNCH OBJECT

NEW IN VERSION 020

RETURNS

BUNCH BUNCH OBJECT WITH THE FOLLOWING ATTRIBUTE

- BUNCHDATA SPARSE MATRIX SHAPE NSAMPLES NFEATURES
- BUNCHTARGET ARRAY SHAPE NSAMPLES
- BUNCHTARGETNAMES A LIST OF CATEGORIES OF THE RETURNED DATA LENGTH NCLASSES
- BUNCHDESCR A DESCRIPTION OF THE DATASET

DATA TARGET TUPLE IFRETURNXY IS TRUE NEW IN VERSION 020

EXAMPLES USING SKLEARNDATASETSFETCH20NEWSGROUPSVECTORIZED

- THE JOHNSONLINDENSTRAUSS BOUND FOR EMBEDDING WITH RANDOM PROJECTIONS
- MODEL COMPLEXITY INFLUENCE
- MULTICLASS SPARSE LOGISITIC REGRESSION ON NEWGROUPS20

SKLEARNDATASETS FETCHCALIFORNIAHOUSING

SKLEARNDATASETS FETCHCALIFORNIAHOUSING DATAHOMENONE DOWNLOADIFMISSINGTRUE

RETURNXYFALSE

LOAD THE CALIFORNIA HOUSING DATASET REGRESSION

SAMPLES TOTAL 20640

DIMENSIONALITY 8

FEATURES REAL

TARGET REAL 015 5

READ MORE IN THE USER GUIDE

PARAMETERS

DATAHOME OPTIONAL DEFAULT NONE SPECIFY ANOTHER DOWNLOAD AND CACHE FOLDER FOR THE DATASETS BY DEFAULT ALL SCIKITLEARN DATA IS STORED IN 'SCIKITLEARNDATA' SUBFOLDERS

68SKLEARNDATASETS DATASETS 1581

SCIKITLEARN USER GUIDE RELEASE 0213

DOWNLOADIFMISSING OPTIONAL DEFAULTTRUE IF FALSE RAISE A IOERROR IF THE DATA IS NOT LOCALLY AVAILABLE INSTEAD OF TRYING TO DOWNLOAD THE DATA FROM THE SOURCE SITE

RETURNXY BOOLEAN DEFAULTFALSE IF TRUE RETURNS DATADATA DATATARGET IN STEAD OF A BUNCH OBJECT

NEW IN VERSION 020

RETURNS

DATASET DICTLIKE OBJECT WITH THE FOLLOWING ATTRIBUTES

DATASETDATA NDARRAY SHAPE 20640 8 EACH ROW CORRESPONDING TO THE 8 FEATURE VALUES IN ORDER

DATSETTARGET NUMPY ARRAY OF SHAPE 20640 EACH VALUE CORRESPONDS TO THE AVERAGE HOUSE VALUE IN UNITS OF 100000

DATASETFEATURENAMES ARRAY OF LENGTH 8 ARRAY OF ORDERED FEATURE NAMES USED IN THE DATASET

DATSETDESCR STRING DESCRIPTION OF THE CALIFORNIA HOUSING DATASET

DATA TARGET TUPLE IFRETURNXY IS TRUE NEW IN VERSION 020

NOTES

THIS DATASET CONSISTS OF 20640 SAMPLES AND 9 FEATURES

EXAMPLES USING SKLEARNDATASETSFETCHCALIFORNIAHOUSING

- IMPUTING MISSING VALUES WITH VARIANTS OF ITERATIVEIMPUTER
- PARTIAL DEPENDENCE PLOTS
- COMPARE THE EFFECT OF DIFFERENT SCALERS ON DATA WITH OUTLIERS

SKLEARNDATASETS FETCHCOVTYPE

SKLEARNDATASETS FETCHCOVTYPE DATAHOMENONE DOWNLOADIFMISSINGTRUE RAN

DOMSTATENONE SHUFFLEFALSE RETURNXYFALSE

LOAD THE COVERTYPE DATASET CLASSIFICATION

DOWNLOAD IT IF NECESSARY

CLASSES 7

SAMPLES TOTAL 581012

DIMENSIONALITY 54

FEATURES INT

READ MORE IN THE USER GUIDE

PARAMETERS

DATASHOME STRING OPTIONAL SPECIFY ANOTHER DOWNLOAD AND CACHE FOLDER FOR THE DATASETS BY

DEFAULT ALL SCIKITLEARN DATA IS STORED IN 'SCIKITLEARNDATA' SUBFOLDERS

DOWNLOADIFMISSING BOOLEAN DEFAULTTRUE IF FALSE RAISE A IOERROR IF THE DATA IS NOT LOCALLY AVAILABLE INSTEAD OF TRYING TO DOWNLOAD THE DATA FROM THE SOURCE SITE

1582 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE DEFAULT DETERMINES RANDOM NUMBER GENERATION FOR DATASET SHUFFLING PASS AN INT FOR REPRODUCIBLE OUTPUT ACROSS MULTIPLE FUNCTION CALLS SEE GLOSSARY

SHUFFLE BOOL DEFAULTFALSE WHETHER TO SHUFFLE DATASET

RETURNXY BOOLEAN DEFAULTFALSE IF TRUE RETURNS DATADATA DATATARGET INSTEAD OF A BUNCH OBJECT

NEW IN VERSION 020

RETURNS

DATASET DICTLIKE OBJECT WITH THE FOLLOWING ATTRIBUTES

DATASETDATA NUMPY ARRAY OF SHAPE 581012 54 EACH ROW CORRESPONDS TO THE 54 FEATURES IN THE DATASET

DATASETTARGET NUMPY ARRAY OF SHAPE 581012 EACH VALUE CORRESPONDS TO ONE OF THE 7 FOREST COVERTYPES WITH VALUES RANGING BETWEEN 1 TO 7

DATASETDESCR STRING DESCRIPTION OF THE FOREST COVERTYPE DATASET

DATA TARGET TUPLE IFRETURNXY IS TRUE NEW IN VERSION 020

SKLEARNDATASETS FETCHKDDCUP99

SKLEARNDATASETS FETCHKDDCUP99 SUBSETNONE DATAHOMENONE SHUFFLEFALSE

RANDOMSTATENONE PERCENT10TRUE DOWN

LOADIFMISSINGTRUE RETURNXYFALSE

LOAD THE KDDCUP99 DATASET CLASSIFICATION

DOWNLOAD IT IF NECESSARY

CLASSES 23

SAMPLES TOTAL 4898431

DIMENSIONALITY 41

FEATURES DISCRETE INT OR CONTINUOUS FLOAT

READ MORE IN THE USER GUIDE

NEW IN VERSION 018

PARAMETERS

SUBSET NONE 'SA' 'SF' 'HTTP' 'SMTP' TO RETURN THE CORRESPONDING CLASSICAL SUBSETS OF KDDCUP 99 IF NONE RETURN THE ENTIRE KDDCUP 99 DATASET

DATAHOME STRING OPTIONAL SPECIFY ANOTHER DOWNLOAD AND CACHE FOLDER FOR THE DATASETS BY DEFAULT ALL SCIKITLEARN DATA IS STORED IN 'SCIKITLEARNDATA' SUBFOLDERS VERSIONADDED 019

SHUFFLE BOOL DEFAULTFALSE WHETHER TO SHUFFLE DATASET

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE DEFAULT DETERMINES RANDOM NUMBER GENERATION FOR DATASET SHUFFLING AND FOR SELECTION OF ABNORMAL SAMPLES IF SUBSETSA PASS AN INT FOR REPRODUCIBLE OUTPUT ACROSS MULTIPLE FUNCTION CALLS SEE GLOSSARY

PERCENT10 BOOL DEFAULTTRUE WHETHER TO LOAD ONLY 10 PERCENT OF THE DATA

DOWNLOADIFMISSING BOOL DEFAULTTRUE IF FALSE RAISE A IOERROR IF THE DATA IS NOT LOCALLY AVAILABLE INSTEAD OF TRYING TO DOWNLOAD THE DATA FROM THE SOURCE SITE

68SKLEARNDATASETS DATASETS 1583

SCIKITLEARN USER GUIDE RELEASE 0213

RETURNXY BOOLEAN DEFAULTFALSE IF TRUE RETURNS DATA TARGET INSTEAD OF A BUNCH  
OBJECT SEE BELOW FOR MORE INFORMATION ABOUT THE DATA ANDTARGET OBJECT  
NEW IN VERSION 020

RETURNS  
DATA BUNCH  
DICTIONARYLIKE OBJECT THE INTERESTING ATTRIBUTES ARE

- ‘DATA’ THE DATA TO LEARN
- ‘TARGET’ THE REGRESSION TARGET FOR EACH SAMPLE
- ‘DESCR’ A DESCRIPTION OF THE DATASET

DATA TARGET TUPLE IFRETURNXY IS TRUE NEW IN VERSION 020

SKLEARNDATASETS FETCHLFWPAIRS  
SKLEARNDATASETS FETCHLFWPAIRS SUBSET‘TRAIN’ DATAHOMENONE FUNNELEDTRUE RE  
SIZE05 COLORFALSE SLICESLICE70 195NONE SLICE78  
172NONE DOWNLOADIFMISSINGTRUE  
LOAD THE LABELED FACES IN THE WILD LFW PAIRS DATASET CLASSIFICATION  
DOWNLOAD IT IF NECESSARY  
CLASSES 5749  
SAMPLES TOTAL 13233  
DIMENSIONALITY 5828  
FEATURES REAL BETWEEN 0 AND 255

IN THE OFFICIAL READMETXT THIS TASK IS DESCRIBED AS THE “RESTRICTED” TASK AS I AM NOT SURE AS TO IMPLEMENT THE  
“UNRESTRICTED” VARIANT CORRECTLY I LEFT IT AS UNSUPPORTED FOR NOW  
THE ORIGINAL IMAGES ARE 250 X 250 PIXELS BUT THE DEFAULT SLICE AND RESIZE ARGUMENTS REDUCE THEM TO 62 X 47  
READ MORE IN THE USER GUIDE

PARAMETERS

SUBSET OPTIONAL DEFAULT ‘TRAIN’ SELECT THE DATASET TO LOAD ‘TRAIN’ FOR THE DEVELOPMENT TRAINING  
SET ‘TEST’ FOR THE DEVELOPMENT TEST SET AND ‘10FOLDS’ FOR THE OFFICIAL EVALUATION SET THAT IS  
MEANT TO BE USED WITH A 10FOLDS CROSS VALIDATION

DATASHOME OPTIONAL DEFAULT NONE SPECIFY ANOTHER DOWNLOAD AND CACHE FOLDER FOR THE  
DATASETS BY DEFAULT ALL SCIKITLEARN DATA IS STORED IN ‘SCIKITLEARNDATA’ SUBFOLDERS

FUNNELED BOOLEAN OPTIONAL DEFAULT TRUE DOWNLOAD AND USE THE FUNNELED VARIANT OF THE  
DATASET

RESIZE FLOAT OPTIONAL DEFAULT 05 RATIO USED TO RESIZE THE EACH FACE PICTURE

COLOR BOOLEAN OPTIONAL DEFAULT FALSE KEEP THE 3 RGB CHANNELS INSTEAD OF AVERAGING THEM TO  
A SINGLE GRAY LEVEL CHANNEL IF COLOR IS TRUE THE SHAPE OF THE DATA HAS ONE MORE DIMENSION  
THAN THE SHAPE WITH COLOR FALSE

SLICE OPTIONAL PROVIDE A CUSTOM 2D SLICE HEIGHT WIDTH TO EXTRACT THE ‘INTERESTING’ PART OF  
THE JPEG FILES AND AVOID USE STATISTICAL CORRELATION FROM THE BACKGROUND

1584 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

DOWNLOADIFMISSING OPTIONAL TRUE BY DEFAULT IF FALSE RAISE A IOERROR IF THE DATA IS NOT  
LOCALLY AVAILABLE INSTEAD OF TRYING TO DOWNLOAD THE DATA FROM THE SOURCE SITE

RETURNS

THE DATA IS RETURNED AS A BUNCH OBJECT WITH THE FOLLOWING ATTRIBUTES

DATA NUMPY ARRAY OF SHAPE 2200 5828 SHAPE DEPENDS ON SUBSET EACH ROW CORRESPONDS  
TO 2 RAVEL'D FACE IMAGES OF ORIGINAL SIZE 62 X 47 PIXELS CHANGING THE SLICE RESIZE OR  
SUBSET PARAMETERS WILL CHANGE THE SHAPE OF THE OUTPUT

PAIRS NUMPY ARRAY OF SHAPE 2200 2 62 47 SHAPE DEPENDS ON SUBSET EACH ROW HAS 2  
FACE IMAGES CORRESPONDING TO SAME OR DIFFERENT PERSON FROM THE DATASET CONTAINING 5749  
PEOPLE CHANGING THE SLICE RESIZE ORSUBSET PARAMETERS WILL CHANGE THE SHAPE OF  
THE OUTPUT

TARGET NUMPY ARRAY OF SHAPE 2200 SHAPE DEPENDS ON SUBSET LABELS ASSOCIATED TO EACH  
PAIR OF IMAGES THE TWO LABEL VALUES BEING DIFFERENT PERSONS OR THE SAME PERSON

DESCR STRING DESCRIPTION OF THE LABELED FACES IN THE WILD LFW DATASET

SKLEARNDATASETS FETCHLFWPEOPLE

SKLEARNDATASETS FETCHLFWPEOPLE DATAHOMENONE FUNNELEDTRUE RESIZE05

MINFACESPERPERSON0 COLORFALSE SLICESLICE70

195 NONE SLICE78 172 NONE DOWN

LOADIFMISSINGTRUE RETURNXYFALSE

LOAD THE LABELED FACES IN THE WILD LFW PEOPLE DATASET CLASSIFICATION

DOWNLOAD IT IF NECESSARY

CLASSES 5749

SAMPLES TOTAL 13233

DIMENSIONALITY 5828

FEATURES REAL BETWEEN 0 AND 255

READ MORE IN THE USER GUIDE

PARAMETERS

DATAHOME OPTIONAL DEFAULT NONE SPECIFY ANOTHER DOWNLOAD AND CACHE FOLDER FOR THE  
DATASETS BY DEFAULT ALL SCIKITLEARN DATA IS STORED IN 'SCIKITLEARNDATA' SUBFOLDERS

FUNNELED BOOLEAN OPTIONAL DEFAULT TRUE DOWNLOAD AND USE THE FUNNELED VARIANT OF THE  
DATASET

RESIZE FLOAT OPTIONAL DEFAULT 05 RATIO USED TO RESIZE THE EACH FACE PICTURE

MINFACESPERPERSON INT OPTIONAL DEFAULT NONE THE EXTRACTED DATASET WILL ONLY RETAIN PIC  
TURES OF PEOPLE THAT HAVE AT LEAST MINFACESPERPERSON DIFFERENT PICTURES

COLOR BOOLEAN OPTIONAL DEFAULT FALSE KEEP THE 3 RGB CHANNELS INSTEAD OF AVERAGING THEM TO  
A SINGLE GRAY LEVEL CHANNEL IF COLOR IS TRUE THE SHAPE OF THE DATA HAS ONE MORE DIMENSION  
THAN THE SHAPE WITH COLOR FALSE

SLICE OPTIONAL PROVIDE A CUSTOM 2D SLICE HEIGHT WIDTH TO EXTRACT THE 'INTERESTING' PART OF  
THE JPEG FILES AND AVOID USE STATISTICAL CORRELATION FROM THE BACKGROUND

DOWNLOADIFMISSING OPTIONAL TRUE BY DEFAULT IF FALSE RAISE A IOERROR IF THE DATA IS NOT  
LOCALLY AVAILABLE INSTEAD OF TRYING TO DOWNLOAD THE DATA FROM THE SOURCE SITE

68SKLEARNDATASETS DATASETS 1585

SCIKITLEARN USER GUIDE RELEASE 0213

RETURNXY BOOLEAN DEFAULTFALSE IF TRUE RETURNS DATASETDATA DATASET  
TARGET INSTEAD OF A BUNCH OBJECT SEE BELOW FOR MORE INFORMATION ABOUT THE DATASET  
DATA ANDDATASETTARGET OBJECT  
NEW IN VERSION 020

RETURNS

DATASETDICTLIKE OBJECT WITH THE FOLLOWING ATTRIBUTES

DATASETDATA NUMPY ARRAY OF SHAPE 13233 2914 EACH ROW CORRESPONDS TO A RAVELLED FACE  
IMAGE OF ORIGINAL SIZE 62 X 47 PIXELS CHANGING THE SLICE OR RESIZE PARAMETERS WILL  
CHANGE THE SHAPE OF THE OUTPUT

DATASETIMAGES NUMPY ARRAY OF SHAPE 13233 62 47 EACH ROW IS A FACE IMAGE CORRESPONDING  
TO ONE OF THE 5749 PEOPLE IN THE DATASET CHANGING THE SLICE OR RESIZE PARAMETERS WILL  
CHANGE THE SHAPE OF THE OUTPUT

DATASETTARGET NUMPY ARRAY OF SHAPE 13233 LABELS ASSOCIATED TO EACH FACE IMAGE THOSE  
LABELS RANGE FROM 05748 AND CORRESPOND TO THE PERSON IDS

DATASETDESCR STRING DESCRIPTION OF THE LABELED FACES IN THE WILD LFW DATASET

DATA TARGET TUPLE IFRETURNXY IS TRUE NEW IN VERSION 020

EXAMPLES USING SKLEARNDATASETSFETCHLFWPEOPLE

- FACES RECOGNITION EXAMPLE USING EIGENFACES AND SVMs

SKLEARNDATASETS FETCHOLIVETTIFACES

SKLEARNDATASETS FETCHOLIVETTIFACES DATAHOMENONE SHUFFLEFALSE RANDOMSTATE0  
DOWNLOADIFMISSINGTRUE

LOAD THE OLIVETTI FACES DATASET FROM ATT CLASSIFICATION

DOWNLOAD IT IF NECESSARY

CLASSES 40

SAMPLES TOTAL 400

DIMENSIONALITY 4096

FEATURES REAL BETWEEN 0 AND 1

READ MORE IN THE USER GUIDE

PARAMETERS

DATASHOME OPTIONAL DEFAULT NONE SPECIFY ANOTHER DOWNLOAD AND CACHE FOLDER FOR THE  
DATASETS BY DEFAULT ALL SCIKITLEARN DATA IS STORED IN 'SCIKITLEARNDATA' SUBFOLDERS

SHUFFLE BOOLEAN OPTIONAL IF TRUE THE ORDER OF THE DATASET IS SHUFFLED TO AVOID HAVING IMAGES OF  
THE SAME PERSON GROUPED

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE DEFAULT0 DETERMINES RANDOM NUMBER  
GENERATION FOR DATASET SHUFFLING PASS AN INT FOR REPRODUCIBLE OUTPUT ACROSS MULTIPLE FUNCTION  
CALLS SEE GLOSSARY

DOWNLOADIFMISSING OPTIONAL TRUE BY DEFAULT IF FALSE RAISE A IOERROR IF THE DATA IS NOT  
LOCALLY AVAILABLE INSTEAD OF TRYING TO DOWNLOAD THE DATA FROM THE SOURCE SITE

1586 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

RETURNS

AN OBJECT WITH THE FOLLOWING ATTRIBUTES

DATA NUMPY ARRAY OF SHAPE 400 4096 EACH ROW CORRESPONDS TO A RAVELLED FACE IMAGE OF ORIGINAL SIZE 64 X 64 PIXELS

IMAGES NUMPY ARRAY OF SHAPE 400 64 64 EACH ROW IS A FACE IMAGE CORRESPONDING TO ONE OF THE 40 SUBJECTS OF THE DATASET

TARGET NUMPY ARRAY OF SHAPE 400 LABELS ASSOCIATED TO EACH FACE IMAGE THOSE LABELS ARE RANGING FROM 039 AND CORRESPOND TO THE SUBJECT IDS

DESCR STRING DESCRIPTION OF THE MODIFIED OLIVETTI FACES DATASET

EXAMPLES USING SKLEARNDATASETSFETCHOLIVETTIFACES

- FACE COMPLETION WITH A MULTIOUTPUT ESTIMATORS
- ONLINE LEARNING OF A DICTIONARY OF PARTS OF FACES
- FACES DATASET DECOMPOSITIONS
- PIXEL IMPORTANCES WITH A PARALLEL FOREST OF TREES

SKLEARNDATASETS FETCHOPENML

SKLEARNDATASETS FETCHOPENML NAMENONE VERSION'ACTIVE' DATAIDNONE DATAHOMENONE

TARGETCOLUMN'DEFAULT'TARGET' CACHETRUE RETURNXYFALSE

FETCH DATASET FROM OPENML BY NAME OR DATASET ID

DATASETS ARE UNIQUELY IDENTIFIED BY EITHER AN INTEGER ID OR BY A COMBINATION OF NAME AND VERSION IE THERE MIGHT BE MULTIPLE VERSIONS OF THE 'IRIS' DATASET PLEASE GIVE EITHER NAME OR DATAID NOT BOTH IN CASE A NAME IS GIVEN A VERSION CAN ALSO BE PROVIDED

READ MORE IN THE USER GUIDE

NOTE EXPERIMENTAL

THE API IS EXPERIMENTAL PARTICULARLY THE RETURN VALUE STRUCTURE AND MIGHT HAVE SMALL BACKWARDINCOMPATIBLE CHANGES IN FUTURE RELEASES

PARAMETERS

NAME STR OR NONE STRING IDENTIFIER OF THE DATASET NOTE THAT OPENML CAN HAVE MULTIPLE

DATASETS WITH THE SAME NAME

VERSION INTEGER OR 'ACTIVE' DEFAULT'ACTIVE' VERSION OF THE DATASET CAN ONLY BE PROVIDED IF

ALSONAME IS GIVEN IF 'ACTIVE' THE OLDEST VERSION THAT'S STILL ACTIVE IS USED SINCE THERE

MAY BE MORE THAN ONE ACTIVE VERSION OF A DATASET AND THOSE VERSIONS MAY FUNDAMENTALLY BE

DIFFERENT FROM ONE ANOTHER SETTING AN EXACT VERSION IS HIGHLY RECOMMENDED

DATAID INT OR NONE OPENML ID OF THE DATASET THE MOST SPECIFIC WAY OF RETRIEVING A DATASET

IF DATAID IS NOT GIVEN NAME AND POTENTIAL VERSION ARE USED TO OBTAIN A DATASET

DATAHOME STRING OR NONE DEFAULT NONE SPECIFY ANOTHER DOWNLOAD AND CACHE FOLDER FOR THE

DATA SETS BY DEFAULT ALL SCIKITLEARN DATA IS STORED IN 'SCIKITLEARNDATA' SUBFOLDERS

68SKLEARNDATASETS DATASETS 1587

SCIKITLEARN USER GUIDE RELEASE 0213

TARGETCOLUMN STRING LIST OR NONE DEFAULT 'DEFAULTTARGET' SPECIFY THE COLUMN NAME IN THE DATA TO USE AS TARGET IF 'DEFAULTTARGET' THE STANDARD TARGET COLUMN A STORED ON THE SERVER IS USED IFNONE ALL COLUMNS ARE RETURNED AS DATA AND THE TARGET IS NONE IF LIST OF STRINGS ALL COLUMNS WITH THESE NAMES ARE RETURNED AS MULTITARGET NOTE NOT ALL SCIKITLEARN CLASSIFIERS CAN HANDLE ALL TYPES OF MULTIOUTPUT COMBINATIONS

CACHE BOOLEAN DEFAULTTRUE WHETHER TO CACHE DOWNLOADED DATASETS USING JOBLIB

RETURNXY BOOLEAN DEFAULTFALSE IF TRUE RETURNS DATA TARGET INSTEAD OF A BUNCH OBJECT SEE BELOW FOR MORE INFORMATION ABOUT THE DATA ANDTARGET OBJECTS

RETURNS

DATA BUNCH DICTIONARYLIKE OBJECT WITH ATTRIBUTES

DATA NPARRAY OR SCIPYSPARSECSRMATRIX OF FLOATS THE FEATURE MATRIX CATEGORICAL FEATURES ARE ENCODED AS ORDINALS

TARGET NPARRAY THE REGRESSION TARGET OR CLASSIFICATION LABELS IF APPLICABLE DTYPE IS FLOAT IF NUMERIC AND OBJECT IF CATEGORICAL

DESCR STR THE FULL DESCRIPTION OF THE DATASET

FEATURENAMES LIST THE NAMES OF THE DATASET COLUMNS

CATEGORIES DICT MAPS EACH CATEGORICAL FEATURE NAME TO A LIST OF VALUES SUCH THAT THE VALUE ENCODED AS I IS ITH IN THE LIST

DETAILS DICT MORE METADATA FROM OPENML

DATA TARGET TUPLE IFRETURNXY IS TRUE

NOTE EXPERIMENTAL

THIS INTERFACE IS EXPERIMENTAL AND SUBSEQUENT RELEASES MAY CHANGE ATTRIBUTES WITHOUT NOTICE ALTHOUGH THERE SHOULD ONLY BE MINOR CHANGES TO DATA ANDTARGET

MISSING VALUES IN THE 'DATA' ARE REPRESENTED AS NAN'S MISSING VALUES IN 'TARGET' ARE REPRESENTED AS NAN'S NUMERICAL TARGET OR NONE CATEGORICAL TARGET

EXAMPLES USING SKLEARNDATASETSFETCHOPENML

- GAUSSIAN PROCESS REGRESSION GPR ON MAUNA LOA CO2 DATA
- MNIST CLASSIFICATION USING MULTINOMIAL LOGISTIC L1
- EARLY STOPPING OF STOCHASTIC GRADIENT DESCENT
- CLASSIFIER CHAIN
- VISUALIZATION OF MLP WEIGHTS ON MNIST

SKLEARNDATASETS FETCHRCV1

SKLEARNDATASETS FETCHRCV1 DATAHOMENONE SUBSET'ALL' DOWNLOADIFMISSINGTRUE RAN

DOMSTATENONE SHUFFLEFALSE RETURNXYFALSE

LOAD THE RCV1 MULTILABEL DATASET CLASSIFICATION

DOWNLOAD IT IF NECESSARY

1588 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213  
VERSION RCV1V2 VECTORS FULL SETS TOPICS MULTILABELS  
CLASSES 103  
SAMPLES TOTAL 804414  
DIMENSIONALITY 47236  
FEATURES REAL BETWEEN 0 AND 1  
READ MORE IN THE USER GUIDE  
NEW IN VERSION 017  
PARAMETERS  
DATAHOME STRING OPTIONAL SPECIFY ANOTHER DOWNLOAD AND CACHE FOLDER FOR THE DATASETS BY  
DEFAULT ALL SCIKITLEARN DATA IS STORED IN 'SCIKITLEARNDATA' SUBFOLDERS  
SUBSET STRING 'TRAIN' 'TEST' OR 'ALL' DEFAULT'ALL' SELECT THE DATASET TO LOAD 'TRAIN' FOR THE  
TRAINING SET 23149 SAMPLES 'TEST' FOR THE TEST SET 781265 SAMPLES 'ALL' FOR BOTH WITH THE  
TRAINING SAMPLES FIRST IF SHUFFLE IS FALSE THIS FOLLOWS THE OFFICIAL LYRL2004 CHRONOLOGICAL  
SPLIT  
DOWNLOADIFMISSING BOOLEAN DEFAULTTRUE IF FALSE RAISE A IOERROR IF THE DATA IS NOT LOCALLY  
AVAILABLE INSTEAD OF TRYING TO DOWNLOAD THE DATA FROM THE SOURCE SITE  
RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE DEFAULT DETERMINES RANDOM NUMBER GEN  
ERATION FOR DATASET SHUFFLING PASS AN INT FOR REPRODUCIBLE OUTPUT ACROSS MULTIPLE FUNCTION  
CALLS SEE GLOSSARY  
SHUFFLE BOOL DEFAULTFALSE WHETHER TO SHUFFLE DATASET  
RETURNXY BOOLEAN DEFAULTFALSE IF TRUE RETURNS DATASETDATA DATASET  
TARGET INSTEAD OF A BUNCH OBJECT SEE BELOW FOR MORE INFORMATION ABOUT THE DATASET  
DATA ANDDATASETTARGET OBJECT  
NEW IN VERSION 020  
RETURNS  
DATASET DICTLIKE OBJECT WITH THE FOLLOWING ATTRIBUTES  
DATASETDATA SCIPY CSR ARRAY DTYPE NPFLOAT64 SHAPE 804414 47236 THE ARRAY HAS 016 OF  
NON ZERO VALUES  
DATASETTARGET SCIPY CSR ARRAY DTYPE NPUINT8 SHAPE 804414 103 EACH SAMPLE HAS A VALUE  
OF 1 IN ITS CATEGORIES AND 0 IN OTHERS THE ARRAY HAS 315 OF NON ZERO VALUES  
DATASETSAMPLEID NUMPY ARRAY DTYPE NPUINT32 SHAPE 804414 IDENTIFICATION NUMBER OF  
EACH SAMPLE AS ORDERED IN DATASETDATA  
DATASETTARGETNAMES NUMPY ARRAY DTYPE OBJECT LENGTH 103 NAMES OF EACH TARGET RCV1  
TOPICS AS ORDERED IN DATASETTARGET  
DATASETDESCR STRING DESCRIPTION OF THE RCV1 DATASET  
DATA TARGET TUPLE IFRETURNXY IS TRUE NEW IN VERSION 020  
SKLEARNDATASETS FETCHSPECIESDISTRIBUTIONS  
SKLEARNDATASETS FETCHSPECIESDISTRIBUTIONS DATAHOMENONE DOWN  
LOADIFMISSINGTRUE  
LOADER FOR SPECIES DISTRIBUTION DATASET FROM PHILLIPS ET AL 2006  
68SKLEARNDATASETS DATASETS 1589

SCIKITLEARN USER GUIDE RELEASE 0213

READ MORE IN THE USER GUIDE

PARAMETERS

DATAHOME OPTIONAL DEFAULT NONE SPECIFY ANOTHER DOWNLOAD AND CACHE FOLDER FOR THE DATASETS BY DEFAULT ALL SCIKITLEARN DATA IS STORED IN 'SCIKITLEARNDATA' SUBFOLDERS  
DOWNLOADIFMISSING OPTIONAL TRUE BY DEFAULT IF FALSE RAISE A IOERROR IF THE DATA IS NOT LOCALLY AVAILABLE INSTEAD OF TRYING TO DOWNLOAD THE DATA FROM THE SOURCE SITE  
RETURNS

THE DATA IS RETURNED AS A BUNCH OBJECT WITH THE FOLLOWING ATTRIBUTES  
COVERAGES ARRAY SHAPE 14 1592 1212 THESE REPRESENT THE 14 FEATURES MEASURED AT EACH POINT OF THE MAP GRID THE LATITUDELONGITUDE VALUES FOR THE GRID ARE DISCUSSED BELOW MISSING DATA IS REPRESENTED BY THE VALUE 9999  
TRAIN RECORD ARRAY SHAPE 1624 THE TRAINING POINTS FOR THE DATA EACH POINT HAS THREE FIELDS  
• TRAIN'SPECIES' IS THE SPECIES NAME  
• TRAIN'DD LONG' IS THE LONGITUDE IN DEGREES  
• TRAIN'DD LAT' IS THE LATITUDE IN DEGREES  
TEST RECORD ARRAY SHAPE 620 THE TEST POINTS FOR THE DATA SAME FORMAT AS THE TRAINING DATA  
NX NY INTEGERS THE NUMBER OF LONGITUDES X AND LATITUDES Y IN THE GRID  
XLEFTLOWERCORNER YLEFTLOWERCORNER FLOATS THE XY POSITION OF THE LOWERLEFT CORNER IN DEGREES  
GRIDSIZE FLOAT THE SPACING BETWEEN POINTS OF THE GRID IN DEGREES

NOTES

THIS DATASET REPRESENTS THE GEOGRAPHIC DISTRIBUTION OF SPECIES THE DATASET IS PROVIDED BY PHILLIPS ET AL 2006 THE TWO SPECIES ARE

- "BRADYPUS VARIEGATUS" THE BROWNTHOATED SLOTH
- "MICRORYZOMYS MINUTUS" ALSO KNOWN AS THE FOREST SMALL RICE RAT A RODENT THAT LIVES IN PERU COLOMBIA ECUADOR PERU AND VENEZUELA

• FOR AN EXAMPLE OF USING THIS DATASET WITH SCIKITLEARN SEE EXAM  
PLESAPPLICATIONSPLOTSPECIESDISTRIBUTIONMODELINGPY

REFERENCES

- "MAXIMUM ENTROPY MODELING OF SPECIES GEOGRAPHIC DISTRIBUTIONS" S J PHILLIPS R P ANDERSON R E SCHAPIRE ECOLOGICAL MODELLING 190231259 2006

EXAMPLES USING SKLEARNDATASETSFETCHSPECIESDISTRIBUTIONS

- SPECIES DISTRIBUTION MODELING
- KERNEL DENSITY ESTIMATE OF SPECIES DISTRIBUTIONS

1590 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

SKLEARNDATASETS GETDATAHOME

SKLEARNDATASETS GETDATAHOME DATAHOMENONE

RETURN THE PATH OF THE SCIKITLEARN DATA DIR

THIS FOLDER IS USED BY SOME LARGE DATASET LOADERS TO AVOID DOWNLOADING THE DATA SEVERAL TIMES

BY DEFAULT THE DATA DIR IS SET TO A FOLDER NAMED 'SCIKITLEARNDATA' IN THE USER HOME FOLDER

ALTERNATIVELY IT CAN BE SET BY THE 'SCIKITLEARNDATA' ENVIRONMENT VARIABLE OR PROGRAMMATICALLY BY GIVING AN EXPLICIT FOLDER PATH THE '' SYMBOL IS EXPANDED TO THE USER HOME FOLDER

IF THE FOLDER DOES NOT ALREADY EXIST IT IS AUTOMATICALLY CREATED

PARAMETERS

DATAHOME STR NONE THE PATH TO SCIKITLEARN DATA DIR

EXAMPLES USING SKLEARNDATASETSGETDATAHOME

- OUTOFCORE CLASSIFICATION OF TEXT DOCUMENTS

SKLEARNDATASETS LOADBOSTON

SKLEARNDATASETS LOADBOSTON RETURNXYFALSE

LOAD AND RETURN THE BOSTON HOUSEPRICES DATASET REGRESSION

SAMPLES TOTAL 506

DIMENSIONALITY 13

FEATURES REAL POSITIVE

TARGETS REAL 5 50

READ MORE IN THE USER GUIDE

PARAMETERS

RETURNXY BOOLEAN DEFAULTFALSE IF TRUE RETURNS DATA TARGET INSTEAD OF A BUNCH

OBJECT SEE BELOW FOR MORE INFORMATION ABOUT THE DATA ANDTARGET OBJECT

NEW IN VERSION 018

RETURNS

DATA BUNCH DICTIONARYLIKE OBJECT THE INTERESTING ATTRIBUTES ARE 'DATA' THE DATA TO LEARN

'TARGET' THE REGRESSION TARGETS 'DESCR' THE FULL DESCRIPTION OF THE DATASET AND 'FILENAME'

THE PHYSICAL LOCATION OF BOSTON CSV DATASET ADDED IN VERSION 020

DATA TARGET TUPLE IFRETURNXY IS TRUE NEW IN VERSION 018

NOTES

CHANGED IN VERSION 020 FIXED A WRONG DATA POINT AT 445 0

68SKLEARNDATASETS DATASETS 1591

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES

FROM SKLEARNDATASETS IMPORT LOADBOSTON

BOSTON LOADBOSTON

PRINTBOSTONDATA SHAPE

506 13

EXAMPLES USING SKLEARNDATASETSLOADBOSTON

- OUTLIER DETECTION ON A REAL DATA SET
- MODEL COMPLEXITY INFLUENCE
- EFFECT OF TRANSFORMING THE TARGETS IN REGRESSION MODEL
- PLOT INDIVIDUAL AND VOTING REGRESSION PREDICTIONS
- GRADIENT BOOSTING REGRESSION
- FEATURE SELECTION USING SELECTFROMMODEL AND LASSOCV
- IMPUTING MISSING VALUES BEFORE BUILDING AN ESTIMATOR
- PLOTTING CROSSVALIDATED PREDICTIONS

SKLEARNDATASETS LOADBREASTCANCER

SKLEARNDATASETS LOADBREASTCANCER RETURNXYFALSE

LOAD AND RETURN THE BREAST CANCER WISCONSIN DATASET CLASSIFICATION

THE BREAST CANCER DATASET IS A CLASSIC AND VERY EASY BINARY CLASSIFICATION DATASET

CLASSES 2

SAMPLES PER CLASS 212M357B

SAMPLES TOTAL 569

DIMENSIONALITY 30

FEATURES REAL POSITIVE

READ MORE IN THE USER GUIDE

PARAMETERS

RETURNXY BOOLEAN DEFAULTFALSE IF TRUE RETURNS DATA TARGET INSTEAD OF A BUNCH

OBJECT SEE BELOW FOR MORE INFORMATION ABOUT THE DATA ANDTARGET OBJECT

NEW IN VERSION 018

RETURNS

DATA BUNCH DICTIONARYLIKE OBJECT THE INTERESTING ATTRIBUTES ARE 'DATA' THE DATA TO LEARN

'TARGET' THE CLASSIFICATION LABELS 'TARGETNAMES' THE MEANING OF THE LABELS 'FEATURENAMES'

THE MEANING OF THE FEATURES AND 'DESCR' THE FULL DESCRIPTION OF THE DATASET 'FILENAME' THE

PHYSICAL LOCATION OF BREAST CANCER CSV DATASET ADDED IN VERSION 020

DATA TARGET TUPLE IFRETURNXY IS TRUE NEW IN VERSION 018

THE COPY OF UCI ML BREAST CANCER WISCONSIN DIAGNOSTIC DATASET IS

DOWNLOADED FROM

1592 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213  
HTTPSGOOGLELU2UWZ2  
EXAMPLES  
LET'S SAY YOU ARE INTERESTED IN THE SAMPLES 10 50 AND 85 AND WANT TO KNOW THEIR CLASS NAME  
FROM SKLEARNDATASETS IMPORT LOADBREASTCANCER  
DATA LOADBREASTCANCER  
DATATARGET10 50 85  
ARRAY0 1 0  
LISTDATATARGETNAMES  
MALIGNANT BENIGN  
SKLEARNDATASETS LOADDIABETES  
SKLEARNDATASETS LOADDIABETES RETURNXYFALSE  
LOAD AND RETURN THE DIABETES DATASET REGRESSION  
SAMPLES TOTAL 442  
DIMENSIONALITY 10  
FEATURES REAL 2 X 2  
TARGETS INTEGER 25 346  
READ MORE IN THE USER GUIDE  
PARAMETERS  
RETURNXY BOOLEAN DEFAULTFALSE IF TRUE RETURNS DATA TARGET INSTEAD OF A BUNCH  
OBJECT SEE BELOW FOR MORE INFORMATION ABOUT THE DATA ANDTARGET OBJECT  
NEW IN VERSION 018  
RETURNS  
DATA BUNCH DICTIONARYLIKE OBJECT THE INTERESTING ATTRIBUTES ARE 'DATA' THE DATA TO LEARN 'TARGET' THE REGRESSION TARGET FOR EACH SAMPLE 'DATAFILENAME' THE PHYSICAL LOCATION OF DIABETES  
DATA CSV DATASET AND 'TARGETFILENAME' THE PHYSICAL LOCATION OF DIABETES TARGETS CSV DATASET  
ADDED IN VERSION 020  
DATA TARGET TUPLE IFRETURNXY IS TRUE NEW IN VERSION 018  
EXAMPLES USING SKLEARNDATASETSLOADDIABETES  
•CROSSVALIDATION ON DIABETES DATASET EXERCISE  
•IMPUTING MISSING VALUES BEFORE BUILDING AN ESTIMATOR  
•LASSO PATH USING LARS  
•LINEAR REGRESSION EXAMPLE  
•SPARSITY EXAMPLE FITTING ONLY FEATURES 1 AND 2  
•LASSO AND ELASTIC NET  
•LASSO MODEL SELECTION CROSSVALIDATION AIC BIC  
68SKLEARNDATASETS DATASETS 1593

SCIKITLEARN USER GUIDE RELEASE 0213

SKLEARNDATASETS LOADDIGITS

SKLEARNDATASETS LOADDIGITS NCLASS10 RETURNXYFALSE

LOAD AND RETURN THE DIGITS DATASET CLASSIFICATION

EACH DATAPOINT IS A 8X8 IMAGE OF A DIGIT

CLASSES 10

SAMPLES PER CLASS 180

SAMPLES TOTAL 1797

DIMENSIONALITY 64

FEATURES INTEGERS 016

READ MORE IN THE USER GUIDE

PARAMETERS

NCLASS INTEGER BETWEEN 0 AND 10 OPTIONAL DEFAULT10 THE NUMBER OF CLASSES TO RETURN

RETURNXY BOOLEAN DEFAULTFALSE IF TRUE RETURNS DATA TARGET INSTEAD OF A BUNCH

OBJECT SEE BELOW FOR MORE INFORMATION ABOUT THE DATA ANDTARGET OBJECT

NEW IN VERSION 018

RETURNS

DATA BUNCH DICTIONARYLIKE OBJECT THE INTERESTING ATTRIBUTES ARE 'DATA' THE DATA TO LEARN

'IMAGES' THE IMAGES CORRESPONDING TO EACH SAMPLE 'TARGET' THE CLASSIFICATION LABELS FOR

EACH SAMPLE 'TARGETNAMES' THE MEANING OF THE LABELS AND 'DESCR' THE FULL DESCRIPTION

OF THE DATASET

DATA TARGET TUPLE IFRETURNXY IS TRUE NEW IN VERSION 018

THIS IS A COPY OF THE TEST SET OF THE UCI ML HANDWRITTEN DIGITS DATASETS

[HTTPSARCHIVEICSUCIEDUMLDATASETSOPTICALRECOGNITIONOFHANDWRITTENDIGITS](https://archive.ics.uci.edu/ml/dataset/optical_recognition_of_handwritten_digits)

EXAMPLES

TO LOAD THE DATA AND VISUALIZE THE IMAGES

FROM SKLEARNDATASETS IMPORT LOADDIGITS

DIGITS LOADDIGITS

PRINTDIGITSDATASHAPE

1797 64

IMPORT MATPLOTLIBPYPLOT AS PLT

PLTGRAY

PLTMATSHOWDIGITSIMAGES0

PLTSHOW

EXAMPLES USING SKLEARNDATASETSLOADDIGITS

- THE JOHNSONLINDENSTRAUSS BOUND FOR EMBEDDING WITH RANDOM PROJECTIONS
- EXPLICIT FEATURE MAP APPROXIMATION FOR RBF KERNELS
- RECOGNIZING HANDWRITTEN DIGITS

1594 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

- FEATURE AGGLOMERATION
- VARIOUS AGGLOMERATIVE CLUSTERING ON A 2D EMBEDDING OF DIGITS
- A DEMO OF KMEANS CLUSTERING ON THE HANDWRITTEN DIGITS DATA
- PIPELINING CHAINING A PCA AND A LOGISTIC REGRESSION
- SELECTING DIMENSIONALITY REDUCTION WITH PIPELINE AND GRIDSEARCHCV
- THE DIGIT DATASET
- EARLY STOPPING OF GRADIENT BOOSTING
- DIGITS CLASSIFICATION EXERCISE
- CROSSVALIDATION ON DIGITS DATASET EXERCISE
- RECURSIVE FEATURE ELIMINATION
- COMPARING VARIOUS ONLINE SOLVERS
- L1 PENALTY AND SPARSITY IN LOGISTIC REGRESSION
- MANIFOLD LEARNING ON HANDWRITTEN DIGITS LOCALLY LINEAR EMBEDDING ISOMAP
- PLOTING VALIDATION CURVES
- PARAMETER ESTIMATION USING GRID SEARCH WITH CROSSVALIDATION
- COMPARING RANDOMIZED SEARCH AND GRID SEARCH FOR HYPERPARAMETER ESTIMATION
- BALANCE MODEL COMPLEXITY AND CROSSVALIDATED SCORE
- PLOTING LEARNING CURVES
- KERNEL DENSITY ESTIMATION
- DIMENSIONALITY REDUCTION WITH NEIGHBORHOOD COMPONENTS ANALYSIS
- RESTRICTED BOLTZMANN MACHINE FEATURES FOR DIGIT CLASSIFICATION
- COMPARE STOCHASTIC LEARNING STRATEGIES FOR MLPCLASSIFIER
- LABEL PROPAGATION DIGITS DEMONSTRATING PERFORMANCE
- LABEL PROPAGATION DIGITS ACTIVE LEARNING

SKLEARNDATASETS LOADFILES

SKLEARNDATASETS LOADFILES CONTAINERPATH DESCRIPTIONNONE CATEGORIESNONE

LOADCONTENTTRUE SHUFFLETRUE ENCODINGNONE DE

CODEERROR'STRICT' RANDOMSTATE0

LOAD TEXT FILES WITH CATEGORIES AS SUBFOLDER NAMES

INDIVIDUAL SAMPLES ARE ASSUMED TO BE FILES STORED A TWO LEVELS FOLDER STRUCTURE SUCH AS THE FOLLOWING  
CONTAINERFOLDER

CATEGORY1FOLDER FILE1TXT FILE2TXT FILE42TXT

CATEGORY2FOLDER FILE43TXT FILE44TXT

THE FOLDER NAMES ARE USED AS SUPERVISED SIGNAL LABEL NAMES THE INDIVIDUAL FILE NAMES ARE NOT IMPORTANT

THIS FUNCTION DOES NOT TRY TO EXTRACT FEATURES INTO A NUMPY ARRAY OR SCIPY SPARSE MATRIX IN ADDITION IF

LOADCONTENT IS FALSE IT DOES NOT TRY TO LOAD THE FILES IN MEMORY

68SKLEARNDATASETS DATASETS 1595

SCIKITLEARN USER GUIDE RELEASE 0213  
TO USE TEXT FILES IN A SCIKITLEARN CLASSIFICATION OR CLUSTERING ALGORITHM YOU WILL NEED TO USE THE SKLEARN  
FEATUREEXTRACTIONTEXT MODULE TO BUILD A FEATURE EXTRACTION TRANSFORMER THAT SUITS YOUR PROBLEM  
IF YOU SET LOADCONTENTTRUE YOU SHOULD ALSO SPECIFY THE ENCODING OF THE TEXT USING THE 'ENCODING' PARAMETER  
FOR MANY MODERN TEXT FILES 'UTF8' WILL BE THE CORRECT ENCODING IF YOU LEAVE ENCODING EQUAL TO NONE THEN THE  
CONTENT WILL BE MADE OF BYTES INSTEAD OF UNICODE AND YOU WILL NOT BE ABLE TO USE MOST FUNCTIONS IN SKLEARN  
FEATUREEXTRACTIONTEXT  
SIMILAR FEATURE EXTRACTORS SHOULD BE BUILT FOR OTHER KIND OF UNSTRUCTURED DATA INPUT SUCH AS IMAGES AUDIO VIDEO

READ MORE IN THE USER GUIDE

PARAMETERS  
CONTAINERPATH STRING OR UNICODE PATH TO THE MAIN FOLDER HOLDING ONE SUBFOLDER PER CATEGORY  
DESCRIPTION STRING OR UNICODE OPTIONAL DEFAULTNONE A PARAGRAPH DESCRIBING THE CHARACTER  
ISTIC OF THE DATASET ITS SOURCE REFERENCE ETC  
CATEGORIES A COLLECTION OF STRINGS OR NONE OPTIONAL DEFAULTNONE IF NONE DEFAULT LOAD ALL  
THE CATEGORIES IF NOT NONE LIST OF CATEGORY NAMES TO LOAD OTHER CATEGORIES IGNORED  
LOADCONTENT BOOLEAN OPTIONAL DEFAULTTRUE WHETHER TO LOAD OR NOT THE CONTENT OF THE DIF  
FERENT FILES IF TRUE A 'DATA' ATTRIBUTE CONTAINING THE TEXT INFORMATION IS PRESENT IN THE DATA  
STRUCTURE RETURNED IF NOT A FILENAMES ATTRIBUTE GIVES THE PATH TO THE FILES  
SHUFFLE BOOL OPTIONAL DEFAULTTRUE WHETHER OR NOT TO SHUFFLE THE DATA MIGHT BE IMPORTANT  
FOR MODELS THAT MAKE THE ASSUMPTION THAT THE SAMPLES ARE INDEPENDENT AND IDENTICALLY DIS  
TRIBUTED IID SUCH AS STOCHASTIC GRADIENT DESCENT  
ENCODING STRING OR NONE DEFAULT IS NONE IF NONE DO NOT TRY TO DECODE THE CONTENT OF THE FILES  
EG FOR IMAGES OR OTHER NONTEXT CONTENT IF NOT NONE ENCODING TO USE TO DECODE TEXT FILES  
TO UNICODE IF LOADCONTENT IS TRUE  
DECODEERROR 'STRICT' 'IGNORE' 'REPLACE' OPTIONAL INSTRUCTION ON WHAT TO DO IF A BYTE SE  
QUENCE IS GIVEN TO ANALYZE THAT CONTAINS CHARACTERS NOT OF THE GIVEN ENCODING PASSED AS  
KEYWORD ARGUMENT 'ERRORS' TO BYTESDECODE  
RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE DEFAULT0 DETERMINES RANDOM NUMBER  
GENERATION FOR DATASET SHUFFLING PASS AN INT FOR REPRODUCIBLE OUTPUT ACROSS MULTIPLE FUNCTION  
CALLS SEE GLOSSARY  
RETURNS  
DATA BUNCH DICTIONARYLIKE OBJECT THE INTERESTING ATTRIBUTES ARE EITHER DATA THE RAW TEXT DATA  
TO LEARN OR 'FILENAMES' THE FILES HOLDING IT 'TARGET' THE CLASSIFICATION LABELS INTEGER INDEX  
'TARGETNAMES' THE MEANING OF THE LABELS AND 'DESCR' THE FULL DESCRIPTION OF THE DATASET  
SKLEARNDATASETS LOADIRIS  
SKLEARNDATASETS LOADIRIS RETURNXYFALSE  
LOAD AND RETURN THE IRIS DATASET CLASSIFICATION  
THE IRIS DATASET IS A CLASSIC AND VERY EASY MULTICLASS CLASSIFICATION DATASET  
1596 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

CLASSES 3

SAMPLES PER CLASS 50

SAMPLES TOTAL 150

DIMENSIONALITY 4

FEATURES REAL POSITIVE

READ MORE IN THE USER GUIDE

PARAMETERS

RETURNXY BOOLEAN DEFAULTFALSE IF TRUE RETURNS DATA TARGET INSTEAD OF A BUNCH

OBJECT SEE BELOW FOR MORE INFORMATION ABOUT THE DATA ANDTARGET OBJECT

NEW IN VERSION 018

RETURNS

DATA BUNCH DICTIONARYLIKE OBJECT THE INTERESTING ATTRIBUTES ARE 'DATA' THE DATA TO LEARN

'TARGET' THE CLASSIFICATION LABELS 'TARGETNAMES' THE MEANING OF THE LABELS 'FEATURENAMES'

THE MEANING OF THE FEATURES 'DESCR' THE FULL DESCRIPTION OF THE DATASET 'FILENAME' THE

PHYSICAL LOCATION OF IRIS CSV DATASET ADDED IN VERSION 020

DATA TARGET TUPLE IFRETURNXY IS TRUE NEW IN VERSION 018

NOTES

CHANGED IN VERSION 020 FIXED TWO WRONG DATA POINTS ACCORDING TO FISHER'S PAPER THE NEW VERSION IS THE SAME

AS IN R BUT NOT AS IN THE UCI MACHINE LEARNING REPOSITORY

EXAMPLES

LET'S SAY YOU ARE INTERESTED IN THE SAMPLES 10 25 AND 50 AND WANT TO KNOW THEIR CLASS NAME

FROM SKLEARNDATASETS IMPORT LOADIRIS

DATA LOADIRIS

DATATARGET10 25 50

ARRAY0 0 1

LISTDATATARGETNAMES

SETOSA VERSICOLOR VIRGINICA

EXAMPLES USING SKLEARNDATASETSLOADIRIS

- PLOT CLASSIFICATION PROBABILITY
  - KMEANS CLUSTERING
  - CONCATENATING MULTIPLE FEATURE EXTRACTION METHODS
  - THE IRIS DATASET
  - PCA EXAMPLE WITH IRIS DATASET
  - INCREMENTAL PCA
  - COMPARISON OF LDA AND PCA 2D PROJECTION OF IRIS DATASET
  - PLOT THE DECISION BOUNDARIES OF A VOTINGCLASSIFIER
- 68SKLEARNDATASETS DATASETS 1597

SCIKITLEARN USER GUIDE RELEASE 0213

- EARLY STOPPING OF GRADIENT BOOSTING
  - PLOT THE DECISION SURFACES OF ENSEMBLES OF TREES ON THE IRIS DATASET
  - SVM EXERCISE
  - TEST WITH PERMUTATIONS THE SIGNIFICANCE OF A CLASSIFICATION SCORE
  - UNIVARIATE FEATURE SELECTION
  - GAUSSIAN PROCESS CLASSIFICATION GPC ON IRIS DATASET
  - REGULARIZATION PATH OF L1 LOGISTIC REGRESSION
  - LOGISTIC REGRESSION 3CLASS CLASSIFIER
  - PLOT MULTICLASS SGD ON THE IRIS DATASET
  - GMM COVARIANCES
  - RECEIVER OPERATING CHARACTERISTIC ROC WITH CROSS VALIDATION
  - NESTED VERSUS NONNESTED CROSSVALIDATION
  - CONFUSION MATRIX
  - RECEIVER OPERATING CHARACTERISTIC ROC
  - PRECISIONRECALL
  - NEAREST NEIGHBORS CLASSIFICATION
  - NEAREST CENTROID CLASSIFICATION
  - COMPARING NEAREST NEIGHBORS WITH AND WITHOUT NEIGHBORHOOD COMPONENTS ANALYSIS
  - COMPARE STOCHASTIC LEARNING STRATEGIES FOR MLPCLASSIFIER
  - DECISION BOUNDARY OF LABEL PROPAGATION VERSUS SVM ON THE IRIS DATASET
  - SVM WITH CUSTOM KERNEL
  - SVMANOVA SVM WITH UNIVARIATE FEATURE SELECTION
  - PLOT DIFFERENT SVM CLASSIFIERS IN THE IRIS DATASET
  - RBF SVM PARAMETERS
  - PLOT THE DECISION SURFACE OF A DECISION TREE ON THE IRIS DATASET
  - UNDERSTANDING THE DECISION TREE STRUCTURE
- SKLEARNDATASETS LOADLINNERUD  
SKLEARNDATASETS LOADLINNERUD RETURNXYFALSE  
LOAD AND RETURN THE LINNERUD DATASET MULTIVARIATE REGRESSION  
SAMPLES TOTAL 20  
DIMENSIONALITY 3 FOR BOTH DATA AND TARGET  
FEATURES INTEGER  
TARGETS INTEGER  
READ MORE IN THE USER GUIDE  
PARAMETERS  
1598 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

RETURNXY BOOLEAN DEFAULTFALSE IF TRUE RETURNS DATA TARGET INSTEAD OF A BUNCH  
OBJECT SEE BELOW FOR MORE INFORMATION ABOUT THE DATA ANDTARGET OBJECT

NEW IN VERSION 018

RETURNS

DATA BUNCH DICTIONARYLIKE OBJECT THE INTERESTING ATTRIBUTES ARE 'DATA' AND 'TARGET' THE TWO  
MULTIVARIATE DATASETS WITH 'DATA' CORRESPONDING TO THE EXERCISE AND 'TARGET' CORRESPONDING  
TO THE PHYSIOLOGICAL MEASUREMENTS AS WELL AS 'FEATURENAMES' AND 'TARGETNAMES' IN ADDI  
TION YOU WILL ALSO HAVE ACCESS TO 'DATAFILENAME' THE PHYSICAL LOCATION OF LINNERUD DATA CSV  
DATASET AND 'TARGETFILENAME' THE PHYSICAL LOCATION OF LINNERUD TARGETS CSV DATATASET ADDED  
IN VERSION020

DATA TARGET TUPLE IFRETURNXY IS TRUE NEW IN VERSION 018

SKLEARNDATASETS LOADSAMPLEIMAGE

SKLEARNDATASETS LOADSAMPLEIMAGE IMAGENAME

LOAD THE NUMPY ARRAY OF A SINGLE SAMPLE IMAGE

READ MORE IN THE USER GUIDE

PARAMETERS

IMAGENAME CHINA.JPG FLOWER.JPG THE NAME OF THE SAMPLE IMAGE LOADED

RETURNS

IMG 3D ARRAY THE IMAGE AS A NUMPY ARRAY HEIGHT X WIDTH X COLOR

EXAMPLES

FROM SKLEARNDATASETS IMPORT LOADSAMPLEIMAGE

CHINA LOADSAMPLEIMAGECHINA.JPG

CHINADTYPE

DTYPEUINT8

CHINASHAPE

427 640 3

FLOWER LOADSAMPLEIMAGEFLOWER.JPG

FLOWERDTYPE

DTYPEUINT8

FLOWERSHAPE

427 640 3

EXAMPLES USING SKLEARNDATASETSLOADSAMPLEIMAGE

•COLOR QUANTIZATION USING KMEANS

SKLEARNDATASETS LOADSAMPLEIMAGES

SKLEARNDATASETS LOADSAMPLEIMAGES

LOAD SAMPLE IMAGES FOR IMAGE MANIPULATION

LOADS BOTH CHINA ANDFLOWER

68SKLEARNDATASETS DATASETS 1599

SCIKITLEARN USER GUIDE RELEASE 0213

READ MORE IN THE USER GUIDE

RETURNS

DATA BUNCH DICTIONARYLIKE OBJECT WITH THE FOLLOWING ATTRIBUTES ‘IMAGES’ THE TWO SAMPLE IMAGES ‘FILENAMES’ THE FILE NAMES FOR THE IMAGES AND ‘DESCR’ THE FULL DESCRIPTION OF THE DATASET

EXAMPLES

TO LOAD THE DATA AND VISUALIZE THE IMAGES

FROM SKLEARNDATASETS IMPORT LOADSAMPLEIMAGES

DATASET LOADSAMPLEIMAGES

LENDATASETIMAGES

2

FIRSTIMGDATA DATASETIMAGES0

FIRSTIMGDATASHAPE

427 640 3

FIRSTIMGDATADTYPE

DTYPEUINT8

SKLEARNDATASETS LOADSVMLIGHTFILE

SKLEARNDATASETS LOADSVMLIGHTFILE FNFEATURESNONE DTYPESCLASS ‘NUMPYFLOAT64’

MULTILABELFALSE ZEROBASED‘AUTO’ QUERYIDFALSE

OFFSET0 LENGTH1

LOAD DATASETS IN THE SVMLIGHT LIBSVM FORMAT INTO SPARSE CSR MATRIX

THIS FORMAT IS A TEXTBASED FORMAT WITH ONE SAMPLE PER LINE IT DOES NOT STORE ZERO VALUED FEATURES HENCE IS SUITABLE FOR SPARSE DATASET

THE FIRST ELEMENT OF EACH LINE CAN BE USED TO STORE A TARGET VARIABLE TO PREDICT

THIS FORMAT IS USED AS THE DEFAULT FORMAT FOR BOTH SVMLIGHT AND THE LIBSVM COMMAND LINE PROGRAMS

PARSING A TEXT BASED SOURCE CAN BE EXPENSIVE WHEN WORKING ON REPEATEDLY ON THE SAME DATASET IT IS RECOMMENDED TO WRAP THIS LOADER WITH JOBLIBMEMORYCACHE TO STORE A MEMMAPPED BACKUP OF THE CSR RESULTS OF THE FIRST CALL AND BENEFIT FROM THE NEAR INSTANTANEOUS LOADING OF MEMMAPPED STRUCTURES FOR THE SUBSEQUENT CALLS IN CASE THE FILE CONTAINS A PAIRWISE PREFERENCE CONSTRAINT KNOWN AS “QID” IN THE SVMLIGHT FORMAT THESE ARE IGNORED UNLESS THE QUERYID PARAMETER IS SET TO TRUE THESE PAIRWISE PREFERENCE CONSTRAINTS CAN BE USED TO CONSTRAINT THE COMBINATION OF SAMPLES WHEN USING PAIRWISE LOSS FUNCTIONS AS IS THE CASE IN SOME LEARNING TO RANK PROBLEMS SO THAT ONLY PAIRS WITH THE SAME QUERYID VALUE ARE CONSIDERED THIS IMPLEMENTATION IS WRITTEN IN CYTHON AND IS REASONABLY FAST HOWEVER A FASTER APICOMPATIBLE LOADER IS ALSO AVAILABLE AT

HTTPSGITHUBCOMMBLONDELSVMLIGHTLOADER

PARAMETERS

FSTR FILELIKE INT PATH TO A FILE TO LOAD IF A PATH ENDS IN “GZ” OR “BZ2” IT WILL BE UNCOMPRESSED

ON THE FLY IF AN INTEGER IS PASSED IT IS ASSUMED TO BE A FILE DESCRIPTOR A FILELIKE

OR FILE DESCRIPTOR WILL NOT BE CLOSED BY THIS FUNCTION A FILELIKE OBJECT MUST BE OPENED IN

BINARY MODE

1600 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

NFEATURES INT OR NONE THE NUMBER OF FEATURES TO USE IF NONE IT WILL BE INFERRED THIS ARGUMENT IS USEFUL TO LOAD SEVERAL FILES THAT ARE SUBSETS OF A BIGGER SLICED DATASET EACH SUBSET MIGHT NOT HAVE EXAMPLES OF EVERY FEATURE HENCE THE INFERRED SHAPE MIGHT VARY FROM ONE SLICE TO ANOTHER NFEATURES IS ONLY REQUIRED IF OFFSET ORLENGTH ARE PASSED A NON DEFAULT VALUE

DTYPE NUMPY DATA TYPE DEFAULT NPFLOAT64 DATA TYPE OF DATASET TO BE LOADED THIS WILL BE THE DATA TYPE OF THE OUTPUT NUMPY ARRAYS XANDY

MULTILABEL BOOLEAN OPTIONAL DEFAULT FALSE SAMPLES MAY HAVE SEVERAL LABELS EACH SEE [HTTPS WWWCSIENTUEDUTWCJLINLIBSVMTOOLSDATASETSMULTILABELHTML](https://www.cs.ientu.edu.tw/cjlin/libsvm/tools/datasets/multilabel.html)

ZEROBASED BOOLEAN OR "AUTO" OPTIONAL DEFAULT "AUTO" WHETHER COLUMN INDICES IN F ARE ZERO BASED TRUE OR ONEBASED FALSE IF COLUMN INDICES ARE ONEBASED THEY ARE TRANSFORMED TO ZEROBASED TO MATCH PYTHONNUMPY CONVENTIONS IF SET TO "AUTO" A HEURISTIC CHECK IS APPLIED TO DETERMINE THIS FROM THE FILE CONTENTS BOTH KINDS OF FILES OCCUR "IN THE WILD" BUT THEY ARE UNFORTUNATELY NOT SELFIDENTIFYING USING "AUTO" OR TRUE SHOULD ALWAYS BE SAFE WHEN NOOFFSET ORLENGTH IS PASSED IF OFFSET ORLENGTH ARE PASSED THE "AUTO" MODE FALLS BACK TOZEROBASEDTRUE TO AVOID HAVING THE HEURISTIC CHECK YIELD INCONSISTENT RESULTS ON DIFFERENT SEGMENTS OF THE FILE

QUERYID BOOLEAN DEFAULT FALSE IF TRUE WILL RETURN THE QUERYID ARRAY FOR EACH FILE

OFFSET INTEGER OPTIONAL DEFAULT 0 IGNORE THE OFFSET FIRST BYTES BY SEEKING FORWARD THEN DISCARDING THE FOLLOWING BYTES UP UNTIL THE NEXT NEW LINE CHARACTER

LENGTH INTEGER OPTIONAL DEFAULT 1 IF STRICTLY POSITIVE STOP READING ANY NEW LINE OF DATA ONCE THE POSITION IN THE FILE HAS REACHED THE OFFSET LENGTH BYTES THRESHOLD

RETURNS

XSCIPYSPARSE MATRIX OF SHAPE NSAMPLES NFEATURES

YNDARRAY OF SHAPE NSAMPLES OR IN THE MULTILABEL A LIST OF TUPLES OF LENGTH NSAMPLES

QUERYID ARRAY OF SHAPE NSAMPLES QUERYID FOR EACH SAMPLE ONLY RETURNED WHEN QUERYID IS SET TO TRUE

SEE ALSO

LOADSVMLIGHTFILES SIMILAR FUNCTION FOR LOADING MULTIPLE FILES IN THIS FORMAT ENFORCING THE SAME NUMBER OF FEATURES COLUMNS ON ALL OF THEM

EXAMPLES

TO USE JOBLIBMEMORY TO CACHE THE SVMLIGHT FILE

```
FROM JOBLIB IMPORT MEMORY
FROM DATASETS IMPORT LOADSVMLIGHTFILE
MEM MEMORYMYCACHE
MEMCACHE
DEFGETDATA
DATA LOADSVMLIGHTFILEMYSVMLIGHTFILE
RETURNDATA0 DATA1
X Y GETDATA
68SKLEARN DATASETS DATASETS 1601
```

SCIKITLEARN USER GUIDE RELEASE 0213

SKLEARNDATASETS LOADSVMLIGHTFILES

SKLEARNDATASETS LOADSVMLIGHTFILES FILES NFEATURESNONE DTYPECLASS

‘NUMPYFLOAT64’ MULTILABELFALSE

ZEROBASED‘AUTO’ QUERYIDFALSE OFFSET0 LENGTH

1

LOAD DATASET FROM MULTIPLE FILES IN SVMLIGHT FORMAT

THIS FUNCTION IS EQUIVALENT TO MAPPING LOADSVMLIGHTFILE OVER A LIST OF FILES EXCEPT THAT THE RESULTS ARE CONCATENATED INTO A SINGLE FLAT LIST AND THE SAMPLES VECTORS ARE CONSTRAINED TO ALL HAVE THE SAME NUMBER OF FEATURES IN CASE THE FILE CONTAINS A PAIRWISE PREFERENCE CONSTRAINT KNOWN AS “QID” IN THE SVMLIGHT FORMAT THESE ARE IGNORED UNLESS THE QUERYID PARAMETER IS SET TO TRUE THESE PAIRWISE PREFERENCE CONSTRAINTS CAN BE USED TO CONSTRAINT THE COMBINATION OF SAMPLES WHEN USING PAIRWISE LOSS FUNCTIONS AS IS THE CASE IN SOME LEARNING TO RANK PROBLEMS SO THAT ONLY PAIRS WITH THE SAME QUERYID VALUE ARE CONSIDERED

PARAMETERS

FILES ITERABLE OVER STR FILELIKE INT PATHS OF FILES TO LOAD IF A PATH ENDS IN “GZ” OR “BZ2” IT WILL BE UNCOMPRESSED ON THE FLY IF AN INTEGER IS PASSED IT IS ASSUMED TO BE A FILE DESCRIPTOR FILELIKES AND FILE DESCRIPTORS WILL NOT BE CLOSED BY THIS FUNCTION FILELIKE OBJECTS MUST BE OPENED IN BINARY MODE

NFEATURES INT OR NONE THE NUMBER OF FEATURES TO USE IF NONE IT WILL BE INFERRED FROM THE MAXIMUM COLUMN INDEX OCCURRING IN ANY OF THE FILES

THIS CAN BE SET TO A HIGHER VALUE THAN THE ACTUAL NUMBER OF FEATURES IN ANY OF THE INPUT FILES BUT SETTING IT TO A LOWER VALUE WILL CAUSE AN EXCEPTION TO BE RAISED

DTYPE NUMPY DATA TYPE DEFAULT NPFLOAT64 DATA TYPE OF DATASET TO BE LOADED THIS WILL BE THE DATA TYPE OF THE OUTPUT NUMPY ARRAYS XANDY

MULTILABEL BOOLEAN OPTIONAL SAMPLES MAY HAVE SEVERAL LABELS EACH SEE HTTPSWWWCSIENTU

EDUTWCJLINLIBSVMTOOLS DATASETSMULTILABELHTML

ZEROBASED BOOLEAN OR “AUTO” OPTIONAL WHETHER COLUMN INDICES IN F ARE ZEROBASED TRUE OR ONEBASED FALSE IF COLUMN INDICES ARE ONEBASED THEY ARE TRANSFORMED TO ZEROBASED TO MATCH PYTHONNUMPY CONVENTIONS IF SET TO “AUTO” A HEURISTIC CHECK IS APPLIED TO DETERMINE THIS FROM THE FILE CONTENTS BOTH KINDS OF FILES OCCUR “IN THE WILD” BUT THEY ARE UNFORTUNATELY NOT SELFIDENTIFYING USING “AUTO” OR TRUE SHOULD ALWAYS BE SAFE WHEN NO OFFSET OR LENGTH IS PASSED IF OFFSET OR LENGTH ARE PASSED THE “AUTO” MODE FALLS BACK TO ZEROBASEDTRUE TO AVOID HAVING THE HEURISTIC CHECK YIELD INCONSISTENT RESULTS ON DIFFERENT SEGMENTS OF THE FILE

QUERYID BOOLEAN DEFAULTS TO FALSE IF TRUE WILL RETURN THE QUERYID ARRAY FOR EACH FILE

OFFSET INTEGER OPTIONAL DEFAULT 0 IGNORE THE OFFSET FIRST BYTES BY SEEKING FORWARD THEN DISCARDING THE FOLLOWING BYTES UP UNTIL THE NEXT NEW LINE CHARACTER

LENGTH INTEGER OPTIONAL DEFAULT 1 IF STRICTLY POSITIVE STOP READING ANY NEW LINE OF DATA ONCE THE POSITION IN THE FILE HAS REACHED THE OFFSET LENGTH BYTES THRESHOLD

RETURNS

X1 Y1 XN YN

WHERE EACH XI YI PAIR IS THE RESULT FROM LOADSVMLIGHTFILEFILES

IF QUERYID IS SET TO TRUE THIS WILL RETURN INSTEAD X1 Y1 Q1

XN YN QN WHERE XI YI QI IS THE RESULT FROM

LOADSVMLIGHTFILEFILES

SEE ALSO

1602 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

LOADSVMLIGHTFILE

NOTES

WHEN FITTING A MODEL TO A MATRIX XTRAIN AND EVALUATING IT AGAINST A MATRIX XTEST IT IS ESSENTIAL THAT XTRAIN AND XTEST HAVE THE SAME NUMBER OF FEATURES XTRAINSHAPE1 XTESTSHAPE1 THIS MAY NOT BE THE CASE IF YOU LOAD THE FILES INDIVIDUALLY WITH LOADSVMLIGHTFILE

SKLEARNDATASETS LOADWINE

SKLEARNDATASETS LOADWINE RETURNXYFALSE

LOAD AND RETURN THE WINE DATASET CLASSIFICATION

NEW IN VERSION 018

THE WINE DATASET IS A CLASSIC AND VERY EASY MULTICLASS CLASSIFICATION DATASET

CLASSES 3

SAMPLES PER CLASS 597148

SAMPLES TOTAL 178

DIMENSIONALITY 13

FEATURES REAL POSITIVE

READ MORE IN THE USER GUIDE

PARAMETERS

RETURNXY BOOLEAN DEFAULTFALSE IF TRUE RETURNS DATA TARGET INSTEAD OF A BUNCH

OBJECT SEE BELOW FOR MORE INFORMATION ABOUT THE DATA ANDTARGET OBJECT

RETURNS

DATA BUNCH DICTIONARYLIKE OBJECT THE INTERESTING ATTRIBUTES ARE 'DATA' THE DATA TO LEARN

'TARGET' THE CLASSIFICATION LABELS 'TARGETNAMES' THE MEANING OF THE LABELS 'FEATURENAMES'

THE MEANING OF THE FEATURES AND 'DESCR' THE FULL DESCRIPTION OF THE DATASET

DATA TARGET TUPLE IFRETURNXY IS TRUE

THE COPY OF UCI ML WINE DATA SET DATASET IS DOWNLOADED AND MODIFIED TO FIT

STANDARD FORMAT FROM

[HTTPSARCHIVEICSUCIEDUMLMACHINELEARNINGDATABASESWINEWINEDATA](https://archive.ics.uci.edu/ml/machine-learning-databases/wine/wine.data)

EXAMPLES

LET'S SAY YOU ARE INTERESTED IN THE SAMPLES 10 80 AND 140 AND WANT TO KNOW THEIR CLASS NAME

FROM SKLEARNDATASETS IMPORT LOADWINE

DATA LOADWINE

DATATARGET10 80 140

ARRAY0 1 2

LISTDATATARGETNAMES

CLASS0 CLASS1 CLASS2

68SKLEARNDATASETS DATASETS 1603

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES USING SKLEARN DATASETS LOAD WINE

- IMPORTANCE OF FEATURE SCALING

682 SAMPLES GENERATOR

DATASETS MAKE BICLUSTERS SHAPE NCLUSTERS GENERATE AN ARRAY WITH CONSTANT BLOCK DIAGONAL STRUCTURE FOR BICLUSTERING

DATASETS MAKE BLOBS NSAMPLES NFEATURES GENERATE ISOTROPIC GAUSSIAN BLOBS FOR CLUSTERING

DATASETS MAKE CHECKERBOARD SHAPE NCLUSTERS GENERATE AN ARRAY WITH BLOCK CHECKERBOARD STRUCTURE FOR BICLUSTERING

DATASETS MAKE CIRCLES NSAMPLES SHUFFLE MAKE A LARGE CIRCLE CONTAINING A SMALLER CIRCLE IN 2D

DATASETS MAKE CLASSIFICATION NSAMPLES

GENERATE A RANDOM NCLASS CLASSIFICATION PROBLEM

DATASETS MAKE FRIEDMAN1 NSAMPLES GENERATE THE “FRIEDMAN 1” REGRESSION PROBLEM

DATASETS MAKE FRIEDMAN2 NSAMPLES NOISE

GENERATE THE “FRIEDMAN 2” REGRESSION PROBLEM

DATASETS MAKE FRIEDMAN3 NSAMPLES NOISE

GENERATE THE “FRIEDMAN 3” REGRESSION PROBLEM

DATASETS MAKE GAUSSIAN QUANTILES MEAN

GENERATE ISOTROPIC GAUSSIAN AND LABEL SAMPLES BY QUANTILE

DATASETS MAKE HASTIE102 NSAMPLES GENERATES DATA FOR BINARY CLASSIFICATION USED IN HASTIE ET AL

DATASETS MAKE LOW RANK MATRIX NSAMPLES

GENERATE A MOSTLY LOW RANK MATRIX WITH BELL SHAPED SINGULAR VALUES

DATASETS MAKE MOONS NSAMPLES SHUFFLE MAKE TWO INTERLEAVING HALF CIRCLES

DATASETS MAKE MULTILABEL CLASSIFICATION GENERATE A RANDOM MULTILABEL CLASSIFICATION PROBLEM

DATASETS MAKE REGRESSION NSAMPLES GENERATE A RANDOM REGRESSION PROBLEM

DATASETS MAKE S CURVE NSAMPLES NOISE GENERATE AN S CURVE DATASET

DATASETS MAKE SPARSE CODED SIGNAL NSAMPLES

GENERATE A SIGNAL AS A SPARSE COMBINATION OF DICTIONARY ELEMENTS

DATASETS MAKE SPARSE PD MATRIX DIM GENERATE A SPARSE SYMMETRIC DEFINITE POSITIVE MATRIX

DATASETS MAKE SPARSE UNCORRELATED RELATED DESIGN GENERATE A RANDOM REGRESSION PROBLEM WITH SPARSE UNCORRELATED DESIGN

DATASETS MAKE SPD MATRIX NDIM RAN

DOM STATE GENERATE A RANDOM SYMMETRIC POSITIVE DEFINITE MATRIX

DATASETS MAKE SWISS ROLL NSAMPLES NOISE

GENERATE A SWISS ROLL DATASET

SKLEARN DATASETS MAKE BICLUSTERS

SKLEARN DATASETS MAKE BICLUSTERS SHAPE NCLUSTERS NOISE00 MINVAL10 MAXVAL100 SHUFFLE TRUE RANDOM STATE NONE

GENERATE AN ARRAY WITH CONSTANT BLOCK DIAGONAL STRUCTURE FOR BICLUSTERING

READ MORE IN THE USER GUIDE

PARAMETERS

SHAPE ITERABLE NROWS NCOLS THE SHAPE OF THE RESULT

NCLUSTERS INTEGER THE NUMBER OF BICLUSTERS

NOISE FLOAT OPTIONAL DEFAULT 0.0 THE STANDARD DEVIATION OF THE GAUSSIAN NOISE

1604 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

MINVAL INT OPTIONAL DEFAULT10 MINIMUM VALUE OF A BICLUSTER

MAXVAL INT OPTIONAL DEFAULT100 MAXIMUM VALUE OF A BICLUSTER

SHUFFLE BOOLEAN OPTIONAL DEFAULTTRUE SHUFFLE THE SAMPLES

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE DEFAULT DETERMINES RANDOM NUMBER GENERATION FOR DATASET CREATION PASS AN INT FOR REPRODUCIBLE OUTPUT ACROSS MULTIPLE FUNCTION CALLS SEE GLOSSARY

RETURNS

XARRAY OF SHAPE SHAPE THE GENERATED ARRAY

ROWS ARRAY OF SHAPE NCLUSTERS XSHAPE0 THE INDICATORS FOR CLUSTER MEMBERSHIP OF EACH ROW

COLS ARRAY OF SHAPE NCLUSTERS XSHAPE1 THE INDICATORS FOR CLUSTER MEMBERSHIP OF EACH COLUMN

SEE ALSO

MAKECHECKERBOARD

REFERENCES

1

EXAMPLES USING SKLEARNDATASETSMAKEBICLUSTERS

- A DEMO OF THE SPECTRAL COCLUSTERING ALGORITHM

SKLEARNDATASETS MAKEBLOBS

SKLEARNDATASETS MAKEBLOBS NSAMPLES100 NFEATURES2 CENTERSNONE CLUSTERSTD10

CENTERBOX100 100 SHUFFLETRUE RANDOMSTATENONE

GENERATE ISOTROPIC GAUSSIAN BLOBS FOR CLUSTERING

READ MORE IN THE USER GUIDE

PARAMETERS

NSAMPLES INT OR ARRAYLIKE OPTIONAL DEFAULT100 IF INT IT IS THE TOTAL NUMBER OF POINTS EQUALLY DIVIDED AMONG CLUSTERS IF ARRAYLIKE EACH ELEMENT OF THE SEQUENCE INDICATES THE NUMBER OF SAMPLES PER CLUSTER

NFEATURES INT OPTIONAL DEFAULT2 THE NUMBER OF FEATURES FOR EACH SAMPLE

CENTERS INT OR ARRAY OF SHAPE NCENTERS NFEATURES OPTIONAL DEFAULTNONE THE NUMBER OF CENTERS TO GENERATE OR THE FIXED CENTER LOCATIONS IF NSAMPLES IS AN INT AND CENTERS IS NONE 3 CENTERS ARE GENERATED IF NSAMPLES IS ARRAYLIKE CENTERS MUST BE EITHER NONE OR AN ARRAY OF LENGTH EQUAL TO THE LENGTH OF NSAMPLES

CLUSTERSTD FLOAT OR SEQUENCE OF FLOATS OPTIONAL DEFAULT10 THE STANDARD DEVIATION OF THE CLUSTERS

CENTERBOX PAIR OF FLOATS MIN MAX OPTIONAL DEFAULT100 100 THE BOUNDING BOX FOR EACH CLUSTER CENTER WHEN CENTERS ARE GENERATED AT RANDOM

68SKLEARNDATASETS DATASETS 1605

SCIKITLEARN USER GUIDE RELEASE 0213

SHUFFLE BOOLEAN OPTIONAL DEFAULTTRUE SHUFFLE THE SAMPLES

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE DEFAULT DETERMINES RANDOM NUMBER GEN

ERATION FOR DATASET CREATION PASS AN INT FOR REPRODUCIBLE OUTPUT ACROSS MULTIPLE FUNCTION

CALLS SEE GLOSSARY

RETURNS

XARRAY OF SHAPE NSAMPLES NFEATURES THE GENERATED SAMPLES

YARRAY OF SHAPE NSAMPLES THE INTEGER LABELS FOR CLUSTER MEMBERSHIP OF EACH SAMPLE

SEE ALSO

MAKECLASSIFICATION A MORE INTRICATE VARIANT

EXAMPLES

FROM SKLEARNDATASETSAMPLESGENERATOR IMPORT MAKEBLOBS

X Y MAKEBLOBSNSAMPLES10 CENTERS3 NFEATURES2

RANDOMSTATE0

PRINTXSHAPE

10 2

Y

ARRAY0 0 1 0 2 2 2 1 1 0

X Y MAKEBLOBSNSAMPLES3 3 4 CENTERS NONE NFEATURES2

RANDOMSTATE0

PRINTXSHAPE

10 2

Y

ARRAY0 1 2 0 2 2 2 1 1 0

EXAMPLES USING SKLEARNDATASETSMAKEBLOBS

- COMPARING ANOMALY DETECTION ALGORITHMS FOR OUTLIER DETECTION ON TOY DATASETS
- PROBABILITY CALIBRATION OF CLASSIFIERS
- PROBABILITY CALIBRATION FOR 3CLASS CLASSIFICATION
- NORMAL AND SHRINKAGE LINEAR DISCRIMINANT ANALYSIS FOR CLASSIFICATION
- A DEMO OF THE MEANSHIFT CLUSTERING ALGORITHM
- DEMONSTRATION OF KMEANS ASSUMPTIONS
- DEMO OF AFFINITY PROPAGATION CLUSTERING ALGORITHM
- DEMO OF DBSCAN CLUSTERING ALGORITHM
- INDUCTIVE CLUSTERING
- COMPARE BIRCH AND MINIBATCHKMEANS
- COMPARISON OF THE KMEANS AND MINIBATCHKMEANS CLUSTERING ALGORITHMS
- COMPARING DIFFERENT HIERARCHICAL LINKAGE METHODS ON TOY DATASETS
- SELECTING THE NUMBER OF CLUSTERS WITH SILHOUETTE ANALYSIS ON KMEANS CLUSTERING
- COMPARING DIFFERENT CLUSTERING ALGORITHMS ON TOY DATASETS

1606 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

- PLOT RANDOMLY GENERATED CLASSIFICATION DATASET
- SGD MAXIMUM MARGIN SEPARATING HYPERPLANE
- PLOT MULTINOMIAL AND ONEVSREST LOGISTIC REGRESSION
- DEMONSTRATING THE DIFFERENT STRATEGIES OF KBINSDISCRETIZER
- SVM MAXIMUM MARGIN SEPARATING HYPERPLANE
- SVM SEPARATING HYPERPLANE FOR UNBALANCED CLASSES

SKLEARNDATASETS MAKECHECKERBOARD  
SKLEARNDATASETS MAKECHECKERBOARD SHAPE NCLUSTERS NOISE00 MINVAL10 MAXVAL100

SHUFFLETRUE RANDOMSTATENONE  
GENERATE AN ARRAY WITH BLOCK CHECKERBOARD STRUCTURE FOR BICLUSTERING  
READ MORE IN THE USER GUIDE

PARAMETERS  
SHAPE ITERABLE NROWS NCOLS THE SHAPE OF THE RESULT  
NCLUSTERS INTEGER OR ITERABLE NROWCLUSTERS NCOLUMNCLUSTERS THE NUMBER OF ROW AND  
COLUMN CLUSTERS

NOISE FLOAT OPTIONAL DEFAULT00 THE STANDARD DEVIATION OF THE GAUSSIAN NOISE  
MINVAL INT OPTIONAL DEFAULT10 MINIMUM VALUE OF A BICLUSTER  
MAXVAL INT OPTIONAL DEFAULT100 MAXIMUM VALUE OF A BICLUSTER  
SHUFFLE BOOLEAN OPTIONAL DEFAULTTRUE SHUFFLE THE SAMPLES  
RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE DEFAULT DETERMINES RANDOM NUMBER GEN  
ERATION FOR DATASET CREATION PASS AN INT FOR REPRODUCIBLE OUTPUT ACROSS MULTIPLE FUNCTION  
CALLS SEE GLOSSARY

RETURNS  
XARRAY OF SHAPE SHAPE THE GENERATED ARRAY  
ROWS ARRAY OF SHAPE NCLUSTERS XSHAPE0 THE INDICATORS FOR CLUSTER MEMBERSHIP OF EACH  
ROW  
COLS ARRAY OF SHAPE NCLUSTERS XSHAPE1 THE INDICATORS FOR CLUSTER MEMBERSHIP OF EACH  
COLUMN

SEE ALSO  
MAKEBICLUSTERS  
REFERENCES

1  
EXAMPLES USING SKLEARNDATASETSMMAKECHECKERBOARD  
•A DEMO OF THE SPECTRAL BICLUSTERING ALGORITHM  
68SKLEARNDATASETS DATASETS 1607

SCIKITLEARN USER GUIDE RELEASE 0213

SKLEARNDATASETS MAKECIRCLES

SKLEARNDATASETS MAKECIRCLES NSAMPLES100 SHUFFLETRUE NOISENONE RANDOMSTATENONE  
FACTOR08

MAKE A LARGE CIRCLE CONTAINING A SMALLER CIRCLE IN 2D

A SIMPLE TOY DATASET TO VISUALIZE CLUSTERING AND CLASSIFICATION ALGORITHMS

READ MORE IN THE USER GUIDE

PARAMETERS

NSAMPLES INT OPTIONAL DEFAULT100 THE TOTAL NUMBER OF POINTS GENERATED IF ODD THE INNER  
CIRCLE WILL HAVE ONE POINT MORE THAN THE OUTER CIRCLE

SHUFFLE BOOL OPTIONAL DEFAULTTRUE WHETHER TO SHUFFLE THE SAMPLES

NOISE DOUBLE OR NONE DEFAULTNONE STANDARD DEVIATION OF GAUSSIAN NOISE ADDED TO THE DATA

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE DEFAULT DETERMINES RANDOM NUMBER GEN  
ERATION FOR DATASET SHUFFLING AND NOISE PASS AN INT FOR REPRODUCIBLE OUTPUT ACROSS MULTIPLE  
FUNCTION CALLS SEE GLOSSARY

FACTOR 0 DOUBLE 1 DEFAULT8 SCALE FACTOR BETWEEN INNER AND OUTER CIRCLE

RETURNS

XARRAY OF SHAPE NSAMPLES 2 THE GENERATED SAMPLES

YARRAY OF SHAPE NSAMPLES THE INTEGER LABELS 0 OR 1 FOR CLASS MEMBERSHIP OF EACH SAMPLE

EXAMPLES USING SKLEARNDATASETSMAKECIRCLES

- CLASSIFIER COMPARISON
- COMPARING DIFFERENT HIERARCHICAL LINKAGE METHODS ON TOY DATASETS
- COMPARING DIFFERENT CLUSTERING ALGORITHMS ON TOY DATASETS
- KERNEL PCA
- HASHING FEATURE TRANSFORMATION USING TOTALLY RANDOM TREES
- TSNE THE EFFECT OF VARIOUS PERPLEXITY VALUES ON THE SHAPE
- VARYING REGULARIZATION IN MULTILAYER PERCEPTRON
- COMPARE STOCHASTIC LEARNING STRATEGIES FOR MLPCLASSIFIER
- FEATURE DISCRETIZATION
- LABEL PROPAGATION LEARNING A COMPLEX STRUCTURE

SKLEARNDATASETS MAKECLASSIFICATION

SKLEARNDATASETS MAKECLASSIFICATION NSAMPLES100 NFEATURES20 NINFORMATIVE2  
NREDUNDANT2 NREPEATED0 NCLASSES2

NCLUSTERSPERCLASS2 WEIGHTSNONE FLIPY001

CLASSESP10 HYPERCUBETRUE SHIFT00 SCALE10

SHUFFLETRUE RANDOMSTATENONE

GENERATE A RANDOM NCLASS CLASSIFICATION PROBLEM

1608 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

THIS INITIALLY CREATES CLUSTERS OF POINTS NORMALLY DISTRIBUTED STD1 ABOUT VERTICES OF AN NINFORMATIVE DIMENSIONAL HYPERCUBE WITH SIDES OF LENGTH 2CLASSEP AND ASSIGNS AN EQUAL NUMBER OF CLUSTERS TO EACH CLASS IT INTRODUCES INTERDEPENDENCE BETWEEN THESE FEATURES AND ADDS VARIOUS TYPES OF FURTHER NOISE TO THE DATA WITHOUT SHUFFLING XHORIZONTALLY STACKS FEATURES IN THE FOLLOWING ORDER THE PRIMARY NINFORMATIVE FEATURES FOLLOWED BY NREDUNDANT LINEAR COMBINATIONS OF THE INFORMATIVE FEATURES FOLLOWED BY NREPEATED DUPLICATES DRAWN RANDOMLY WITH REPLACEMENT FROM THE INFORMATIVE AND REDUNDANT FEATURES THE REMAINING FEATURES ARE FILLED WITH RANDOM NOISE THUS WITHOUT SHUFFLING ALL USEFUL FEATURES ARE CONTAINED IN THE COLUMNS

X NINFORMATIVE NREDUNDANT NREPEATED

READ MORE IN THE USER GUIDE

PARAMETERS

NSAMPLES INT OPTIONAL DEFAULT100 THE NUMBER OF SAMPLES

NFEATURES INT OPTIONAL DEFAULT20 THE TOTAL NUMBER OF FEATURES THESE COMPRISE NINFORMATIVE INFORMATIVE FEATURES NREDUNDANT REDUNDANT FEATURES NREPEATED DUPLICATED FEATURES AND NFEATURESNINFORMATIVENREDUNDANTNREPEATED USELESS FEATURES DRAWN AT RANDOM

NINFORMATIVE INT OPTIONAL DEFAULT2 THE NUMBER OF INFORMATIVE FEATURES EACH CLASS IS COMPOSED OF A NUMBER OF GAUSSIAN CLUSTERS EACH LOCATED AROUND THE VERTICES OF A HYPERCUBE IN A SUBSPACE OF DIMENSION NINFORMATIVE FOR EACH CLUSTER INFORMATIVE FEATURES ARE DRAWN INDEPENDENTLY FROM NO 1 AND THEN RANDOMLY LINEARLY COMBINED WITHIN EACH CLUSTER IN ORDER TO ADD COVARIANCE THE CLUSTERS ARE THEN PLACED ON THE VERTICES OF THE HYPERCUBE

NREDUNDANT INT OPTIONAL DEFAULT2 THE NUMBER OF REDUNDANT FEATURES THESE FEATURES ARE GENERATED AS RANDOM LINEAR COMBINATIONS OF THE INFORMATIVE FEATURES

NREPEATED INT OPTIONAL DEFAULT0 THE NUMBER OF DUPLICATED FEATURES DRAWN RANDOMLY FROM THE INFORMATIVE AND THE REDUNDANT FEATURES

NCLASSES INT OPTIONAL DEFAULT2 THE NUMBER OF CLASSES OR LABELS OF THE CLASSIFICATION PROBLEM

NCLUSTERSPERCLASS INT OPTIONAL DEFAULT2 THE NUMBER OF CLUSTERS PER CLASS

WEIGHTS LIST OF FLOATS OR NONE DEFAULTNONE THE PROPORTIONS OF SAMPLES ASSIGNED TO EACH CLASS IF NONE THEN CLASSES ARE BALANCED NOTE THAT IF LENWEIGHTS NCLASSES 1 THEN THE LAST CLASS WEIGHT IS AUTOMATICALLY INFERRED MORE THAN NSAMPLES SAMPLES MAY BE RETURNED IF THE SUM OF WEIGHTS EXCEEDS 1

FLIPY FLOAT OPTIONAL DEFAULT001 THE FRACTION OF SAMPLES WHOSE CLASS ARE RANDOMLY EXCHANGED LARGER VALUES INTRODUCE NOISE IN THE LABELS AND MAKE THE CLASSIFICATION TASK HARDER

CLASSEP FLOAT OPTIONAL DEFAULT10 THE FACTOR MULTIPLYING THE HYPERCUBE SIZE LARGER VALUES SPREAD OUT THE CLUSTERSCLASSES AND MAKE THE CLASSIFICATION TASK EASIER

HYPERCUBE BOOLEAN OPTIONAL DEFAULTTRUE IF TRUE THE CLUSTERS ARE PUT ON THE VERTICES OF A HYPERCUBE IF FALSE THE CLUSTERS ARE PUT ON THE VERTICES OF A RANDOM POLYTOPE

SHIFT FLOAT ARRAY OF SHAPE NFEATURES OR NONE OPTIONAL DEFAULT00 SHIFT FEATURES BY THE SPECIFIED VALUE IF NONE THEN FEATURES ARE SHIFTED BY A RANDOM VALUE DRAWN IN CLASSEP CLASSEP

SCALE FLOAT ARRAY OF SHAPE NFEATURES OR NONE OPTIONAL DEFAULT10 MULTIPLY FEATURES BY THE SPECIFIED VALUE IF NONE THEN FEATURES ARE SCALED BY A RANDOM VALUE DRAWN IN 1 100 NOTE THAT SCALING HAPPENS AFTER SHIFTING

SHUFFLE BOOLEAN OPTIONAL DEFAULTTRUE SHUFFLE THE SAMPLES AND THE FEATURES

68SKLEARNDATASETS DATASETS 1609

SCIKITLEARN USER GUIDE RELEASE 0213

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE DEFAULT DETERMINES RANDOM NUMBER GENERATION FOR DATASET CREATION PASS AN INT FOR REPRODUCIBLE OUTPUT ACROSS MULTIPLE FUNCTION CALLS SEE GLOSSARY

RETURNS

XARRAY OF SHAPE NSAMPLES NFEATURES THE GENERATED SAMPLES

YARRAY OF SHAPE NSAMPLES THE INTEGER LABELS FOR CLASS MEMBERSHIP OF EACH SAMPLE

SEE ALSO

MAKEBLOBS SIMPLIFIED VARIANT

MAKEMULTILABELCLASSIFICATION UNRELATED GENERATOR FOR MULTILABEL TASKS

NOTES

THE ALGORITHM IS ADAPTED FROM GUYON 1 AND WAS DESIGNED TO GENERATE THE “MADELON” DATASET

REFERENCES

1

EXAMPLES USING SKLEARNDATASETSMAKECLASSIFICATION

- COMPARISON OF CALIBRATION OF CLASSIFIERS
- PROBABILITY CALIBRATION CURVES
- CLASSIFIER COMPARISON
- PLOT RANDOMLY GENERATED CLASSIFICATION DATASET
- FEATURE IMPORTANCES WITH FORESTS OF TREES
- OOB ERRORS FOR RANDOM FORESTS
- FEATURE TRANSFORMATIONS WITH ENSEMBLES OF TREES
- PIPELINE ANOVA SVM
- RECURSIVE FEATURE ELIMINATION WITH CROSSVALIDATION
- NEIGHBORHOOD COMPONENTS ANALYSIS ILLUSTRATION
- VARYING REGULARIZATION IN MULTILAYER PERCEPTRON
- FEATURE DISCRETIZATION
- SCALING THE REGULARIZATION PARAMETER FOR SVCS

SKLEARNDATASETS MAKEFRIEDMAN1

SKLEARNDATASETS MAKEFRIEDMAN1 NSAMPLES100 NFEATURES10 NOISE00 RAN

DOMSTATENONE

GENERATE THE “FRIEDMAN 1” REGRESSION PROBLEM

THIS DATASET IS DESCRIBED IN FRIEDMAN 1 AND BREIMAN 2

1610 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

INPUTSXARE INDEPENDENT FEATURES UNIFORMLY DISTRIBUTED ON THE INTERVAL 0 1 THE OUTPUT YIS CREATED ACCORDING TO THE FORMULA

$YX = 10 \sin(\pi x_0 x_1 x_2 x_3 x_4)$

↪3 5X 4 NOISE N0 1

OUT OF THE NFEATURES FEATURES ONLY 5 ARE ACTUALLY USED TO COMPUTE Y THE REMAINING FEATURES ARE INDEPENDENT OFY

THE NUMBER OF FEATURES HAS TO BE 5

READ MORE IN THE USER GUIDE

PARAMETERS

NSAMPLES INT OPTIONAL DEFAULT100 THE NUMBER OF SAMPLES

NFEATURES INT OPTIONAL DEFAULT10 THE NUMBER OF FEATURES SHOULD BE AT LEAST 5

NOISE FLOAT OPTIONAL DEFAULT00 THE STANDARD DEVIATION OF THE GAUSSIAN NOISE APPLIED TO THE OUTPUT

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE DEFAULT DETERMINES RANDOM NUMBER GENERATION FOR DATASET NOISE PASS AN INT FOR REPRODUCIBLE OUTPUT ACROSS MULTIPLE FUNCTION CALLS

SEEGLOSSARY

RETURNS

XARRAY OF SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES

YARRAY OF SHAPE NSAMPLES THE OUTPUT VALUES

REFERENCES

12

SKLEARNDATASETS MAKEFRIEDMAN2

SKLEARNDATASETS MAKEFRIEDMAN2 NSAMPLES100 NOISE00 RANDOMSTATENONE

GENERATE THE “FRIEDMAN 2” REGRESSION PROBLEM

THIS DATASET IS DESCRIBED IN FRIEDMAN 1 AND BREIMAN 2

INPUTSXARE 4 INDEPENDENT FEATURES UNIFORMLY DISTRIBUTED ON THE INTERVALS

$0 \leq x_0 \leq 100$

$40\pi \leq x_1 \leq 560\pi$

$0 \leq x_2 \leq 1$

$1 \leq x_3 \leq 11$

THE OUTPUT YIS CREATED ACCORDING TO THE FORMULA

$YX = x_0 x_2 x_1 x_2 x_1 x_1 x_3 x_3 x_0$

↪5 NOISE N0 1

READ MORE IN THE USER GUIDE

PARAMETERS

NSAMPLES INT OPTIONAL DEFAULT100 THE NUMBER OF SAMPLES

68SKLEARNDATASETS DATASETS 1611

SCIKITLEARN USER GUIDE RELEASE 0213  
 NOISE FLOAT OPTIONAL DEFAULT00 THE STANDARD DEVIATION OF THE GAUSSIAN NOISE APPLIED TO THE OUTPUT  
 RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE DEFAULT DETERMINES RANDOM NUMBER GENERATION FOR DATASET NOISE PASS AN INT FOR REPRODUCIBLE OUTPUT ACROSS MULTIPLE FUNCTION CALLS  
 SEEGLOSSARY  
 RETURNS  
 XARRAY OF SHAPE NSAMPLES 4 THE INPUT SAMPLES  
 YARRAY OF SHAPE NSAMPLES THE OUTPUT VALUES  
 REFERENCES  
 12  
 SKLEARNDATASETS MAKEFRIEDMAN3  
 SKLEARNDATASETS MAKEFRIEDMAN3 NSAMPLES100 NOISE00 RANDOMSTATENONE  
 GENERATE THE “FRIEDMAN 3” REGRESSION PROBLEM  
 THIS DATASET IS DESCRIBED IN FRIEDMAN 1 AND BREIMAN 2  
 INPUTSXARE 4 INDEPENDENT FEATURES UNIFORMLY DISTRIBUTED ON THE INTERVALS  
 0 X 0 100  
 40PI X 1 560 PI  
 0 X 2 1  
 1 X 3 11  
 THE OUTPUT YIS CREATED ACCORDING TO THE FORMULA  

$$YX = \arctan(X_1 X_2 - 1 X_1 X_3 - X_0 \text{ NOISE})$$
 ↪N0 1  
 READ MORE IN THE USER GUIDE  
 PARAMETERS  
 NSAMPLES INT OPTIONAL DEFAULT100 THE NUMBER OF SAMPLES  
 NOISE FLOAT OPTIONAL DEFAULT00 THE STANDARD DEVIATION OF THE GAUSSIAN NOISE APPLIED TO THE OUTPUT  
 RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE DEFAULT DETERMINES RANDOM NUMBER GENERATION FOR DATASET NOISE PASS AN INT FOR REPRODUCIBLE OUTPUT ACROSS MULTIPLE FUNCTION CALLS  
 SEEGLOSSARY  
 RETURNS  
 XARRAY OF SHAPE NSAMPLES 4 THE INPUT SAMPLES  
 YARRAY OF SHAPE NSAMPLES THE OUTPUT VALUES  
 REFERENCES  
 12  
 1612 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

SKLEARNDATASETS MAKEGAUSSIANQUANTILES

SKLEARNDATASETS MAKEGAUSSIANQUANTILES MEANNONE COV10 NSAMPLES100

NFEATURES2 NCLASSES3 SHUFFLETRUE

RANDOMSTATENONE

GENERATE ISOTROPIC GAUSSIAN AND LABEL SAMPLES BY QUANTILE

THIS CLASSIFICATION DATASET IS CONSTRUCTED BY TAKING A MULTIDIMENSIONAL STANDARD NORMAL DISTRIBUTION AND DEFINING CLASSES SEPARATED BY NESTED CONCENTRIC MULTIDIMENSIONAL SPHERES SUCH THAT ROUGHLY EQUAL NUMBERS OF SAMPLES ARE IN EACH CLASS QUANTILES OF THE  $\chi^2$ DISTRIBUTION

READ MORE IN THE USER GUIDE

PARAMETERS

MEAN ARRAY OF SHAPE NFEATURES OPTIONAL DEFAULTNONE THE MEAN OF THE MULTI DIMENSIONAL NORMAL DISTRIBUTION IF NONE THEN USE THE ORIGIN 0 0

COV FLOAT OPTIONAL DEFAULT1 THE COVARIANCE MATRIX WILL BE THIS VALUE TIMES THE UNIT MATRIX

THIS DATASET ONLY PRODUCES SYMMETRIC NORMAL DISTRIBUTIONS

NSAMPLES INT OPTIONAL DEFAULT100 THE TOTAL NUMBER OF POINTS EQUALLY DIVIDED AMONG CLASSES

NFEATURES INT OPTIONAL DEFAULT2 THE NUMBER OF FEATURES FOR EACH SAMPLE

NCLASSES INT OPTIONAL DEFAULT3 THE NUMBER OF CLASSES

SHUFFLE BOOLEAN OPTIONAL DEFAULTTRUE SHUFFLE THE SAMPLES

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE DEFAULT DETERMINES RANDOM NUMBER GENERATION FOR DATASET CREATION PASS AN INT FOR REPRODUCIBLE OUTPUT ACROSS MULTIPLE FUNCTION CALLS SEE GLOSSARY

RETURNS

XARRAY OF SHAPE NSAMPLES NFEATURES THE GENERATED SAMPLES

YARRAY OF SHAPE NSAMPLES THE INTEGER LABELS FOR QUANTILE MEMBERSHIP OF EACH SAMPLE

NOTES

THE DATASET IS FROM ZHU ET AL 1

REFERENCES

1

EXAMPLES USING SKLEARNDATASETSMAKEGAUSSIANQUANTILES

- PLOT RANDOMLY GENERATED CLASSIFICATION DATASET
- TWOCLASS ADABOOST
- MULTICLASS ADABOOSTED DECISION TREES

68SKLEARNDATASETS DATASETS 1613

SCIKITLEARN USER GUIDE RELEASE 0213

SKLEARNDATASETS MAKEHASTIE102

SKLEARNDATASETS MAKEHASTIE102 NSAMPLES12000 RANDOMSTATENONE

GENERATES DATA FOR BINARY CLASSIFICATION USED IN HASTIE ET AL 2009 EXAMPLE 102

THE TEN FEATURES ARE STANDARD INDEPENDENT GAUSSIAN AND THE TARGET YIS DEFINED BY

YI 1 IFNPSUMXI 2 934 ELSE1

READ MORE IN THE USER GUIDE

PARAMETERS

NSAMPLES INT OPTIONAL DEFAULT12000 THE NUMBER OF SAMPLES

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE DEFAULT DETERMINES RANDOM NUMBER GEN

ERATION FOR DATASET CREATION PASS AN INT FOR REPRODUCIBLE OUTPUT ACROSS MULTIPLE FUNCTION

CALLS SEE GLOSSARY

RETURNS

XARRAY OF SHAPE NSAMPLES 10 THE INPUT SAMPLES

YARRAY OF SHAPE NSAMPLES THE OUTPUT VALUES

SEE ALSO

MAKEGAUSSIANQUANTILES A GENERALIZATION OF THIS DATASET APPROACH

REFERENCES

1

EXAMPLES USING SKLEARNDATASETSMAKEHASTIE102

- GRADIENT BOOSTING REGULARIZATION
- DISCRETE VERSUS REAL ADABOOST
- EARLY STOPPING OF GRADIENT BOOSTING
- DEMONSTRATION OF MULTIMETRIC EVALUATION ON CROSSVALSCORE AND GRIDSEARCHCV

SKLEARNDATASETS MAKELOWRANKMATRIX

SKLEARNDATASETS MAKELOWRANKMATRIX NSAMPLES100 NFEATURES100 EFFECTIVERANK10

TAILSTRENGTH05 RANDOMSTATENONE

GENERATE A MOSTLY LOW RANK MATRIX WITH BELLSHAPED SINGULAR VALUES

MOST OF THE VARIANCE CAN BE EXPLAINED BY A BELLSHAPED CURVE OF WIDTH EFFECTIVERANK THE LOW RANK PART OF THE

SINGULAR VALUES PROFILE IS

1 TAILSTRENGTH EXP10 I EFFECTIVERANK 2

THE REMAINING SINGULAR VALUES' TAIL IS FAT DECREASING AS

TAILSTRENGTH EXP01 I EFFECTIVERANK

1614 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

THE LOW RANK PART OF THE PROFILE CAN BE CONSIDERED THE STRUCTURED SIGNAL PART OF THE DATA WHILE THE TAIL CAN BE CONSIDERED THE NOISY PART OF THE DATA THAT CANNOT BE SUMMARIZED BY A LOW NUMBER OF LINEAR COMPONENTS SINGULAR VECTORS

THIS KIND OF SINGULAR PROFILES IS OFTEN SEEN IN PRACTICE FOR INSTANCE

- GRAY LEVEL PICTURES OF FACES
- TFIDF VECTORS OF TEXT DOCUMENTS CRAWLED FROM THE WEB

READ MORE IN THE USER GUIDE

PARAMETERS

NSAMPLES INT OPTIONAL DEFAULT100 THE NUMBER OF SAMPLES

NFEATURES INT OPTIONAL DEFAULT100 THE NUMBER OF FEATURES

EFFECTIVERANK INT OPTIONAL DEFAULT10 THE APPROXIMATE NUMBER OF SINGULAR VECTORS REQUIRED TO EXPLAIN MOST OF THE DATA BY LINEAR COMBINATIONS

TAILSTRENGTH FLOAT BETWEEN 00 AND 10 OPTIONAL DEFAULT05 THE RELATIVE IMPORTANCE OF THE FAT NOISY TAIL OF THE SINGULAR VALUES PROFILE

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE DEFAULT DETERMINES RANDOM NUMBER GENERATION FOR DATASET CREATION PASS AN INT FOR REPRODUCIBLE OUTPUT ACROSS MULTIPLE FUNCTION CALLS SEE GLOSSARY

RETURNS

XARRAY OF SHAPE NSAMPLES NFEATURES THE MATRIX

SKLEARNDATASETS MAKEMOONS

SKLEARNDATASETS MAKEMOONS NSAMPLES100 SHUFFLETRUE NOISENONE RANDOMSTATENONE

MAKE TWO INTERLEAVING HALF CIRCLES

A SIMPLE TOY DATASET TO VISUALIZE CLUSTERING AND CLASSIFICATION ALGORITHMS READ MORE IN THE USER GUIDE

PARAMETERS

NSAMPLES INT OPTIONAL DEFAULT100 THE TOTAL NUMBER OF POINTS GENERATED

SHUFFLE BOOL OPTIONAL DEFAULTTRUE WHETHER TO SHUFFLE THE SAMPLES

NOISE DOUBLE OR NONE DEFAULTNONE STANDARD DEVIATION OF GAUSSIAN NOISE ADDED TO THE DATA

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE DEFAULT DETERMINES RANDOM NUMBER GENERATION FOR DATASET SHUFFLING AND NOISE PASS AN INT FOR REPRODUCIBLE OUTPUT ACROSS MULTIPLE FUNCTION CALLS SEE GLOSSARY

RETURNS

XARRAY OF SHAPE NSAMPLES 2 THE GENERATED SAMPLES

YARRAY OF SHAPE NSAMPLES THE INTEGER LABELS 0 OR 1 FOR CLASS MEMBERSHIP OF EACH SAMPLE

EXAMPLES USING SKLEARNDATASETSMAKEMOONS

- COMPARING ANOMALY DETECTION ALGORITHMS FOR OUTLIER DETECTION ON TOY DATASETS
- CLASSIFIER COMPARISON
- COMPARING DIFFERENT HIERARCHICAL LINKAGE METHODS ON TOY DATASETS

68SKLEARNDATASETS DATASETS 1615

SCIKITLEARN USER GUIDE RELEASE 0213

- COMPARING DIFFERENT CLUSTERING ALGORITHMS ON TOY DATASETS
- VARYING REGULARIZATION IN MULTILAYER PERCEPTRON
- COMPARE STOCHASTIC LEARNING STRATEGIES FOR MLPCLASSIFIER
- FEATURE DISCRETIZATION

SKLEARN DATASETS MAKE MULTILABEL CLASSIFICATION

SKLEARN DATASETS MAKE MULTILABEL CLASSIFICATION NSAMPLES100 NFEATURES20  
NCLASSES5 NLABELS2 LENGTH50  
ALLOW UNLABELED TRUE SPARSE FALSE  
RETURN INDICATOR 'DENSE' RETURN DISTRIBUTIONS FALSE RANDOM STATE NONE

GENERATE A RANDOM MULTILABEL CLASSIFICATION PROBLEM

FOR EACH SAMPLE THE GENERATIVE PROCESS IS

- PICK THE NUMBER OF LABELS  $N \sim \text{POISSON}(\text{NLABELS})$
- $N$  TIMES CHOOSE A CLASS  $C \sim \text{MULTINOMIAL}(\text{THETA})$
- PICK THE DOCUMENT LENGTH  $K \sim \text{POISSON}(\text{LENGTH})$
- $K$  TIMES CHOOSE A WORD  $W \sim \text{MULTINOMIAL}(\text{THETA}_C)$

IN THE ABOVE PROCESS REJECTION SAMPLING IS USED TO MAKE SURE THAT  $N$  IS NEVER ZERO OR MORE THAN  $N_{\text{CLASSES}}$  AND THAT THE DOCUMENT LENGTH IS NEVER ZERO LIKEWISE WE REJECT CLASSES WHICH HAVE ALREADY BEEN CHOSEN

READ MORE IN THE USER GUIDE

PARAMETERS

NSAMPLES INT OPTIONAL DEFAULT 100 THE NUMBER OF SAMPLES

NFEATURES INT OPTIONAL DEFAULT 20 THE TOTAL NUMBER OF FEATURES

NCLASSES INT OPTIONAL DEFAULT 5 THE NUMBER OF CLASSES OF THE CLASSIFICATION PROBLEM

NLABELS INT OPTIONAL DEFAULT 2 THE AVERAGE NUMBER OF LABELS PER INSTANCE MORE PRECISELY THE NUMBER OF LABELS PER SAMPLE IS DRAWN FROM A POISSON DISTRIBUTION WITH NLABELS AS ITS EXPECTED VALUE BUT SAMPLES ARE BOUNDED USING REJECTION SAMPLING BY NCLASSES AND MUST BE NONZERO IF ALLOW UNLABELED IS FALSE

LENGTH INT OPTIONAL DEFAULT 50 THE SUM OF THE FEATURES NUMBER OF WORDS IF DOCUMENTS IS DRAWN FROM A POISSON DISTRIBUTION WITH THIS EXPECTED VALUE

ALLOW UNLABELED BOOL OPTIONAL DEFAULT TRUE IF TRUE SOME INSTANCES MIGHT NOT BELONG TO ANY CLASS

SPARSE BOOL OPTIONAL DEFAULT FALSE IF TRUE RETURN A SPARSE FEATURE MATRIX

NEW IN VERSION 0.17 PARAMETER TO ALLOW SPARSE OUTPUT

RETURN INDICATOR 'DENSE' DEFAULT 'SPARSE' FALSE IF DENSE RETURN YES IN THE DENSE BINARY INDICATOR FORMAT IF SPARSE RETURN YES IN THE SPARSE BINARY INDICATOR FORMAT FALSE

RETURNS A LIST OF LISTS OF LABELS

RETURN DISTRIBUTIONS BOOL OPTIONAL DEFAULT FALSE IF TRUE RETURN THE PRIOR CLASS PROBABILITY AND CONDITIONAL PROBABILITIES OF FEATURES GIVEN CLASSES FROM WHICH THE DATA WAS DRAWN

1616 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE DEFAULT DETERMINES RANDOM NUMBER GENERATION FOR DATASET CREATION PASS AN INT FOR REPRODUCIBLE OUTPUT ACROSS MULTIPLE FUNCTION CALLS SEE GLOSSARY

RETURNS

XARRAY OF SHAPE NSAMPLES NFEATURES THE GENERATED SAMPLES

YARRAY OR SPARSE CSR MATRIX OF SHAPE NSAMPLES NCLASSES THE LABEL SETS

PC ARRAY SHAPE NCLASSES THE PROBABILITY OF EACH CLASS BEING DRAWN ONLY RETURNED IF RETURNNDISTRIBUTIONSTRUE

PWC ARRAY SHAPE NFEATURES NCLASSES THE PROBABILITY OF EACH FEATURE BEING DRAWN GIVEN EACH CLASS ONLY RETURNED IF RETURNNDISTRIBUTIONSTRUE

EXAMPLES USING SKLEARNDATASETSMAKEMULTILABELCLASSIFICATION

- MULTILABEL CLASSIFICATION
- PLOT RANDOMLY GENERATED MULTILABEL DATASET

SKLEARNDATASETS MAKEREGRESSION

SKLEARNDATASETS MAKEREGRESSION NSAMPLES100 NFEATURES100 NINFORMATIVE10 NTARGETS1 BIAS00 EFFECTIVERANKNONE

TAILSTRENGTH05 NOISE00 SHUFFLETRUE COEFFALSE

RANDOMSTATENONE

GENERATE A RANDOM REGRESSION PROBLEM

THE INPUT SET CAN EITHER BE WELL CONDITIONED BY DEFAULT OR HAVE A LOW RANKFAT TAIL SINGULAR PROFILE SEE MAKELOWRANKMATRIX FOR MORE DETAILS

THE OUTPUT IS GENERATED BY APPLYING A POTENTIALLY BIASED RANDOM LINEAR REGRESSION MODEL WITH NINFORMATIVE NONZERO REGRESSORS TO THE PREVIOUSLY GENERATED INPUT AND SOME GAUSSIAN CENTERED NOISE WITH SOME ADJUSTABLE SCALE

READ MORE IN THE USER GUIDE

PARAMETERS

NSAMPLES INT OPTIONAL DEFAULT100 THE NUMBER OF SAMPLES

NFEATURES INT OPTIONAL DEFAULT100 THE NUMBER OF FEATURES

NINFORMATIVE INT OPTIONAL DEFAULT10 THE NUMBER OF INFORMATIVE FEATURES IE THE NUMBER OF FEATURES USED TO BUILD THE LINEAR MODEL USED TO GENERATE THE OUTPUT

NTARGETS INT OPTIONAL DEFAULT1 THE NUMBER OF REGRESSION TARGETS IE THE DIMENSION OF THE Y OUTPUT VECTOR ASSOCIATED WITH A SAMPLE BY DEFAULT THE OUTPUT IS A SCALAR

BIAS FLOAT OPTIONAL DEFAULT00 THE BIAS TERM IN THE UNDERLYING LINEAR MODEL

EFFECTIVERANK INT OR NONE OPTIONAL DEFAULTNONE

IF NOT NONE THE APPROXIMATE NUMBER OF SINGULAR VECTORS REQUIRED TO EXPLAIN MOST OF THE INPUT DATA BY LINEAR COMBINATIONS USING THIS KIND OF SINGULAR SPECTRUM IN THE INPUT ALLOWS THE GENERATOR TO REPRODUCE THE CORRELATIONS OFTEN OBSERVED IN PRACTICE

IF NONE THE INPUT SET IS WELL CONDITIONED CENTERED AND GAUSSIAN WITH UNIT VARIANCE

68SKLEARNDATASETS DATASETS 1617

SCIKITLEARN USER GUIDE RELEASE 0213

TAILSTRENGTH FLOAT BETWEEN 00 AND 10 OPTIONAL DEFAULT05 THE RELATIVE IMPORTANCE OF THE FAT NOISY TAIL OF THE SINGULAR VALUES PROFILE IF EFFECTIVERANK IS NOT NONE

NOISE FLOAT OPTIONAL DEFAULT00 THE STANDARD DEVIATION OF THE GAUSSIAN NOISE APPLIED TO THE OUTPUT

SHUFFLE BOOLEAN OPTIONAL DEFAULTTRUE SHUFFLE THE SAMPLES AND THE FEATURES

COEF BOOLEAN OPTIONAL DEFAULTFALSE IF TRUE THE COEFFICIENTS OF THE UNDERLYING LINEAR MODEL ARE RETURNED

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE DEFAULT DETERMINES RANDOM NUMBER GENERATION FOR DATASET CREATION PASS AN INT FOR REPRODUCIBLE OUTPUT ACROSS MULTIPLE FUNCTION CALLS SEE GLOSSARY

RETURNS

XARRAY OF SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES

YARRAY OF SHAPE NSAMPLES OR NSAMPLES NTARGETS THE OUTPUT VALUES

COEF ARRAY OF SHAPE NFEATURES OR NFEATURES NTARGETS OPTIONAL THE COEFFICIENT OF THE UNDERLYING LINEAR MODEL IT IS RETURNED ONLY IF COEF IS TRUE

EXAMPLES USING SKLEARNDATASETSMAKEREgression

- PREDICTION LATENCY
- EFFECT OF TRANSFORMING THE TARGETS IN REGRESSION MODEL
- PLOT RIDGE COEFFICIENTS AS A FUNCTION OF THE L2 REGULARIZATION
- ROBUST LINEAR MODEL ESTIMATION USING RANSAC
- LASSO ON DENSE AND SPARSE DATA
- HUBERREGRESSOR VS RIDGE ON DATASET WITH STRONG OUTLIERS

SKLEARNDATASETS MAKESCURVE

SKLEARNDATASETS MAKESCURVE NSAMPLES100 NOISE00 RANDOMSTATENONE

GENERATE AN S CURVE DATASET

READ MORE IN THE USER GUIDE

PARAMETERS

NSAMPLES INT OPTIONAL DEFAULT100 THE NUMBER OF SAMPLE POINTS ON THE S CURVE

NOISE FLOAT OPTIONAL DEFAULT00 THE STANDARD DEVIATION OF THE GAUSSIAN NOISE

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE DEFAULT DETERMINES RANDOM NUMBER GENERATION FOR DATASET CREATION PASS AN INT FOR REPRODUCIBLE OUTPUT ACROSS MULTIPLE FUNCTION CALLS SEE GLOSSARY

RETURNS

XARRAY OF SHAPE NSAMPLES 3 THE POINTS

TARRAY OF SHAPE NSAMPLES THE UNIVARIATE POSITION OF THE SAMPLE ACCORDING TO THE MAIN DIMENSION OF THE POINTS IN THE MANIFOLD

1618 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES USING SKLEARNDATASETSMAKESCURVE

- TSNE THE EFFECT OF VARIOUS PERPLEXITY VALUES ON THE SHAPE
- COMPARISON OF MANIFOLD LEARNING METHODS

SKLEARNDATASETS MAKESPARSECODEDSIGNAL

SKLEARNDATASETS MAKESPARSECODEDSIGNAL NSAMPLES NCOMPONENTS NFEATURES

NNONZEROCOEF RANDOMSTATENONE

GENERATE A SIGNAL AS A SPARSE COMBINATION OF DICTIONARY ELEMENTS

RETURNS A MATRIX Y DX SUCH AS D IS NFEATURES NCOMPONENTS X IS NCOMPONENTS NSAMPLES AND EACH COLUMN OF X HAS EXACTLY NNONZEROCOEF NONZERO ELEMENTS

READ MORE IN THE USER GUIDE

PARAMETERS

NSAMPLES INT NUMBER OF SAMPLES TO GENERATE

NCOMPONENTS INT NUMBER OF COMPONENTS IN THE DICTIONARY

NFEATURES INT NUMBER OF FEATURES OF THE DATASET TO GENERATE

NNONZEROCOEF INT NUMBER OF ACTIVE NONZERO COEFFICIENTS IN EACH SAMPLE

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE DEFAULT DETERMINES RANDOM NUMBER GENERATION FOR DATASET CREATION PASS AN INT FOR REPRODUCIBLE OUTPUT ACROSS MULTIPLE FUNCTION CALLS SEE GLOSSARY

RETURNS

DATA ARRAY OF SHAPE NFEATURES NSAMPLES THE ENCODED SIGNAL Y

DICTIONARY ARRAY OF SHAPE NFEATURES NCOMPONENTS THE DICTIONARY WITH NORMALIZED COMPONENTS D

CODE ARRAY OF SHAPE NCOMPONENTS NSAMPLES THE SPARSE CODE SUCH THAT EACH COLUMN OF THIS MATRIX HAS EXACTLY NNONZEROCOEF NONZERO ITEMS X

EXAMPLES USING SKLEARNDATASETSMAKESPARSECODEDSIGNAL

- ORTHOGONAL MATCHING PURSUIT

SKLEARNDATASETS MAKESPARSESPDMATRIX

SKLEARNDATASETS MAKESPARSESPDMATRIX DIM1 ALPHA095 NORMDIAGFALSE

SMALLESTCOEF01 LARGESTCOEF09 RAN

DOMSTATENONE

GENERATE A SPARSE SYMMETRIC DEFINITE POSITIVE MATRIX

READ MORE IN THE USER GUIDE

PARAMETERS

DIM INTEGER OPTIONAL DEFAULT1 THE SIZE OF THE RANDOM MATRIX TO GENERATE

ALPHA FLOAT BETWEEN 0 AND 1 OPTIONAL DEFAULT095 THE PROBABILITY THAT A COEFFICIENT IS ZERO

SEE NOTES LARGER VALUES ENFORCE MORE SPARSITY

68SKLEARNDATASETS DATASETS 1619

SCIKITLEARN USER GUIDE RELEASE 0213

NORMDIAG BOOLEAN OPTIONAL DEFAULTFALSE WHETHER TO NORMALIZE THE OUTPUT MATRIX TO MAKE THE LEADING DIAGONAL ELEMENTS ALL 1

SMALLESTCOEF FLOAT BETWEEN 0 AND 1 OPTIONAL DEFAULT01 THE VALUE OF THE SMALLEST COEFFICIENT

LARGESTCOEF FLOAT BETWEEN 0 AND 1 OPTIONAL DEFAULT09 THE VALUE OF THE LARGEST COEFFICIENT

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE DEFAULT DETERMINES RANDOM NUMBER GENERATION FOR DATASET CREATION PASS AN INT FOR REPRODUCIBLE OUTPUT ACROSS MULTIPLE FUNCTION CALLS SEE GLOSSARY

RETURNS

PREC SPARSE MATRIX OF SHAPE DIM DIM THE GENERATED MATRIX

SEE ALSO

MAKESPDMATRIX

NOTES

THE SPARSITY IS ACTUALLY IMPOSED ON THE CHOLESKY FACTOR OF THE MATRIX THUS ALPHA DOES NOT TRANSLATE DIRECTLY INTO THE FILLING FRACTION OF THE MATRIX ITSELF

EXAMPLES USING SKLEARNDATASETSMAKESPARSESPDMATRIX

- SPARSE INVERSE COVARIANCE ESTIMATION

SKLEARNDATASETS MAKESPARSEUNCORRELATED

SKLEARNDATASETS MAKESPARSEUNCORRELATED NSAMPLES100 NFEATURES10 RANDOMSTATENONE

GENERATE A RANDOM REGRESSION PROBLEM WITH SPARSE UNCORRELATED DESIGN

THIS DATASET IS DESCRIBED IN CELEUX ET AL 1 AS

$X = \begin{bmatrix} 0 & 1 \\ 0 & 2 \\ 1 & 1 \\ 2 & 2 \end{bmatrix}$   $Y = \begin{bmatrix} 15 \\ 3 \end{bmatrix}$

ONLY THE FIRST 4 FEATURES ARE INFORMATIVE THE REMAINING FEATURES ARE USELESS

READ MORE IN THE USER GUIDE

PARAMETERS

NSAMPLES INT OPTIONAL DEFAULT100 THE NUMBER OF SAMPLES

NFEATURES INT OPTIONAL DEFAULT10 THE NUMBER OF FEATURES

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE DEFAULT DETERMINES RANDOM NUMBER GENERATION FOR DATASET CREATION PASS AN INT FOR REPRODUCIBLE OUTPUT ACROSS MULTIPLE FUNCTION CALLS SEE GLOSSARY

RETURNS

XARRAY OF SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES

YARRAY OF SHAPE NSAMPLES THE OUTPUT VALUES

1620 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

REFERENCES

1

SKLEARNDATASETS MAKESPDMATRIX

SKLEARNDATASETS MAKESPDMATRIX NDIM RANDOMSTATENONE

GENERATE A RANDOM SYMMETRIC POSITIVEDEFINITE MATRIX

READ MORE IN THE USER GUIDE

PARAMETERS

NDIM INT THE MATRIX DIMENSION

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE DEFAULT DETERMINES RANDOM NUMBER GEN

ERATION FOR DATASET CREATION PASS AN INT FOR REPRODUCIBLE OUTPUT ACROSS MULTIPLE FUNCTION

CALLS SEE GLOSSARY

RETURNS

XARRAY OF SHAPE NDIM NDIM THE RANDOM SYMMETRIC POSITIVEDEFINITE MATRIX

SEE ALSO

MAKESPARSESPDMATRIX

SKLEARNDATASETS MAKESWISSROLL

SKLEARNDATASETS MAKESWISSROLL NSAMPLES100 NOISE00 RANDOMSTATENONE

GENERATE A SWISS ROLL DATASET

READ MORE IN THE USER GUIDE

PARAMETERS

NSAMPLES INT OPTIONAL DEFAULT100 THE NUMBER OF SAMPLE POINTS ON THE S CURVE

NOISE FLOAT OPTIONAL DEFAULT00 THE STANDARD DEVIATION OF THE GAUSSIAN NOISE

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE DEFAULT DETERMINES RANDOM NUMBER GEN

ERATION FOR DATASET CREATION PASS AN INT FOR REPRODUCIBLE OUTPUT ACROSS MULTIPLE FUNCTION

CALLS SEE GLOSSARY

RETURNS

XARRAY OF SHAPE NSAMPLES 3 THE POINTS

TARRAY OF SHAPE NSAMPLES THE UNIVARIATE POSITION OF THE SAMPLE ACCORDING TO THE MAIN

DIMENSION OF THE POINTS IN THE MANIFOLD

NOTES

THE ALGORITHM IS FROM MARSLAND 1

REFERENCES

1

68SKLEARNDATASETS DATASETS 1621

SCIKITLEARN USER GUIDE RELEASE 0213  
EXAMPLES USING SKLEARNDATASETSMAKESWISSROLL  
•HIERARCHICAL CLUSTERING STRUCTURED VS UNSTRUCTURED WARD  
•SWISS ROLL REDUCTION WITH LLE  
69SKLEARNDECOMPOSITION MATRIX DECOMPOSITION  
THESKLEARNDECOMPOSITION MODULE INCLUDES MATRIX DECOMPOSITION ALGORITHMS INCLUDING AMONG OTHERS PCA  
NMF OR ICA MOST OF THE ALGORITHMS OF THIS MODULE CAN BE REGARDED AS DIMENSIONALITY REDUCTION TECHNIQUES  
USER GUIDE SEE THE DECOMPOSING SIGNALS IN COMPONENTS MATRIX FACTORIZATION PROBLEMS SECTION FOR FURTHER DETAILS  
DECOMPOSITIONDICTIONARYLEARNING DICTIONARY LEARNING  
DECOMPOSITIONFACTORANALYSIS NCOMPONENTS  
FACTOR ANALYSIS FA  
DECOMPOSITIONFASTICA NCOMPONENTS FASTICA A FAST ALGORITHM FOR INDEPENDENT COMPONENT  
ANALYSIS  
DECOMPOSITIONINCREMENTALPCA NCOMPONENTS  
INCREMENTAL PRINCIPAL COMPONENTS ANALYSIS IPCA  
DECOMPOSITIONKERNELPCA NCOMPONENTS KERNEL PRINCIPAL COMPONENT ANALYSIS KPCA  
DECOMPOSITIONLATENTDIRICHLETALLOCATION LATENT DIRICHLET ALLOCATION WITH ONLINE VARIATIONAL BAYES  
ALGORITHM  
DECOMPOSITIONMINIBATCHDICTIONARYLEARNING MINIBATCH DICTIONARY LEARNING  
DECOMPOSITIONMINIBATCHSPARSEPCA MINIBATCH SPARSE PRINCIPAL COMPONENTS ANALYSIS  
DECOMPOSITIONNMF NCOMPONENTS INIT NONNEGATIVE MATRIX FACTORIZATION NMF  
DECOMPOSITIONPCA NCOMPONENTS COPY PRINCIPAL COMPONENT ANALYSIS PCA  
DECOMPOSITIONSPARSEPCA NCOMPONENTS SPARSE PRINCIPAL COMPONENTS ANALYSIS SPARSEPCA  
DECOMPOSITIONSPARSECODER DICTIONARY SPARSE CODING  
DECOMPOSITIONTRUNCATEDSVD NCOMPONENTS  
DIMENSIONALITY REDUCTION USING TRUNCATED SVD AKA LSA  
691SKLEARNDECOMPOSITION DICTIONARYLEARNING  
CLASSSSKLEARNDECOMPOSITION DICTIONARYLEARNING NCOMPONENTSNONE ALPHA1  
MAXITER1000 TOL1E08  
FITALGORITHM'LARS' TRANS  
FORMALGORITHM'OMP' TRANS  
FORMNNONZEROCOEFSSNONE TRANS  
FORMALPHANONE NJOBSNONE  
CODEINITNONE DICTINITNONE VER  
BOSEFALSE SPLITSIGNFALSE RAN  
DOMSTATENONE POSITIVECODEFALSE  
POSITIVEDICTFALSE  
DICTIONARY LEARNING  
FINDS A DICTIONARY A SET OF ATOMS THAT CAN BEST BE USED TO REPRESENT DATA USING A SPARSE CODE  
SOLVES THE OPTIMIZATION PROBLEM  
UV ARGMIN 05 Y U V 22 ALPHA U 1  
UV  
WITH VK 2 1 FORALL 0 K NCOMPONENTS  
1622 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213  
READ MORE IN THE USER GUIDE  
PARAMETERS  
NCOMPONENTS INT NUMBER OF DICTIONARY ELEMENTS TO EXTRACT  
ALPHA FLOAT SPARSITY CONTROLLING PARAMETER  
MAXITER INT MAXIMUM NUMBER OF ITERATIONS TO PERFORM  
TOLFLOAT TOLERANCE FOR NUMERICAL ERROR  
FITALGORITHM 'LARS' 'CD' LARS USES THE LEAST ANGLE REGRESSION METHOD TO SOLVE THE LASSO  
PROBLEM LINEARMODELLARSPATH CD USES THE COORDINATE DESCENT METHOD TO COMPUTE THE  
LASSO SOLUTION LINEARMODELLASSO LARS WILL BE FASTER IF THE ESTIMATED COMPONENTS ARE  
SPARSE  
NEW IN VERSION 017 CDCOORDINATE DESCENT METHOD TO IMPROVE SPEED  
TRANSFORMALGORITHM 'LASSOLARS' 'LASSOCD' 'LARS' 'OMP' 'THRESHOLD' ALGORITHM USED  
TO TRANSFORM THE DATA LARS USES THE LEAST ANGLE REGRESSION METHOD LINEARMODELLARSPATH  
LASSOLARS USES LARS TO COMPUTE THE LASSO SOLUTION LASSOCD USES THE COORDINATE DESCENT  
METHOD TO COMPUTE THE LASSO SOLUTION LINEARMODELLASSO LASSOLARS WILL BE FASTER IF  
THE ESTIMATED COMPONENTS ARE SPARSE OMP USES ORTHOGONAL MATCHING PURSUIT TO ESTIMATE  
THE SPARSE SOLUTION THRESHOLD SQUASHES TO ZERO ALL COEFFICIENTS LESS THAN ALPHA FROM THE  
PROJECTIONDICTIONARY X  
NEW IN VERSION 017 LASSOCD COORDINATE DESCENT METHOD TO IMPROVE SPEED  
TRANSFORMNNNONZEROCOEF5 INTO1NFEATURES BY DEFAULT NUMBER OF NONZERO  
COEFFICIENTS TO TARGET IN EACH COLUMN OF THE SOLUTION THIS IS ONLY USED BY  
ALGORITHMMLARS ANDALGORITHMOMP AND IS OVERRIDDEN BY ALPHA IN THE OR  
THOGONAL MATCHING PURSUIT OMP CASE  
TRANSFORMALPHA FLOAT 1 BY DEFAULT IF ALGORITHMMLASSOLARS OR  
ALGORITHMMLASSOCD ALPHA IS THE PENALTY APPLIED TO THE L1 NORM IF  
ALGORITHMTHRESHOLD ALPHA IS THE ABSOLUTE VALUE OF THE THRESHOLD BELOW  
WHICH COEFFICIENTS WILL BE SQUASHED TO ZERO IF ALGORITHMOMP ALPHA IS THE  
TOLERANCE PARAMETER THE VALUE OF THE RECONSTRUCTION ERROR TARGETED IN THIS CASE IT OVERRIDES  
NNNONZEROCOEF5  
NJOBS INT OR NONE OPTIONAL DEFAULTNONE NUMBER OF PARALLEL JOBS TO RUN NONE MEANS 1  
UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL PROCESSORS SEE  
GLOSSARY FOR MORE DETAILS  
CODEINIT ARRAY OF SHAPE NSAMPLES NCOMPONENTS INITIAL VALUE FOR THE CODE FOR WARM  
RESTART  
DICTINIT ARRAY OF SHAPE NCOMPONENTS NFEATURES INITIAL VALUES FOR THE DICTIONARY FOR WARM  
RESTART  
VERBOSE BOOL OPTIONAL DEFAULT FALSE TO CONTROL THE VERBOSITY OF THE PROCEDURE  
SPLITSIGN BOOL FALSE BY DEFAULT WHETHER TO SPLIT THE SPARSE FEATURE VECTOR INTO THE CONCATENATION  
OF ITS NEGATIVE PART AND ITS POSITIVE PART THIS CAN IMPROVE THE PERFORMANCE OF DOWN  
STREAM CLASSIFIERS  
RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RANDOM  
DOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE  
RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS  
THE RANDOMSTATE INSTANCE USED BY NRANDOM  
695KLEARNDECOMPOSITION MATRIX DECOMPOSITION 1623

SCIKITLEARN USER GUIDE RELEASE 0213

POSITIVECODE BOOL WHETHER TO ENFORCE POSITIVITY WHEN FINDING THE CODE  
NEW IN VERSION 020

POSITIVEDICT BOOL WHETHER TO ENFORCE POSITIVITY WHEN FINDING THE DICTIONARY  
NEW IN VERSION 020

ATTRIBUTES

COMPONENTS ARRAY NCOMPONENTS NFEATURES DICTIONARY ATOMS EXTRACTED FROM THE DATA

ERROR ARRAY VECTOR OF ERRORS AT EACH ITERATION

NITER INT NUMBER OF ITERATIONS RUN

SEE ALSO

SPARSECODER

MINIBATCHDICTIONARYLEARNING

SPARSEPCA

MINIBATCHSPARSEPCA

NOTES

REFERENCES

J MAIRAL F BACH J PONCE G SAPIRO 2009 ONLINE DICTIONARY LEARNING FOR SPARSE CODING [HTTPSWWWDIENS](https://www.di.ens.fr/~sierrapdf/sicml09.pdf)  
FRSIERRAPDFSICML09PDF

METHODS

FITSELF X Y FIT THE MODEL FROM DATA IN X

FITTRANSFORM SELF X Y FIT TO DATA THEN TRANSFORM IT

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF X ENCODE THE DATA AS A SPARSE COMBINATION OF THE DICTIO  
NARY ATOMS

INIT SELFNCOMPONENTSNONE ALPHA1 MAXITER1000 TOL1E08 FITALGORITHM'LARS'

TRANSFORMALGORITHM'OMP' TRANSFORMMNNONZEROCOEFNONE TRANSFORMALPHANONE

NJOBSNONE CODEINITNONE DICTINITNONE VERBOSEFALSE SPLITSIGNFALSE RAN

DOMSTATENONE POSITIVECODEFALSE POSITIVEDICTFALSE

FITSELFXYNONE

FIT THE MODEL FROM DATA IN X

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING VECTOR WHERE NSAMPLES IN THE NUM  
BER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

YIGNORED

RETURNS

SELF OBJECT RETURNS THE OBJECT ITSELF

1624 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

`fit_transform(self, X=None, fit_params=None)`  
Fit to data then transform it

`fit_transformer(X, y, fit_params=None)`  
Fits transformer to X and y with optional parameters fit\_params and returns a transformed version of X

`parameters`  
X: Numpy array of shape (n\_samples, n\_features) Training set  
y: Numpy array of shape (n\_samples) Target values

`returns`  
X: New Numpy array of shape (n\_samples, n\_features) New transformed array

`get_params(self, deep=True)`  
Get parameters for this estimator

`parameters`  
deep: Boolean optional. If True will return the parameters for this estimator and contained subobjects that are estimators

`returns`  
params: Mapping of string to any parameter names mapped to their values

`set_params(self, **params)`  
Set the parameters of this estimator

The method works on simple estimators as well as on nested objects such as pipelines. The latter have parameters of the form `component__parameter` so that it's possible to update each component of a nested object

`returns`  
self

`transform(self, X)`  
Encode the data as a sparse combination of the dictionary atoms

Coding method is determined by the object parameter `transform_algorithm`

`parameters`  
X: Array of shape (n\_samples, n\_features) Test data to be transformed must have the same number of features as the data used to train the model

`returns`  
X: New array shape (n\_samples, n\_components) Transformed data

6925 `sklearn.decomposition.FactorAnalysis`

CLASS `sklearn.decomposition.FactorAnalysis` n\_components: int, None, 'tol' 0.01 copy: bool

MAX\_ITER: 1000 noise\_variance: float, None

`SVDMethod` 'RANDOMIZED' iterated\_power: 3

`RANDOMSTATE` 0

`FactorAnalysis` FA

A simple linear generative model with Gaussian latent variables

The observations are assumed to be caused by a linear transformation of lower dimensional latent factors and added Gaussian noise. Without loss of generality the factors are distributed according to a Gaussian with zero mean and unit covariance. The noise is also zero mean and has an arbitrary diagonal covariance matrix

695 `sklearn.decomposition.MatrixDecomposition` 1625

SCIKITLEARN USER GUIDE RELEASE 0213

IF WE WOULD RESTRICT THE MODEL FURTHER BY ASSUMING THAT THE GAUSSIAN NOISE IS EVEN ISOTROPIC ALL DIAGONAL ENTRIES ARE THE SAME WE WOULD OBTAIN PPCA

FACTORANALYSIS PERFORMS A MAXIMUM LIKELIHOOD ESTIMATE OF THE SOCALLED LOADING MATRIX THE TRANSFORMATION OF THE LATENT VARIABLES TO THE OBSERVED ONES USING EXPECTATIONMAXIMIZATION EM

READ MORE IN THE USER GUIDE

PARAMETERS

NCOMPONENTS INT NONE DIMENSIONALITY OF LATENT SPACE THE NUMBER OF COMPONENTS OF XTHAT ARE OBTAINED AFTER TRANSFORM IF NONE NCOMPONENTS IS SET TO THE NUMBER OF FEATURES

TOLFLOAT STOPPING TOLERANCE FOR EM ALGORITHM

COPY BOOL WHETHER TO MAKE A COPY OF X IF FALSE THE INPUT X GETS OVERWRITTEN DURING FITTING

MAXITER INT MAXIMUM NUMBER OF ITERATIONS

NOISEVARIANCEINIT NONE ARRAY SHAPENFEATURES THE INITIAL GUESS OF THE NOISE VARIANCE FOR EACH FEATURE IF NONE IT DEFAULTS TO NPONESNFEATURES

SVDMETHOD 'LAPACK' 'RANDOMIZED' WHICH SVD METHOD TO USE IF 'LAPACK' USE STANDARD SVD FROM SCIPYLINALG IF 'RANDOMIZED' USE FAST RANDOMIZEDSVD FUNCTION DEFAULTS TO 'RANDOMIZED' FOR MOST APPLICATIONS 'RANDOMIZED' WILL BE SUFFICIENTLY PRECISE WHILE PROVIDING SIGNIFICANT SPEED GAINS ACCURACY CAN ALSO BE IMPROVED BY SETTING HIGHER VALUES FOR ITERATEDPOWER IF THIS IS NOT SUFFICIENT FOR MAXIMUM PRECISION YOU SHOULD CHOOSE 'LAPACK'

ITERATEDPOWER INT OPTIONAL NUMBER OF ITERATIONS FOR THE POWER METHOD 3 BY DEFAULT ONLY USED IFSVDMETHOD EQUALS 'RANDOMIZED'

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULT0 IF INT RANDOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NPRANDOM ONLY USED WHEN SVDMETHOD EQUALS 'RANDOMIZED'

ATTRIBUTES

COMPONENTS ARRAY NCOMPONENTS NFEATURES COMPONENTS WITH MAXIMUM VARIANCE

LOGLIKE LIST NITERATIONS THE LOG LIKELIHOOD AT EACH ITERATION

NOISEVARIANCE ARRAY SHAPENFEATURES THE ESTIMATED NOISE VARIANCE FOR EACH FEATURE

NITER INT NUMBER OF ITERATIONS RUN

SEE ALSO

PCA PRINCIPAL COMPONENT ANALYSIS IS ALSO A LATENT LINEAR VARIABLE MODEL WHICH HOWEVER ASSUMES EQUAL NOISE VARIANCE FOR EACH FEATURE THIS EXTRA ASSUMPTION MAKES PROBABILISTIC PCA FASTER AS IT CAN BE COMPUTED IN CLOSED FORM

FASTICA INDEPENDENT COMPONENT ANALYSIS A LATENT VARIABLE MODEL WITH NONGAUSSIAN LATENT VARIABLES

REFERENCES

EXAMPLES

1626 CHAPTER 6 API REFERENCE



```
SCIKITLEARN USER GUIDE RELEASE 0213
FROM SKLEARNDATASETS IMPORT LOADDIGITS
FROM SKLEARNDECOMPOSITION IMPORT FACTORANALYSIS
X LOADDIGITSRETURNXY TRUE
TRANSFORMER FACTORANALYSISNCOMPONENTS7 RANDOMSTATE0
XTRANSFORMED TRANSFORMERFITTRANSFORMX
XTRANSFORMEDSHAPE
1797 7
METHODS
FITSELF X Y FIT THE FACTORANALYSIS MODEL TO X USING EM
FITTRANSFORM SELF X Y FIT TO DATA THEN TRANSFORM IT
GETCOVARIANCE SELF COMPUTE DATA COVARIANCE WITH THE FACTORANALYSIS
MODEL
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR
GETPRECISION SELF COMPUTE DATA PRECISION MATRIX WITH THE FACTORANALYSIS
MODEL
SCORE SELF X Y COMPUTE THE AVERAGE LOGLIKELIHOOD OF THE SAMPLES
SCORESAMPLES SELF X COMPUTE THE LOGLIKELIHOOD OF EACH SAMPLE
SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR
TRANSFORM SELF X APPLY DIMENSIONALITY REDUCTION TO X USING THE MODEL
INIT SELF NCOMPONENTSNONE TOL001 COPYTRUE MAXITER1000
NOISEVARIANCEINITNONE SVDMETHOD'RANDOMIZED' ITERATEDPOWER3 RAN
DOMSTATE0
FITSELFXYNONE
FIT THE FACTORANALYSIS MODEL TO X USING EM
PARAMETERS
XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA
YIGNORED
RETURNS
SELF
FITTRANSFORM SELFXYNONE FITPARAMS
FIT TO DATA THEN TRANSFORM IT
FITS TRANSFORMER TO X AND Y WITH OPTIONAL PARAMETERS FITPARAMS AND RETURNS A TRANSFORMED VERSION OF X
PARAMETERS
XNUMPY ARRAY OF SHAPE NSAMPLES NFEATURES TRAINING SET
YNUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES
RETURNS
XNEW NUMPY ARRAY OF SHAPE NSAMPLES NFEATURESNEW TRANSFORMED ARRAY
GETCOVARIANCE SELF
COMPUTE DATA COVARIANCE WITH THE FACTORANALYSIS MODEL
COV COMPONENTSTST COMPONENTS DIAGNOISEVARIANCE
69SKLEARNDECOMPOSITION MATRIX DECOMPOSITION 1627
```

SCIKITLEARN USER GUIDE RELEASE 0213

RETURNS

COV ARRAY SHAPE NFEATURES NFEATURES ESTIMATED COVARIANCE OF DATA

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

GETPRECISION SELF

COMPUTE DATA PRECISION MATRIX WITH THE FACTORANALYSIS MODEL

RETURNS

PRECISION ARRAY SHAPE NFEATURES NFEATURES ESTIMATED PRECISION OF DATA

SCORESELFXYNONE

COMPUTE THE AVERAGE LOGLIKELIHOOD OF THE SAMPLES

PARAMETERS

XARRAY SHAPE NSAMPLES NFEATURES THE DATA

YIGNORED

RETURNS

LLFLOAT AVERAGE LOGLIKELIHOOD OF THE SAMPLES UNDER THE CURRENT MODEL

SCORESAMPLES SELFXY

COMPUTE THE LOGLIKELIHOOD OF EACH SAMPLE

PARAMETERS

XARRAY SHAPE NSAMPLES NFEATURES THE DATA

RETURNS

LLARRAY SHAPE NSAMPLES LOGLIKELIHOOD OF EACH SAMPLE UNDER THE CURRENT MODEL

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE

PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT

OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELFXY

APPLY DIMENSIONALITY REDUCTION TO X USING THE MODEL

COMPUTE THE EXPECTED MEAN OF THE LATENT VARIABLES SEE BARBER 21233 OR BISHOP 1266

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA

RETURNS

1628 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

XNEW ARRAYLIKE SHAPE NSAMPLES NCOMPONENTS THE LATENT VARIABLES OF X

EXAMPLES USING SKLEARNDECOMPOSITIONFACTORANALYSIS

- MODEL SELECTION WITH PROBABILISTIC PCA AND FACTOR ANALYSIS FA
- FACES DATASET DECOMPOSITIONS

693SKLEARNDECOMPOSITION FASTICA

CLASSSSKLEARNDECOMPOSITION FASTICANCOMPONENTSNONE ALGORITHM'PARALLEL' WHITENTRUE

FUN'LOGCOSH' FUNARGSNONE MAXITER200 TOL00001

WINITNONE RANDOMSTATENONE

FASTICA A FAST ALGORITHM FOR INDEPENDENT COMPONENT ANALYSIS

READ MORE IN THE USER GUIDE

PARAMETERS

NCOMPONENTS INT OPTIONAL NUMBER OF COMPONENTS TO USE IF NONE IS PASSED ALL ARE USED

ALGORITHM 'PARALLEL' 'DEFLATION' APPLY PARALLEL OR DEFLATIONAL ALGORITHM FOR FASTICA

WHITEN BOOLEAN OPTIONAL IF WHITEN IS FALSE THE DATA IS ALREADY CONSIDERED TO BE WHITENED AND NO WHITENING IS PERFORMED

FUN STRING OR FUNCTION OPTIONAL DEFAULT 'LOGCOSH' THE FUNCTIONAL FORM OF THE G FUNCTION USED IN THE APPROXIMATION TO NEGENTROPY COULD BE EITHER 'LOGCOSH' 'EXP' OR 'CUBE' YOU CAN ALSO PROVIDE YOUR OWN FUNCTION IT SHOULD RETURN A TUPLE CONTAINING THE VALUE OF THE FUNCTION AND OF ITS DERIVATIVE IN THE POINT EXAMPLE

DEF MYGX RETURN X 3 3 X 2MEANAXIS1

FUNARGS DICTIONARY OPTIONAL ARGUMENTS TO SEND TO THE FUNCTIONAL FORM IF EMPTY AND IF FUN'LOGCOSH' FUNARGS WILL TAKE VALUE 'ALPHA' 10

MAXITER INT OPTIONAL MAXIMUM NUMBER OF ITERATIONS DURING FIT

TOLFLOAT OPTIONAL TOLERANCE ON UPDATE AT EACH ITERATION

WINIT NONE OF AN NCOMPONENTS NCOMPONENTS NDARRAY THE MIXING MATRIX TO BE USED TO INITIALIZE THE ALGORITHM

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RAN

DOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE

RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NPRANDOM

ATTRIBUTES

COMPONENTS 2D ARRAY SHAPE NCOMPONENTS NFEATURES THE UNMIXING MATRIX

MIXING ARRAY SHAPE NFEATURES NCOMPONENTS THE MIXING MATRIX

NITER INT IF THE ALGORITHM IS "DEFLATION" NITER IS THE MAXIMUM NUMBER OF ITERATIONS RUN ACROSS ALL COMPONENTS ELSE THEY ARE JUST THE NUMBER OF ITERATIONS TAKEN TO CONVERGE

69SKLEARNDECOMPOSITION MATRIX DECOMPOSITION 1629

SCIKITLEARN USER GUIDE RELEASE 0213

NOTES

IMPLEMENTATION BASED ON A HYVARINEN AND E OJA INDEPENDENT COMPONENT ANALYSIS ALGORITHMS AND APPLI  
CATIONS NEURAL NETWORKS 1345 2000 PP 411430

EXAMPLES

FROM SKLEARNDATASETS IMPORT LOADDIGITS  
FROM SKLEARNDECOMPOSITION IMPORT FASTICA  
X LOADDIGITSRETURNXY TRUE  
TRANSFORMER FASTICANCOMPONENTS7  
RANDOMSTATE0  
XTRANSFORMED TRANSFORMERFITTRANSFORMX  
XTRANSFORMEDSHAPE  
1797 7

METHODS

FITSELF X Y FIT THE MODEL TO X  
FITTRANSFORM SELF X Y FIT THE MODEL AND RECOVER THE SOURCES FROM X  
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR  
INVERSETRANSFORM SELF X COPY TRANSFORM THE SOURCES BACK TO THE MIXED DATA APPLY  
MIXING MATRIX  
SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR  
TRANSFORM SELF X COPY RECOVER THE SOURCES FROM X APPLY THE UNMIXING MA  
TRIX  
INIT SELFNCOMPONENTSNONE ALGORITHM'PARALLEL' WHITENTRUE FUN'LOGCOSH'  
FUNARGSNONE MAXITER200 TOL00001 WINITNONE RANDOMSTATENONE  
FITSELFXYNONE  
FIT THE MODEL TO X  
PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA WHERE NSAMPLES IS THE NUMBER  
OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES  
YIGNORED  
RETURNS  
SELF  
FITTRANSFORM SELFXYNONE  
FIT THE MODEL AND RECOVER THE SOURCES FROM X  
PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA WHERE NSAMPLES IS THE NUMBER  
OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES  
YIGNORED  
RETURNS

1630 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

XNEW ARRAYLIKE SHAPE NSAMPLES NCOMPONENTS

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

INVERSETRANSFORM SELFXCOPYTRUE

TRANSFORM THE SOURCES BACK TO THE MIXED DATA APPLY MIXING MATRIX

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NCOMPONENTS SOURCES WHERE NSAMPLES IS THE NUMBER

OF SAMPLES AND NCOMPONENTS IS THE NUMBER OF COMPONENTS

COPY BOOL OPTIONAL IF FALSE DATA PASSED TO FIT ARE OVERWRITTEN DEFAULTS TO TRUE

RETURNS

XNEW ARRAYLIKE SHAPE NSAMPLES NFEATURES

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE

PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT

OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELFXCOPYTRUE

RECOVER THE SOURCES FROM X APPLY THE UNMIXING MATRIX

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES DATA TO TRANSFORM WHERE NSAMPLES IS THE

NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

COPY BOOL OPTIONAL IF FALSE DATA PASSED TO FIT ARE OVERWRITTEN DEFAULTS TO TRUE

RETURNS

XNEW ARRAYLIKE SHAPE NSAMPLES NCOMPONENTS

EXAMPLES USING SKLEARNDECOMPOSITIONFASTICA

- BLIND SOURCE SEPARATION USING FASTICA
- FASTICA ON 2D POINT CLOUDS
- FACES DATASET DECOMPOSITIONS

69SKLEARNDECOMPOSITION MATRIX DECOMPOSITION 1631

SCIKITLEARN USER GUIDE RELEASE 0213

6945KLEARNDECOMPOSITION INCREMENTALPCA

CLASSSSKLEARNDECOMPOSITION INCREMENTALPCA NCOMPONENTSNONE WHITENFALSE COPYTRUE

BATCHSIZENONE

INCREMENTAL PRINCIPAL COMPONENTS ANALYSIS IPCA

LINEAR DIMENSIONALITY REDUCTION USING SINGULAR VALUE DECOMPOSITION OF THE DATA KEEPING ONLY THE MOST SIGNIF

ICANT SINGULAR VECTORS TO PROJECT THE DATA TO A LOWER DIMENSIONAL SPACE THE INPUT DATA IS CENTERED BUT NOT SCALED

FOR EACH FEATURE BEFORE APPLYING THE SVD

DEPENDING ON THE SIZE OF THE INPUT DATA THIS ALGORITHM CAN BE MUCH MORE MEMORY EFFICIENT THAN A PCA

THIS ALGORITHM HAS CONSTANT MEMORY COMPLEXITY ON THE ORDER OF BATCHSIZE ENABLING USE OF NPMEMMAP

FILES WITHOUT LOADING THE ENTIRE FILE INTO MEMORY

THE COMPUTATIONAL OVERHEAD OF EACH SVD IS OBATCHSIZE NFEATURES 2 BUT ONLY 2

BATCHSIZE SAMPLES REMAIN IN MEMORY AT A TIME THERE WILL BE NSAMPLES BATCHSIZE SVD COMPU

TATIONS TO GET THE PRINCIPAL COMPONENTS VERSUS 1 LARGE SVD OF COMPLEXITY ONSAMPLES NFEATURES

2FOR PCA

READ MORE IN THE USER GUIDE

PARAMETERS

NCOMPONENTS INT OR NONE DEFAULTNONE NUMBER OF COMPONENTS TO KEEP IF

NCOMPONENTS IS NONE THENNCOMPONENTS IS SET TOMINNSAMPLES

NFEATURES

WHITEN BOOL OPTIONAL WHEN TRUE FALSE BY DEFAULT THE COMPONENTS VECTORS ARE DI

VIDED BYNSAMPLES TIMESCOMPONENTS TO ENSURE UNCORRELATED OUTPUTS WITH UNIT

COMPONENTWISE VARIANCES

WHITENING WILL REMOVE SOME INFORMATION FROM THE TRANSFORMED SIGNAL THE RELATIVE VARIANCE

SCALES OF THE COMPONENTS BUT CAN SOMETIMES IMPROVE THE PREDICTIVE ACCURACY OF THE DOWN

STREAM ESTIMATORS BY MAKING DATA RESPECT SOME HARDWIRED ASSUMPTIONS

COPY BOOL DEFAULTTRUE IF FALSE X WILL BE OVERWRITTEN COPYFALSE CAN BE USED TO SAVE

MEMORY BUT IS UNSAFE FOR GENERAL USE

BATCHSIZE INT OR NONE DEFAULTNONE THE NUMBER OF SAMPLES TO USE FOR EACH BATCH ONLY

USED WHEN CALLING FIT IFBATCHSIZE ISNONE THENBATCHSIZE IS INFERRED FROM THE

DATA AND SET TO 5NFEATURES TO PROVIDE A BALANCE BETWEEN APPROXIMATION ACCURACY

AND MEMORY CONSUMPTION

ATTRIBUTES

COMPONENTS ARRAY SHAPE NCOMPONENTS NFEATURES COMPONENTS WITH MAXIMUM VARIANCE

EXPLAINEDVARIANCE ARRAY SHAPE NCOMPONENTS VARIANCE EXPLAINED BY EACH OF THE SE

LECTED COMPONENTS

EXPLAINEDVARIANCERATIO ARRAY SHAPE NCOMPONENTS PERCENTAGE OF VARIANCE EXPLAINED

BY EACH OF THE SELECTED COMPONENTS IF ALL COMPONENTS ARE STORED THE SUM OF EXPLAINED

VARIANCES IS EQUAL TO 10

SINGULARVALUES ARRAY SHAPE NCOMPONENTS THE SINGULAR VALUES CORRESPONDING TO EACH

OF THE SELECTED COMPONENTS THE SINGULAR VALUES ARE EQUAL TO THE 2NORMS OF THE

NCOMPONENTS VARIABLES IN THE LOWERDIMENSIONAL SPACE

MEAN ARRAY SHAPE NFEATURES PERFEATURE EMPIRICAL MEAN AGGREGATE OVER CALLS TO

PARTIALFIT

1632 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

VAR ARRAY SHAPE NFEATURES PERFEATURE EMPIRICAL VARIANCE AGGREGATE OVER CALLS TO PARTIALFIT

NOISEVARIANCE FLOAT THE ESTIMATED NOISE COVARIANCE FOLLOWING THE PROBABILISTIC PCA MODEL FROM TIPPING AND BISHOP 1999 SEE “PATTERN RECOGNITION AND MACHINE LEARNING” BY C BISHOP 1221 P 574 OR [HTTPWWWMIKETIPPINGCOMPAPERSMETMPPCAPDF](http://www.miketipping.com/papers/metmppca.pdf)

NCOMPONENTS INT THE ESTIMATED NUMBER OF COMPONENTS RELEVANT WHEN NCOMPONENTSNONE

NSAMPLESSEEN INT THE NUMBER OF SAMPLES PROCESSED BY THE ESTIMATOR WILL BE RESET ON NEW CALLS TO FIT BUT INCREMENTS ACROSS PARTIALFIT CALLS

SEE ALSO

PCA

KERNELPCA

SPARSEPCA

TRUNCATEDSVD

NOTES

IMPLEMENTS THE INCREMENTAL PCA MODEL FROM D ROSS J LIM R LIN M YANG INCREMENTAL LEARNING FOR ROBUST VISUAL TRACKING INTERNATIONAL JOURNAL OF COMPUTER VISION VOLUME 77 ISSUE 13 PP 125141 MAY 2008 SEE [HTTPSWWWCSSTORONTOEDUDROSSIVTROSSLIMLINYANGIJCVPDF](https://www.cs.toronto.edu/drossiv/trosslimlinyangijcv.pdf)

THIS MODEL IS AN EXTENSION OF THE SEQUENTIAL KARHUNENLOEVE TRANSFORM FROM A LEVY AND M LINDENBAUM SEQUENTIAL KARHUNENLOEVE BASIS EXTRACTION AND ITS APPLICATION TO IMAGES IEEE TRANSACTIONS ON IMAGE PROCESSING VOLUME 9 NUMBER 8 PP 13711374 AUGUST 2000 SEE [HTTPSWWWCSTECHNIONACILMICDOCSKLIPPDF](http://www.cse.cmu.edu/~cmv/papers/levy_lindenbaum_klpp.pdf)

WE HAVE SPECIFICALLY ABSTAINED FROM AN OPTIMIZATION USED BY AUTHORS OF BOTH PAPERS A QR DECOMPOSITION USED IN SPECIFIC SITUATIONS TO REDUCE THE ALGORITHMIC COMPLEXITY OF THE SVD THE SOURCE FOR THIS TECHNIQUE IS MATRIX COMPUTATIONS THIRD EDITION G GOLUB AND C VAN LOAN CHAPTER 5 SECTION 544 PP 252253 THIS TECHNIQUE HAS BEEN OMITTED BECAUSE IT IS ADVANTAGEOUS ONLY WHEN DECOMPOSING A MATRIX WITH NSAMPLES ROWS 53 NFEATURES COLUMNS AND HURTS THE READABILITY OF THE IMPLEMENTED ALGORITHM THIS WOULD BE A GOOD OPPORTUNITY FOR FUTURE OPTIMIZATION IF IT IS DEEMED NECESSARY

REFERENCES

D ROSS J LIM R LIN M YANG INCREMENTAL LEARNING FOR ROBUST VISUAL TRACKING INTERNATIONAL JOURNAL OF COMPUTER VISION V OLUME 77 ISSUE 13 PP 125141 MAY 2008

G GOLUB AND C VAN LOAN MATRIX COMPUTATIONS THIRD EDITION CHAPTER 5 SECTION 544 PP 252253

EXAMPLES

FROM SKLEARNDATASETS IMPORT LOADDIGITS

FROM SKLEARNDECOMPOSITION IMPORT INCREMENTALPCA

X LOADDIGITSRETURNXY TRUE

TRANSFORMER INCREMENTALPCANCOMPONENTS7 BATCHSIZE200

EITHER PARTIALLY FIT ON SMALLER BATCHES OF DATA

TRANSFORMERPARTIALFITX100

INCREMENTALPCABATCHSIZE200 COPYTRUE NCOMPONENTS7 WHITENFALSE

69SKLEARNDECOMPOSITION MATRIX DECOMPOSITION 1633

SCIKITLEARN USER GUIDE RELEASE 0213

OR LET THE FIT FUNCTION ITSELF DIVIDE THE DATA INTO BATCHES

XTRANSFORMED TRANSFORMERFITTRANSFORMX

XTRANSFORMEDSHAPE

1797 7

METHODS

FITSELF X Y FIT THE MODEL WITH X USING MINIBATCHES OF SIZE

BATCHSIZE

FITTRANSFORM SELF X Y FIT TO DATA THEN TRANSFORM IT

GETCOVARIANCE SELF COMPUTE DATA COVARIANCE WITH THE GENERATIVE MODEL

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

GETPRECISION SELF COMPUTE DATA PRECISION MATRIX WITH THE GENERATIVE

MODEL

INVERSETRANSFORM SELF X TRANSFORM DATA BACK TO ITS ORIGINAL SPACE

PARTIALFIT SELF X Y CHECKINPUT INCREMENTAL FIT WITH X

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF X APPLY DIMENSIONALITY REDUCTION TO X

INIT SELFNCOMPONENTSNONE WHITENFALSE COPYTRUE BATCHSIZENONE

FITSELFXYNONE

FIT THE MODEL WITH X USING MINIBATCHES OF SIZE BATCHSIZE

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA WHERE NSAMPLES IS THE NUMBER

OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

YIGNORED

RETURNS

SELF OBJECT RETURNS THE INSTANCE ITSELF

FITTRANSFORM SELFXYNONE FITPARAMS

FIT TO DATA THEN TRANSFORM IT

FITS TRANSFORMER TO X AND Y WITH OPTIONAL PARAMETERS FITPARAMS AND RETURNS A TRANSFORMED VERSION OF X

PARAMETERS

XNUMPY ARRAY OF SHAPE NSAMPLES NFEATURES TRAINING SET

YNUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES

RETURNS

XNEW NUMPY ARRAY OF SHAPE NSAMPLES NFEATURESNEW TRANSFORMED ARRAY

GETCOVARIANCE SELF

COMPUTE DATA COVARIANCE WITH THE GENERATIVE MODEL

COV COMPONENTST S2COMPONENTS SIGMA2 EYENFEATURES

WHERE S2 CONTAINS THE EXPLAINED VARIANCES AND SIGMA2 CONTAINS THE NOISE VARIANCES

RETURNS

COV ARRAY SHAPENFEATURES NFEATURES ESTIMATED COVARIANCE OF DATA

1634 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

GETPARAMS SELFDEEPTURE  
GET PARAMETERS FOR THIS ESTIMATOR  
PARAMETERS  
DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBOBJECTS THAT ARE ESTIMATORS  
RETURNS  
PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES  
GETPRECISION SELF  
COMPUTE DATA PRECISION MATRIX WITH THE GENERATIVE MODEL  
EQUALS THE INVERSE OF THE COVARIANCE BUT COMPUTED WITH THE MATRIX INVERSION LEMMA FOR EFFICIENCY  
RETURNS  
PRECISION ARRAY SHAPENFEATURES NFEATURES ESTIMATED PRECISION OF DATA  
INVERSETRANSFORM SELF  
TRANSFORM DATA BACK TO ITS ORIGINAL SPACE  
IN OTHER WORDS RETURN AN INPUT XORIGINAL WHOSE TRANSFORM WOULD BE X  
PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NCOMPONENTS NEW DATA WHERE NSAMPLES IS THE NUMBER  
OF SAMPLES AND NCOMPONENTS IS THE NUMBER OF COMPONENTS  
RETURNS  
XORIGINAL ARRAYLIKE SHAPE NSAMPLES NFEATURES  
NOTES  
IF WHITENING IS ENABLED INVERSETRANSFORM WILL COMPUTE THE EXACT INVERSE OPERATION WHICH INCLUDES RE  
VERSING WHITENING  
PARTIALFIT SELFXYNONE CHECKINPUTTRUE  
INCREMENTAL FIT WITH X ALL OF X IS PROCESSED AS A SINGLE BATCH  
PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA WHERE NSAMPLES IS THE NUMBER  
OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES  
CHECKINPUT BOOL RUN CHECKARRAY ON X  
YIGNORED  
RETURNS  
SELF OBJECT RETURNS THE INSTANCE ITSELF  
SETPARAMS SELFPARAMS  
SET THE PARAMETERS OF THIS ESTIMATOR  
THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE  
PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT  
OF A NESTED OBJECT  
RETURNS  
69SKLEARNDECOMPOSITION MATRIX DECOMPOSITION 1635

SCIKITLEARN USER GUIDE RELEASE 0213

SELF

TRANSFORM SELF

APPLY DIMENSIONALITY REDUCTION TO X

X IS PROJECTED ON THE FIRST PRINCIPAL COMPONENTS PREVIOUSLY EXTRACTED FROM A TRAINING SET

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES NEW DATA WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

RETURNS

XNEW ARRAYLIKE SHAPE NSAMPLES NCOMPONENTS

EXAMPLES

```
import numpy as np
from sklearn.decomposition import IncrementalPCA
X = np.array([1, 2, 1, 3, 2, 1, 1, 2, 1, 3, 2])
ipca = IncrementalPCA(n_components=2, batch_size=3)
ipca.fit(X)
IncrementalPCABatchSize3(copy=True, n_components=2, whiten=False)
ipca.transform(X)
```

EXAMPLES USING SKLEARNDECOMPOSITIONINCREMENTALPCA

- INCREMENTAL PCA

```
from sklearn.decomposition import KernelPCA
class SKLearnDecomposition(KernelPCA):
    n_components = None
    kernel = 'linear'
    gamma = None
    degree = 3
    coef0 = 1
    kernel_params = None
    alpha = 10
    fit_inverse_transform = False
    eigensolver = 'auto'
    tol = 0
    max_iter = None
    remove_zero_eig = False
    random_state = None
    dom = 'stat'
    copy_x = True
    n_jobs = None
    kernel_principal_component_analysis = 'kPCA'
    nonlinear_dimensionality_reduction_through_the_use_of_kernels = True
    pairwise_metrics_affinities_and_kernels = True
```

READ MORE IN THE USER GUIDE

PARAMETERS

NCOMPONENTS INT DEFAULTNONE NUMBER OF COMPONENTS IF NONE ALL NONZERO COMPONENTS ARE KEPT

KERNEL “LINEAR” “POLY” “RBF” “SIGMOID” “COSINE” “PRECOMPUTED” KERNEL DE FAULT”LINEAR”

GAMMA FLOAT DEFAULT1NFEATURES KERNEL COEFFICIENT FOR RBF POLY AND SIGMOID KERNELS IGNORED BY OTHER KERNELS

DEGREE INT DEFAULT3 DEGREE FOR POLY KERNELS IGNORED BY OTHER KERNELS

COEF0 FLOAT DEFAULT1 INDEPENDENT TERM IN POLY AND SIGMOID KERNELS IGNORED BY OTHER KERNELS

1636 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

KERNELPARAMS MAPPING OF STRING TO ANY DEFAULTNONE PARAMETERS KEYWORD ARGUMENTS AND VALUES FOR KERNEL PASSED AS CALLABLE OBJECT IGNORED BY OTHER KERNELS

ALPHA INT DEFAULT10 HYPERPARAMETER OF THE RIDGE REGRESSION THAT LEARNS THE INVERSE TRANSFORM WHEN FITINVERSESTRANSFORMTRUE

FITINVERSESTRANSFORM BOOL DEFAULTFALSE LEARN THE INVERSE TRANSFORM FOR NONPRECOMPUTED KERNELS IE LEARN TO FIND THE PREIMAGE OF A POINT

EIGENSOLVER STRING 'AUTO''DENSE''ARPACK' DEFAULT'AUTO' SELECT EIGENSOLVER TO USE IF NCOMPONENTS IS MUCH LESS THAN THE NUMBER OF TRAINING SAMPLES ARPACK MAY BE MORE EFFICIENT THAN THE DENSE EIGENSOLVER

TOLFLOAT DEFAULT0 CONVERGENCE TOLERANCE FOR ARPACK IF 0 OPTIMAL VALUE WILL BE CHOSEN BY ARPACK

MAXITER INT DEFAULTNONE MAXIMUM NUMBER OF ITERATIONS FOR ARPACK IF NONE OPTIMAL VALUE WILL BE CHOSEN BY ARPACK

REMOVEZEROEIG BOOLEAN DEFAULTFALSE IF TRUE THEN ALL COMPONENTS WITH ZERO EIGENVALUES ARE REMOVED SO THAT THE NUMBER OF COMPONENTS IN THE OUTPUT MAY BE NCOMPONENTS AND SOMETIMES EVEN ZERO DUE TO NUMERICAL INSTABILITY WHEN NCOMPONENTS IS NONE THIS PARAMETER IS IGNORED AND COMPONENTS WITH ZERO EIGENVALUES ARE REMOVED REGARDLESS

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RANDOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NPRANDOM USED WHEN EIGENSOLVER 'ARPACK'

NEW IN VERSION 018

COPYX BOOLEAN DEFAULTTRUE IF TRUE INPUT X IS COPIED AND STORED BY THE MODEL IN THE XFIT ATTRIBUTE IF NO FURTHER CHANGES WILL BE DONE TO X SETTING COPYXFALSE SAVES MEMORY BY STORING A REFERENCE

NEW IN VERSION 018

NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF PARALLEL JOBS TO RUN NONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

NEW IN VERSION 018

ATTRIBUTES

LAMBDAS ARRAY NCOMPONENTS EIGENVALUES OF THE CENTERED KERNEL MATRIX IN DECREASING ORDER IF NCOMPONENTS ANDREMOVEZEROEIG ARE NOT SET THEN ALL VALUES ARE STORED

ALPHAS ARRAY NSAMPLES NCOMPONENTS EIGENVECTORS OF THE CENTERED KERNEL MATRIX IF NCOMPONENTS ANDREMOVEZEROEIG ARE NOT SET THEN ALL COMPONENTS ARE STORED

DUALCOEF ARRAY NSAMPLES NFEATURES INVERSE TRANSFORM MATRIX ONLY AVAILABLE WHEN FITINVERSESTRANSFORM IS TRUE

XTRANSFORMEDFIT ARRAY NSAMPLES NCOMPONENTS PROJECTION OF THE FITTED DATA ON THE KERNEL PRINCIPAL COMPONENTS ONLY AVAILABLE WHEN FITINVERSESTRANSFORM IS TRUE

XFIT NSAMPLES NFEATURES THE DATA USED TO FIT THE MODEL IF COPYXFALSE THEN XFIT IS A REFERENCE THIS ATTRIBUTE IS USED FOR THE CALLS TO TRANSFORM

69SKLEARNDECOMPOSITION MATRIX DECOMPOSITION 1637

SCIKITLEARN USER GUIDE RELEASE 0213

REFERENCES

KERNEL PCA WAS INTRODUCED IN BERNHARD SCHOELKOPF ALEXANDER J SMOLA AND KLAUSROBERT MUELLER 1999

KERNEL PRINCIPAL COMPONENT ANALYSIS IN ADVANCES IN KERNEL METHODS MIT PRESS CAMBRIDGE MA USA

327352

EXAMPLES

FROM SKLEARNDATASETS IMPORT LOADDIGITS

FROM SKLEARNDECOMPOSITION IMPORT KERNELPCA

X LOADDIGITSRETURNXY TRUE

TRANSFORMER KERNELPCANCOMPONENTS7 KERNELLINEAR

XTRANSFORMED TRANSFORMERFITTRANSFORMX

XTRANSFORMEDSHAPE

1797 7

METHODS

FITSELF X Y FIT THE MODEL FROM DATA IN X

FITTRANSFORM SELF X Y FIT THE MODEL FROM DATA IN X AND TRANSFORM X

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

INVERSETRANSFORM SELF X TRANSFORM X BACK TO ORIGINAL SPACE

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF X TRANSFORM X

INIT SELFNCOMPONENTSNONE KERNEL'LINEAR' GAMMANONE DEGREE3 COEF01 KER

NELPARAMSNONE ALPHA10 FITINVERSETRANSFORMFALSE EIGENSOLVER'AUTO'

TOLO MAXITERNONE REMOVEZEROEIGFALSE RANDOMSTATENONE COPYXTRUE

NJOBSNONE

FITSELFXYNONE

FIT THE MODEL FROM DATA IN X

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING VECTOR WHERE NSAMPLES IN THE NUM

BER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

RETURNS

SELF OBJECT RETURNS THE INSTANCE ITSELF

FITTRANSFORM SELFXYNONE PARAMS

FIT THE MODEL FROM DATA IN X AND TRANSFORM X

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING VECTOR WHERE NSAMPLES IN THE NUM

BER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

RETURNS

XNEW ARRAYLIKE SHAPE NSAMPLES NCOMPONENTS

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

1638 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

INVERSETRANSFORM SELF

TRANSFORM X BACK TO ORIGINAL SPACE

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NCOMPONENTS

RETURNS

XNEW ARRAYLIKE SHAPE NSAMPLES NFEATURES

REFERENCES

“LEARNING TO FIND PREIMAGES” G BAKIR ET AL 2004

SETPARAMS SELF

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELF

TRANSFORM X

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES

RETURNS

XNEW ARRAYLIKE SHAPE NSAMPLES NCOMPONENTS

EXAMPLES USING SKLEARNDECOMPOSITIONKERNELPCA

- KERNEL PCA

69SKLEARNDECOMPOSITION MATRIX DECOMPOSITION 1639

SCIKITLEARN USER GUIDE RELEASE 0213  
696SKLEARNDECOMPOSITION LATENTDIRICHLETALLOCATION  
CLASSSSKLEARNDECOMPOSITION LATENTDIRICHLETALLOCATION NCOMPONENTS10  
DOCTOPICPRIORNONE  
TOPICWORDPRIORNONE  
LEARNINGMETHOD'BATCH'  
LEARNINGDECAY07  
LEARNINGOFFSET100  
MAXITER10 BATCHSIZE128  
EVALUATEEVERY1 TO  
TALSAMPLES10000000  
PERPTOL01  
MEANCHANGETOL0001  
MAXDOCUPDATEITER100  
NJOBSNONE VERBOSE0 RAN  
DOMSTATENONE  
LATENT DIRICHLET ALLOCATION WITH ONLINE VARIATIONAL BAYES ALGORITHM  
NEW IN VERSION 017  
READ MORE IN THE USER GUIDE  
PARAMETERS  
NCOMPONENTS INT OPTIONAL DEFAULT10 NUMBER OF TOPICS  
DOCTOPICPRIOR FLOAT OPTIONAL DEFAULTNONE PRIOR OF DOCUMENT TOPIC DISTRIBUTION THETA  
IF THE VALUE IS NONE DEFAULTS TO 1 NCOMPONENTS IN1 THIS IS CALLED ALPHA  
TOPICWORDPRIOR FLOAT OPTIONAL DEFAULTNONE PRIOR OF TOPIC WORD DISTRIBUTION BETA IF  
THE VALUE IS NONE DEFAULTS TO 1 NCOMPONENTS IN1 THIS IS CALLED ETA  
LEARNINGMETHOD 'BATCH' 'ONLINE' DEFAULT'BATCH' METHOD USED TO UPDATE COMPONENT  
ONLY USED IN FIT METHOD IN GENERAL IF THE DATA SIZE IS LARGE THE ONLINE UPDATE WILL BE  
MUCH FASTER THAN THE BATCH UPDATE  
VALID OPTIONS  
BATCH BATCH VARIATIONAL BAYES METHOD USE ALL TRAINING DATA IN  
EACH EM UPDATE  
OLD COMPONENTS WILL BE OVERWRITTEN IN EACH ITERATION  
ONLINE ONLINE VARIATIONAL BAYES METHOD IN EACH EM UPDATE USE  
MINIBATCH OF TRAINING DATA TO UPDATE THE COMPONENTS  
VARIABLE INCREMENTALLY THE LEARNING RATE IS CONTROLLED BY THE  
LEARNINGDECAY AND THE LEARNINGOFFSET PARAMETERS  
CHANGED IN VERSION 020 THE DEFAULT LEARNING METHOD IS NOW BATCH  
LEARNINGDECAY FLOAT OPTIONAL DEFAULT07 IT IS A PARAMETER THAT CONTROL LEARNING RATE IN THE  
ONLINE LEARNING METHOD THE VALUE SHOULD BE SET BETWEEN 05 10 TO GUARANTEE ASYMPTOTIC  
CONVERGENCE WHEN THE VALUE IS 00 AND BATCHSIZE IS NSAMPLES THE UPDATE METHOD IS  
SAME AS BATCH LEARNING IN THE LITERATURE THIS IS CALLED KAPPA  
LEARNINGOFFSET FLOAT OPTIONAL DEFAULT10 A POSITIVE PARAMETER THAT DOWNWEIGHTS EARLY  
ITERATIONS IN ONLINE LEARNING IT SHOULD BE GREATER THAN 10 IN THE LITERATURE THIS IS CALLED  
TAU0  
MAXITER INTEGER OPTIONAL DEFAULT10 THE MAXIMUM NUMBER OF ITERATIONS  
BATCHSIZE INT OPTIONAL DEFAULT128 NUMBER OF DOCUMENTS TO USE IN EACH EM ITERATION  
ONLY USED IN ONLINE LEARNING  
1640 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

EVALUATEEVERY INT OPTIONAL DEFAULT0 HOW OFTEN TO EVALUATE PERPLEXITY ONLY USED IN FIT METHOD SET IT TO 0 OR NEGATIVE NUMBER TO NOT EVALUTE PERPLEXITY IN TRAINING AT ALL EVALUATING PERPLEXITY CAN HELP YOU CHECK CONVERGENCE IN TRAINING PROCESS BUT IT WILL ALSO INCREASE TOTAL TRAINING TIME EVALUATING PERPLEXITY IN EVERY ITERATION MIGHT INCREASE TRAINING TIME UP TO TWOFOLD

TOTALSAMPLES INT OPTIONAL DEFAULT1E6 TOTAL NUMBER OF DOCUMENTS ONLY USED IN THE PARTIALFIT METHOD

PERTOL FLOAT OPTIONAL DEFAULT1E1 PERPLEXITY TOLERANCE IN BATCH LEARNING ONLY USED WHEN EVALUATEEVERY IS GREATER THAN 0

MEANCHANGETOL FLOAT OPTIONAL DEFAULT1E3 STOPPING TOLERANCE FOR UPDATING DOCUMENT TOPIC DISTRIBUTION IN ESTEP

MAXDOCUPDATEITER INT DEFAULT100 MAX NUMBER OF ITERATIONS FOR UPDATING DOCUMENT TOPIC DISTRIBUTION IN THE ESTEP

NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF JOBS TO USE IN THE ESTEP NONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL PROCES

SORS SEE GLOSSARY FOR MORE DETAILS

VERBOSE INT OPTIONAL DEFAULT0 VERBOSITY LEVEL

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RAN DOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NPRANDOM

ATTRIBUTES

COMPONENTS ARRAY NCOMPONENTS NFEATURES VARIATIONAL PARAMETERS FOR TOPIC WORD DIS TRIBUTION SINCE THE COMPLETE CONDITIONAL FOR TOPIC WORD DISTRIBUTION IS A DIRICHLET COMPONENTSI J CAN BE VIEWED AS PSEUDOCOUNT THAT REPRESENTS THE NUMBER OF TIMES WORDJWAS ASSIGNED TO TOPIC I IT CAN ALSO BE VIEWED AS DISTRIBUTION OVER THE WORDS FOR EACH TOPIC AFTER NORMALIZATION MODELCOMPONENTS MODELCOMPONENTS

SUMAXIS1 NPNEWAXIS

NBATCHITER INT NUMBER OF ITERATIONS OF THE EM STEP

NITER INT NUMBER OF PASSES OVER THE DATASET

REFERENCES

1 “ONLINE LEARNING FOR LATENT DIRICHLET ALLOCATION” MATTHEW D HOFFMAN DAVID M BLEI FRANCIS BACH 2010

2 “STOCHASTIC VARIATIONAL INFERENCE” MATTHEW D HOFFMAN DAVID M BLEI CHONG WANG JOHN PAISLEY 2013

3 MATTHEW D HOFFMAN’S ONLINELDAVB CODE LINK HTTPSGITHUBCOMBLEILABONLINELDAVB

EXAMPLES

FROM SKLEARNDECOMPOSITION IMPORT LATENTDIRICHLETALLOCATION

FROM SKLEARNDATASETS IMPORT MAKEMULTILABELCLASSIFICATION

THIS PRODUCES A FEATURE MATRIX OF TOKEN COUNTS SIMILAR TO WHAT COUNTVECTORIZER WOULD PRODUCE ON TEXT

69SKLEARNDECOMPOSITION MATRIX DECOMPOSITION 1641

SCIKITLEARN USER GUIDE RELEASE 0213

X MAKEMULTILABELCLASSIFICATIONRANDOMSTATE0

LDA LATENTDIRICHLETALLOCATIONNNCOMPONENTS5

RANDOMSTATE0

LDAFITX

LATENTDIRICHLETALLOCATION

GET TOPICS FOR SOME GIVEN SAMPLES

LDATRANSFORMX2

ARRAY000360392 025499205 00036211 064236448 009541846

015297572 000362644 044412786 039568399 0003586

METHODS

FITSELF X Y LEARN MODEL FOR THE DATA X WITH VARIATIONAL BAYES

METHOD

FITTRANSFORM SELF X Y FIT TO DATA THEN TRANSFORM IT

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PARTIALFIT SELF X Y ONLINE VB WITH MINIBATCH UPDATE

PERPLEXITY SELF X SUBSAMPLING CALCULATE APPROXIMATE PERPLEXITY FOR DATA X

SCORE SELF X Y CALCULATE APPROXIMATE LOGLIKELIHOOD AS SCORE

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF X TRANSFORM DATA X ACCORDING TO THE FITTED MODEL

INIT SELFNNCOMPONENTS10 DOCTOPICPRIORNONE TOPICWORDPRIORNONE LEARN

INGMETHOD'BATCH' LEARNINGDECAY07 LEARNINGOFFSET100 MAXITER10

BATCHSIZE128 EVALUATEEVERY1 TOTALSAMPLES10000000 PERPTOL01

MEANCHANGETOL0001 MAXDOCUPDATEITER100 NJOBSNONE VERBOSE0 RAN

DOMSTATENONE

FITSELFXYNONE

LEARN MODEL FOR THE DATA X WITH VARIATIONAL BAYES METHOD

WHENLEARNINGMETHOD IS 'ONLINE' USE MINIBATCH UPDATE OTHERWISE USE BATCH UPDATE

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPENSAMPLES NFEATURES DOCUMENT WORD MATRIX

YIGNORED

RETURNS

SELF

FITTRANSFORM SELFXYNONE FITPARAMS

FIT TO DATA THEN TRANSFORM IT

FITS TRANSFORMER TO X AND Y WITH OPTIONAL PARAMETERS FITPARAMS AND RETURNS A TRANSFORMED VERSION OF X

PARAMETERS

XNUMPY ARRAY OF SHAPE NSAMPLES NFEATURES TRAINING SET

YNUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES

RETURNS

XNEW NUMPY ARRAY OF SHAPE NSAMPLES NFEATURESNEW TRANSFORMED ARRAY

1642 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

GETPARAMS SELFDEEPTTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PARTIALFIT SELFXYNONE

ONLINE VB WITH MINIBATCH UPDATE

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPENSAMPLES NFEATURES DOCUMENT WORD MATRIX

YIGNORED

RETURNS

SELF

PERPLEXITY SELFXYSUBSAMPLINGFALSE

CALCULATE APPROXIMATE PERPLEXITY FOR DATA X

PERPLEXITY IS DEFINED AS  $\exp(-\text{LOGLIKELIHOOD PER WORD})$

CHANGED IN VERSION 019 DOCTOPICDISTR ARGUMENT HAS BEEN DEPRECATED AND IS IGNORED BECAUSE USER NO

LONGER HAS ACCESS TO UNNORMALIZED DISTRIBUTION

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX NSAMPLES NFEATURES DOCUMENT WORD MATRIX

SUBSAMPLING BOOL DO SUBSAMPLING OR NOT

RETURNS

SCORE FLOAT PERPLEXITY SCORE

SCORESELFXYNONE

CALCULATE APPROXIMATE LOGLIKELIHOOD AS SCORE

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPENSAMPLES NFEATURES DOCUMENT WORD MATRIX

YIGNORED

RETURNS

SCORE FLOAT USE APPROXIMATE BOUND AS SCORE

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE

PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT

OF A NESTED OBJECT

RETURNS

SELF

69SKLEARNDECOMPOSITION MATRIX DECOMPOSITION 1643

SCIKITLEARN USER GUIDE RELEASE 0213

TRANSFORM SELF X

TRANSFORM DATA X ACCORDING TO THE FITTED MODEL

CHANGED IN VERSION 018 DOCTOPICDISTR IS NOW NORMALIZED

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPENSAMPLES NFEATURES DOCUMENT WORD MATRIX

RETURNS

DOCTOPICDISTR SHAPENSAMPLES NCOMPONENTS DOCUMENT TOPIC DISTRIBUTION FOR X

EXAMPLES USING SKLEARNDECOMPOSITIONLATENTDIRICHLETALLOCATION

- TOPIC EXTRACTION WITH NONNEGATIVE MATRIX FACTORIZATION AND LATENT DIRICHLET ALLOCATION

697SKLEARNDECOMPOSITION MINIBATCHDICTIONARYLEARNING

CLASSSSKLEARNDECOMPOSITION MINIBATCHDICTIONARYLEARNING NCOMPONENTSNONE AL

PHA1 NITER1000

FITALGORITHM'LARS'

NJOBSNONE

BATCHSIZE3 SHUFFLETRUE

DICTINITNONE TRANS

FORMALGORITHM'OMP' TRANS

FORMNNONZEROCOEFNONE

TRANSFORMALPHANONE VER

BOSEFALSE SPLITSIGNFALSE

RANDOMSTATENONE POS

ITIVECODEFALSE POSI

TIVEDICTFALSE

MINIBATCH DICTIONARY LEARNING

FINDS A DICTIONARY A SET OF ATOMS THAT CAN BEST BE USED TO REPRESENT DATA USING A SPARSE CODE

SOLVES THE OPTIMIZATION PROBLEM

UV ARGMIN 05 Y U V 22 ALPHA U 1

UV

WITH VK 2 1 FORALL 0 K NCOMPONENTS

READ MORE IN THE USER GUIDE

PARAMETERS

NCOMPONENTS INT NUMBER OF DICTIONARY ELEMENTS TO EXTRACT

ALPHA FLOAT SPARSITY CONTROLLING PARAMETER

NITER INT TOTAL NUMBER OF ITERATIONS TO PERFORM

FITALGORITHM 'LARS' 'CD' LARS USES THE LEAST ANGLE REGRESSION METHOD TO SOLVE THE LASSO

PROBLEM LINEARMODELLARSPATH CD USES THE COORDINATE DESCENT METHOD TO COMPUTE THE

LASSO SOLUTION LINEARMODELLASSO LARS WILL BE FASTER IF THE ESTIMATED COMPONENTS ARE

SPARSE

1644 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

NJOBS INT OR NONE OPTIONAL DEFAULTNONE NUMBER OF PARALLEL JOBS TO RUN NONE MEANS 1  
UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL PROCESSORS SEE  
GLOSSARY FOR MORE DETAILS

BATCHSIZE INT NUMBER OF SAMPLES IN EACH MINIBATCH

SHUFFLE BOOL WHETHER TO SHUFFLE THE SAMPLES BEFORE FORMING BATCHES

DICTINIT ARRAY OF SHAPE NCOMPONENTS NFEATURES INITIAL VALUE OF THE DICTIONARY FOR WARM  
RESTART SCENARIOS

TRANSFORMALGORITHM ‘LASSOLARS’ ‘LASSOCD’ ‘LARS’ ‘OMP’ ‘THRESHOLD’ ALGORITHM USED TO  
TRANSFORM THE DATA LARS USES THE LEAST ANGLE REGRESSION METHOD LINEARMODELLARSPATH  
LASSOLARS USES LARS TO COMPUTE THE LASSO SOLUTION LASSOCD USES THE COORDINATE DESCENT  
METHOD TO COMPUTE THE LASSO SOLUTION LINEARMODELLASSO LASSOLARS WILL BE FASTER IF  
THE ESTIMATED COMPONENTS ARE SPARSE OMP USES ORTHOGONAL MATCHING PURSUIT TO ESTIMATE  
THE SPARSE SOLUTION THRESHOLD SQUASHES TO ZERO ALL COEFFICIENTS LESS THAN ALPHA FROM THE  
PROJECTION DICTIONARY ‘X’

TRANSFORMMNNONZEROCOEF5 INTO1NFEATURES BY DEFAULT NUMBER OF NONZERO  
COEFFICIENTS TO TARGET IN EACH COLUMN OF THE SOLUTION THIS IS ONLY USED BY  
ALGORITHMMLARS ANDALGORITHMOMP AND IS OVERRIDDEN BY ALPHA IN THE OR  
THOGONAL MATCHING PURSUIT OMP CASE

TRANSFORMALPHA FLOAT 1 BY DEFAULT IF ALGORITHMMLASSOLARS OR  
ALGORITHMMLASSOCD ALPHA IS THE PENALTY APPLIED TO THE L1 NORM IF  
ALGORITHMTHRESHOLD ALPHA IS THE ABSOLUTE VALUE OF THE THRESHOLD BELOW  
WHICH COEFFICIENTS WILL BE SQUASHED TO ZERO IF ALGORITHMOMP ALPHA IS THE  
TOLERANCE PARAMETER THE VALUE OF THE RECONSTRUCTION ERROR TARGETED IN THIS CASE IT OVERRIDES  
NNONZEROCOEF5

VERBOSE BOOL OPTIONAL DEFAULT FALSE TO CONTROL THE VERBOSITY OF THE PROCEDURE

SPLITSIGN BOOL FALSE BY DEFAULT WHETHER TO SPLIT THE SPARSE FEATURE VECTOR INTO THE CONCATENATION  
OF ITS NEGATIVE PART AND ITS POSITIVE PART THIS CAN IMPROVE THE PERFORMANCE OF DOWN  
STREAM CLASSIFIERS

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RANDOMSTATE  
IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE  
RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS  
THE RANDOMSTATE INSTANCE USED BY NRANDOM

POSITIVECODE BOOL WHETHER TO ENFORCE POSITIVITY WHEN FINDING THE CODE  
NEW IN VERSION 020

POSITIVEDICT BOOL WHETHER TO ENFORCE POSITIVITY WHEN FINDING THE DICTIONARY  
NEW IN VERSION 020

ATTRIBUTES

COMPONENTS ARRAY NCOMPONENTS NFEATURES COMPONENTS EXTRACTED FROM THE DATA

INNERSTATS TUPLE OF A B NDARRAYS INTERNAL SUFFICIENT STATISTICS THAT ARE KEPT BY THE ALGORITHM  
KEEPING THEM IS USEFUL IN ONLINE SETTINGS TO AVOID LOOSING THE HISTORY OF THE EVOLUTION  
BUT THEY SHOULDN’T HAVE ANY USE FOR THE END USER A NCOMPONENTS NCOMPONENTS IS  
THE DICTIONARY COVARIANCE MATRIX B NFEATURES NCOMPONENTS IS THE DATA APPROXIMATION  
MATRIX

NITER INT NUMBER OF ITERATIONS RUN

SEE ALSO

69SKLEARNDECOMPOSITION MATRIX DECOMPOSITION 1645

SCIKITLEARN USER GUIDE RELEASE 0213

SPARSECODER

DICTIONARYLEARNING

SPARSEPCA

MINIBATCHSPARSEPCA

NOTES

REFERENCES

J MAIRAL F BACH J PONCE G SAPIRO 2009 ONLINE DICTIONARY LEARNING FOR SPARSE CODING [HTTPSWWWDIENS](https://www.di.enscm.fr/~sierra/pdfs/icml09.pdf)

FRSIERRAPDFSICML09PDF

METHODS

FITSELF X Y FIT THE MODEL FROM DATA IN X

FITTRANSFORM SELF X Y FIT TO DATA THEN TRANSFORM IT

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PARTIALFIT SELF X Y ITEROFFSET UPDATES THE MODEL USING THE DATA IN X AS A MINIBATCH

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF X ENCODE THE DATA AS A SPARSE COMBINATION OF THE DICTIO

NARY ATOMS

INIT SELFNCOMPONENTSNONE ALPHA1 NITER1000 FITALGORITHM'LARS' NJOBSNONE

BATCHSIZE3 SHUFFLETRUE DICTINITNONE TRANSFORMALGORITHM'OMP' TRANS

FORMNNONZEROCOEFNONE TRANSFORMALPHANONE VERBOSEFALSE SPLITSIGNFALSE

RANDOMSTATENONE POSITIVECODEFALSE POSITIVEDICTFALSE

FITSELFXYNONE

FIT THE MODEL FROM DATA IN X

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING VECTOR WHERE NSAMPLES IN THE NUM

BER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

YIGNORED

RETURNS

SELF OBJECT RETURNS THE INSTANCE ITSELF

FITTRANSFORM SELFXYNONE FITPARAMS

FIT TO DATA THEN TRANSFORM IT

FITS TRANSFORMER TO X AND Y WITH OPTIONAL PARAMETERS FITPARAMS AND RETURNS A TRANSFORMED VERSION OF X

PARAMETERS

XNUMPY ARRAY OF SHAPE NSAMPLES NFEATURES TRAINING SET

YNUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES

RETURNS

XNEW NUMPY ARRAY OF SHAPE NSAMPLES NFEATURESNEW TRANSFORMED ARRAY

1646 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

GETPARAMS SELFDEEPTTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PARTIALFIT SELFXYNONE ITEROFFSETNONE

UPDATES THE MODEL USING THE DATA IN X AS A MINIBATCH

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING VECTOR WHERE NSAMPLES IN THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

YIGNORED

ITEROFFSET INTEGER OPTIONAL THE NUMBER OF ITERATION ON DATA BATCHES THAT HAS BEEN PERFORMED BEFORE THIS CALL TO PARTIALFIT THIS IS OPTIONAL IF NO NUMBER IS PASSED THE MEMORY OF THE OBJECT IS USED

RETURNS

SELF OBJECT RETURNS THE INSTANCE ITSELF

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELFXY

ENCODE THE DATA AS A SPARSE COMBINATION OF THE DICTIONARY ATOMS

CODING METHOD IS DETERMINED BY THE OBJECT PARAMETER TRANSFORMALGORITHM

PARAMETERS

XARRAY OF SHAPE NSAMPLES NFEATURES TEST DATA TO BE TRANSFORMED MUST HAVE THE SAME NUMBER OF FEATURES AS THE DATA USED TO TRAIN THE MODEL

RETURNS

XNEW ARRAY SHAPE NSAMPLES NCOMPONENTS TRANSFORMED DATA

EXAMPLES USING SKLEARNDECOMPOSITIONMINIBATCHDICTIONARYLEARNING

- IMAGE DENOISING USING DICTIONARY LEARNING
- FACES DATASET DECOMPOSITIONS

69SKLEARNDECOMPOSITION MATRIX DECOMPOSITION 1647

SCIKITLEARN USER GUIDE RELEASE 0213

698SKLEARNDECOMPOSITION MINIBATCHSPARSEPCA

CLASSSSKLEARNDECOMPOSITION MINIBATCHSPARSEPCA NCOMPONENTSNONE ALPHA1

RIDGEALPHA001 NITER100 CALL

BACKNONE BATCHSIZE3 VER

BOSEFALSE SHUFFLETRUE NJOBSNONE

METHOD'LARS' RANDOMSTATENONE

NORMALIZECOMPONENTSFALSE

MINIBATCH SPARSE PRINCIPAL COMPONENTS ANALYSIS

FINDS THE SET OF SPARSE COMPONENTS THAT CAN OPTIMALLY RECONSTRUCT THE DATA THE AMOUNT OF SPARSENESS IS CONTROL

LABLE BY THE COEFFICIENT OF THE L1 PENALTY GIVEN BY THE PARAMETER ALPHA

READ MORE IN THE USER GUIDE

PARAMETERS

NCOMPONENTS INT NUMBER OF SPARSE ATOMS TO EXTRACT

ALPHA INT SPARSITY CONTROLLING PARAMETER HIGHER VALUES LEAD TO SPARSER COMPONENTS

RIDGEALPHA FLOAT AMOUNT OF RIDGE SHRINKAGE TO APPLY IN ORDER TO IMPROVE CONDITIONING WHEN

CALLING THE TRANSFORM METHOD

NITER INT NUMBER OF ITERATIONS TO PERFORM FOR EACH MINI BATCH

CALLBACK CALLABLE OR NONE OPTIONAL DEFAULT NONE CALLABLE THAT GETS INVOKED EVERY FIVE ITER

ATIONS

BATCHSIZE INT THE NUMBER OF FEATURES TO TAKE IN EACH MINI BATCH

VERBOSE INT CONTROLS THE VERBOSITY THE HIGHER THE MORE MESSAGES DEFAULTS TO 0

SHUFFLE BOOLEAN WHETHER TO SHUFFLE THE DATA BEFORE SPLITTING IT IN BATCHES

NJOBS INT OR NONE OPTIONAL DEFAULTNONE NUMBER OF PARALLEL JOBS TO RUN NONE MEANS 1

UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL PROCESSORS SEE

GLOSSARY FOR MORE DETAILS

METHOD 'LARS' 'CD' LARS USES THE LEAST ANGLE REGRESSION METHOD TO SOLVE THE LASSO PROB

LEM LINEARMODELLARSPATH CD USES THE COORDINATE DESCENT METHOD TO COMPUTE THE LASSO

SOLUTION LINEARMODELLASSO LARS WILL BE FASTER IF THE ESTIMATED COMPONENTS ARE SPARSE

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RAN

DOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE

RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS

THE RANDOMSTATE INSTANCE USED BY NPRANDOM

NORMALIZECOMPONENTS BOOLEAN OPTIONAL DEFAULTFALSE

- IF FALSE USE A VERSION OF SPARSE PCA WITHOUT COMPONENTS NORMALIZATION AND WITHOUT DATA

CENTERING THIS IS LIKELY A BUG AND EVEN THOUGH IT'S THE DEFAULT FOR BACKWARD COMPATIBILITY

THIS SHOULD NOT BE USED

- IF TRUE USE A VERSION OF SPARSE PCA WITH COMPONENTS NORMALIZATION AND DATA CENTERING

NEW IN VERSION 020

DEPRECATED SINCE VERSION 022 NORMALIZECOMPONENTS WAS ADDED AND SET TO FALSE

FOR BACKWARD COMPATIBILITY IT WOULD BE SET TO TRUE FROM 022 ONWARDS

ATTRIBUTES

COMPONENTS ARRAY NCOMPONENTS NFEATURES SPARSE COMPONENTS EXTRACTED FROM THE DATA

1648 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213  
NITER INT NUMBER OF ITERATIONS RUN  
MEAN ARRAY SHAPE NFEATURES PERFEATURE EMPIRICAL MEAN ESTIMATED FROM THE TRAINING SET  
EQUAL TOXMEANAXISO  
SEE ALSO  
PCA  
SPARSEPCA  
DICTIONARYLEARNING  
EXAMPLES  
IMPORT NUMPY AS NP  
FROM SKLEARNDATASETS IMPORT MAKEFRIEDMAN1  
FROM SKLEARNDECOMPOSITION IMPORT MINIBATCHSPARSEPCA  
X MAKEFRIEDMAN1NSAMPLES200 NFEATURES30 RANDOMSTATE0  
TRANSFORMER MINIBATCHSPARSEPCANCOMPONENTS5  
BATCHSIZE50  
NORMALIZECOMPONENTS TRUE  
RANDOMSTATE0  
TRANSFORMERFITX  
MINIBATCHSPARSEPCA  
XTRANSFORMED TRANSFORMERTRANSFORMX  
XTRANSFORMEDSHAPE  
200 5  
MOST VALUES IN THE COMPONENTS ARE ZERO SPARSITY  
NPMEANTRANSFORMERCOMPONENTS 0  
094  
METHODS  
FITSELF X Y FIT THE MODEL FROM DATA IN X  
FITTRANSFORM SELF X Y FIT TO DATA THEN TRANSFORM IT  
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR  
SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR  
TRANSFORM SELF X LEAST SQUARES PROJECTION OF THE DATA ONTO THE SPARSE  
COMPONENTS  
INIT SELFNCOMPONENTSNONE ALPHA1 RIDGEALPHA001 NITER100 CALLBACKNONE  
BATCHSIZE3 VERBOSEFALSE SHUFFLETRUE NJOBSNONE METHOD'LARS' RAN  
DOMSTATENONE NORMALIZECOMPONENTSFALSE  
FITSELFXYNONE  
FIT THE MODEL FROM DATA IN X  
PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING VECTOR WHERE NSAMPLES IN THE NUM  
BER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES  
YIGNORED  
RETURNS  
69SKLEARNDECOMPOSITION MATRIX DECOMPOSITION 1649

SCIKITLEARN USER GUIDE RELEASE 0213

SELF OBJECT RETURNS THE INSTANCE ITSELF

FITTRANSFORM SELFXYNONE FITPARAMS

FIT TO DATA THEN TRANSFORM IT

FITS TRANSFORMER TO X AND Y WITH OPTIONAL PARAMETERS FITPARAMS AND RETURNS A TRANSFORMED VERSION OF X

PARAMETERS

XNUMPY ARRAY OF SHAPE NSAMPLES NFEATURES TRAINING SET

YNUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES

RETURNS

XNEW NUMPY ARRAY OF SHAPE NSAMPLES NFEATURESNEW TRANSFORMED ARRAY

GETPARAMS SELFDEEPTTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE

PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT

OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELFX

LEAST SQUARES PROJECTION OF THE DATA ONTO THE SPARSE COMPONENTS

TO AVOID INSTABILITY ISSUES IN CASE THE SYSTEM IS UNDERDETERMINED REGULARIZATION CAN BE APPLIED RIDGE

REGRESSION VIA THE RIDGEALPHA PARAMETER

NOTE THAT SPARSE PCA COMPONENTS ORTHOGONALITY IS NOT ENFORCED AS IN PCA HENCE ONE CANNOT USE A SIMPLE

LINEAR PROJECTION

PARAMETERS

XARRAY OF SHAPE NSAMPLES NFEATURES TEST DATA TO BE TRANSFORMED MUST HAVE THE SAME

NUMBER OF FEATURES AS THE DATA USED TO TRAIN THE MODEL

RETURNS

XNEW ARRAY SHAPE NSAMPLES NCOMPONENTS TRANSFORMED DATA

EXAMPLES USING SKLEARNDECOMPOSITIONMINIBATCHSPARSEPCA

- FACES DATASET DECOMPOSITIONS

1650 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213  
699SKLEARNDECOMPOSITION NMF  
CLASSSSKLEARNDECOMPOSITION NMFNCOMPONENTSNONE INITNONE SOLVER'CD'  
BETALOSS'FROBENIUS' TOL00001 MAXITER200 RAN  
DOMSTATENONE ALPHA00 L1RATIO00 VERBOSE0 SHUF  
FLEFALSE  
NONNEGATIVE MATRIX FACTORIZATION NMF  
FIND TWO NONNEGATIVE MATRICES W H WHOSE PRODUCT APPROXIMATES THE NON NEGATIVE MATRIX X THIS FACTORIZA  
TION CAN BE USED FOR EXAMPLE FOR DIMENSIONALITY REDUCTION SOURCE SEPARATION OR TOPIC EXTRACTION  
THE OBJECTIVE FUNCTION IS  
05X WHFRO2  
ALPHA L1RATIO VECW1  
ALPHA L1RATIO VECH1  
05ALPHA1 L1RATIO WFRO2  
05ALPHA1 L1RATIO HFRO2  
WHERE  
AFRO2 SUMIJ AIJ2 FROBENIUS NORM  
VECA1 SUMIJ ABSAIJ ELEMENTWISE L1 NORM  
FOR MULTIPLICATIVEUPDATE 'MU' SOLVER THE FROBENIUS NORM 05 X WHFRO2 CAN BE CHANGED INTO  
ANOTHER BETADIVERGENCE LOSS BY CHANGING THE BETALOSS PARAMETER  
THE OBJECTIVE FUNCTION IS MINIMIZED WITH AN ALTERNATING MINIMIZATION OF W AND H  
READ MORE IN THE USER GUIDE  
PARAMETERS  
NCOMPONENTS INT OR NONE NUMBER OF COMPONENTS IF NCOMPONENTS IS NOT SET ALL FEATURES  
ARE KEPT  
INIT NONE 'RANDOM' 'NNDSDVD' 'NNDSDVDA' 'NNDSDVDAR' 'CUSTOM' METHOD USED TO INITIALIZE  
THE PROCEDURE DEFAULT NONE VALID OPTIONS  
•NONE 'NNDSDVD' IF NCOMPONENTS MINNSAMPLES NFEATURES OTHERWISE RANDOM  
•'RANDOM' NONNEGATIVE RANDOM MATRICES SCALED WITH SQRTXMEAN  
NCOMPONENTS  
•'NNDSDVD' NONNEGATIVE DOUBLE SINGULAR VALUE DECOMPOSITION NNDSDVD  
INITIALIZATION BETTER FOR SPARSENESS  
•'NNDSDVDA' NNDSDVD WITH ZEROS FILLED WITH THE AVERAGE OF X BETTER WHEN SPARSITY IS  
NOT DESIRED  
•'NNDSDVDAR' NNDSDVD WITH ZEROS FILLED WITH SMALL RANDOM VALUES GENERALLY FASTER  
LESS ACCURATE ALTERNATIVE TO NNDSDVDA FOR WHEN SPARSITY IS NOT DESIRED  
• 'CUSTOM' USE CUSTOM MATRICES W AND H  
SOLVER 'CD' 'MU' NUMERICAL SOLVER TO USE 'CD' IS A COORDINATE DESCENT SOLVER 'MU' IS A  
MULTIPLICATIVE UPDATE SOLVER  
NEW IN VERSION 017 COORDINATE DESCENT SOLVER  
NEW IN VERSION 019 MULTIPLICATIVE UPDATE SOLVER  
69SKLEARNDECOMPOSITION MATRIX DECOMPOSITION 1651

SCIKITLEARN USER GUIDE RELEASE 0213

BETALOSS FLOAT OR STRING DEFAULT 'FROBENIUS' STRING MUST BE IN 'FROBENIUS' 'KULLBACKLEIBLER' 'ITAKURASAITO' BETA DIVERGENCE TO BE MINIMIZED MEASURING THE DISTANCE BETWEEN X AND THE DOT PRODUCT WH NOTE THAT VALUES DIFFERENT FROM 'FROBENIUS' OR 2 AND 'KULLBACKLEIBLER' OR 1 LEAD TO SIGNIFICANTLY SLOWER FITS NOTE THAT FOR BETALOSS 0 OR 'ITAKURASAITO' THE INPUT MATRIX X CANNOT CONTAIN ZEROS USED ONLY IN 'MU' SOLVER

NEW IN VERSION 019

TOLFLOAT DEFAULT 1E4 TOLERANCE OF THE STOPPING CONDITION

MAXITER INTEGER DEFAULT 200 MAXIMUM NUMBER OF ITERATIONS BEFORE TIMING OUT

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULT NONE IF INT RAN

DOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE

RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NPRANDOM

ALPHA DOUBLE DEFAULT 0 CONSTANT THAT MULTIPLIES THE REGULARIZATION TERMS SET IT TO ZERO TO HAVE NO REGULARIZATION

NEW IN VERSION 017 ALPHA USED IN THE COORDINATE DESCENT SOLVER

L1RATIO DOUBLE DEFAULT 0 THE REGULARIZATION MIXING PARAMETER WITH 0 L1RATIO 1 FOR L1RATIO 0 THE PENALTY IS AN ELEMENTWISE L2 PENALTY AKA FROBENIUS NORM FOR L1RATIO 1 IT IS AN ELEMENTWISE L1 PENALTY FOR 0 L1RATIO 1 THE PENALTY IS A COMBINATION OF L1 AND L2

NEW IN VERSION 017 REGULARIZATION PARAMETER L1RATIO USED IN THE COORDINATE DESCENT SOLVER

VERBOSE BOOL DEFAULTFALSE WHETHER TO BE VERBOSE

SHUFFLE BOOLEAN DEFAULT FALSE IF TRUE RANDOMIZE THE ORDER OF COORDINATES IN THE CD SOLVER

NEW IN VERSION 017 SHUFFLE PARAMETER USED IN THE COORDINATE DESCENT SOLVER

ATTRIBUTES

COMPONENTS ARRAY NCOMPONENTS NFEATURES FACTORIZATION MATRIX SOMETIMES CALLED 'DICTIONARY'

RECONSTRUCTIONERR NUMBER FROBENIUS NORM OF THE MATRIX DIFFERENCE OR BETADIVERGENCE BETWEEN THE TRAINING DATA XAND THE RECONSTRUCTED DATA WHFROM THE FITTED MODEL

NITER INT ACTUAL NUMBER OF ITERATIONS

REFERENCES

CICHOCKI ANDRZEJ AND P H A N ANHHUY "FAST LOCAL ALGORITHMS FOR LARGE SCALE NONNEGATIVE MATRIX AND TENSOR FACTORIZATIONS" IEICE TRANSACTIONS ON FUNDAMENTALS OF ELECTRONICS COMMUNICATIONS AND COMPUTER SCIENCES 923 708721 2009

FEVOTTE C IDIER J 2011 ALGORITHMS FOR NONNEGATIVE MATRIX FACTORIZATION WITH THE BETADIVERGENCE NEURAL COMPUTATION 239

EXAMPLES

1652 CHAPTER 6 API REFERENCE

```
SCIKITLEARN USER GUIDE RELEASE 0213
import numpy as np
X = np.array([1, 2, 1, 3, 12, 4, 1, 5, 08, 6, 1])
from sklearn.decomposition import NMF
model = NMF(n_components=2, init='random', random_state=0)
W = model.fit_transform(X)
H = model.components_

METHODS
fit(self, X, Y=None) Learn a NMF model for the data X
fit_transform(self, X, Y=None, W=None, H=None) Learn a NMF model for the data X and returns the transformed data
get_params(self, deep=True) Get parameters for this estimator
inverse_transform(self, W) Transform data back to its original space
set_params(self, **kwargs) Set the parameters of this estimator
transform(self, X) Transform the data X according to the fitted NMF model
init(self, n_components=None, init='random', solver='cd', beta_loss='frobenius', tol=0.0001, max_iter=200, random_state=None, alpha=0.0, l1_ratio=0.0, verbose=0, shuffle=False)
fit(self, X, Y=None, W=None, H=None) Learn a NMF model for the data X
parameters_ Parameters
X_arraylike_shape nsamples nfeatures Data matrix to be decomposed
Y_ignored Returns
self fit_transform(self, X, Y=None, W=None, H=None) Learn a NMF model for the data X and returns the transformed data
this_is_more_efficient_than_calling_fit_followed_by_transform Parameters
X_arraylike_shape nsamples nfeatures Data matrix to be decomposed
Y_ignored Returns
W_arraylike_shape nsamples ncomponents If init='custom' it is used as initial guess for the solution
H_arraylike_shape ncomponents nfeatures If init='custom' it is used as initial guess for the solution
Returns
W_array_shape nsamples ncomponents Transformed data
get_params(self, deep=True) Get parameters for this estimator
parameters_ Parameters
69sklearn.decomposition.MatrixDecomposition 1653
```

SCIKITLEARN USER GUIDE RELEASE 0213

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

INVERSETRANSFORM SELF

TRANSFORM DATA BACK TO ITS ORIGINAL SPACE

PARAMETERS

WARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NCOMPONENTS TRANSFORMED DATA MATRIX

RETURNS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES DATA MATRIX OF ORIGINAL SHAPE

NEW IN VERSION 018

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELF

TRANSFORM THE DATA X ACCORDING TO THE FITTED NMF MODEL

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES DATA MATRIX TO BE TRANSFORMED BY THE MODEL

RETURNS

WARRAY SHAPE NSAMPLES NCOMPONENTS TRANSFORMED DATA

EXAMPLES USING SKLEARNDECOMPOSITIONNMF

- TOPIC EXTRACTION WITH NONNEGATIVE MATRIX FACTORIZATION AND LATENT DIRICHLET ALLOCATION
- SELECTING DIMENSIONALITY REDUCTION WITH PIPELINE AND GRIDSEARCHCV
- FACES DATASET DECOMPOSITIONS

6910SKLEARNDECOMPOSITION PCA

CLASSSKLEARNDECOMPOSITION PCANCOMPONENTSNONE COPYTRUE WHITENFALSE

SVDSOLVER'AUTO' TOL00 ITERATEDPOWER'AUTO' RAN

DOMSTATENONE

PRINCIPAL COMPONENT ANALYSIS PCA

LINEAR DIMENSIONALITY REDUCTION USING SINGULAR VALUE DECOMPOSITION OF THE DATA TO PROJECT IT TO A LOWER DIMENSIONAL SPACE THE INPUT DATA IS CENTERED BUT NOT SCALED FOR EACH FEATURE BEFORE APPLYING THE SVD

IT USES THE LAPACK IMPLEMENTATION OF THE FULL SVD OR A RANDOMIZED TRUNCATED SVD BY THE METHOD OF HALKO ET AL 2009 DEPENDING ON THE SHAPE OF THE INPUT DATA AND THE NUMBER OF COMPONENTS TO EXTRACT

1654 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

IT CAN ALSO USE THE SCIPYSPARSELINALG ARPACK IMPLEMENTATION OF THE TRUNCATED SVD  
NOTICE THAT THIS CLASS DOES NOT SUPPORT SPARSE INPUT SEE TRUNCATEDSVD FOR AN ALTERNATIVE WITH SPARSE DATA  
READ MORE IN THE USER GUIDE

PARAMETERS

NCOMPONENTS INT FLOAT NONE OR STRING NUMBER OF COMPONENTS TO KEEP IF NCOMPONENTS IS  
NOT SET ALL COMPONENTS ARE KEPT

NCOMPONENTS MINNSAMPLES NFEATURES

IFNCOMPONENTS MLE ANDSVDSOLVER FULL MINKA’S MLE IS USED

TO GUESS THE DIMENSION USE OF NCOMPONENTS MLE WILL INTERPRET SVDSOLVER  
AUTO ASSVDSOLVER FULL

IF0 NCOMPONENTS 1 ANDSVDSOLVER FULL SELECT THE NUMBER OF  
COMPONENTS SUCH THAT THE AMOUNT OF VARIANCE THAT NEEDS TO BE EXPLAINED IS GREATER THAN THE  
PERCENTAGE SPECIFIED BY NCOMPONENTS

IFSVDSOLVER ARPACK THE NUMBER OF COMPONENTS MUST BE STRICTLY LESS THAN THE  
MINIMUM OF NFEATURES AND NSAMPLES

HENCE THE NONE CASE RESULTS IN

NCOMPONENTS MINNSAMPLES NFEATURES 1

COPY BOOL DEFAULT TRUE IF FALSE DATA PASSED TO FIT ARE OVERWRITTEN AND RUNNING  
FITXTRANSFORMX WILL NOT YIELD THE EXPECTED RESULTS USE FITTRANSFORMX INSTEAD

WHITEN BOOL OPTIONAL DEFAULT FALSE WHEN TRUE FALSE BY DEFAULT THE COMPONENTS VECTORS  
ARE MULTIPLIED BY THE SQUARE ROOT OF NSAMPLES AND THEN DIVIDED BY THE SINGULAR VALUES TO  
ENSURE UNCORRELATED OUTPUTS WITH UNIT COMPONENTWISE VARIANCES

WHITENING WILL REMOVE SOME INFORMATION FROM THE TRANSFORMED SIGNAL THE RELATIVE VARIANCE  
SCALES OF THE COMPONENTS BUT CAN SOMETIME IMPROVE THE PREDICTIVE ACCURACY OF THE DOWN  
STREAM ESTIMATORS BY MAKING THEIR DATA RESPECT SOME HARDWIRED ASSUMPTIONS

SVDSOLVER STRING ‘AUTO’ ‘FULL’ ‘ARPACK’ ‘RANDOMIZED’

AUTO THE SOLVER IS SELECTED BY A DEFAULT POLICY BASED ON XSHAPE ANDNCOMPONENTS IF  
THE INPUT DATA IS LARGER THAN 500X500 AND THE NUMBER OF COMPONENTS TO EXTRACT IS LOWER  
THAN 80 OF THE SMALLEST DIMENSION OF THE DATA THEN THE MORE EFFICIENT ‘RANDOMIZED’  
METHOD IS ENABLED OTHERWISE THE EXACT FULL SVD IS COMPUTED AND OPTIONALLY TRUNCATED  
AFTERWARDS

FULL RUN EXACT FULL SVD CALLING THE STANDARD LAPACK SOLVER VIA SCIPYLINALGSVD  
AND SELECT THE COMPONENTS BY POSTPROCESSING

ARPACK RUN SVD TRUNCATED TO NCOMPONENTS CALLING ARPACK SOLVER VIA SCIPY

SPARSELINALGSVDS IT REQUIRES STRICTLY 0 NCOMPONENTS MINXSHAPE

RANDOMIZED RUN RANDOMIZED SVD BY THE METHOD OF HALKO ET AL

NEW IN VERSION 0180

TOLFLOAT 0 OPTIONAL DEFAULT 0 TOLERANCE FOR SINGULAR VALUES COMPUTED BY SVDSOLVER  
‘ARPACK’

NEW IN VERSION 0180

69SKLEARNDECOMPOSITION MATRIX DECOMPOSITION 1655

SCIKITLEARN USER GUIDE RELEASE 0213  
ITERATEDPOWER INT 0 OR 'AUTO' DEFAULT 'AUTO' NUMBER OF ITERATIONS FOR THE POWER METHOD  
COMPUTED BY SVDSOLVER 'RANDOMIZED'  
NEW IN VERSION 0180  
RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULT NONE IF INT RAN  
DOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE  
RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS  
THE RANDOMSTATE INSTANCE USED BY NPRANDOM USED WHEN SVDSOLVER 'ARPACK' OR  
'RANDOMIZED'  
NEW IN VERSION 0180  
ATTRIBUTES  
COMPONENTS ARRAY SHAPE NCOMPONENTS NFEATURES PRINCIPAL AXES IN FEATURE SPACE REP  
RESENTING THE DIRECTIONS OF MAXIMUM VARIANCE IN THE DATA THE COMPONENTS ARE SORTED BY  
EXPLAINEDVARIANCE  
EXPLAINEDVARIANCE ARRAY SHAPE NCOMPONENTS THE AMOUNT OF VARIANCE EXPLAINED BY  
EACH OF THE SELECTED COMPONENTS  
EQUAL TO NCOMPONENTS LARGEST EIGENVALUES OF THE COVARIANCE MATRIX OF X  
NEW IN VERSION 018  
EXPLAINEDVARIANCERATIO ARRAY SHAPE NCOMPONENTS PERCENTAGE OF VARIANCE EXPLAINED  
BY EACH OF THE SELECTED COMPONENTS  
IFNCOMPONENTS IS NOT SET THEN ALL COMPONENTS ARE STORED AND THE SUM OF THE RATIOS IS  
EQUAL TO 10  
SINGULARVALUES ARRAY SHAPE NCOMPONENTS THE SINGULAR VALUES CORRESPONDING TO EACH  
OF THE SELECTED COMPONENTS THE SINGULAR VALUES ARE EQUAL TO THE 2NORMS OF THE  
NCOMPONENTS VARIABLES IN THE LOWERDIMENSIONAL SPACE  
NEW IN VERSION 019  
MEAN ARRAY SHAPE NFEATURES PERFEATURE EMPIRICAL MEAN ESTIMATED FROM THE TRAINING SET  
EQUAL TOXMEANAXIS0  
NCOMPONENTS INT THE ESTIMATED NUMBER OF COMPONENTS WHEN NCOMPONENTS IS SET TO  
'MLE' OR A NUMBER BETWEEN 0 AND 1 WITH SVDSOLVER 'FULL' THIS NUMBER IS ESTIMATED  
FROM INPUT DATA OTHERWISE IT EQUALS THE PARAMETER NCOMPONENTS OR THE LESSER VALUE OF  
NFEATURES AND NSAMPLES IF NCOMPONENTS IS NONE  
NOISEVARIANCE FLOAT THE ESTIMATED NOISE COVARIANCE FOLLOWING THE PROBABILISTIC PCA MODEL  
FROM TIPPING AND BISHOP 1999 SEE "PATTERN RECOGNITION AND MACHINE LEARNING" BY C  
BISHOP 1221 P 574 OR HTTPWWWMIKETIPPINGCOMPAPERSMETMPPCAPDF IT IS REQUIRED  
TO COMPUTE THE ESTIMATED DATA COVARIANCE AND SCORE SAMPLES  
EQUAL TO THE AVERAGE OF MINNFEATURES NSAMPLES NCOMPONENTS SMALLEST EIGENVALUES  
OF THE COVARIANCE MATRIX OF X  
SEE ALSO  
KERNELPCA  
SPARSEPCA  
TRUNCATEDSVD  
INCREMENTALPCA  
1656 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

REFERENCES

FOR NCOMPONENTS ‘MLE’ THIS CLASS USES THE METHOD OF MINKA T P “AUTOMATIC CHOICE OF DIMENSIONALITY FOR PCA” IN NIPS PP 598604

IMPLEMENTS THE PROBABILISTIC PCA MODEL FROM TIPPING M E AND BISHOP C M 1999 “PROBABILISTIC PRINCIPAL COMPONENT ANALYSIS” JOURNAL OF THE ROYAL STATISTICAL SOCIETY SERIES B STATISTICAL METHODOLOGY 613 611622 VIA THE SCORE AND SCORESAMPLES METHODS SEE [HTTPWWWMIKETIPPINGCOMPAPERSMETMPPCA PDF](http://www.miketipping.com/papers/METMPPCA.pdf)

FOR SVDSOLVER ‘ARPACK’ REFER TO SCIPYSPARSELINALGSVDS

FOR SVDSOLVER ‘RANDOMIZED’ SEE HALKO N MARTINSSON P G AND TROPP J A 2011 “FINDING STRUCTURE WITH RANDOMNESS PROBABILISTIC ALGORITHMS FOR CONSTRUCTING APPROXIMATE MATRIX DECOMPOSITIONS” SIAM REVIEW 532 217288 AND ALSO MARTINSSON P G ROKHLIN V AND TYGERT M 2011 “A RANDOMIZED ALGORITHM FOR THE DECOMPOSITION OF MATRICES” APPLIED AND COMPUTATIONAL HARMONIC ANALYSIS 301 4768

EXAMPLES

```
import numpy as np
from sklearn.decomposition import PCA
X = np.array([1, 2, 1, 3, 2, 1, 1, 2, 1, 3, 2])
PCA(n_components=2).fit(X)
print(PCA.explained_variance_ratio_)
0.9924 0.0075
print(PCA.singular_values_)
630061 0.54980
PCA(n_components=2, svd_solver='full').fit(X)
print(PCA.explained_variance_ratio_)
0.9924 0.0075
print(PCA.singular_values_)
630061 0.54980
PCA(n_components=1, svd_solver='arpack').fit(X)
print(PCA.explained_variance_ratio_)
0.99244
print(PCA.singular_values_)
630061
METHODS
fit(self, X, Y=None) Fit the model with X
Continued on next page
95SKLEARNDECOMPOSITION MATRIX DECOMPOSITION 1657
```

SCIKITLEARN USER GUIDE RELEASE 0213

TABLE 659 – CONTINUED FROM PREVIOUS PAGE

FITTRANSFORM SELF X Y FIT THE MODEL WITH X AND APPLY THE DIMENSIONALITY REDUCTION ON X

GETCOVARIANCE SELF COMPUTE DATA COVARIANCE WITH THE GENERATIVE MODEL

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

GETPRECISION SELF COMPUTE DATA PRECISION MATRIX WITH THE GENERATIVE MODEL

INVERSETRANSFORM SELF X TRANSFORM DATA BACK TO ITS ORIGINAL SPACE

SCORE SELF X Y RETURN THE AVERAGE LOGLIKELIHOOD OF ALL SAMPLES

SCORESAMPLES SELF X RETURN THE LOGLIKELIHOOD OF EACH SAMPLE

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF X APPLY DIMENSIONALITY REDUCTION TO X

INIT SELF NCOMPONENTS NONE COPY TRUE WHITEN FALSE SVDSOLVER 'AUTO' TOL 0.0 ITERATED POWER 'AUTO' RANDOM STATE NONE

FIT SELF X Y NONE

FIT THE MODEL WITH X

PARAMETERS

X ARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

Y IGNORED

RETURNS

SELF OBJECT RETURNS THE INSTANCE ITSELF

FITTRANSFORM SELF X Y NONE

FIT THE MODEL WITH X AND APPLY THE DIMENSIONALITY REDUCTION ON X

PARAMETERS

X ARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

Y IGNORED

RETURNS

X NEW ARRAYLIKE SHAPE NSAMPLES NCOMPONENTS

GETCOVARIANCE SELF

COMPUTE DATA COVARIANCE WITH THE GENERATIVE MODEL

COV COMPONENTS S2COMPONENTS SIGMA2 EYE NFEATURES

WHERE S2 CONTAINS THE EXPLAINED VARIANCES AND SIGMA2 CONTAINS THE NOISE VARIANCES

RETURNS

COV ARRAY SHAPENFEATURES NFEATURES ESTIMATED COVARIANCE OF DATA

GETPARAMS SELF DEEP TRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

1658 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

GETPRECISION SELF

COMPUTE DATA PRECISION MATRIX WITH THE GENERATIVE MODEL

EQUALS THE INVERSE OF THE COVARIANCE BUT COMPUTED WITH THE MATRIX INVERSION LEMMA FOR EFFICIENCY

RETURNS

PRECISION ARRAY SHAPENFEATURES NFEATURES ESTIMATED PRECISION OF DATA

INVERSETRANSFORM SELF

TRANSFORM DATA BACK TO ITS ORIGINAL SPACE

IN OTHER WORDS RETURN AN INPUT XORIGINAL WHOSE TRANSFORM WOULD BE X

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NCOMPONENTS NEW DATA WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NCOMPONENTS IS THE NUMBER OF COMPONENTS

RETURNS

XORIGINAL ARRAYLIKE SHAPE NSAMPLES NFEATURES

NOTES

IF WHITENING IS ENABLED INVERSETRANSFORM WILL COMPUTE THE EXACT INVERSE OPERATION WHICH INCLUDES REVERSING WHITENING

SCORESELFXYNONE

RETURN THE AVERAGE LOGLIKELIHOOD OF ALL SAMPLES

SEE "PATTERN RECOGNITION AND MACHINE LEARNING" BY C BISHOP 1221 P 574 OR [HTTPWWWMIKETIPPING.COMPAPERSMETMPPCAPDF](http://www.miketipping.com/papers/metmppcapdf)

PARAMETERS

XARRAY SHAPENSAMPLES NFEATURES THE DATA

YIGNORED

RETURNS

LLFLOAT AVERAGE LOGLIKELIHOOD OF THE SAMPLES UNDER THE CURRENT MODEL

SCORESAMPLES SELF

RETURN THE LOGLIKELIHOOD OF EACH SAMPLE

SEE "PATTERN RECOGNITION AND MACHINE LEARNING" BY C BISHOP 1221 P 574 OR [HTTPWWWMIKETIPPING.COMPAPERSMETMPPCAPDF](http://www.miketipping.com/papers/metmppcapdf)

PARAMETERS

XARRAY SHAPENSAMPLES NFEATURES THE DATA

RETURNS

LLARRAY SHAPE NSAMPLES LOGLIKELIHOOD OF EACH SAMPLE UNDER THE CURRENT MODEL

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

69SKLEARNDECOMPOSITION MATRIX DECOMPOSITION 1659

SCIKITLEARN USER GUIDE RELEASE 0213

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS  
SELF

TRANSFORM SELF  
X

APPLY DIMENSIONALITY REDUCTION TO X  
X IS PROJECTED ON THE FIRST PRINCIPAL COMPONENTS PREVIOUSLY EXTRACTED FROM A TRAINING SET

PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES NEW DATA WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

RETURNS  
XNEW ARRAYLIKE SHAPE NSAMPLES NCOMPONENTS

EXAMPLES  
IMPORT NUMPY AS NP  
FROM SKLEARNDECOMPOSITION IMPORT INCREMENTALPCA  
X NPARRAY1 1 2 1 3 2 1 1 2 1 3 2  
IPCA INCREMENTALPCANCOMPONENTS2 BATCHSIZE3  
IPCAFITX  
INCREMENTALPCABATCHSIZE3 COPYTRUE NCOMPONENTS2 WHITENFALSE  
IPCATRANSFORMX

EXAMPLES USING SKLEARNDECOMPOSITIONPCA

- MULTILABEL CLASSIFICATION
- EXPLICIT FEATURE MAP APPROXIMATION FOR RBF KERNELS
- FACES RECOGNITION EXAMPLE USING EIGENFACES AND SVMs
- A DEMO OF KMEANS CLUSTERING ON THE HANDWRITTEN DIGITS DATA
- CONCATENATING MULTIPLE FEATURE EXTRACTION METHODS
- PIPELINING CHAINING A PCA AND A LOGISTIC REGRESSION
- SELECTING DIMENSIONALITY REDUCTION WITH PIPELINE AND GRIDSEARCHCV
- THE IRIS DATASET
- PCA EXAMPLE WITH IRIS DATASET
- INCREMENTAL PCA
- COMPARISON OF LDA AND PCA 2D PROJECTION OF IRIS DATASET
- BLIND SOURCE SEPARATION USING FASTICA
- PRINCIPAL COMPONENTS ANALYSIS PCA
- FASTICA ON 2D POINT CLOUDS
- KERNEL PCA

1660 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

- MODEL SELECTION WITH PROBABILISTIC PCA AND FACTOR ANALYSIS FA
- FACES DATASET DECOMPOSITIONS
- MULTIDIMENSIONAL SCALING
- BALANCE MODEL COMPLEXITY AND CROSSVALIDATED SCORE
- KERNEL DENSITY ESTIMATION
- DIMENSIONALITY REDUCTION WITH NEIGHBORHOOD COMPONENTS ANALYSIS
- USING FUNCTIONTRANSFORMER TO SELECT COLUMNS
- IMPORTANCE OF FEATURE SCALING

6911SKLEARNDECOMPOSITION SPARSEPCA  
CLASSSSKLEARNDECOMPOSITION SPARSEPCA NCOMPONENTSNONE ALPHA1 RIDGEALPHA001

MAXITER1000 TOL1E08 METHOD'LARS'  
NJOBSNONE UINITNONE VINITNONE  
VERBOSEFALSE RANDOMSTATENONE NORMAL  
IZECOMPONENTSFALSE

SPARSE PRINCIPAL COMPONENTS ANALYSIS SPARSEPCA  
FINDS THE SET OF SPARSE COMPONENTS THAT CAN OPTIMALLY RECONSTRUCT THE DATA THE AMOUNT OF SPARSENESS IS CONTROL  
LABLE BY THE COEFFICIENT OF THE L1 PENALTY GIVEN BY THE PARAMETER ALPHA  
READ MORE IN THE USER GUIDE

PARAMETERS  
NCOMPONENTS INT NUMBER OF SPARSE ATOMS TO EXTRACT  
ALPHA FLOAT SPARSITY CONTROLLING PARAMETER HIGHER VALUES LEAD TO SPARSER COMPONENTS  
RIDGEALPHA FLOAT AMOUNT OF RIDGE SHRINKAGE TO APPLY IN ORDER TO IMPROVE CONDITIONING WHEN  
CALLING THE TRANSFORM METHOD  
MAXITER INT MAXIMUM NUMBER OF ITERATIONS TO PERFORM  
TOLFLOAT TOLERANCE FOR THE STOPPING CONDITION  
METHOD 'LARS' 'CD' LARS USES THE LEAST ANGLE REGRESSION METHOD TO SOLVE THE LASSO PROB  
LEM LINEARMODELLARSPATH CD USES THE COORDINATE DESCENT METHOD TO COMPUTE THE LASSO  
SOLUTION LINEARMODELLASSO LARS WILL BE FASTER IF THE ESTIMATED COMPONENTS ARE SPARSE  
NJOBS INT OR NONE OPTIONAL DEFAULTNONE NUMBER OF PARALLEL JOBS TO RUN NONE MEANS 1  
UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL PROCESSORS SEE

GLOSSARY FOR MORE DETAILS  
UINIT ARRAY OF SHAPE NSAMPLES NCOMPONENTS INITIAL VALUES FOR THE LOADINGS FOR WARM  
RESTART SCENARIOS  
VINIT ARRAY OF SHAPE NCOMPONENTS NFEATURES INITIAL VALUES FOR THE COMPONENTS FOR WARM  
RESTART SCENARIOS  
VERBOSE INT CONTROLS THE VERBOSITY THE HIGHER THE MORE MESSAGES DEFAULTS TO 0  
RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RAN  
DOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE  
RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS  
THE RANDOMSTATE INSTANCE USED BY NPRANDOM  
69SKLEARNDECOMPOSITION MATRIX DECOMPOSITION 1661

SCIKITLEARN USER GUIDE RELEASE 0213

NORMALIZECOMPONENTS BOOLEAN OPTIONAL DEFAULTFALSE

- IF FALSE USE A VERSION OF SPARSE PCA WITHOUT COMPONENTS NORMALIZATION AND WITHOUT DATA CENTERING THIS IS LIKELY A BUG AND EVEN THOUGH IT’S THE DEFAULT FOR BACKWARD COMPATIBILITY THIS SHOULD NOT BE USED
- IF TRUE USE A VERSION OF SPARSE PCA WITH COMPONENTS NORMALIZATION AND DATA CENTERING NEW IN VERSION 020

DEPRECATED SINCE VERSION 022 NORMALIZECOMPONENTS WAS ADDED AND SET TO FALSE FOR BACKWARD COMPATIBILITY IT WOULD BE SET TO TRUE FROM 022 ONWARDS

ATTRIBUTES

COMPONENTS ARRAY NCOMPONENTS NFEATURES SPARSE COMPONENTS EXTRACTED FROM THE DATA

ERROR ARRAY VECTOR OF ERRORS AT EACH ITERATION

NITER INT NUMBER OF ITERATIONS RUN

MEAN ARRAY SHAPE NFEATURES PERFEATURE EMPIRICAL MEAN ESTIMATED FROM THE TRAINING SET

EQUAL TOXMEANAXIS0

SEE ALSO

PCA

MINIBATCHSPARSEPCA

DICTIONARYLEARNING

EXAMPLES

```
import numpy as np
from sklearn.datasets import make_friedman1
from sklearn.decomposition import sparse_pca
X, make_friedman1(n_samples=200, n_features=30, random_state=0)
transformer = sparse_pca(n_components=5,
                          normalize_components=True,
                          random_state=0)
transformer.fit(X)
sparse_pca
X_transformed = transformer.transform(X)
X_transformed.shape
200 5
# Most values in the components are zero
np.mean(transformer.components_)
0.9666
```

METHODS

fit(X, Y) fit the model from data in X

fit\_transform(X, Y) fit to data then transform it

get\_params() self deep get parameters for this estimator

set\_params(self, params) set the parameters of this estimator

CONTINUED ON NEXT PAGE

1662 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

TABLE 660 – CONTINUED FROM PREVIOUS PAGE

TRANSFORM SELF X LEAST SQUARES PROJECTION OF THE DATA ONTO THE SPARSE COMPONENTS

INIT SELF NCOMPONENTS NONE ALPHA1 RIDGE ALPHA001 MAXITER 1000 TOL 1E-08 METHOD ‘LARS’ NJOBS NONE UINIT NONE VINIT NONE VERBOSE FALSE RAN DOMSTAT NONE NORMALIZE COMPONENTS FALSE

FIT SELF X Y NONE

FIT THE MODEL FROM DATA IN X

PARAMETERS

X ARRAY LIKE SHAPE NSAMPLES NFEATURES TRAINING VECTOR WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

Y IGNORED

RETURNS

SELF OBJECT RETURNS THE INSTANCE ITSELF

FIT TRANSFORM SELF X Y NONE FIT PARAMS

FIT TO DATA THEN TRANSFORM IT

FITS TRANSFORMER TO X AND Y WITH OPTIONAL PARAMETERS FIT PARAMS AND RETURNS A TRANSFORMED VERSION OF X

PARAMETERS

X NUMPY ARRAY OF SHAPE NSAMPLES NFEATURES TRAINING SET

Y NUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES

RETURNS

X NEW NUMPY ARRAY OF SHAPE NSAMPLES NFEATURES NEW TRANSFORMED ARRAY

GET PARAMS SELF DEEP TRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

SET PARAMS SELF PARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENT PARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELF X

LEAST SQUARES PROJECTION OF THE DATA ONTO THE SPARSE COMPONENTS

TO AVOID INSTABILITY ISSUES IN CASE THE SYSTEM IS UNDERDETERMINED REGULARIZATION CAN BE APPLIED RIDGE REGRESSION VIA THE RIDGE ALPHA PARAMETER

69 SKLEARN DECOMPOSITION MATRIX DECOMPOSITION 1663

SCIKITLEARN USER GUIDE RELEASE 0213

NOTE THAT SPARSE PCA COMPONENTS ORTHOGONALITY IS NOT ENFORCED AS IN PCA HENCE ONE CANNOT USE A SIMPLE LINEAR PROJECTION

PARAMETERS

XARRAY OF SHAPE NSAMPLES NFEATURES TEST DATA TO BE TRANSFORMED MUST HAVE THE SAME NUMBER OF FEATURES AS THE DATA USED TO TRAIN THE MODEL

RETURNS

XNEW ARRAY SHAPE NSAMPLES NCOMPONENTS TRANSFORMED DATA

6912SKLEARNDECOMPOSITION SPARSECODER

CLASSSKLEARNDECOMPOSITION SPARSECODER DICTIONARY TRANSFORMALGORITHM'OMP'

TRANSFORMNNONZEROCOEFSSNONE TRANS

FORMALPHANONE SPLITSIGNFALSE NJOBSNONE

POSITIVECODEFALSE

SPARSE CODING

FINDS A SPARSE REPRESENTATION OF DATA AGAINST A FIXED PRECOMPUTED DICTIONARY

EACH ROW OF THE RESULT IS THE SOLUTION TO A SPARSE CODING PROBLEM THE GOAL IS TO FIND A SPARSE ARRAY CODE SUCH THAT

X CODE DICTIONARY

READ MORE IN THE USER GUIDE

PARAMETERS

DICTIONARY ARRAY NCOMPONENTS NFEATURES THE DICTIONARY ATOMS USED FOR SPARSE CODING

LINEAS ARE ASSUMED TO BE NORMALIZED TO UNIT NORM

TRANSFORMALGORITHM 'LASSOLARS' 'LASSOCD' 'LARS' 'OMP' 'THRESHOLD' ALGORITHM USED TO TRANSFORM THE DATA LARS USES THE LEAST ANGLE REGRESSION METHOD LINEARMODELLARSPATH LASSOLARS USES LARS TO COMPUTE THE LASSO SOLUTION LASSOCD USES THE COORDINATE DESCENT METHOD TO COMPUTE THE LASSO SOLUTION LINEARMODELLASSO LASSOLARS WILL BE FASTER IF THE ESTIMATED COMPONENTS ARE SPARSE OMP USES ORTHOGONAL MATCHING PURSUIT TO ESTIMATE THE SPARSE SOLUTION THRESHOLD SQUASHES TO ZERO ALL COEFFICIENTS LESS THAN ALPHA FROM THE PROJECTIONDICTIONARY X

TRANSFORMNNONZEROCOEFSS INTO1NFEATURES BY DEFAULT NUMBER OF NONZERO COEFFICIENTS TO TARGET IN EACH COLUMN OF THE SOLUTION THIS IS ONLY USED BY ALGORITHMMLARS ANDALGORITHMOMP AND IS OVERRIDDEN BY ALPHA IN THE OR THOGONAL MATCHING PURSUIT OMP CASE

TRANSFORMALPHA FLOAT 1 BY DEFAULT IF ALGORITHMMLASSOLARS OR ALGORITHMMLASSOCD ALPHA IS THE PENALTY APPLIED TO THE L1 NORM IF ALGORITHMTHRESHOLD ALPHA IS THE ABSOLUTE VALUE OF THE THRESHOLD BELOW WHICH COEFFICIENTS WILL BE SQUASHED TO ZERO IF ALGORITHMOMP ALPHA IS THE TOLERANCE PARAMETER THE VALUE OF THE RECONSTRUCTION ERROR TARGETED IN THIS CASE IT OVERRIDES NNONZEROCOEFSS

SPLITSIGN BOOL FALSE BY DEFAULT WHETHER TO SPLIT THE SPARSE FEATURE VECTOR INTO THE CONCATENATION OF ITS NEGATIVE PART AND ITS POSITIVE PART THIS CAN IMPROVE THE PERFORMANCE OF DOWN STREAM CLASSIFIERS

1664 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213  
NJOBS INT OR NONE OPTIONAL DEFAULTNONE NUMBER OF PARALLEL JOBS TO RUN NONE MEANS 1  
UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL PROCESSORS SEE  
GLOSSARY FOR MORE DETAILS  
POSITIVECODE BOOL WHETHER TO ENFORCE POSITIVITY WHEN FINDING THE CODE  
NEW IN VERSION 020  
ATTRIBUTES  
COMPONENTS ARRAY NCOMPONENTS NFEATURES THE UNCHANGED DICTIONARY ATOMS  
SEE ALSO  
DICTIONARYLEARNING  
MINIBATCHDICTIONARYLEARNING  
SPARSEPCA  
MINIBATCHSPARSEPCA  
SPARSEENCODE  
METHODS  
FITSELF X Y DO NOTHING AND RETURN THE ESTIMATOR UNCHANGED  
FITTRANSFORM SELF X Y FIT TO DATA THEN TRANSFORM IT  
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR  
SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR  
TRANSFORM SELF X ENCODE THE DATA AS A SPARSE COMBINATION OF THE DICTIO  
NARY ATOMS  
INIT SELFDICTIONARY TRANSFORMALGORITHM'OMP' TRANSFORMMNNONZEROCOEFSSNONE TRANS  
FORMALPHANONE SPLITSIGNFALSE NJOBSNONE POSITIVECODEFALSE  
FITSELFXYNONE  
DO NOTHING AND RETURN THE ESTIMATOR UNCHANGED  
THIS METHOD IS JUST THERE TO IMPLEMENT THE USUAL API AND HENCE WORK IN PIPELINES  
PARAMETERS  
XIGNORED  
YIGNORED  
RETURNS  
SELF OBJECT RETURNS THE OBJECT ITSELF  
FITTRANSFORM SELFXYNONE FITPARAMS  
FIT TO DATA THEN TRANSFORM IT  
FITS TRANSFORMER TO X AND Y WITH OPTIONAL PARAMETERS FITPARAMS AND RETURNS A TRANSFORMED VERSION OF X  
PARAMETERS  
XNUMPY ARRAY OF SHAPE NSAMPLES NFEATURES TRAINING SET  
YNUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES  
RETURNS  
69SKLEARNDECOMPOSITION MATRIX DECOMPOSITION 1665

SCIKITLEARN USER GUIDE RELEASE 0213

XNEW NUMPY ARRAY OF SHAPE NSAMPLES NFEATURESNEW TRANSFORMED ARRAY

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE

PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT

OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELF

ENCODE THE DATA AS A SPARSE COMBINATION OF THE DICTIONARY ATOMS

CODING METHOD IS DETERMINED BY THE OBJECT PARAMETER TRANSFORMALGORITHM

PARAMETERS

XARRAY OF SHAPE NSAMPLES NFEATURES TEST DATA TO BE TRANSFORMED MUST HAVE THE SAME

NUMBER OF FEATURES AS THE DATA USED TO TRAIN THE MODEL

RETURNS

XNEW ARRAY SHAPE NSAMPLES NCOMPONENTS TRANSFORMED DATA

EXAMPLES USING SKLEARNDECOMPOSITIONSPARSECODER

- SPARSE CODING WITH A PRECOMPUTED DICTIONARY

6913SKLEARNDECOMPOSITION TRUNCATEDSVD

CLASSSSKLEARNDECOMPOSITION TRUNCATEDSVD NCOMPONENTS2 ALGORITHM’RANDOMIZED’

NITER5 RANDOMSTATENONE TOL00

DIMENSIONALITY REDUCTION USING TRUNCATED SVD AKA LSA

THIS TRANSFORMER PERFORMS LINEAR DIMENSIONALITY REDUCTION BY MEANS OF TRUNCATED SINGULAR VALUE DECOMPOSITION

SVD CONTRARY TO PCA THIS ESTIMATOR DOES NOT CENTER THE DATA BEFORE COMPUTING THE SINGULAR VALUE DECOMPO

SITION THIS MEANS IT CAN WORK WITH SCIPYSPARSE MATRICES EFFICIENTLY

IN PARTICULAR TRUNCATED SVD WORKS ON TERM COUNTTFIDF MATRICES AS RETURNED BY THE VECTORIZERS IN

SKLEARNFEATUREEXTRACTIONTEXT IN THAT CONTEXT IT IS KNOWN AS LATENT SEMANTIC ANALYSIS LSA

THIS ESTIMATOR SUPPORTS TWO ALGORITHMS A FAST RANDOMIZED SVD SOLVER AND A “NAIVE” ALGORITHM THAT USES

ARPACK AS AN EIGENSOLVER ON X XT OR XT X WHICHEVER IS MORE EFFICIENT

READ MORE IN THE USER GUIDE

PARAMETERS

1666 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

NCOMPONENTS INT DEFAULT 2 DESIRED DIMENSIONALITY OF OUTPUT DATA MUST BE STRICTLY LESS THAN THE NUMBER OF FEATURES THE DEFAULT VALUE IS USEFUL FOR VISUALISATION FOR LSA A VALUE OF 100 IS RECOMMENDED

ALGORITHM STRING DEFAULT "RANDOMIZED" SVD SOLVER TO USE EITHER "ARPACK" FOR THE ARPACK WRAPPER IN SCIPY SCIPYSPARSELINALGSVDS OR "RANDOMIZED" FOR THE RANDOMIZED ALGORITHM DUE TO HALKO 2009

NITER INT OPTIONAL DEFAULT 5 NUMBER OF ITERATIONS FOR RANDOMIZED SVD SOLVER NOT USED BY ARPACK THE DEFAULT IS LARGER THAN THE DEFAULT IN RANDOMIZEDSVD TO HANDLE SPARSE MATRICES THAT MAY HAVE LARGE SLOWLY DECAYING SPECTRUM

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULT NONE IF INT RANDOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NPRANDOM

TOLFLOAT OPTIONAL TOLERANCE FOR ARPACK 0 MEANS MACHINE PRECISION IGNORED BY RANDOMIZED SVD SOLVER

ATTRIBUTES

COMPONENTS ARRAY SHAPE NCOMPONENTS NFEATURES

EXPLAINEDVARIANCE ARRAY SHAPE NCOMPONENTS THE VARIANCE OF THE TRAINING SAMPLES TRANSFORMED BY A PROJECTION TO EACH COMPONENT

EXPLAINEDVARIANCERATIO ARRAY SHAPE NCOMPONENTS PERCENTAGE OF VARIANCE EXPLAINED BY EACH OF THE SELECTED COMPONENTS

SINGULARVALUES ARRAY SHAPE NCOMPONENTS THE SINGULAR VALUES CORRESPONDING TO EACH OF THE SELECTED COMPONENTS THE SINGULAR VALUES ARE EQUAL TO THE 2NORMS OF THE NCOMPONENTS VARIABLES IN THE LOWERDIMENSIONAL SPACE

SEE ALSO

PCA

NOTES

SVD SUFFERS FROM A PROBLEM CALLED "SIGN INDETERMINACY" WHICH MEANS THE SIGN OF THE COMPONENTS AND THE OUTPUT FROM TRANSFORM DEPEND ON THE ALGORITHM AND RANDOM STATE TO WORK AROUND THIS FIT INSTANCES OF THIS CLASS TO DATA ONCE THEN KEEP THE INSTANCE AROUND TO DO TRANSFORMATIONS

REFERENCES

FINDING STRUCTURE WITH RANDOMNESS STOCHASTIC ALGORITHMS FOR CONSTRUCTING APPROXIMATE MATRIX DECOMPOSITIONS HALKO ET AL 2009 ARXIV909 [HTTPSARXIVORGPDF09094061PDF](https://arxiv.org/pdf/0909.4061.pdf)

EXAMPLES

```
FROM SKLEARNDECOMPOSITION IMPORT TRUNCATEDSVD
FROM SKLEARNRANDOMPROJECTION IMPORT SPARSERANDOMMATRIX
X SPARSERANDOMMATRIX100 100 DENSITY001 RANDOMSTATE42
SVD TRUNCATEDSVDNCOMPONENTS5 NITER7 RANDOMSTATE42
SVDFITX
69SKLEARNDECOMPOSITION MATRIX DECOMPOSITION 1667
```

SCIKITLEARN USER GUIDE RELEASE 0213  
TRUNCATEDSVDALGORITHMRANDOMIZED NCOMPONENTS5 NITER7  
RANDOMSTATE42 TOL00  
PRINTSVDEXPLAINEDVARIANCERATIO  
00606 00584 00497 00434 00372  
PRINTSVDEXPLAINEDVARIANCERATIOSUM  
0249  
PRINTSVDSINGULARVALUES  
25841 25245 23201 21753 20443  
METHODS  
FITSELF X Y FIT LSI MODEL ON TRAINING DATA X  
FITTRANSFORM SELF X Y FIT LSI MODEL TO X AND PERFORM DIMENSIONALITY REDUC  
TION ON X  
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR  
INVERSETRANSFORM SELF X TRANSFORM X BACK TO ITS ORIGINAL SPACE  
SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR  
TRANSFORM SELF X PERFORM DIMENSIONALITY REDUCTION ON X  
INIT SELF NCOMPONENTS2 ALGORITHM'RANDOMIZED' NITER5 RANDOMSTATENONE TOL00  
FITSELFXYNONE  
FIT LSI MODEL ON TRAINING DATA X  
PARAMETERS  
XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES TRAINING DATA  
YIGNORED  
RETURNS  
SELF OBJECT RETURNS THE TRANSFORMER OBJECT  
FITTRANSFORM SELFXYNONE  
FIT LSI MODEL TO X AND PERFORM DIMENSIONALITY REDUCTION ON X  
PARAMETERS  
XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES TRAINING DATA  
YIGNORED  
RETURNS  
XNEW ARRAY SHAPE NSAMPLES NCOMPONENTS REDUCED VERSION OF X THIS WILL ALWAYS  
BE A DENSE ARRAY  
GETPARAMS SELFDEEPTTRUE  
GET PARAMETERS FOR THIS ESTIMATOR  
PARAMETERS  
DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBOBJECTS THAT ARE ESTIMATORS  
RETURNS  
PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES  
1668 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

INVERSETRANSFORM SELF  
X  
TRANSFORM X BACK TO ITS ORIGINAL SPACE  
RETURNS AN ARRAY XORIGINAL WHOSE TRANSFORM WOULD BE X  
PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NCOMPONENTS NEW DATA  
RETURNS  
XORIGINAL ARRAY SHAPE NSAMPLES NFEATURES NOTE THAT THIS IS ALWAYS A DENSE ARRAY  
SETPARAMS SELFPARAMS  
SET THE PARAMETERS OF THIS ESTIMATOR  
THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE  
PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT  
OF A NESTED OBJECT  
RETURNS  
SELF  
TRANSFORM SELF  
X  
PERFORM DIMENSIONALITY REDUCTION ON X  
PARAMETERS  
XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES NEW DATA  
RETURNS  
XNEW ARRAY SHAPE NSAMPLES NCOMPONENTS REDUCED VERSION OF X THIS WILL ALWAYS  
BE A DENSE ARRAY  
EXAMPLES USING SKLEARNDECOMPOSITIONTRUNCATEDSVD  
•COLUMN TRANSFORMER WITH HETEROGENEOUS DATA SOURCES  
•HASHING FEATURE TRANSFORMATION USING TOTALLY RANDOM TREES  
•MANIFOLD LEARNING ON HANDWRITTEN DIGITS LOCALLY LINEAR EMBEDDING ISOMAP  
•CLUSTERING TEXT DOCUMENTS USING KMEANS  
DECOMPOSITIONDICTLEARNING X  
NCOMPONENTS SOLVES A DICTIONARY LEARNING MATRIX FACTORIZATION PROBLEM  
DECOMPOSITIONDICTLEARNINGONLINE X  
SOLVES A DICTIONARY LEARNING MATRIX FACTORIZATION PROBLEM  
ONLINE  
DECOMPOSITIONFASTICA X NCOMPONENTS PERFORM FAST INDEPENDENT COMPONENT ANALYSIS  
DECOMPOSITIONNONNEGATIVEFACTORIZATION X COMPUTE NONNEGATIVE MATRIX FACTORIZATION NMF  
DECOMPOSITIONSPARSEENCODE X DICTIONARY  
SPARSE CODING  
69SKLEARNDECOMPOSITION MATRIX DECOMPOSITION 1669

SCIKITLEARN USER GUIDE RELEASE 0213  
69145KLEARNDECOMPOSITION DICTLEARNING  
SKLEARNDECOMPOSITION DICTLEARNING XNCOMPONENTS ALPHA MAXITER100 TOL1E  
08METHOD'LARS' NJOBSNONE DICTINITNONE  
CODEINITNONE CALLBACKNONE VERBOSEFALSE  
RANDOMSTATENONE RETURNNITERFALSE POSI  
TIVEDICTFALSE POSITIVECODEFALSE  
SOLVES A DICTIONARY LEARNING MATRIX FACTORIZATION PROBLEM  
FINDS THE BEST DICTIONARY AND THE CORRESPONDING SPARSE CODE FOR APPROXIMATING THE DATA MATRIX X BY SOLVING  
U V ARGMIN 05 X U V 22 ALPHA U 1  
UV  
WITH VK 2 1 FORALL 0 K NCOMPONENTS  
WHERE V IS THE DICTIONARY AND U IS THE SPARSE CODE  
READ MORE IN THE USER GUIDE  
PARAMETERS  
XARRAY OF SHAPE NSAMPLES NFEATURES DATA MATRIX  
NCOMPONENTS INT NUMBER OF DICTIONARY ATOMS TO EXTRACT  
ALPHA INT SPARSITY CONTROLLING PARAMETER  
MAXITER INT MAXIMUM NUMBER OF ITERATIONS TO PERFORM  
TOLFLOAT TOLERANCE FOR THE STOPPING CONDITION  
METHOD 'LARS' 'CD' LARS USES THE LEAST ANGLE REGRESSION METHOD TO SOLVE THE LASSO PROB  
LEM LINEARMODELLARSPATH CD USES THE COORDINATE DESCENT METHOD TO COMPUTE THE LASSO  
SOLUTION LINEARMODELLASSO LARS WILL BE FASTER IF THE ESTIMATED COMPONENTS ARE SPARSE  
NJOBS INT OR NONE OPTIONAL DEFAULTNONE NUMBER OF PARALLEL JOBS TO RUN NONE MEANS 1  
UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL PROCESSORS SEE  
GLOSSARY FOR MORE DETAILS  
DICTINIT ARRAY OF SHAPE NCOMPONENTS NFEATURES INITIAL VALUE FOR THE DICTIONARY FOR WARM  
RESTART SCENARIOS  
CODEINIT ARRAY OF SHAPE NSAMPLES NCOMPONENTS INITIAL VALUE FOR THE SPARSE CODE FOR  
WARM RESTART SCENARIOS  
CALLBACK CALLABLE OR NONE OPTIONAL DEFAULT NONE CALLABLE THAT GETS INVOKED EVERY FIVE ITER  
ATIONS  
VERBOSE BOOL OPTIONAL DEFAULT FALSE TO CONTROL THE VERBOSITY OF THE PROCEDURE  
RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RAN  
DOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE  
RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS  
THE RANDOMSTATE INSTANCE USED BY NPRANDOM  
RETURNNITER BOOL WHETHER OR NOT TO RETURN THE NUMBER OF ITERATIONS  
POSITIVEDICT BOOL WHETHER TO ENFORCE POSITIVITY WHEN FINDING THE DICTIONARY  
NEW IN VERSION 020  
POSITIVECODE BOOL WHETHER TO ENFORCE POSITIVITY WHEN FINDING THE CODE  
NEW IN VERSION 020  
1670 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

RETURNS

CODE ARRAY OF SHAPE NSAMPLES NCOMPONENTS THE SPARSE CODE FACTOR IN THE MATRIX FACTORIZATION

DICTIONARY ARRAY OF SHAPE NCOMPONENTS NFEATURES THE DICTIONARY FACTOR IN THE MATRIX FACTORIZATION

ERRORS ARRAY VECTOR OF ERRORS AT EACH ITERATION

NITER INT NUMBER OF ITERATIONS RUN RETURNED ONLY IF RETURNNITER IS SET TO TRUE

SEE ALSO

DICTLEARNINGONLINE

DICTIONARYLEARNING

MINIBATCHDICTIONARYLEARNING

SPARSEPCA

MINIBATCHSPARSEPCA

6915SKLEARNDECOMPOSITION DICTLEARNINGONLINE

SKLEARNDECOMPOSITION DICTLEARNINGONLINE XNCOMPONENTS2 ALPHA1 NITER100

RETURNCODETRUE DICTINITNONE CALL

BACKNONE BATCHSIZE3 VERBOSEFALSE

SHUFFLETRUE NJOBSNONE METHOD'LARS'

ITEROFFSET0 RANDOMSTATENONE RE

TURNINNERSTATSFALSE INNERSTATSNONE

RETURNNITERFALSE POSITIVEDICTFALSE

POSITIVECODEFALSE

SOLVES A DICTIONARY LEARNING MATRIX FACTORIZATION PROBLEM ONLINE

FINDS THE BEST DICTIONARY AND THE CORRESPONDING SPARSE CODE FOR APPROXIMATING THE DATA MATRIX X BY SOLVING

$U V \argmin_{U, V} \|X - UV\|_2^2$  ALPHA U 1

UV

WITH VK 2 1 FORALL 0 K NCOMPONENTS

WHERE V IS THE DICTIONARY AND U IS THE SPARSE CODE THIS IS ACCOMPLISHED BY REPEATEDLY ITERATING OVER MINI

BATCHES BY SLICING THE INPUT DATA

READ MORE IN THE USER GUIDE

PARAMETERS

XARRAY OF SHAPE NSAMPLES NFEATURES DATA MATRIX

NCOMPONENTS INT NUMBER OF DICTIONARY ATOMS TO EXTRACT

ALPHA FLOAT SPARSITY CONTROLLING PARAMETER

NITER INT NUMBER OF ITERATIONS TO PERFORM

RETURNCODE BOOLEAN WHETHER TO ALSO RETURN THE CODE U OR JUST THE DICTIONARY V

DICTINIT ARRAY OF SHAPE NCOMPONENTS NFEATURES INITIAL VALUE FOR THE DICTIONARY FOR WARM

RESTART SCENARIOS

69SKLEARNDECOMPOSITION MATRIX DECOMPOSITION 1671

SCIKITLEARN USER GUIDE RELEASE 0213

CALLBACK CALLABLE OR NONE OPTIONAL DEFAULT NONE CALLABLE THAT GETS INVOKED EVERY FIVE ITERATIONS

BATCHSIZE INT THE NUMBER OF SAMPLES TO TAKE IN EACH BATCH

VERBOSE BOOL OPTIONAL DEFAULT FALSE TO CONTROL THE VERBOSITY OF THE PROCEDURE

SHUFFLE BOOLEAN WHETHER TO SHUFFLE THE DATA BEFORE SPLITTING IT IN BATCHES

NJOBS INT OR NONE OPTIONAL DEFAULTNONE NUMBER OF PARALLEL JOBS TO RUN NONE MEANS 1 UNLESS IN A JOBLIBPARALLELBKEND CONTEXT1MEANS USING ALL PROCESSORS SEE

GLOSSARY FOR MORE DETAILS

METHOD ‘LARS’ ‘CD’ LARS USES THE LEAST ANGLE REGRESSION METHOD TO SOLVE THE LASSO PROBLEM LINEARMODELLARSPATH CD USES THE COORDINATE DESCENT METHOD TO COMPUTE THE LASSO SOLUTION LINEARMODELLASSO LARS WILL BE FASTER IF THE ESTIMATED COMPONENTS ARE SPARSE

ITEROFFSET INT DEFAULT 0 NUMBER OF PREVIOUS ITERATIONS COMPLETED ON THE DICTIONARY USED FOR INITIALIZATION

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RANDOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NPRANDOM

RETURNINNERSTATS BOOLEAN OPTIONAL RETURN THE INNER STATISTICS A DICTIONARY COVARIANCE AND B DATA APPROXIMATION USEFUL TO RESTART THE ALGORITHM IN AN ONLINE SETTING IF RETURNINNERSTATS IS TRUE RETURNCODE IS IGNORED

INNERSTATS TUPLE OF A B NDARRAYS INNER SUFFICIENT STATISTICS THAT ARE KEPT BY THE ALGORITHM PASSING THEM AT INITIALIZATION IS USEFUL IN ONLINE SETTINGS TO AVOID LOOSING THE HISTORY OF THE EVOLUTION A NCOMPONENTS NCOMPONENTS IS THE DICTIONARY COVARIANCE MATRIX B NFEATURES NCOMPONENTS IS THE DATA APPROXIMATION MATRIX

RETURNNITER BOOL WHETHER OR NOT TO RETURN THE NUMBER OF ITERATIONS

POSITIVEDICT BOOL WHETHER TO ENFORCE POSITIVITY WHEN FINDING THE DICTIONARY

NEW IN VERSION 020

POSITIVECODE BOOL WHETHER TO ENFORCE POSITIVITY WHEN FINDING THE CODE

NEW IN VERSION 020

RETURNS

CODE ARRAY OF SHAPE NSAMPLES NCOMPONENTS THE SPARSE CODE ONLY RETURNED IF RETURNCODETRUE

DICTIONARY ARRAY OF SHAPE NCOMPONENTS NFEATURES THE SOLUTIONS TO THE DICTIONARY LEARNING PROBLEM

NITER INT NUMBER OF ITERATIONS RUN RETURNED ONLY IF RETURNNITER IS SET TOTRUE

SEE ALSO

DICTLEARNING

DICTIONARYLEARNING

MINIBATCHDICTIONARYLEARNING

SPARSEPCA

MINIBATCHSPARSEPCA

1672 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213  
6916SKLEARNDECOMPOSITION FASTICA  
SKLEARNDECOMPOSITION FASTICAXNCOMPONENTSNONE ALGORITHM'PARALLEL' WHITENTRUE  
FUN'LOGCOSH' FUNARGSNONE MAXITER200 TOL00001  
WINITNONE RANDOMSTATENONE RETURNXMEANFALSE  
COMPUTESOURCETRUE RETURNNITERFALSE  
PERFORM FAST INDEPENDENT COMPONENT ANALYSIS  
READ MORE IN THE USER GUIDE  
PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING VECTOR WHERE NSAMPLES IS THE NUMBER  
OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES  
NCOMPONENTS INT OPTIONAL NUMBER OF COMPONENTS TO EXTRACT IF NONE NO DIMENSION REDUC  
TION IS PERFORMED  
ALGORITHM 'PARALLEL' 'DEFLATION' OPTIONAL APPLY A PARALLEL OR DEFLATIONAL FASTICA ALGO  
RITHM  
WHITEN BOOLEAN OPTIONAL IF TRUE PERFORM AN INITIAL WHITENING OF THE DATA IF FALSE THE DATA  
IS ASSUMED TO HAVE ALREADY BEEN PREPROCESSED IT SHOULD BE CENTERED NORMED AND WHITE  
OTHERWISE YOU WILL GET INCORRECT RESULTS IN THIS CASE THE PARAMETER NCOMPONENTS WILL BE  
IGNORED  
FUN STRING OR FUNCTION OPTIONAL DEFAULT 'LOGCOSH' THE FUNCTIONAL FORM OF THE G FUNCTION  
USED IN THE APPROXIMATION TO NEGENTROPY COULD BE EITHER 'LOGCOSH' 'EXP' OR 'CUBE' YOU  
CAN ALSO PROVIDE YOUR OWN FUNCTION IT SHOULD RETURN A TUPLE CONTAINING THE VALUE OF THE  
FUNCTION AND OF ITS DERIVATIVE IN THE POINT THE DERIVATIVE SHOULD BE AVERAGED ALONG ITS LAST  
DIMENSION EXAMPLE  
DEF MYGX RETURN X 3 NPMEAN3 X 2 AXIS1  
FUNARGS DICTIONARY OPTIONAL ARGUMENTS TO SEND TO THE FUNCTIONAL FORM IF EMPTY OR NONE  
AND IF FUN'LOGCOSH' FUNARGS WILL TAKE VALUE 'ALPHA' 10  
MAXITER INT OPTIONAL MAXIMUM NUMBER OF ITERATIONS TO PERFORM  
TOLFLOAT OPTIONAL A POSITIVE SCALAR GIVING THE TOLERANCE AT WHICH THE UNMIXING MATRIX IS  
CONSIDERED TO HAVE CONVERGED  
WINIT NCOMPONENTS NCOMPONENTS ARRAY OPTIONAL INITIAL UNMIXING ARRAY OF DIMENSION  
NCOMPNCOMP IF NONE DEFAULT THEN AN ARRAY OF NORMAL RV'S IS USED  
RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RAN  
DOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE  
RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS  
THE RANDOMSTATE INSTANCE USED BY NPRANDOM  
RETURNXMEAN BOOL OPTIONAL IF TRUE XMEAN IS RETURNED TOO  
COMPUTESOURCES BOOL OPTIONAL IF FALSE SOURCES ARE NOT COMPUTED BUT ONLY THE ROTATION  
MATRIX THIS CAN SAVE MEMORY WHEN WORKING WITH BIG DATA DEFAULTS TO TRUE  
RETURNNITER BOOL OPTIONAL WHETHER OR NOT TO RETURN THE NUMBER OF ITERATIONS  
RETURNS  
KARRAY SHAPE NCOMPONENTS NFEATURES NONE IF WHITEN IS 'TRUE' K IS THE PREWHITENING  
MATRIX THAT PROJECTS DATA ONTO THE FIRST NCOMPONENTS PRINCIPAL COMPONENTS IF WHITEN IS  
'FALSE' K IS 'NONE'  
69SKLEARNDECOMPOSITION MATRIX DECOMPOSITION 1673

SCIKITLEARN USER GUIDE RELEASE 0213

WARRAY SHAPE NCOMPONENTS NCOMPONENTS ESTIMATED UNMIXING MATRIX THE MIXING  
MATRIX CAN BE OBTAINED BY

W NPDOTW KT

A WT WWTI

SARRAY SHAPE NSAMPLES NCOMPONENTS NONE ESTIMATED SOURCE MATRIX

XMEAN ARRAY SHAPE NFEATURES THE MEAN OVER FEATURES RETURNED ONLY IF RETURNXMEAN  
IS TRUE

NITER INT IF THE ALGORITHM IS “DEFLATION” NITER IS THE MAXIMUM NUMBER OF ITERATIONS RUN  
ACROSS ALL COMPONENTS ELSE THEY ARE JUST THE NUMBER OF ITERATIONS TAKEN TO CONVERGE THIS IS  
RETURNED ONLY WHEN RETURNNITER IS SET TO TRUE

NOTES

THE DATA MATRIX X IS CONSIDERED TO BE A LINEAR COMBINATION OF NONGAUSSIAN INDEPENDENT COMPONENTS IE X  
AS WHERE COLUMNS OF S CONTAIN THE INDEPENDENT COMPONENTS AND A IS A LINEAR MIXING MATRIX IN SHORT ICA  
ATTEMPTS TO UNMIX THE DATA BY ESTIMATING AN UNMIXING MATRIX W WHERE S W  
K X

THIS IMPLEMENTATION WAS ORIGINALLY MADE FOR DATA OF SHAPE NFEATURES NSAMPLES NOW THE INPUT IS TRANSPOSED  
BEFORE THE ALGORITHM IS APPLIED THIS MAKES IT SLIGHTLY FASTER FOR FORTRANORDERED INPUT  
IMPLEMENTED USING FASTICA A HYVARINEN AND E OJA INDEPENDENT COMPONENT ANALYSIS ALGORITHMS AND  
APPLICATIONS NEURAL NETWORKS 1345 2000 PP 411430  
6917SKLEARNDECOMPOSITION NONNEGATIVEFACTORIZATION  
SKLEARNDECOMPOSITION NONNEGATIVEFACTORIZATION X WNONE HNONE  
NCOMPONENTSNONE INIT‘WARN’  
UPDATEHTRUE SOLVER‘CD’  
BETALOSS‘FROBENIUS’ TOL00001  
MAXITER200 ALPHA00  
L1RATIO00 REGULARIZATIONNONE  
RANDOMSTATENONE VERBOSE0  
SHUFFLEFALSE

COMPUTE NONNEGATIVE MATRIX FACTORIZATION NMF

FIND TWO NONNEGATIVE MATRICES W H WHOSE PRODUCT APPROXIMATES THE NON NEGATIVE MATRIX X THIS FACTORIZA  
TION CAN BE USED FOR EXAMPLE FOR DIMENSIONALITY REDUCTION SOURCE SEPARATION OR TOPIC EXTRACTION

THE OBJECTIVE FUNCTION IS

05X WHFRO2

ALPHA L1RATIO VECW1

ALPHA L1RATIO VECH1

05ALPHA1 L1RATIO WFRO2

05ALPHA1 L1RATIO HFRO2

WHERE

AFRO2 SUMIJ AIJ2 FROBENIUS NORM

VECA1 SUMIJ ABSAIJ ELEMENTWISE L1 NORM

1674 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

FOR MULTIPLICATIVEUPDATE 'MU' SOLVER THE FROBENIUS NORM 05 X WHFRO2 CAN BE CHANGED INTO ANOTHER BETADIVERGENCE LOSS BY CHANGING THE BETALOSS PARAMETER

THE OBJECTIVE FUNCTION IS MINIMIZED WITH AN ALTERNATING MINIMIZATION OF W AND H IF H IS GIVEN AND UP DATEHFALSE IT SOLVES FOR W ONLY

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES CONSTANT MATRIX

WARRAYLIKE SHAPE NSAMPLES NCOMPONENTS IF INIT'CUSTOM' IT IS USED AS INITIAL GUESS FOR THE SOLUTION

HARRAYLIKE SHAPE NCOMPONENTS NFEATURES IF INIT'CUSTOM' IT IS USED AS INITIAL GUESS FOR THE SOLUTION IF UPDATEHFALSE IT IS USED AS A CONSTANT TO SOLVE FOR W ONLY

NCOMPONENTS INTEGER NUMBER OF COMPONENTS IF NCOMPONENTS IS NOT SET ALL FEATURES ARE KEPT

INIT NONE 'RANDOM' 'NNDSDVD' 'NNDSDVDA' 'NNDSDVDAR' 'CUSTOM' METHOD USED TO INITIALIZE THE PROCEDURE DEFAULT 'RANDOM'

THE DEFAULT VALUE WILL CHANGE FROM 'RANDOM' TO NONE IN VERSION 023 TO MAKE IT CONSISTENT WITH DECOMPOSITIONNMF

VALID OPTIONS

- NONE 'NNDSDVD' IF NCOMPONENTS NFEATURES OTHERWISE 'RANDOM'
- 'RANDOM' NONNEGATIVE RANDOM MATRICES SCALED WITH SQRTXMEAN

NCOMPONENTS

- 'NNDSDVD' NONNEGATIVE DOUBLE SINGULAR VALUE DECOMPOSITION NNDSDVD

INITIALIZATION BETTER FOR SPARSENESS

- 'NNDSDVDA' NNDSDVD WITH ZEROS FILLED WITH THE AVERAGE OF X BETTER WHEN SPARSITY IS NOT DESIRED
- 'NNDSDVDAR' NNDSDVD WITH ZEROS FILLED WITH SMALL RANDOM VALUES GENERALLY FASTER

LESS ACCURATE ALTERNATIVE TO NNDSDVDA FOR WHEN SPARSITY IS NOT DESIRED

- 'CUSTOM' USE CUSTOM MATRICES W AND H

UPDATEH BOOLEAN DEFAULT TRUE SET TO TRUE BOTH W AND H WILL BE ESTIMATED FROM INITIAL GUESSES SET TO FALSE ONLY W WILL BE ESTIMATED

SOLVER 'CD' 'MU' NUMERICAL SOLVER TO USE 'CD' IS A COORDINATE DESCENT SOLVER THAT USES FAST HIERARCHICAL

ALTERNATING LEAST SQUARES FAST HALS

'MU' IS A MULTIPLICATIVE UPDATE SOLVER

NEW IN VERSION 017 COORDINATE DESCENT SOLVER

NEW IN VERSION 019 MULTIPLICATIVE UPDATE SOLVER

BETALOSS FLOAT OR STRING DEFAULT 'FROBENIUS' STRING MUST BE IN 'FROBENIUS' 'KULLBACKLEIBLER' 'ITAKURASAITO' BETA DIVERGENCE TO BE MINIMIZED MEASURING THE DISTANCE BETWEEN X AND THE DOT PRODUCT WH NOTE THAT VALUES DIFFERENT FROM 'FROBENIUS' OR 2 AND 'KULLBACKLEIBLER' OR 1 LEAD TO SIGNIFICANTLY SLOWER FITS NOTE THAT FOR BETALOSS 0 OR 'ITAKURASAITO' THE INPUT MATRIX X CANNOT CONTAIN ZEROS USED ONLY IN 'MU' SOLVER

NEW IN VERSION 019

TOLFLOAT DEFAULT 1E4 TOLERANCE OF THE STOPPING CONDITION

69SKLEARNDECOMPOSITION MATRIX DECOMPOSITION 1675

SCIKITLEARN USER GUIDE RELEASE 0213

MAXITER INTEGER DEFAULT 200 MAXIMUM NUMBER OF ITERATIONS BEFORE TIMING OUT

ALPHA DOUBLE DEFAULT 0 CONSTANT THAT MULTIPLIES THE REGULARIZATION TERMS

L1RATIO DOUBLE DEFAULT 0 THE REGULARIZATION MIXING PARAMETER WITH 0 L1RATIO 1 FOR L1RATIO 0 THE PENALTY IS AN ELEMENTWISE L2 PENALTY AKA FROBENIUS NORM FOR L1RATIO 1 IT IS AN ELEMENTWISE L1 PENALTY FOR 0 L1RATIO 1 THE PENALTY IS A COMBINATION OF L1 AND L2

REGULARIZATION 'BOTH' 'COMPONENTS' 'TRANSFORMATION' NONE SELECT WHETHER THE REGULARIZATION AFFECTS THE COMPONENTS H THE TRANSFORMATION W BOTH OR NONE OF THEM

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULT NONE IF INT RANDOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NPRANDOM

VERBOSE INTEGER DEFAULT 0 THE VERBOSITY LEVEL

SHUFFLE BOOLEAN DEFAULT FALSE IF TRUE RANDOMIZE THE ORDER OF COORDINATES IN THE CD SOLVER

RETURNS

WARRAYLIKE SHAPE NSAMPLES NCOMPONENTS SOLUTION TO THE NONNEGATIVE LEAST SQUARES PROBLEM

HARRAYLIKE SHAPE NCOMPONENTS NFEATURES SOLUTION TO THE NONNEGATIVE LEAST SQUARES PROBLEM

NITER INT ACTUAL NUMBER OF ITERATIONS

REFERENCES

CICHOCKI ANDRZEJ AND P H A N ANHHUY "FAST LOCAL ALGORITHMS FOR LARGE SCALE NONNEGATIVE MATRIX AND TENSOR FACTORIZATIONS" IEICE TRANSACTIONS ON FUNDAMENTALS OF ELECTRONICS COMMUNICATIONS AND COMPUTER SCIENCES 923 708721 2009

FEVOTTE C IDIER J 2011 ALGORITHMS FOR NONNEGATIVE MATRIX FACTORIZATION WITH THE BETADIVERGENCE NEURAL COMPUTATION 239

EXAMPLES

```
import numpy as np
X = np.array([[1, 2, 1, 3, 12, 4, 1, 5, 0, 8, 6, 1],
              [1, 2, 1, 3, 12, 4, 1, 5, 0, 8, 6, 1]])
from sklearn.decomposition import NonnegativeFactorization
w, h = NITER NonnegativeFactorization(X, n_components=2)
init_random, random_state = 0
6918sklearn.decomposition.sparse_encode
sklearn.decomposition.sparse_encode(x, dictionary, gram=None, cov=None, algo='RITHM', lasso_lars='NNONZERO', coeffs=None, alphanone=True, copy_cov=True, init=None, max_iter=1000, n_jobs=None, check_input=True, verbose=0, positive=False)
sparse_coding
1676 CHAPTER 6 API REFERENCE
```

SCIKITLEARN USER GUIDE RELEASE 0213

EACH ROW OF THE RESULT IS THE SOLUTION TO A SPARSE CODING PROBLEM THE GOAL IS TO FIND A SPARSE ARRAY CODE SUCH THAT

X CODE DICTIONARY

READ MORE IN THE USER GUIDE

PARAMETERS

XARRAY OF SHAPE NSAMPLES NFEATURES DATA MATRIX

DICTIONARY ARRAY OF SHAPE NCOMPONENTS NFEATURES THE DICTIONARY MATRIX AGAINST WHICH TO SOLVE THE SPARSE CODING OF THE DATA SOME OF THE ALGORITHMS ASSUME NORMALIZED ROWS FOR MEANINGFUL OUTPUT

GRAM ARRAY SHAPENCOMPONENTS NCOMPONENTS PRECOMPUTED GRAM MATRIX DICTIONARY

DICTIONARY'

COV ARRAY SHAPENCOMPONENTS NSAMPLES PRECOMPUTED COVARIANCE DICTIONARY' X

ALGORITHM 'LASSOLARS' 'LASSOCD' 'LARS' 'OMP' 'THRESHOLD' LARS USES THE LEAST ANGLE REGRESSION METHOD LINEARMODELLARSPATH LASSOLARS USES LARS TO COMPUTE THE LASSO SOLUTION LASSOCD USES THE COORDINATE DESCENT METHOD TO COMPUTE THE LASSO SOLUTION LINEARMODELLASSO LASSOLARS WILL BE FASTER IF THE ESTIMATED COMPONENTS ARE SPARSE OMP USES ORTHOGONAL MATCHING PURSUIT TO ESTIMATE THE SPARSE SOLUTION THRESHOLD SQUASHES TO ZERO ALL COEFFICIENTS LESS THAN ALPHA FROM THE PROJECTION DICTIONARY X'

NNONZEROCOEF INT 01 NFEATURES BY DEFAULT NUMBER OF NONZERO COEFFICIENTS TO TARGET IN EACH COLUMN OF THE SOLUTION THIS IS ONLY USED BY ALGORITHM LARS AND ALGORITHM OMP AND IS OVERRIDDEN BY ALPHA IN THE ORTHOGONAL MATCHING PURSUIT OMP CASE

ALPHA FLOAT 1 BY DEFAULT IF ALGORITHM LASSOLARS OR ALGORITHM LASSOCD ALPHA IS THE PENALTY APPLIED TO THE L1 NORM IF ALGORITHM THRESHOLD ALPHA IS THE ABSOLUTE VALUE OF THE THRESHOLD BELOW WHICH COEFFICIENTS WILL BE SQUASHED TO ZERO IF ALGORITHM OMP ALPHA IS THE TOLERANCE PARAMETER THE VALUE OF THE RECONSTRUCTION ERROR TARGETED IN THIS CASE IT OVERRIDES NNONZEROCOEF

COPYCOV BOOLEAN OPTIONAL WHETHER TO COPY THE PRECOMPUTED COVARIANCE MATRIX IF FALSE IT MAY BE OVERRITTEN

INIT ARRAY OF SHAPE NSAMPLES NCOMPONENTS INITIALIZATION VALUE OF THE SPARSE CODES ONLY USED IF ALGORITHM LASSOCD

MAXITER INT 1000 BY DEFAULT MAXIMUM NUMBER OF ITERATIONS TO PERFORM IF ALGORITHM LASSOCD

NJOBS INT OR NONE OPTIONAL DEFAULT NONE NUMBER OF PARALLEL JOBS TO RUN NONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT 1 MEANS USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

CHECKINPUT BOOLEAN OPTIONAL IF FALSE THE INPUT ARRAYS X AND DICTIONARY WILL NOT BE CHECKED

VERBOSE INT OPTIONAL CONTROLS THE VERBOSITY THE HIGHER THE MORE MESSAGES DEFAULTS TO 0

POSITIVE BOOLEAN OPTIONAL WHETHER TO ENFORCE POSITIVITY WHEN FINDING THE ENCODING

NEW IN VERSION 020

RETURNS

CODE ARRAY OF SHAPE NSAMPLES NCOMPONENTS THE SPARSE CODES

69SKLEARNDECOMPOSITION MATRIX DECOMPOSITION 1677

SCIKITLEARN USER GUIDE RELEASE 0213  
SEE ALSO  
SKLEARNLINEARMODELLARSPATH  
SKLEARNLINEARMODELORTHOGONALMP  
SKLEARNLINEARMODELLASSO  
SPARSECODER  
610SKLEARNDISCRIMINANTANALYSIS DISCRIMINANT ANALYSIS  
LINEAR DISCRIMINANT ANALYSIS AND QUADRATIC DISCRIMINANT ANALYSIS  
USER GUIDE SEE THE LINEAR AND QUADRATIC DISCRIMINANT ANALYSIS SECTION FOR FURTHER DETAILS  
DISCRIMINANTANALYSIS  
LINEARDISCRIMINANTANALYSIS LINEAR DISCRIMINANT ANALYSIS  
DISCRIMINANTANALYSIS  
QUADRATICDISCRIMINANTANALYSIS QUADRATIC DISCRIMINANT ANALYSIS  
6101SKLEARNDISCRIMINANTANALYSIS LINEARDISCRIMINANTANALYSIS  
CLASSSKLEARNDISCRIMINANTANALYSIS LINEARDISCRIMINANTANALYSIS SOLVER'SVD'  
SHRINKAGENONE  
PRIORSNONE  
NCOMPONENTSNONE  
STORECOVARIANCEFALSE  
TOL00001  
LINEAR DISCRIMINANT ANALYSIS  
A CLASSIFIER WITH A LINEAR DECISION BOUNDARY GENERATED BY FITTING CLASS CONDITIONAL DENSITIES TO THE DATA AND USING  
BAYES' RULE  
THE MODEL FITS A GAUSSIAN DENSITY TO EACH CLASS ASSUMING THAT ALL CLASSES SHARE THE SAME COVARIANCE MATRIX  
THE FITTED MODEL CAN ALSO BE USED TO REDUCE THE DIMENSIONALITY OF THE INPUT BY PROJECTING IT TO THE MOST DISCRIM  
INATIVE DIRECTIONS  
NEW IN VERSION 017 LINEARDISCRIMINANTANALYSIS  
READ MORE IN THE USER GUIDE  
PARAMETERS  
SOLVER STRING OPTIONAL  
SOLVER TO USE POSSIBLE VALUES  
• 'SVD' SINGULAR VALUE DECOMPOSITION DEFAULT DOES NOT COMPUTE THE COVARIANCE MA  
TRIX THEREFORE THIS SOLVER IS RECOMMENDED FOR DATA WITH A LARGE NUMBER OF FEATURES  
• 'LSQR' LEAST SQUARES SOLUTION CAN BE COMBINED WITH SHRINKAGE  
• 'EIGEN' EIGENVALUE DECOMPOSITION CAN BE COMBINED WITH SHRINKAGE  
SHRINKAGE STRING OR FLOAT OPTIONAL  
SHRINKAGE PARAMETER POSSIBLE VALUES  
1678 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

- NONE NO SHRINKAGE DEFAULT
- ‘AUTO’ AUTOMATIC SHRINKAGE USING THE LEDOITWOLF LEMMA
- FLOAT BETWEEN 0 AND 1 FIXED SHRINKAGE PARAMETER

NOTE THAT SHRINKAGE WORKS ONLY WITH ‘LSQR’ AND ‘EIGEN’ SOLVERS

PRIORS ARRAY OPTIONAL SHAPE NCLASSES CLASS PRIORS

NCOMPONENTS INT OPTIONAL DEFAULTNONE NUMBER OF COMPONENTS MINNCLASSES

1 NFEATURES FOR DIMENSIONALITY REDUCTION IF NONE WILL BE SET TO MINNCLASSES 1

NFEATURES

STORECOVARIANCE BOOL OPTIONAL ADDITIONALLY COMPUTE CLASS COVARIANCE MATRIX DEFAULT

FALSE USED ONLY IN ‘SVD’ SOLVER

NEW IN VERSION 017

TOLFLOAT OPTIONAL DEFAULT 10E4 THRESHOLD USED FOR RANK ESTIMATION IN SVD SOLVER

NEW IN VERSION 017

ATTRIBUTES

COEF ARRAY SHAPE NFEATURES OR NCLASSES NFEATURES WEIGHT VECTORS

INTERCEPT ARRAY SHAPE NFEATURES INTERCEPT TERM

COVARIANCE ARRAYLIKE SHAPE NFEATURES NFEATURES COVARIANCE MATRIX SHARED BY ALL CLASSES

EXPLAINEDVARIANCERATIO ARRAY SHAPE NCOMPONENTS PERCENTAGE OF VARIANCE EXPLAINED

BY EACH OF THE SELECTED COMPONENTS IF NCOMPONENTS IS NOT SET THEN ALL COMPONENTS ARE

STORED AND THE SUM OF EXPLAINED VARIANCES IS EQUAL TO 10 ONLY AVAILABLE WHEN EIGEN OR SVD SOLVER IS USED

MEANS ARRAYLIKE SHAPE NCLASSES NFEATURES CLASS MEANS

PRIORS ARRAYLIKE SHAPE NCLASSES CLASS PRIORS SUM TO 1

SCALINGS ARRAYLIKE SHAPE RANK NCLASSES 1 SCALING OF THE FEATURES IN THE SPACE SPANNED

BY THE CLASS CENTROIDS

XBAR ARRAYLIKE SHAPE NFEATURES OVERALL MEAN

CLASSES ARRAYLIKE SHAPE NCLASSES UNIQUE CLASS LABELS

SEE ALSO

SKLEARNDISCRIMINANTANALYSISQUADRATICDISCRIMINANTANALYSIS QUADRATIC DISCRIMI

NANT ANALYSIS

NOTES

THE DEFAULT SOLVER IS ‘SVD’ IT CAN PERFORM BOTH CLASSIFICATION AND TRANSFORM AND IT DOES NOT RELY ON THE CALCULATION OF THE COVARIANCE MATRIX THIS CAN BE AN ADVANTAGE IN SITUATIONS WHERE THE NUMBER OF FEATURES IS LARGE

HOWEVER THE ‘SVD’ SOLVER CANNOT BE USED WITH SHRINKAGE

THE ‘LSQR’ SOLVER IS AN EFFICIENT ALGORITHM THAT ONLY WORKS FOR CLASSIFICATION IT SUPPORTS SHRINKAGE

THE ‘EIGEN’ SOLVER IS BASED ON THE OPTIMIZATION OF THE BETWEEN CLASS SCATTER TO WITHIN CLASS SCATTER RATIO IT CAN BE USED FOR BOTH CLASSIFICATION AND TRANSFORM AND IT SUPPORTS SHRINKAGE HOWEVER THE ‘EIGEN’ SOLVER NEEDS TO

COMPUTE THE COVARIANCE MATRIX SO IT MIGHT NOT BE SUITABLE FOR SITUATIONS WITH A HIGH NUMBER OF FEATURES

610SKLEARNDISCRIMINANTANALYSIS DISCRIMINANT ANALYSIS 1679

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES

```
import numpy as np
from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
X = np.array([1, 2, 1, 3, 2, 1, 1, 2, 1, 3, 2])
Y = np.array([1, 1, 2, 2, 2])
clf = LinearDiscriminantAnalysis()
clf.fit(X, Y)
lda = LinearDiscriminantAnalysis(n_components=None, prior=None, shrinkage=None,
solver='svd', store_covariance=False, tol=0.0001)
print(clf.predict(0.8, 1))
```

METHODS

```
decision_function(self, X) Predict confidence scores for samples
fit(self, X, Y) Fit LinearDiscriminantAnalysis model according to the
given training data and parameters
fit_transform(self, X, Y) Fit to data then transform it
get_params(self, deep=True) Get parameters for this estimator
predict(self, X) Predict class labels for samples in X
predict_log_proba(self, X) Estimate log probability
predict_proba(self, X) Estimate probability
score(self, X, Y, sample_weight=None) Returns the mean accuracy on the given test data and
labels
set_params(self, **kwargs) Set the parameters of this estimator
transform(self, X) Project data to maximize class separation
init(self, solver='svd', shrinkage=None, prior=None, n_components=None,
store_covariance=False, tol=0.0001)
decision_function(self, X) Predict confidence scores for samples
the confidence score for a sample is the signed distance of that sample to the hyperplane
parameters
X: array-like or sparse matrix, shape (n_samples, n_features) Samples
Returns
array, shape (n_samples, n_classes - 1) If n_classes == 2 else (n_samples, n_classes) Confidence
scores per sample class combination in the binary case Confidence score for
self.classes_[1] where 0 means this class would be predicted
fit(self, X, Y) Fit LinearDiscriminantAnalysis model according to the given training data and parameters
Changed in version 0.19: store_covariance has been moved to main constructor
Changed in version 0.19: tol has been moved to main constructor
Parameters
X: array-like, shape (n_samples, n_features) Training data
1680 CHAPTER 6 API REFERENCE
```

SCIKITLEARN USER GUIDE RELEASE 0213  
YARRAY SHAPE NSAMPLES TARGET VALUES  
FITTRANSFORM SELFXYNONE FITPARAMS  
FIT TO DATA THEN TRANSFORM IT  
FITS TRANSFORMER TO X AND Y WITH OPTIONAL PARAMETERS FITPARAMS AND RETURNS A TRANSFORMED VERSION OF X  
PARAMETERS  
XNUMPY ARRAY OF SHAPE NSAMPLES NFEATURES TRAINING SET  
YNUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES  
RETURNS  
XNEW NUMPY ARRAY OF SHAPE NSAMPLES NFEATURESNEW TRANSFORMED ARRAY  
GETPARAMS SELFDEEPTTRUE  
GET PARAMETERS FOR THIS ESTIMATOR  
PARAMETERS  
DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBOBJECTS THAT ARE ESTIMATORS  
RETURNS  
PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES  
PREDICTSELF  
PREDICT CLASS LABELS FOR SAMPLES IN X  
PARAMETERS  
XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES SAMPLES  
RETURNS  
CARRAY SHAPE NSAMPLES PREDICTED CLASS LABEL PER SAMPLE  
PREDICTLOGPROBA SELF  
ESTIMATE LOG PROBABILITY  
PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES INPUT DATA  
RETURNS  
CARRAY SHAPE NSAMPLES NCLASSES ESTIMATED LOG PROBABILITIES  
PREDICTPROBA SELF  
ESTIMATE PROBABILITY  
PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES INPUT DATA  
RETURNS  
CARRAY SHAPE NSAMPLES NCLASSES ESTIMATED PROBABILITIES  
SCORESELFXYSAMPLEWEIGHTNONE  
RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS  
IN MULTILABEL CLASSIFICATION THIS IS THE SUBSET ACCURACY WHICH IS A HARSH METRIC SINCE YOU REQUIRE FOR EACH  
SAMPLE THAT EACH LABEL SET BE CORRECTLY PREDICTED  
610SKLEARNDISCRIMINANTANALYSIS DISCRIMINANT ANALYSIS 1681

SCIKITLEARN USER GUIDE RELEASE 0213

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE LABELS FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT MEAN ACCURACY OF SELF PREDICTX WRT Y

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELF X

PROJECT DATA TO MAXIMIZE CLASS SEPARATION

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES INPUT DATA

RETURNS

XNEW ARRAY SHAPE NSAMPLES NCOMPONENTS TRANSFORMED DATA

EXAMPLES USING SKLEARN DISCRIMINANT ANALYSIS LINEAR DISCRIMINANT ANALYSIS

- NORMAL AND SHRINKAGE LINEAR DISCRIMINANT ANALYSIS FOR CLASSIFICATION
- LINEAR AND QUADRATIC DISCRIMINANT ANALYSIS WITH COVARIANCE ELLIPSOID
- COMPARISON OF LDA AND PCA 2D PROJECTION OF IRIS DATASET
- MANIFOLD LEARNING ON HANDWRITTEN DIGITS LOCALLY LINEAR EMBEDDING ISOMAP
- DIMENSIONALITY REDUCTION WITH NEIGHBORHOOD COMPONENTS ANALYSIS

6102 SKLEARN DISCRIMINANT ANALYSIS QUADRATIC DISCRIMINANT ANALYSIS

CLASS SKLEARN DISCRIMINANT ANALYSIS QUADRATIC DISCRIMINANT ANALYSIS PRIORS NONE

REGPARAM 00

STORE COVARIANCE FALSE

TOL 00001

QUADRATIC DISCRIMINANT ANALYSIS

A CLASSIFIER WITH A QUADRATIC DECISION BOUNDARY GENERATED BY FITTING CLASS CONDITIONAL DENSITIES TO THE DATA AND USING BAYES' RULE

THE MODEL FITS A GAUSSIAN DENSITY TO EACH CLASS

NEW IN VERSION 017 QUADRATIC DISCRIMINANT ANALYSIS

READ MORE IN THE USER GUIDE

PARAMETERS

1682 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213  
PRIORS ARRAY OPTIONAL SHAPE NCLASSES PRIORS ON CLASSES  
REGPARAM FLOAT OPTIONAL REGULARIZES THE COVARIANCE ESTIMATE AS  
1REGPARAM SIGMA REGPARAM NPEYENFEATURES  
STORECOVARIANCE BOOLEAN IF TRUE THE COVARIANCE MATRICES ARE COMPUTED AND STORED IN THE  
SELF COVARIANCE ATTRIBUTE  
NEW IN VERSION 017  
TOLFLOAT OPTIONAL DEFAULT 10E4 THRESHOLD USED FOR RANK ESTIMATION  
NEW IN VERSION 017  
ATTRIBUTES  
COVARIANCE LIST OF ARRAYLIKE SHAPE NFEATURES NFEATURES COVARIANCE MATRICES OF EACH  
CLASS  
MEANS ARRAYLIKE SHAPE NCLASSES NFEATURES CLASS MEANS  
PRIORS ARRAYLIKE SHAPE NCLASSES CLASS PRIORS SUM TO 1  
ROTATIONS LIST OF ARRAYS FOR EACH CLASS K AN ARRAY OF SHAPE NFEATURES NK WITH NK  
MINNFEATURES NUMBER OF ELEMENTS IN CLASS K IT IS THE ROTATION OF THE  
GAUSSIAN DISTRIBUTION IE ITS PRINCIPAL AXIS  
SCALINGS LIST OF ARRAYS FOR EACH CLASS K AN ARRAY OF SHAPE NK IT CONTAINS THE SCALING OF THE  
GAUSSIAN DISTRIBUTIONS ALONG ITS PRINCIPAL AXES IE THE VARIANCE IN THE ROTATED COORDINATE  
SYSTEM  
SEE ALSO  
SKLEARNDISCRIMINANTANALYSISLINEARDISCRIMINANTANALYSIS LINEAR DISCRIMINANT  
ANALYSIS  
EXAMPLES  
FROM SKLEARNDISCRIMINANTANALYSIS IMPORT QUADRATICDISCRIMINANTANALYSIS  
IMPORT NUMPY AS NP  
X NPARRAY1 1 2 1 3 2 1 1 2 1 3 2  
Y NPARRAY1 1 1 2 2 2  
CLF QUADRATICDISCRIMINANTANALYSIS  
CLFFITX Y  
  
QUADRATICDISCRIMINANTANALYSISPRIORSNONE REGPARAM00  
STORECOVARIANCEFALSE TOL00001  
PRINTCLFPREDICT08 1  
1  
METHODS  
DECISIONFUNCTION SELF X APPLY DECISION FUNCTION TO AN ARRAY OF SAMPLES  
FITSELF X Y FIT THE MODEL ACCORDING TO THE GIVEN TRAINING DATA AND  
PARAMETERS  
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR  
CONTINUED ON NEXT PAGE  
610SKLEARNDISCRIMINANTANALYSIS DISCRIMINANT ANALYSIS 1683

SCIKITLEARN USER GUIDE RELEASE 0213

TABLE 666 – CONTINUED FROM PREVIOUS PAGE

PREDICT SELF X PERFORM CLASSIFICATION ON AN ARRAY OF TEST VECTORS X

PREDICTLOGPROBA SELF X RETURN POSTERIOR PROBABILITIES OF CLASSIFICATION

PREDICTPROBA SELF X RETURN POSTERIOR PROBABILITIES OF CLASSIFICATION

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELFPRIORSNONE REGPARAM00 STORECOVARIANCEFALSE TOL00001

DECISIONFUNCTION SELF X

APPLY DECISION FUNCTION TO AN ARRAY OF SAMPLES

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES ARRAY OF SAMPLES TEST VECTORS

RETURNS

CARRAY SHAPE NSAMPLES NCLASSES OR NSAMPLES DECISION FUNCTION VALUES RELATED TO EACH CLASS PER SAMPLE IN THE TWOCLASS CASE THE SHAPE IS NSAMPLES GIVING THE LOG LIKELIHOOD RATIO OF THE POSITIVE CLASS

FITSELFXY

FIT THE MODEL ACCORDING TO THE GIVEN TRAINING DATA AND PARAMETERS

CHANGED IN VERSION 019 STORECOVARIANCES HAS BEEN MOVED TO MAIN CONSTRUCTOR AS STORECOVARIANCE

CHANGED IN VERSION 019 TOL HAS BEEN MOVED TO MAIN CONSTRUCTOR

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING VECTOR WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

YARRAY SHAPE NSAMPLES TARGET VALUES INTEGERS

GETPARAMS SELFDEEPTTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELF X

PERFORM CLASSIFICATION ON AN ARRAY OF TEST VECTORS X

THE PREDICTED CLASS C FOR EACH SAMPLE IN X IS RETURNED

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES

RETURNS

CARRAY SHAPE NSAMPLES

1684 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

PREDICTLOGPROBA SELF  
RETURN POSTERIOR PROBABILITIES OF CLASSIFICATION  
PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES ARRAY OF SAMPLETEST VECTORS  
RETURNS  
CARRAY SHAPE NSAMPLES NCLASSES POSTERIOR LOGPROBABILITIES OF CLASSIFICATION PER CLASS  
PREDICTPROBA SELF  
RETURN POSTERIOR PROBABILITIES OF CLASSIFICATION  
PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES ARRAY OF SAMPLETEST VECTORS  
RETURNS  
CARRAY SHAPE NSAMPLES NCLASSES POSTERIOR PROBABILITIES OF CLASSIFICATION PER CLASS  
SCORESELFXYSAMPLEWEIGHTNONE  
RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS  
IN MULTILABEL CLASSIFICATION THIS IS THE SUBSET ACCURACY WHICH IS A HARSH METRIC SINCE YOU REQUIRE FOR EACH SAMPLE THAT EACH LABEL SET BE CORRECTLY PREDICTED  
PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES  
YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE LABELS FOR X  
SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS  
RETURNS  
SCORE FLOAT MEAN ACCURACY OF SELF  
SETPARAMS SELF  
SET THE PARAMETERS OF THIS ESTIMATOR  
THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT  
RETURNS  
SELF  
EXAMPLES USING SKLEARN  
DISCRIMINANTANALYSIS  
QUADRATICDISCRIMINANTANALYSIS  
•CLASSIFIER COMPARISON  
•LINEAR AND QUADRATIC DISCRIMINANT ANALYSIS WITH COVARIANCE ELLIPSOID  
611SKLEARN  
DUMMY ESTIMATORS  
USER GUIDE SEE THE MODEL EVALUATION QUANTIFYING THE QUALITY OF PREDICTIONS SECTION FOR FURTHER DETAILS  
611SKLEARN  
DUMMY ESTIMATORS 1685

SCIKITLEARN USER GUIDE RELEASE 0213

DUMMYDUMMYCLASSIFIER STRATEGY DUMMYCLASSIFIER IS A CLASSIFIER THAT MAKES PREDICTIONS USING SIMPLE RULES

DUMMYDUMMYREGRESSOR STRATEGY CONSTANT DUMMYREGRESSOR IS A REGRESSOR THAT MAKES PREDICTIONS USING SIMPLE RULES

6111SKLEARNDUMMY DUMMYCLASSIFIER

CLASSSKLEARNDUMMY DUMMYCLASSIFIER STRATEGY‘STRATIFIED’ RANDOMSTATENONE CONSTANTNONE

DUMMYCLASSIFIER IS A CLASSIFIER THAT MAKES PREDICTIONS USING SIMPLE RULES

THIS CLASSIFIER IS USEFUL AS A SIMPLE BASELINE TO COMPARE WITH OTHER REAL CLASSIFIERS DO NOT USE IT FOR REAL PROBLEMS

READ MORE IN THE USER GUIDE

PARAMETERS

STRATEGY STR DEFAULT“STRATIFIED” STRATEGY TO USE TO GENERATE PREDICTIONS

- “STRATIFIED” GENERATES PREDICTIONS BY RESPECTING THE TRAINING SET’S CLASS DISTRIBUTION
- “MOSTFREQUENT” ALWAYS PREDICTS THE MOST FREQUENT LABEL IN THE TRAINING SET
- “PRIOR” ALWAYS PREDICTS THE CLASS THAT MAXIMIZES THE CLASS PRIOR LIKE “MOSTFREQUENT”

ANDPREDICTPROBA RETURNS THE CLASS PRIOR

- “UNIFORM” GENERATES PREDICTIONS UNIFORMLY AT RANDOM
- “CONSTANT” ALWAYS PREDICTS A CONSTANT LABEL THAT IS PROVIDED BY THE USER THIS IS USEFUL FOR METRICS THAT EVALUATE A NONMAJORITY CLASS

NEW IN VERSION 017 DUMMY CLASSIFIER NOW SUPPORTS PRIOR FITTING STRATEGY USING PARAMETERPRIOR

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RANDOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NRANDOM

CONSTANT INT OR STR OR ARRAY OF SHAPE NOUTPUTS THE EXPLICIT CONSTANT AS PREDICTED BY THE “CONSTANT” STRATEGY THIS PARAMETER IS USEFUL ONLY FOR THE “CONSTANT” STRATEGY

ATTRIBUTES

CLASSES ARRAY OR LIST OF ARRAY OF SHAPE NCLASSES CLASS LABELS FOR EACH OUTPUT

NCLASSES ARRAY OR LIST OF ARRAY OF SHAPE NCLASSES NUMBER OF LABEL FOR EACH OUTPUT

CLASSPRIOR ARRAY OR LIST OF ARRAY OF SHAPE NCLASSES PROBABILITY OF EACH CLASS FOR EACH OUTPUT

NOUTPUTS INT NUMBER OF OUTPUTS

SPARSEOUTPUT BOOL TRUE IF THE ARRAY RETURNED FROM PREDICT IS TO BE IN SPARSE CSC FORMAT IS AUTOMATICALLY SET TO TRUE IF THE INPUT Y IS PASSED IN SPARSE FORMAT

METHODS

1686 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

FITSELF X Y SAMPLEWEIGHT FIT THE RANDOM CLASSIFIER

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PREDICT SELF X PERFORM CLASSIFICATION ON TEST VECTORS X

PREDICTLOGPROBA SELF X RETURN LOG PROBABILITY ESTIMATES FOR THE TEST VECTORS X

PREDICTPROBA SELF X RETURN PROBABILITY ESTIMATES FOR THE TEST VECTORS X

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELFSTRATEGY'STRATIFIED' RANDOMSTATENONE CONSTANTNONE

FITSELFXYSAMPLEWEIGHTNONE

FIT THE RANDOM CLASSIFIER

PARAMETERS

XARRAYLIKE OBJECT WITH FINITE LENGTH OR SHAPE TRAINING DATA REQUIRES LENGTH NSAMPLES

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TARGET VALUES

SAMPLEWEIGHT ARRAYLIKE OF SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SELF OBJECT

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELF

PERFORM CLASSIFICATION ON TEST VECTORS X

PARAMETERS

XARRAYLIKE OBJECT WITH FINITE LENGTH OR SHAPE TRAINING DATA REQUIRES LENGTH NSAMPLES

RETURNS

YARRAY SHAPE NSAMPLES OR NSAMPLES NOUTPUTS PREDICTED TARGET VALUES FOR X

PREDICTLOGPROBA SELF

RETURN LOG PROBABILITY ESTIMATES FOR THE TEST VECTORS X

PARAMETERS

XARRAYLIKE OBJECT WITH FINITE LENGTH OR SHAPE TRAINING DATA REQUIRES LENGTH NSAMPLES

RETURNS

PARRAYLIKE OR LIST OF ARRAYLIKE OF SHAPE NSAMPLES NCLASSES RETURNS THE LOG PROBABILITY OF THE SAMPLE FOR EACH CLASS IN THE MODEL WHERE CLASSES ARE ORDERED ARITHMETICALLY FOR EACH OUTPUT

611SKLEARNDUMMY DUMMY ESTIMATORS 1687

SCIKITLEARN USER GUIDE RELEASE 0213

PREDICTPROBA SELF  
RETURN PROBABILITY ESTIMATES FOR THE TEST VECTORS X

PARAMETERS  
XARRAYLIKE OBJECT WITH FINITE LENGTH OR SHAPE TRAINING DATA REQUIRES LENGTH  
NSAMPLES

RETURNS  
PARRAYLIKE OR LIST OF ARRAYLIKE OF SHAPE NSAMPLES NCLASSES RETURNS THE PROBABILITY  
OF THE SAMPLE FOR EACH CLASS IN THE MODEL WHERE CLASSES ARE ORDERED ARITHMETICALLY FOR  
EACH OUTPUT

SCORESELFXYSAMPLEWEIGHTNONE  
RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS  
IN MULTILABEL CLASSIFICATION THIS IS THE SUBSET ACCURACY WHICH IS A HARSH METRIC SINCE YOU REQUIRE FOR EACH  
SAMPLE THAT EACH LABEL SET BE CORRECTLY PREDICTED

PARAMETERS  
XARRAYLIKE NONE TEST SAMPLES WITH SHAPE NSAMPLES NFEATURES OR NONE PASSING  
NONE AS TEST SAMPLES GIVES THE SAME RESULT AS PASSING REAL TEST SAMPLES SINCE DUMMYCLAS  
SIFIER OPERATES INDEPENDENTLY OF THE SAMPLED OBSERVATIONS  
YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE LABELS FOR X  
SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS  
SCORE FLOAT MEAN ACCURACY OF SELF PREDICTX WRT Y

SETPARAMS SELFPARAMS  
SET THE PARAMETERS OF THIS ESTIMATOR  
THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE  
PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT  
OF A NESTED OBJECT

RETURNS  
SELF

6112SKLEARNDUMMY DUMMYREGRESSOR  
CLASSSKLEARNDUMMY DUMMYREGRESSOR STRATEGY'MEAN' CONSTANTNONE QUANTILENONE  
DUMMYREGRESSOR IS A REGRESSOR THAT MAKES PREDICTIONS USING SIMPLE RULES  
THIS REGRESSOR IS USEFUL AS A SIMPLE BASELINE TO COMPARE WITH OTHER REAL REGRESSORS DO NOT USE IT FOR REAL  
PROBLEMS

READ MORE IN THE USER GUIDE

PARAMETERS  
STRATEGY STR STRATEGY TO USE TO GENERATE PREDICTIONS

- "MEAN" ALWAYS PREDICTS THE MEAN OF THE TRAINING SET
- "MEDIAN" ALWAYS PREDICTS THE MEDIAN OF THE TRAINING SET
- "QUANTILE" ALWAYS PREDICTS A SPECIFIED QUANTILE OF THE TRAINING SET PROVIDED WITH THE  
QUANTILE PARAMETER

1688 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

- “CONSTANT” ALWAYS PREDICTS A CONSTANT VALUE THAT IS PROVIDED BY THE USER  
CONSTANT INT OR FLOAT OR ARRAY OF SHAPE NOUTPUTS THE EXPLICIT CONSTANT AS PREDICTED BY THE  
“CONSTANT” STRATEGY THIS PARAMETER IS USEFUL ONLY FOR THE “CONSTANT” STRATEGY  
QUANTILE FLOAT IN 00 10 THE QUANTILE TO PREDICT USING THE “QUANTILE” STRATEGY A QUANTILE OF  
05 CORRESPONDS TO THE MEDIAN WHILE 00 TO THE MINIMUM AND 10 TO THE MAXIMUM

ATTRIBUTES

CONSTANT FLOAT OR ARRAY OF SHAPE NOUTPUTS MEAN OR MEDIAN OR QUANTILE OF THE TRAINING TARGETS OR CONSTANT VALUE GIVEN BY THE USER

NOUTPUTS INT NUMBER OF OUTPUTS

METHODS

FITSELF X Y SAMPLEWEIGHT FIT THE RANDOM REGRESSOR

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PREDICT SELF X RETURNSTD PERFORM CLASSIFICATION ON TEST VECTORS X

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELFSTRATEGY‘MEAN’ CONSTANTNONE QUANTILENONE

FITSELFXYSAMPLEWEIGHTNONE

FIT THE RANDOM REGRESSOR

PARAMETERS

XARRAYLIKE OBJECT WITH FINITE LENGTH OR SHAPE TRAINING DATA REQUIRES LENGTH

NSAMPLES

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TARGET VALUES

SAMPLEWEIGHT ARRAYLIKE OF SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SELF OBJECT

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELFXYRETURNSTDFALSE

PERFORM CLASSIFICATION ON TEST VECTORS X

PARAMETERS

XARRAYLIKE OBJECT WITH FINITE LENGTH OR SHAPE TRAINING DATA REQUIRES LENGTH

NSAMPLES

611SKLEARNDUMMY DUMMY ESTIMATORS 1689

SCIKITLEARN USER GUIDE RELEASE 0213

RETURNSTD BOOLEAN OPTIONAL WHETHER TO RETURN THE STANDARD DEVIATION OF POSTERIOR PREDICTION ALL ZEROS IN THIS CASE

RETURNS

YARRAY SHAPE NSAMPLES OR NSAMPLES NOUTPUTS PREDICTED TARGET VALUES FOR X

YSTD ARRAY SHAPE NSAMPLES OR NSAMPLES NOUTPUTS STANDARD DEVIATION OF PREDICTIVE DISTRIBUTION OF QUERY POINTS

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{U}{V}$  WHERE U IS THE RESIDUAL SUM OF SQUARES  $Y_{TRUE} - Y_{PRED}$  2SUM AND V IS THE TOTAL SUM OF SQUARES  $Y_{TRUE} - Y_{TRUEMEAN}$  2SUM THE BEST POSSIBLE SCORE IS 10 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 00

PARAMETERS

XARRAYLIKE NONE TEST SAMPLES WITH SHAPE NSAMPLES NFEATURES OR NONE FOR SOME ESTIMATORS THIS MAY BE A PRECOMPUTED KERNEL MATRIX INSTEAD SHAPE NSAMPLES NSAMPLESFITTED WHERE NSAMPLESFITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR PASSING NONE AS TEST SAMPLES GIVES THE SAME RESULT AS PASSING REAL TEST SAMPLES SINCE DUMMYREGRESSOR OPERATES INDEPENDENTLY OF THE SAMPLED OBSERVATIONS

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE VALUES FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT R2 OF SELF PREDICTX WRT Y

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

612SKLEARNENSEMBLE ENSEMBLE METHODS

THESKLEARNENSEMBLE MODULE INCLUDES ENSEMBLEBASED METHODS FOR CLASSIFICATION REGRESSION AND ANOMALY DETECTION

USER GUIDE SEE THE ENSEMBLE METHODS SECTION FOR FURTHER DETAILS

ENSEMBLEADABOOSTCLASSIFIER AN ADABOOST CLASSIFIER

ENSEMBLEADABOOSTREGRESSOR BASEESTIMATOR

AN ADABOOST REGRESSOR

ENSEMBLEBAGGINGCLASSIFIER BASEESTIMATOR

A BAGGING CLASSIFIER

CONTINUED ON NEXT PAGE

1690 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

TABLE 671 – CONTINUED FROM PREVIOUS PAGE

ENSEMBLEBAGGINGREGRESSOR BASEESTIMATOR

    A BAGGING REGRESSOR

ENSEMBLEEXTRATREESCLASSIFIER    AN EXTRATREES CLASSIFIER

ENSEMBLEEXTRATREESREGRESSOR NESTIMATORS

    AN EXTRATREES REGRESSOR

ENSEMBLEGRADIENTBOOSTINGCLASSIFIER LOSS

    GRADIENT BOOSTING FOR CLASSIFICATION

ENSEMBLEGRADIENTBOOSTINGREGRESSOR LOSS

    GRADIENT BOOSTING FOR REGRESSION

ENSEMBLEISOLATIONFOREST NESTIMATORS   ISOLATION FOREST ALGORITHM

ENSEMBLERANDOMFORESTCLASSIFIER    A RANDOM FOREST CLASSIFIER

ENSEMBLERANDOMFORESTREGRESSOR    A RANDOM FOREST REGRESSOR

ENSEMBLERANDOMTREEEMBEDDING    AN ENSEMBLE OF TOTALLY RANDOM TREES

ENSEMBLEVOTINGCLASSIFIER ESTIMATORS   SOFT V OTINGMAJORITY RULE CLASSIFIER FOR UNFITTED ESTIMATORS

ENSEMBLEVOTINGREGRESSOR ESTIMATORS   PREDICTION VOTING REGRESSOR FOR UNFITTED ESTIMATORS

ENSEMBLEHISTGRADIENTBOOSTINGREGRESSOR   HISTOGRAMBASED GRADIENT BOOSTING REGRESSION TREE

ENSEMBLEHISTGRADIENTBOOSTINGCLASSIFIER   HISTOGRAMBASED GRADIENT BOOSTING CLASSIFICATION TREE

6121SKLEARNENSEMBLE ADABOOSTCLASSIFIER

CLASSSSKLEARNENSEMBLE ADABOOSTCLASSIFIER BASEESTIMATORNONE NESTIMATORS50 LEARN

INGRATE10 ALGORITHM’SAMMER’ RAN

DOMSTATENONE

AN ADABOOST CLASSIFIER

AN ADABOOST 1 CLASSIFIER IS A METAESTIMATOR THAT BEGINS BY FITTING A CLASSIFIER ON THE ORIGINAL DATASET AND THEN FITS ADDITIONAL COPIES OF THE CLASSIFIER ON THE SAME DATASET BUT WHERE THE WEIGHTS OF INCORRECTLY CLASSIFIED INSTANCES ARE ADJUSTED SUCH THAT SUBSEQUENT CLASSIFIERS FOCUS MORE ON DIFFICULT CASES

THIS CLASS IMPLEMENTS THE ALGORITHM KNOWN AS ADABOOSTSAMME 2

READ MORE IN THE USER GUIDE

PARAMETERS

BASEESTIMATOR OBJECT OPTIONAL DEFAULTNONE THE BASE ESTIMATOR FROM WHICH THE BOOSTED

ENSEMBLE IS BUILT SUPPORT FOR SAMPLE WEIGHTING IS REQUIRED AS WELL AS PROPER

CLASSES ANDNCLASSES ATTRIBUTES IF NONE THEN THE BASE ESTIMATOR IS

DECISIONTREECLASSIFIERMAXDEPTH1

NESTIMATORS INTEGER OPTIONAL DEFAULT50 THE MAXIMUM NUMBER OF ESTIMATORS AT WHICH

BOOSTING IS TERMINATED IN CASE OF PERFECT FIT THE LEARNING PROCEDURE IS STOPPED EARLY

LEARNINGRATE FLOAT OPTIONAL DEFAULT1 LEARNING RATE SHRINKS THE CONTRIBUTION OF EACH

CLASSIFIER BY LEARNINGRATE THERE IS A TRADEOFF BETWEEN LEARNINGRATE AND

NESTIMATORS

ALGORITHM ‘SAMME’ ‘SAMMER’ OPTIONAL DEFAULT’SAMMER’ IF ‘SAMMER’ THEN

USE THE SAMMER REAL BOOSTING ALGORITHM BASEESTIMATOR MUST SUPPORT CALCULATION

OF CLASS PROBABILITIES IF ‘SAMME’ THEN USE THE SAMME DISCRETE BOOSTING ALGORITHM THE

SAMMER ALGORITHM TYPICALLY CONVERGES FASTER THAN SAMME ACHIEVING A LOWER TEST ERROR

WITH FEWER BOOSTING ITERATIONS

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RAN

DOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE

612SKLEARNENSEMBLE ENSEMBLE METHODS 1691

SCIKITLEARN USER GUIDE RELEASE 0213

RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NPRANDOM

ATTRIBUTES

ESTIMATORS LIST OF CLASSIFIERS THE COLLECTION OF FITTED SUBESTIMATORS

CLASSES ARRAY OF SHAPE NCLASSES THE CLASSES LABELS

NCLASSES INT THE NUMBER OF CLASSES

ESTIMATORWEIGHTS ARRAY OF FLOATS WEIGHTS FOR EACH ESTIMATOR IN THE BOOSTED ENSEMBLE

ESTIMATORERRORS ARRAY OF FLOATS CLASSIFICATION ERROR FOR EACH ESTIMATOR IN THE BOOSTED ENSEMBLE

FEATUREIMPORTANCES ARRAY OF SHAPE NFEATURES RETURN THE FEATURE IMPORTANCES THE HIGHER THE MORE IMPORTANT THE FEATURE

SEE ALSO

ADABOOSTREGRESSOR GRADIENTBOOSTINGCLASSIFIER

SKLEARNTREEDECISIONTREECLASSIFIER

REFERENCES

R33E4EC8C4AD51 R33E4EC8C4AD52

EXAMPLES

```
FROM SKLEARNENSEMBLE IMPORT ADABOOSTCLASSIFIER
FROM SKLEARNDATASETS IMPORT MAKECLASSIFICATION
X Y MAKECLASSIFICATIONNNSAMPLES1000 NFEATURES4
NINFORMATIVE2 NREDUNDANT0
RANDOMSTATE0 SHUFFLE FALSE
CLF ADABOOSTCLASSIFIERNESTIMATORS100 RANDOMSTATE0
CLFFITX Y
ADABOOSTCLASSIFIERALGORITHMSAMMER BASEESTIMATORNONE
LEARNINGRATE10 NESTIMATORS100 RANDOMSTATE0
CLFFEATUREIMPORTANCES
ARRAY028 042 014 016
CLFPREDICT0 0 0 0
ARRAY1
CLFSCOREX Y
0983
```

METHODS

DECISIONFUNCTION SELF X COMPUTE THE DECISION FUNCTION OF X

FITSELF X Y SAMPLEWEIGHT BUILD A BOOSTED CLASSIFIER FROM THE TRAINING SET X Y

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PREDICT SELF X PREDICT CLASSES FOR X

PREDICTLOGPROBA SELF X PREDICT CLASS LOGPROBABILITIES FOR X

CONTINUED ON NEXT PAGE

1692 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

TABLE 672 – CONTINUED FROM PREVIOUS PAGE

PREDICTPROBA SELF X PREDICT CLASS PROBABILITIES FOR X

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

STAGEDDECISIONFUNCTION SELF X COMPUTE DECISION FUNCTION OF XFOR EACH BOOSTING ITERATION

STAGEDPREDICT SELF X RETURN STAGED PREDICTIONS FOR X

STAGEDPREDICTPROBA SELF X PREDICT CLASS PROBABILITIES FOR X

STAGEDSCORE SELF X Y SAMPLEWEIGHT RETURN STAGED SCORES FOR X Y

INIT SELFBASEESTIMATORNONE NESTIMATORS50 LEARNINGRATE10 ALGORITHM’SAMMER’

RANDOMSTATENONE

DECISIONFUNCTION SELF X

COMPUTE THE DECISION FUNCTION OF X

PARAMETERS

XARRAYLIKE SPARSE MATRIX OF SHAPE NSAMPLES NFEATURES THE TRAINING INPUT SAMPLES SPARSE MATRIX CAN BE CSC CSR COO DOK OR LIL COO DOK AND LIL ARE CONVERTED TO CSR

RETURNS

SCORE ARRAY SHAPE NSAMPLES K THE DECISION FUNCTION OF THE INPUT SAMPLES THE ORDER OF OUTPUTS IS THE SAME OF THAT OF THE CLASSES ATTRIBUTE BINARY CLASSIFICATION IS A SPECIAL CASES WITHK 1 OTHERWISE KNCLASSES FOR BINARY CLASSIFICATION VALUES CLOSER TO 1 OR 1 MEAN MORE LIKE THE FIRST OR SECOND CLASS IN CLASSES RESPECTIVELY

FEATUREIMPORTANCES

RETURN THE FEATURE IMPORTANCES THE HIGHER THE MORE IMPORTANT THE FEATURE

RETURNS

FEATUREIMPORTANCES ARRAY SHAPE NFEATURES

FITSELFXYSAMPLEWEIGHTNONE

BUILD A BOOSTED CLASSIFIER FROM THE TRAINING SET X Y

PARAMETERS

XARRAYLIKE SPARSE MATRIX OF SHAPE NSAMPLES NFEATURES THE TRAINING INPUT SAMPLES SPARSE MATRIX CAN BE CSC CSR COO DOK OR LIL COO DOK AND LIL ARE CONVERTED TO CSR

YARRAYLIKE OF SHAPE NSAMPLES THE TARGET VALUES CLASS LABELS

SAMPLEWEIGHT ARRAYLIKE OF SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS IF NONE THE SAMPLE WEIGHTS ARE INITIALIZED TO 1 NSAMPLES

RETURNS

SELF OBJECT

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

612SKLEARNENSEMBLE ENSEMBLE METHODS 1693

SCIKITLEARN USER GUIDE RELEASE 0213

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELF X

PREDICT CLASSES FOR X

THE PREDICTED CLASS OF AN INPUT SAMPLE IS COMPUTED AS THE WEIGHTED MEAN PREDICTION OF THE CLASSIFIERS IN THE ENSEMBLE

PARAMETERS

XARRAYLIKE SPARSE MATRIX OF SHAPE NSAMPLES NFEATURES THE TRAINING INPUT SAMPLES SPARSE MATRIX CAN BE CSC CSR COO DOK OR LIL COO DOK AND LIL ARE CONVERTED TO CSR

RETURNS

YARRAY OF SHAPE NSAMPLES THE PREDICTED CLASSES

PREDICTLOGPROBA SELF X

PREDICT CLASS LOGPROBABILITIES FOR X

THE PREDICTED CLASS LOGPROBABILITIES OF AN INPUT SAMPLE IS COMPUTED AS THE WEIGHTED MEAN PREDICTED CLASS LOGPROBABILITIES OF THE CLASSIFIERS IN THE ENSEMBLE

PARAMETERS

XARRAYLIKE SPARSE MATRIX OF SHAPE NSAMPLES NFEATURES THE TRAINING INPUT SAMPLES SPARSE MATRIX CAN BE CSC CSR COO DOK OR LIL COO DOK AND LIL ARE CONVERTED TO CSR

RETURNS

PARRAY OF SHAPE NSAMPLES NCLASSES THE CLASS PROBABILITIES OF THE INPUT SAMPLES THE ORDER OF OUTPUTS IS THE SAME OF THAT OF THE CLASSES ATTRIBUTE

PREDICTPROBA SELF X

PREDICT CLASS PROBABILITIES FOR X

THE PREDICTED CLASS PROBABILITIES OF AN INPUT SAMPLE IS COMPUTED AS THE WEIGHTED MEAN PREDICTED CLASS PROBABILITIES OF THE CLASSIFIERS IN THE ENSEMBLE

PARAMETERS

XARRAYLIKE SPARSE MATRIX OF SHAPE NSAMPLES NFEATURES THE TRAINING INPUT SAMPLES SPARSE MATRIX CAN BE CSC CSR COO DOK OR LIL COO DOK AND LIL ARE CONVERTED TO CSR

RETURNS

PARRAY OF SHAPE NSAMPLES NCLASSES THE CLASS PROBABILITIES OF THE INPUT SAMPLES THE ORDER OF OUTPUTS IS THE SAME OF THAT OF THE CLASSES ATTRIBUTE

SCORESELF X Y SAMPLEWEIGHT NONE

RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS

IN MULTILABEL CLASSIFICATION THIS IS THE SUBSET ACCURACY WHICH IS A HARSH METRIC SINCE YOU REQUIRE FOR EACH SAMPLE THAT EACH LABEL SET BE CORRECTLY PREDICTED

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES

1694 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE LABELS FOR X  
SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT MEAN ACCURACY OF SELF PREDICT X WRT Y

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE  
PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT  
OF A NESTED OBJECT

RETURNS

SELF

STAGEDDECISIONFUNCTION SELF X

COMPUTE DECISION FUNCTION OF X FOR EACH BOOSTING ITERATION

THIS METHOD ALLOWS MONITORING IE DETERMINE ERROR ON TESTING SET AFTER EACH BOOSTING ITERATION

PARAMETERS

XARRAYLIKE SPARSE MATRIX OF SHAPE NSAMPLES NFEATURES THE TRAINING INPUT SAMPLES  
SPARSE MATRIX CAN BE CSC CSR COO DOK OR LIL COO DOK AND LIL ARE  
CONVERTED TO CSR

RETURNS

SCORE GENERATOR OF ARRAY SHAPE NSAMPLES K THE DECISION FUNCTION OF THE INPUT SAMPLES  
THE ORDER OF OUTPUTS IS THE SAME OF THAT OF THE CLASSES ATTRIBUTE BINARY CLASSIFICATION  
IS A SPECIAL CASES WITH K 1 OTHERWISE KNCASSES FOR BINARY CLASSIFICATION  
VALUES CLOSER TO 1 OR -1 MEAN MORE LIKE THE FIRST OR SECOND CLASS IN CLASSES RESPECTIVELY

STAGEDPREDICT SELF X

RETURN STAGED PREDICTIONS FOR X

THE PREDICTED CLASS OF AN INPUT SAMPLE IS COMPUTED AS THE WEIGHTED MEAN PREDICTION OF THE CLASSIFIERS IN  
THE ENSEMBLE

THIS GENERATOR METHOD YIELDS THE ENSEMBLE PREDICTION AFTER EACH ITERATION OF BOOSTING AND THEREFORE ALLOWS  
MONITORING SUCH AS TO DETERMINE THE PREDICTION ON A TEST SET AFTER EACH BOOST

PARAMETERS

XARRAYLIKE OF SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES SPARSE MATRIX CAN BE  
CSC CSR COO DOK OR LIL COO DOK AND LIL ARE CONVERTED TO CSR

RETURNS

YGENERATOR OF ARRAY SHAPE NSAMPLES THE PREDICTED CLASSES

STAGEDPREDICTPROBA SELF X

PREDICT CLASS PROBABILITIES FOR X

THE PREDICTED CLASS PROBABILITIES OF AN INPUT SAMPLE IS COMPUTED AS THE WEIGHTED MEAN PREDICTED CLASS  
PROBABILITIES OF THE CLASSIFIERS IN THE ENSEMBLE

THIS GENERATOR METHOD YIELDS THE ENSEMBLE PREDICTED CLASS PROBABILITIES AFTER EACH ITERATION OF BOOSTING  
AND THEREFORE ALLOWS MONITORING SUCH AS TO DETERMINE THE PREDICTED CLASS PROBABILITIES ON A TEST SET AFTER  
EACH BOOST

PARAMETERS

612SKLEARNENSEMBLE ENSEMBLE METHODS 1695

SCIKITLEARN USER GUIDE RELEASE 0213

XARRAYLIKE SPARSE MATRIX OF SHAPE NSAMPLES NFEATURES THE TRAINING INPUT SAMPLES SPARSE MATRIX CAN BE CSC CSR COO DOK OR LIL COO DOK AND LIL ARE CONVERTED TO CSR

RETURNS

PGENERATOR OF ARRAY SHAPE NSAMPLES THE CLASS PROBABILITIES OF THE INPUT SAMPLES THE ORDER OF OUTPUTS IS THE SAME OF THAT OF THE CLASSES ATTRIBUTE

STAGEDSCORE SELFXYSAMPLEWEIGHTNONE

RETURN STAGED SCORES FOR X Y

THIS GENERATOR METHOD YIELDS THE ENSEMBLE SCORE AFTER EACH ITERATION OF BOOSTING AND THEREFORE ALLOWS MONITORING SUCH AS TO DETERMINE THE SCORE ON A TEST SET AFTER EACH BOOST

PARAMETERS

XARRAYLIKE SPARSE MATRIX OF SHAPE NSAMPLES NFEATURES THE TRAINING INPUT SAMPLES SPARSE MATRIX CAN BE CSC CSR COO DOK OR LIL COO DOK AND LIL ARE CONVERTED TO CSR

YARRAYLIKE SHAPE NSAMPLES LABELS FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

ZFLOAT

EXAMPLES USING SKLEARNENSEMBLEADABOOSTCLASSIFIER

- CLASSIFIER COMPARISON
- TWOCLASS ADABOOST
- MULTICLASS ADABOOSTED DECISION TREES
- DISCRETE VERSUS REAL ADABOOST
- PLOT THE DECISION SURFACES OF ENSEMBLES OF TREES ON THE IRIS DATASET

6122SKLEARNENSEMBLE ADABOOSTREGRESSOR

CLASSSKLEARNENSEMBLE ADABOOSTREGRESSOR BASEESTIMATORNONE NESTIMATORS50 LEARN INGRATE10 LOSS'LINEAR' RANDOMSTATENONE

AN ADABOOST REGRESSOR

AN ADABOOST 1 REGRESSOR IS A METAESTIMATOR THAT BEGINS BY FITTING A REGRESSOR ON THE ORIGINAL DATASET AND THEN FITS ADDITIONAL COPIES OF THE REGRESSOR ON THE SAME DATASET BUT WHERE THE WEIGHTS OF INSTANCES ARE ADJUSTED ACCORDING TO THE ERROR OF THE CURRENT PREDICTION AS SUCH SUBSEQUENT REGRESSORS FOCUS MORE ON DIFFICULT CASES

THIS CLASS IMPLEMENTS THE ALGORITHM KNOWN AS ADABOOSTR2 2

READ MORE IN THE USER GUIDE

PARAMETERS

BASEESTIMATOR OBJECT OPTIONAL DEFAULTNONE THE BASE ESTIMATOR FROM WHICH THE BOOSTED ENSEMBLE IS BUILT SUPPORT FOR SAMPLE WEIGHTING IS REQUIRED IF NONE THEN THE BASE ESTIMATOR ISDECISIONTREEREGRESSORMAXDEPTH3

NESTIMATORS INTEGER OPTIONAL DEFAULT50 THE MAXIMUM NUMBER OF ESTIMATORS AT WHICH BOOSTING IS TERMINATED IN CASE OF PERFECT FIT THE LEARNING PROCEDURE IS STOPPED EARLY

1696 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

LEARNINGRATE FLOAT OPTIONAL DEFAULT1 LEARNING RATE SHRINKS THE CONTRIBUTION OF EACH REGRESSOR BY LEARNINGRATE THERE IS A TRADEOFF BETWEEN LEARNINGRATE AND NESTIMATORS

LOSS 'LINEAR' 'SQUARE' 'EXPONENTIAL' OPTIONAL DEFAULT'LINEAR' THE LOSS FUNCTION TO USE WHEN UPDATING THE WEIGHTS AFTER EACH BOOSTING ITERATION

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RAN DOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NPRANDOM

ATTRIBUTES

ESTIMATORS LIST OF CLASSIFIERS THE COLLECTION OF FITTED SUBESTIMATORS

ESTIMATORWEIGHTS ARRAY OF FLOATS WEIGHTS FOR EACH ESTIMATOR IN THE BOOSTED ENSEMBLE

ESTIMATORERRORS ARRAY OF FLOATS REGRESSION ERROR FOR EACH ESTIMATOR IN THE BOOSTED ENSEMBLE

FEATUREIMPORTANCES ARRAY OF SHAPE NFEATURES RETURN THE FEATURE IMPORTANCES THE HIGHER THE MORE IMPORTANT THE FEATURE

SEE ALSO

ADABOOSTCLASSIFIER GRADIENTBOOSTINGREGRESSOR

SKLEARNTREEDECISIONTREEREGRESSOR

REFERENCES

R0C261B7DEE9D1 R0C261B7DEE9D2

EXAMPLES

FROM SKLEARNENSEMBLE IMPORT ADABOOSTREGRESSOR

FROM SKLEARNDATASETS IMPORT MAKEREGRESSION

X Y MAKEREGRESSIONNFEATURES4 NINFORMATIVE2

RANDOMSTATE0 SHUFFLE FALSE

REGR ADABOOSTREGRESSORRANDOMSTATE0 NESTIMATORS100

REGRFITX Y

ADABOOSTREGRESSORBASEESTIMATORNONE LEARNINGRATE10 LOSSLINEAR

NESTIMATORS100 RANDOMSTATE0

REGRFEATUREIMPORTANCES

ARRAY02788 07109 00065 00036

REGRPREDICT0 0 0 0

ARRAY47972

REGRSCOREX Y

09771

METHODS

FITSELF X Y SAMPLEWEIGHT BUILD A BOOSTED REGRESSOR FROM THE TRAINING SET X Y

CONTINUED ON NEXT PAGE

612SKLEARNENSEMBLE ENSEMBLE METHODS 1697

SCIKITLEARN USER GUIDE RELEASE 0213

TABLE 673 – CONTINUED FROM PREVIOUS PAGE

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PREDICT SELF X PREDICT REGRESSION VALUE FOR X

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

STAGEDPREDICT SELF X RETURN STAGED PREDICTIONS FOR X

STAGEDSCORE SELF X Y SAMPLEWEIGHT RETURN STAGED SCORES FOR X Y

INIT SELFBASEESTIMATORNONE NESTIMATORS50 LEARNINGRATE10 LOSS'LINEAR' RANDOMSTATENONE

FEATUREIMPORTANCES

RETURN THE FEATURE IMPORTANCES THE HIGHER THE MORE IMPORTANT THE FEATURE RETURNS

FEATUREIMPORTANCES ARRAY SHAPE NFEATURES

FITSELFXYSAMPLEWEIGHTNONE

BUILD A BOOSTED REGRESSOR FROM THE TRAINING SET X Y

PARAMETERS

XARRAYLIKE SPARSE MATRIX OF SHAPE NSAMPLES NFEATURES THE TRAINING INPUT SAMPLES SPARSE MATRIX CAN BE CSC CSR COO DOK OR LIL COO DOK AND LIL ARE CONVERTED TO CSR

YARRAYLIKE OF SHAPE NSAMPLES THE TARGET VALUES REAL NUMBERS

SAMPLEWEIGHT ARRAYLIKE OF SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS IF NONE THE SAMPLE WEIGHTS ARE INITIALIZED TO 1 NSAMPLES

RETURNS

SELF OBJECT

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELF

PREDICT REGRESSION VALUE FOR X

THE PREDICTED REGRESSION VALUE OF AN INPUT SAMPLE IS COMPUTED AS THE WEIGHTED MEDIAN PREDICTION OF THE CLASSIFIERS IN THE ENSEMBLE

PARAMETERS

XARRAYLIKE SPARSE MATRIX OF SHAPE NSAMPLES NFEATURES THE TRAINING INPUT SAMPLES SPARSE MATRIX CAN BE CSC CSR COO DOK OR LIL COO DOK AND LIL ARE CONVERTED TO CSR

RETURNS

1698 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

YARRAY OF SHAPE (NSAMPLES, THE PREDICTED REGRESSION VALUES)

SCORESELFXY(SAMPLEWEIGHT=None)

RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{U}{V}$  WHERE U IS THE RESIDUAL SUM OF SQUARES  $\sum (y_{true} - y_{pred})^2$  AND V IS THE TOTAL SUM OF SQUARES  $\sum (y_{true} - y_{true\_mean})^2$  THE BEST POSSIBLE SCORE IS 1.0 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 0.0

PARAMETERS

X: YARRAY-LIKE SHAPE (NSAMPLES, NFEATURES) TEST SAMPLES FOR SOME ESTIMATORS THIS MAY BE A PRECOMPUTED KERNEL MATRIX INSTEAD

SHAPE (NSAMPLES, NSAMPLESFITTED) WHERE NSAMPLESFITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR

Y: YARRAY-LIKE SHAPE (NSAMPLES) OR NSAMPLES NO OUTPUTS TRUE VALUES FOR X

SAMPLEWEIGHT: YARRAY-LIKE SHAPE (NSAMPLES) OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE: FLOAT R2 OF SELF.PREDICT(X) WRT Y

NOTES

THE R2 SCORE USED WHEN CALLING SCORE ON A REGRESSOR WILL USE MULTIOUTPUTUNIFORMAVERAGE FROM VERSION 0.23 TO KEEP CONSISTENT WITH METRICS.R2 SCORE THIS WILL INFLUENCE THE SCORE METHOD OF ALL THE MULTIOUTPUT REGRESSORS EXCEPT FOR MULTIOUTPUTMULTIOUTPUTREGRESSOR TO SPECIFY THE DEFAULT VALUE MANUALLY AND AVOID THE WARNING PLEASE EITHER CALL METRICS.R2 SCORE DIRECTLY OR MAKE A CUSTOM SCORER WITH METRICS.MAKESCORER THE BUILT-IN SCORER R2 USES MULTIOUTPUTUNIFORMAVERAGE

SETPARAMS(self, params)

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENT.PARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

STAGED.PREDICT(X)

RETURN STAGED PREDICTIONS FOR X

THE PREDICTED REGRESSION VALUE OF AN INPUT SAMPLE IS COMPUTED AS THE WEIGHTED MEDIAN PREDICTION OF THE CLASSIFIERS IN THE ENSEMBLE

THIS GENERATOR METHOD YIELDS THE ENSEMBLE PREDICTION AFTER EACH ITERATION OF BOOSTING AND THEREFORE ALLOWS MONITORING SUCH AS TO DETERMINE THE PREDICTION ON A TEST SET AFTER EACH BOOST

PARAMETERS

X: YARRAY-LIKE SPARSE MATRIX OF SHAPE (NSAMPLES, NFEATURES) THE TRAINING INPUT SAMPLES

RETURNS

Y: GENERATOR OF ARRAY SHAPE (NSAMPLES, THE PREDICTED REGRESSION VALUES)

6125KLEARN.ENSEMBLE.ENSEMBLE METHODS 1699

SCIKITLEARN USER GUIDE RELEASE 0213

STAGEDSCORE SELFXYSAMPLEWEIGHTNONE

RETURN STAGED SCORES FOR X Y

THIS GENERATOR METHOD YIELDS THE ENSEMBLE SCORE AFTER EACH ITERATION OF BOOSTING AND THEREFORE ALLOWS MONITORING SUCH AS TO DETERMINE THE SCORE ON A TEST SET AFTER EACH BOOST

PARAMETERS

XARRAYLIKE SPARSE MATRIX OF SHAPE NSAMPLES NFEATURES THE TRAINING INPUT SAMPLES SPARSE MATRIX CAN BE CSC CSR COO DOK OR LIL COO DOK AND LIL ARE CONVERTED TO CSR

YARRAYLIKE SHAPE NSAMPLES LABELS FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

ZFLOAT

EXAMPLES USING SKLEARNENSEMBLEADABOOSTREGRESSOR

- DECISION TREE REGRESSION WITH ADABOOST

6123SKLEARNENSEMBLE BAGGINGCLASSIFIER

CLASSSKLEARNENSEMBLE BAGGINGCLASSIFIER BASEESTIMATORNONE NESTIMATORS10

MAXSAMPLES10 MAXFEATURES10 BOOT

STRAPTRUE BOOTSTRAPFEATURESFALSE

OOBSCOREFALSE WARMSTARTFALSE NJOBSNONE

RANDOMSTATENONE VERBOSE0

A BAGGING CLASSIFIER

A BAGGING CLASSIFIER IS AN ENSEMBLE METAESTIMATOR THAT FITS BASE CLASSIFIERS EACH ON RANDOM SUBSETS OF THE ORIGINAL DATASET AND THEN AGGREGATE THEIR INDIVIDUAL PREDICTIONS EITHER BY VOTING OR BY AVERAGING TO FORM A FINAL PREDICTION SUCH A METAESTIMATOR CAN TYPICALLY BE USED AS A WAY TO REDUCE THE VARIANCE OF A BLACKBOX ESTIMATOR EG A DECISION TREE BY INTRODUCING RANDOMIZATION INTO ITS CONSTRUCTION PROCEDURE AND THEN MAKING AN ENSEMBLE OUT OF IT

THIS ALGORITHM ENCOMPASSES SEVERAL WORKS FROM THE LITERATURE WHEN RANDOM SUBSETS OF THE DATASET ARE DRAWN AS RANDOM SUBSETS OF THE SAMPLES THEN THIS ALGORITHM IS KNOWN AS PASTING RB1846455D0E51 IF SAMPLES ARE DRAWN WITH REPLACEMENT THEN THE METHOD IS KNOWN AS BAGGING RB1846455D0E52 WHEN RANDOM SUBSETS OF THE DATASET ARE DRAWN AS RANDOM SUBSETS OF THE FEATURES THEN THE METHOD IS KNOWN AS RANDOM SUBSPACES RB1846455D0E53 FINALLY WHEN BASE ESTIMATORS ARE BUILT ON SUBSETS OF BOTH SAMPLES AND FEATURES THEN THE METHOD IS KNOWN AS RANDOM PATCHES RB1846455D0E54

READ MORE IN THE USER GUIDE

PARAMETERS

BASEESTIMATOR OBJECT OR NONE OPTIONAL DEFAULTNONE THE BASE ESTIMATOR TO FIT ON RANDOM SUBSETS OF THE DATASET IF NONE THEN THE BASE ESTIMATOR IS A DECISION TREE

NESTIMATORS INT OPTIONAL DEFAULT10 THE NUMBER OF BASE ESTIMATORS IN THE ENSEMBLE

MAXSAMPLES INT OR FLOAT OPTIONAL DEFAULT10 THE NUMBER OF SAMPLES TO DRAW FROM X TO TRAIN EACH BASE ESTIMATOR

- IF INT THEN DRAW MAXSAMPLES SAMPLES
- IF FLOAT THEN DRAW MAXSAMPLES XSHAPE0 SAMPLES

1700 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

MAXFEATURES INT OR FLOAT OPTIONAL DEFAULT10 THE NUMBER OF FEATURES TO DRAW FROM X TO TRAIN EACH BASE ESTIMATOR

- IF INT THEN DRAW MAXFEATURES FEATURES
- IF FLOAT THEN DRAW MAXFEATURES XSHAPE1 FEATURES

BOOTSTRAP BOOLEAN OPTIONAL DEFAULTTRUE WHETHER SAMPLES ARE DRAWN WITH REPLACEMENT IF FALSE SAMPLING WITHOUT REPLACEMENT IS PERFORMED

BOOTSTRAPFEATURES BOOLEAN OPTIONAL DEFAULTFALSE WHETHER FEATURES ARE DRAWN WITH REPLACEMENT

OOBSCORE BOOL OPTIONAL DEFAULTFALSE WHETHER TO USE OUTFBAG SAMPLES TO ESTIMATE THE GENERALIZATION ERROR

WARMSTART BOOL OPTIONAL DEFAULTFALSE WHEN SET TO TRUE REUSE THE SOLUTION OF THE PREVIOUS CALL TO FIT AND ADD MORE ESTIMATORS TO THE ENSEMBLE OTHERWISE JUST FIT A WHOLE NEW ENSEMBLE SEE THE GLOSSARY

NEW IN VERSION 017 WARMSTART CONSTRUCTOR PARAMETER

NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF JOBS TO RUN IN PARALLEL FOR BOTH FIT ANDPREDICT NONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT

1MEANS USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RAN

DOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE

RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NPRANDOM

VERBOSE INT OPTIONAL DEFAULT0 CONTROLS THE VERBOSITY WHEN FITTING AND PREDICTING ATTRIBUTES

BASEESTIMATOR ESTIMATOR THE BASE ESTIMATOR FROM WHICH THE ENSEMBLE IS GROWN

ESTIMATORS LIST OF ESTIMATORS THE COLLECTION OF FITTED BASE ESTIMATORS

ESTIMATORSSAMPLES LIST OF ARRAYS THE SUBSET OF DRAWN SAMPLES FOR EACH BASE ESTIMATOR

ESTIMATORSFEATURES LIST OF ARRAYS THE SUBSET OF DRAWN FEATURES FOR EACH BASE ESTIMATOR

CLASSES ARRAY OF SHAPE NCLASSES THE CLASSES LABELS

NCLASSES INT OR LIST THE NUMBER OF CLASSES

OOBSCORE FLOAT SCORE OF THE TRAINING DATASET OBTAINED USING AN OUTFBAG ESTIMATE

OOBDECISIONFUNCTION ARRAY OF SHAPE NSAMPLES NCLASSES DECISION FUNCTION COM

PUTED WITH OUTFBAG ESTIMATE ON THE TRAINING SET IF NESTIMATORS IS SMALL IT MIGHT BE POSSIBLE THAT A DATA POINT WAS NEVER LEFT OUT DURING THE BOOTSTRAP IN THIS CASE

OOBDECISIONFUNCTION MIGHT CONTAIN NAN

REFERENCES

RB1846455D0E51 RB1846455D0E52 RB1846455D0E53 RB1846455D0E54

METHODS

612SKLEARNENSEMBLE ENSEMBLE METHODS 1701

SCIKITLEARN USER GUIDE RELEASE 0213

DECISIONFUNCTION SELF X AVERAGE OF THE DECISION FUNCTIONS OF THE BASE CLASSIFIERS

FITSELF X Y SAMPLEWEIGHT BUILD A BAGGING ENSEMBLE OF ESTIMATORS FROM THE TRAINING SET X Y

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PREDICT SELF X PREDICT CLASS FOR X

PREDICTLOGPROBA SELF X PREDICT CLASS LOGPROBABILITIES FOR X

PREDICTPROBA SELF X PREDICT CLASS PROBABILITIES FOR X

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELFBASEESTIMATORNONE NESTIMATORS10 MAXSAMPLES10 MAXFEATURES10 BOOTSTRAPTRUE BOOTSTRAPFEATURESFALSE OOBSCOREFALSE WARMSTARTFALSE NJOBSNONE

RANDOMSTATENONE VERBOSE0

DECISIONFUNCTION SELF X

AVERAGE OF THE DECISION FUNCTIONS OF THE BASE CLASSIFIERS

PARAMETERS

XARRAYLIKE SPARSE MATRIX OF SHAPE NSAMPLES NFEATURES THE TRAINING INPUT SAMPLES SPARSE MATRICES ARE ACCEPTED ONLY IF THEY ARE SUPPORTED BY THE BASE ESTIMATOR

RETURNS

SCORE ARRAY SHAPE NSAMPLES K THE DECISION FUNCTION OF THE INPUT SAMPLES THE COLUMNS CORRESPOND TO THE CLASSES IN SORTED ORDER AS THEY APPEAR IN THE ATTRIBUTE CLASSES REGRESSION AND BINARY CLASSIFICATION ARE SPECIAL CASES WITH K 1 OTHERWISEKNCLASSES

ESTIMATORSSAMPLES

THE SUBSET OF DRAWN SAMPLES FOR EACH BASE ESTIMATOR

RETURNS A DYNAMICALLY GENERATED LIST OF INDICES IDENTIFYING THE SAMPLES USED FOR FITTING EACH MEMBER OF THE ENSEMBLE IE THE INBAG SAMPLES

NOTE THE LIST IS RECREATED AT EACH CALL TO THE PROPERTY IN ORDER TO REDUCE THE OBJECT MEMORY FOOTPRINT BY NOT STORING THE SAMPLING DATA THUS FETCHING THE PROPERTY MAY BE SLOWER THAN EXPECTED

FITSELFXYSAMPLEWEIGHTNONE

BUILD A BAGGING ENSEMBLE OF ESTIMATORS FROM THE TRAINING SET X Y

PARAMETERS

XARRAYLIKE SPARSE MATRIX OF SHAPE NSAMPLES NFEATURES THE TRAINING INPUT SAMPLES SPARSE MATRICES ARE ACCEPTED ONLY IF THEY ARE SUPPORTED BY THE BASE ESTIMATOR

YARRAYLIKE SHAPE NSAMPLES THE TARGET VALUES CLASS LABELS IN CLASSIFICATION REAL NUMBERS IN REGRESSION

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OR NONE SAMPLE WEIGHTS IF NONE THEN SAMPLES ARE EQUALLY WEIGHTED NOTE THAT THIS IS SUPPORTED ONLY IF THE BASE ESTIMATOR SUPPORTS SAMPLE WEIGHTING

RETURNS

SELF OBJECT

1702 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

GETPARAMS SELFDEEPTURE  
GET PARAMETERS FOR THIS ESTIMATOR  
PARAMETERS  
DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBOBJECTS THAT ARE ESTIMATORS  
RETURNS  
PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES  
PREDICTSELF  
PREDICT CLASS FOR X  
THE PREDICTED CLASS OF AN INPUT SAMPLE IS COMPUTED AS THE CLASS WITH THE HIGHEST MEAN PREDICTED PROBABILITY  
IF BASE ESTIMATORS DO NOT IMPLEMENT A PREDICTPROBA METHOD THEN IT RESORTS TO VOTING  
PARAMETERS  
XARRAYLIKE SPARSE MATRIX OF SHAPE NSAMPLES NFEATURES THE TRAINING INPUT SAM  
PLES SPARSE MATRICES ARE ACCEPTED ONLY IF THEY ARE SUPPORTED BY THE BASE ESTIMATOR  
RETURNS  
YARRAY OF SHAPE NSAMPLES THE PREDICTED CLASSES  
PREDICTLOGPROBA SELF  
PREDICT CLASS LOGPROBABILITIES FOR X  
THE PREDICTED CLASS LOGPROBABILITIES OF AN INPUT SAMPLE IS COMPUTED AS THE LOG OF THE MEAN PREDICTED CLASS  
PROBABILITIES OF THE BASE ESTIMATORS IN THE ENSEMBLE  
PARAMETERS  
XARRAYLIKE SPARSE MATRIX OF SHAPE NSAMPLES NFEATURES THE TRAINING INPUT SAM  
PLES SPARSE MATRICES ARE ACCEPTED ONLY IF THEY ARE SUPPORTED BY THE BASE ESTIMATOR  
RETURNS  
PARRAY OF SHAPE NSAMPLES NCLASSES THE CLASS LOGPROBABILITIES OF THE INPUT SAMPLES  
THE ORDER OF THE CLASSES CORRESPONDS TO THAT IN THE ATTRIBUTE CLASSES  
PREDICTPROBA SELF  
PREDICT CLASS PROBABILITIES FOR X  
THE PREDICTED CLASS PROBABILITIES OF AN INPUT SAMPLE IS COMPUTED AS THE MEAN PREDICTED CLASS PROBABILITIES  
OF THE BASE ESTIMATORS IN THE ENSEMBLE IF BASE ESTIMATORS DO NOT IMPLEMENT A PREDICTPROBA METHOD  
THEN IT RESORTS TO VOTING AND THE PREDICTED CLASS PROBABILITIES OF AN INPUT SAMPLE REPRESENTS THE PROPORTION  
OF ESTIMATORS PREDICTING EACH CLASS  
PARAMETERS  
XARRAYLIKE SPARSE MATRIX OF SHAPE NSAMPLES NFEATURES THE TRAINING INPUT SAM  
PLES SPARSE MATRICES ARE ACCEPTED ONLY IF THEY ARE SUPPORTED BY THE BASE ESTIMATOR  
RETURNS  
PARRAY OF SHAPE NSAMPLES NCLASSES THE CLASS PROBABILITIES OF THE INPUT SAMPLES THE  
ORDER OF THE CLASSES CORRESPONDS TO THAT IN THE ATTRIBUTE CLASSES  
SCORESELFXYSAMPLEWEIGHTNONE  
RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS  
IN MULTILABEL CLASSIFICATION THIS IS THE SUBSET ACCURACY WHICH IS A HARSH METRIC SINCE YOU REQUIRE FOR EACH  
SAMPLE THAT EACH LABEL SET BE CORRECTLY PREDICTED  
612SKLEARNENSEMBLE ENSEMBLE METHODS 1703

SCIKITLEARN USER GUIDE RELEASE 0213

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE LABELS FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT MEAN ACCURACY OF SELF PREDICTX WRT Y

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT

OF A NESTED OBJECT

RETURNS

SELF

61245KLEARNENSEMBLE BAGGINGREGRESSOR

CLASSSSKLEARNENSEMBLE BAGGINGREGRESSOR BASEESTIMATORNONE NESTIMATORS10

MAXSAMPLES10 MAXFEATURES10 BOOT

STRAPTRUE BOOTSTRAPFEATURESFALSE

OOBSCOREFALSE WARMSTARTFALSE NJOBSNONE

RANDOMSTATENONE VERBOSE0

A BAGGING REGRESSOR

A BAGGING REGRESSOR IS AN ENSEMBLE METAESTIMATOR THAT FITS BASE REGRESSORS EACH ON RANDOM SUBSETS OF THE ORIGINAL DATASET AND THEN AGGREGATE THEIR INDIVIDUAL PREDICTIONS EITHER BY VOTING OR BY AVERAGING TO FORM A FINAL PREDICTION SUCH A METAESTIMATOR CAN TYPICALLY BE USED AS A WAY TO REDUCE THE VARIANCE OF A BLACKBOX ESTIMATOR EG A DECISION TREE BY INTRODUCING RANDOMIZATION INTO ITS CONSTRUCTION PROCEDURE AND THEN MAKING AN ENSEMBLE OUT OF IT

THIS ALGORITHM ENCOMPASSES SEVERAL WORKS FROM THE LITERATURE WHEN RANDOM SUBSETS OF THE DATASET ARE DRAWN AS RANDOM SUBSETS OF THE SAMPLES THEN THIS ALGORITHM IS KNOWN AS PASTING R4D113BA76FC01 IF SAMPLES ARE DRAWN WITH REPLACEMENT THEN THE METHOD IS KNOWN AS BAGGING R4D113BA76FC02 WHEN RANDOM SUBSETS OF THE DATASET ARE DRAWN AS RANDOM SUBSETS OF THE FEATURES THEN THE METHOD IS KNOWN AS RANDOM SUBSPACES R4D113BA76FC03 FINALLY WHEN BASE ESTIMATORS ARE BUILT ON SUBSETS OF BOTH SAMPLES AND FEATURES THEN THE METHOD IS KNOWN AS RANDOM PATCHES R4D113BA76FC04

READ MORE IN THE USER GUIDE

PARAMETERS

BASEESTIMATOR OBJECT OR NONE OPTIONAL DEFAULTNONE THE BASE ESTIMATOR TO FIT ON RANDOM SUBSETS OF THE DATASET IF NONE THEN THE BASE ESTIMATOR IS A DECISION TREE

NESTIMATORS INT OPTIONAL DEFAULT10 THE NUMBER OF BASE ESTIMATORS IN THE ENSEMBLE

MAXSAMPLES INT OR FLOAT OPTIONAL DEFAULT10 THE NUMBER OF SAMPLES TO DRAW FROM X TO TRAIN EACH BASE ESTIMATOR

- IF INT THEN DRAW MAXSAMPLES SAMPLES

- IF FLOAT THEN DRAW MAXSAMPLES XSHAPE0 SAMPLES

MAXFEATURES INT OR FLOAT OPTIONAL DEFAULT10 THE NUMBER OF FEATURES TO DRAW FROM X TO TRAIN EACH BASE ESTIMATOR

1704 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

- IF INT THEN DRAW MAXFEATURES FEATURES
- IF FLOAT THEN DRAW MAXFEATURES XSHAPE1 FEATURES

BOOTSTRAP BOOLEAN OPTIONAL DEFAULTTRUE WHETHER SAMPLES ARE DRAWN WITH REPLACEMENT IF FALSE SAMPLING WITHOUT REPLACEMENT IS PERFORMED

BOOTSTRAPFEATURES BOOLEAN OPTIONAL DEFAULTFALSE WHETHER FEATURES ARE DRAWN WITH REPLACEMENT

OOBSCORE BOOL WHETHER TO USE OUTFBAG SAMPLES TO ESTIMATE THE GENERALIZATION ERROR

WARMSTART BOOL OPTIONAL DEFAULTFALSE WHEN SET TO TRUE REUSE THE SOLUTION OF THE PREVIOUS CALL TO FIT AND ADD MORE ESTIMATORS TO THE ENSEMBLE OTHERWISE JUST FIT A WHOLE NEW ENSEMBLE SEE THE GLOSSARY

NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF JOBS TO RUN IN PARALLEL FOR BOTH FIT ANDPREDICT NONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT

1MEANS USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RAN

DOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE

RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NPRANDOM

VERBOSE INT OPTIONAL DEFAULT0 CONTROLS THE VERBOSITY WHEN FITTING AND PREDICTING ATTRIBUTES

ESTIMATORS LIST OF ESTIMATORS THE COLLECTION OF FITTED SUBESTIMATORS

ESTIMATORSSAMPLES LIST OF ARRAYS THE SUBSET OF DRAWN SAMPLES FOR EACH BASE ESTIMATOR

ESTIMATORSFEATURES LIST OF ARRAYS THE SUBSET OF DRAWN FEATURES FOR EACH BASE ESTIMATOR

OOBSCORE FLOAT SCORE OF THE TRAINING DATASET OBTAINED USING AN OUTFBAG ESTIMATE

OOBPREDICTION ARRAY OF SHAPE NSAMPLES PREDICTION COMPUTED WITH OUTFBAG ESTIMATE

ON THE TRAINING SET IF NESTIMATORS IS SMALL IT MIGHT BE POSSIBLE THAT A DATA POINT WAS NEVER LEFT OUT DURING THE BOOTSTRAP IN THIS CASE OOBPREDICTION MIGHT CONTAIN NAN

REFERENCES

R4D113BA76FC01 R4D113BA76FC02 R4D113BA76FC03 R4D113BA76FC04

METHODS

FITSELF X Y SAMPLEWEIGHT BUILD A BAGGING ENSEMBLE OF ESTIMATORS FROM THE TRAINING SET X Y

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PREDICT SELF X PREDICT REGRESSION TARGET FOR X

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

612SKLEARNENSEMBLE ENSEMBLE METHODS 1705

SCIKITLEARN USER GUIDE RELEASE 0213

INIT SELFBASEESTIMATORNONE NESTIMATORS10 MAXSAMPLES10 MAXFEATURES10 BOOTSTRAPTRUE BOOTSTRAPFEATURESFALSE OOBSCOREFALSE WARMSTARTFALSE NJOBSNONE

RANDOMSTATENONE VERBOSE0

ESTIMATORSSAMPLES

THE SUBSET OF DRAWN SAMPLES FOR EACH BASE ESTIMATOR

RETURNS A DYNAMICALLY GENERATED LIST OF INDICES IDENTIFYING THE SAMPLES USED FOR FITTING EACH MEMBER OF THE ENSEMBLE IE THE INBAG SAMPLES

NOTE THE LIST IS RECREATED AT EACH CALL TO THE PROPERTY IN ORDER TO REDUCE THE OBJECT MEMORY FOOTPRINT BY NOT STORING THE SAMPLING DATA THUS FETCHING THE PROPERTY MAY BE SLOWER THAN EXPECTED

FITSELFXYSAMPLEWEIGHTNONE

BUILD A BAGGING ENSEMBLE OF ESTIMATORS FROM THE TRAINING SET X Y

PARAMETERS

XARRAYLIKE SPARSE MATRIX OF SHAPE NSAMPLES NFEATURES THE TRAINING INPUT SAMPLES SPARSE MATRICES ARE ACCEPTED ONLY IF THEY ARE SUPPORTED BY THE BASE ESTIMATOR

YARRAYLIKE SHAPE NSAMPLES THE TARGET VALUES CLASS LABELS IN CLASSIFICATION REAL NUMBERS IN REGRESSION

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OR NONE SAMPLE WEIGHTS IF NONE THEN SAMPLES ARE EQUALLY WEIGHTED NOTE THAT THIS IS SUPPORTED ONLY IF THE BASE ESTIMATOR SUPPORTS SAMPLE WEIGHTING

RETURNS

SELF OBJECT

GETPARAMS SELFDEEPTREE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELFXY

PREDICT REGRESSION TARGET FOR X

THE PREDICTED REGRESSION TARGET OF AN INPUT SAMPLE IS COMPUTED AS THE MEAN PREDICTED REGRESSION TARGETS OF THE ESTIMATORS IN THE ENSEMBLE

PARAMETERS

XARRAYLIKE SPARSE MATRIX OF SHAPE NSAMPLES NFEATURES THE TRAINING INPUT SAMPLES SPARSE MATRICES ARE ACCEPTED ONLY IF THEY ARE SUPPORTED BY THE BASE ESTIMATOR

RETURNS

YARRAY OF SHAPE NSAMPLES THE PREDICTED VALUES

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{U}{V}$  WHERE U IS THE RESIDUAL SUM OF SQUARES YTRUE YPRED 2SUM AND V IS THE TOTAL SUM OF SQUARES YTRUE YTRUEMEAN 2SUM THE BEST POSSIBLE SCORE

1706 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

IS 10 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 00

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES FOR SOME ESTIMATORS THIS MAY BE A PRECOMPUTED KERNEL MATRIX INSTEAD SHAPE NSAMPLES NSAMPLESFITTED WHERE NSAMPLESFITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE VALUES FOR X SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT R2 OF SELF PREDICT X WRT Y

NOTES

THE R2 SCORE USED WHEN CALLING SCORE ON A REGRESSOR WILL USE MULTIOUTPUT UNIFORM AVERAGE FROM VERSION 023 TO KEEP CONSISTENT WITH METRIC SR2 SCORE THIS WILL INFLUENCE THE SCORE METHOD OF ALL THE MULTIOUTPUT REGRESSORS EXCEPT FOR MULTIOUTPUT MULTIOUTPUT REGRESSOR TO SPECIFY THE DEFAULT VALUE MANUALLY AND AVOID THE WARNING PLEASE EITHER CALL METRIC SR2 SCORE DIRECTLY OR MAKE A CUSTOM SCORER WITH METRIC S MAKE SCORER THE BUILT IN SCORER R2 USES MULTIOUTPUT UNIFORM AVERAGE

SET PARAMS SELF PARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENT PARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

EXAMPLES USING SKLEARN ENSEMBLE BAGGING REGRESSOR

- SINGLE ESTIMATOR VERSUS BAGGING BIAS VARIANCE DECOMPOSITION

6125 SKLEARN ENSEMBLE ISOLATION FOREST

CLASS SKLEARN ENSEMBLE ISOLATION FOREST NESTIMATORS 100 MAX SAMPLES 'AUTO' CONTAMINATION 'LEGACY' MAX FEATURES 10 BOOTSTRAP FALSE

NJOBS NONE BEHAVIOUR 'OLD' RANDOM STATE NONE

VERBOSE 0 WARM START FALSE

ISOLATION FOREST ALGORITHM

RETURN THE ANOMALY SCORE OF EACH SAMPLE USING THE ISOLATION FOREST ALGORITHM

THE ISOLATION FOREST 'ISOLATES' OBSERVATIONS BY RANDOMLY SELECTING A FEATURE AND THEN RANDOMLY SELECTING A SPLIT VALUE BETWEEN THE MAXIMUM AND MINIMUM VALUES OF THE SELECTED FEATURE

SINCE RECURSIVE PARTITIONING CAN BE REPRESENTED BY A TREE STRUCTURE THE NUMBER OF SPLITTINGS REQUIRED TO ISOLATE A SAMPLE IS EQUIVALENT TO THE PATH LENGTH FROM THE ROOT NODE TO THE TERMINATING NODE

THIS PATH LENGTH AVERAGED OVER A FOREST OF SUCH RANDOM TREES IS A MEASURE OF NORMALITY AND OUR DECISION FUNCTION

6125 SKLEARN ENSEMBLE ENSEMBLE METHODS 1707

SCIKITLEARN USER GUIDE RELEASE 0213

RANDOM PARTITIONING PRODUCES NOTICEABLY SHORTER PATHS FOR ANOMALIES HENCE WHEN A FOREST OF RANDOM TREES COLLECTIVELY PRODUCE SHORTER PATH LENGTHS FOR PARTICULAR SAMPLES THEY ARE HIGHLY LIKELY TO BE ANOMALIES

READ MORE IN THE USER GUIDE

NEW IN VERSION 018

PARAMETERS

NESTIMATORS INT OPTIONAL DEFAULT100 THE NUMBER OF BASE ESTIMATORS IN THE ENSEMBLE

MAXSAMPLES INT OR FLOAT OPTIONAL DEFAULT" AUTO"

THE NUMBER OF SAMPLES TO DRAW FROM X TO TRAIN EACH BASE ESTIMATOR

- IF INT THEN DRAW MAXSAMPLES SAMPLES
- IF FLOAT THEN DRAW MAXSAMPLES XSHAPE0 SAMPLES
- IF " AUTO" THEN MAXSAMPLESMIN256 NSAMPLES

IF MAXSAMPLES IS LARGER THAN THE NUMBER OF SAMPLES PROVIDED ALL SAMPLES WILL BE USED FOR ALL TREES NO SAMPLING

CONTAMINATION FLOAT IN 0 05 OPTIONAL DEFAULT01 THE AMOUNT OF CONTAMINATION OF THE DATA SET IE THE PROPORTION OF OUTLIERS IN THE DATA SET USED WHEN FITTING TO DEFINE THE THRESHOLD ON THE DECISION FUNCTION IF ' AUTO' THE DECISION FUNCTION THRESHOLD IS DETERMINED AS IN THE ORIGINAL PAPER

CHANGED IN VERSION 020 THE DEFAULT VALUE OF CONTAMINATION WILL CHANGE FROM 01 IN 020 TOAUTO IN 022

MAXFEATURES INT OR FLOAT OPTIONAL DEFAULT10 THE NUMBER OF FEATURES TO DRAW FROM X TO TRAIN EACH BASE ESTIMATOR

- IF INT THEN DRAW MAXFEATURES FEATURES
- IF FLOAT THEN DRAW MAXFEATURES XSHAPE1 FEATURES

BOOTSTRAP BOOLEAN OPTIONAL DEFAULTFALSE IF TRUE INDIVIDUAL TREES ARE FIT ON RANDOM SUBSETS OF THE TRAINING DATA SAMPLED WITH REPLACEMENT IF FALSE SAMPLING WITHOUT REPLACEMENT IS PERFORMED

NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF JOBS TO RUN IN PARALLEL FOR BOTH FIT ANDPREDICT NONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT

1MEANS USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

BEHAVIOUR STR DEFAULT' OLD' BEHAVIOUR OF THE DECISIONFUNCTION WHICH CAN BE EITHER ' OLD' OR ' NEW' PASSING BEHAVIOURNEW MAKES THE DECISIONFUNCTION CHANGE TO MATCH OTHER ANOMALY DETECTION ALGORITHM API WHICH WILL BE THE DEFAULT BEHAVIOUR IN THE FUTURE AS EXPLAINED IN DETAILS IN THE OFFSET ATTRIBUTE DOCUMENTATION

THEDECISIONFUNCTION BECOMES DEPENDENT ON THE CONTAMINATION PARAMETER IN SUCH A WAY THAT 0 BECOMES ITS NATURAL THRESHOLD TO DETECT OUTLIERS

NEW IN VERSION 020 BEHAVIOUR IS ADDED IN 020 FOR BACKCOMPATIBILITY PURPOSE

DEPRECATED SINCE VERSION 020 BEHAVIOUROLD IS DEPRECATED IN 020 AND WILL NOT BE POSSIBLE IN 022

DEPRECATED SINCE VERSION 022 BEHAVIOUR PARAMETER WILL BE DEPRECATED IN 022 AND REMOVED IN 024

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RANDOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE

1708 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NPRANDOM  
VERBOSE INT OPTIONAL DEFAULT0 CONTROLS THE VERBOSITY OF THE TREE BUILDING PROCESS  
WARMSTART BOOL OPTIONAL DEFAULTFALSE WHEN SET TO TRUE REUSE THE SOLUTION OF THE PREVIOUS CALL TO FIT AND ADD MORE ESTIMATORS TO THE ENSEMBLE OTHERWISE JUST FIT A WHOLE NEW FOREST SEE THE GLOSSARY  
NEW IN VERSION 021

ATTRIBUTES  
ESTIMATORS LIST OF DECISIONTREECLASSIFIER THE COLLECTION OF FITTED SUBESTIMATORS  
ESTIMATORSSAMPLES LIST OF ARRAYS THE SUBSET OF DRAWN SAMPLES FOR EACH BASE ESTIMATOR

MAXSAMPLES INTEGER THE ACTUAL NUMBER OF SAMPLES  
OFFSET FLOAT OFFSET USED TO DEFINE THE DECISION FUNCTION FROM THE RAW SCORES WE HAVE THE RELATIONDECISIONFUNCTION SCORESAMPLES OFFSET ASSUMING BEHAVIOUR 'NEW' OFFSET IS DEFINED AS FOLLOWS WHEN THE CONTAMINATION PARAMETER IS SET TO "AUTO" THE OFFSET IS EQUAL TO 0.5 AS THE SCORES OF INLIERS ARE CLOSE TO 0 AND THE SCORES OF OUTLIERS ARE CLOSE TO 1 WHEN A CONTAMINATION PARAMETER DIFFERENT THAN "AUTO" IS PROVIDED THE OFFSET IS DEFINED IN SUCH A WAY WE OBTAIN THE EXPECTED NUMBER OF OUTLIERS SAMPLES WITH DECISION FUNCTION 0 IN TRAINING ASSUMING THE BEHAVIOUR PARAMETER IS SET TO 'OLD' WE ALWAYS HAVE OFFSET 0.5 MAKING THE DECISION FUNCTION INDEPENDENT FROM THE CONTAMINATION PARAMETER

NOTES  
THE IMPLEMENTATION IS BASED ON AN ENSEMBLE OF EXTRATREEREgressor THE MAXIMUM DEPTH OF EACH TREE IS SET TO CEILLOG2N WHERE n IS THE NUMBER OF SAMPLES USED TO BUILD THE TREE SEE LIU ET AL 2008 FOR MORE

DETAILS

REFERENCES  
RD7AE0A2AE6881 RD7AE0A2AE6882

METHODS  
DECISIONFUNCTION SELF X AVERAGE ANOMALY SCORE OF X OF THE BASE CLASSIFIERS  
FITSELF X Y SAMPLEWEIGHT FIT ESTIMATOR  
FITPREDICT SELF X Y PERFORMS FIT ON X AND RETURNS LABELS FOR X  
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR  
PREDICT SELF X PREDICT IF A PARTICULAR SAMPLE IS AN OUTLIER OR NOT  
SCORESAMPLES SELF X OPPOSITE OF THE ANOMALY SCORE DEFINED IN THE ORIGINAL PAPER  
SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR  
INIT SELFNESTIMATORS100 MAXSAMPLES'AUTO' CONTAMINATION'LEGACY' MAXFEATURES10  
BOOTSTRAPFALSE NJOBSNONE BEHAVIOUR'OLD' RANDOMSTATENONE VERBOSE0  
WARMSTARTFALSE  
612SKLEARNENSEMBLE ENSEMBLE METHODS 1709

SCIKITLEARN USER GUIDE RELEASE 0213

DECISIONFUNCTION SELF\_X

AVERAGE ANOMALY SCORE OF X OF THE BASE CLASSIFIERS

THE ANOMALY SCORE OF AN INPUT SAMPLE IS COMPUTED AS THE MEAN ANOMALY SCORE OF THE TREES IN THE FOREST

THE MEASURE OF NORMALITY OF AN OBSERVATION GIVEN A TREE IS THE DEPTH OF THE LEAF CONTAINING THIS OBSERVATION

WHICH IS EQUIVALENT TO THE NUMBER OF SPLITTINGS REQUIRED TO ISOLATE THIS POINT IN CASE OF SEVERAL OBSERVATIONS

NLEFT IN THE LEAF THE AVERAGE PATH LENGTH OF A NLEFT SAMPLES ISOLATION TREE IS ADDED

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTERNALLY

IT WILL BE CONVERTED TO DTYPE NPFLOAT32 AND IF A SPARSE MATRIX IS PROVIDED TO A SPARSE CSR MATRIX

RETURNS

SCORES ARRAY SHAPE NSAMPLES THE ANOMALY SCORE OF THE INPUT SAMPLES THE LOWER THE MORE ABNORMAL NEGATIVE SCORES REPRESENT OUTLIERS POSITIVE SCORES REPRESENT INLIERS

ESTIMATOR SSAMPLES

THE SUBSET OF DRAWN SAMPLES FOR EACH BASE ESTIMATOR

RETURNS A DYNAMICALLY GENERATED LIST OF INDICES IDENTIFYING THE SAMPLES USED FOR FITTING EACH MEMBER OF THE ENSEMBLE IE THE INBAG SAMPLES

NOTE THE LIST IS RECREATED AT EACH CALL TO THE PROPERTY IN ORDER TO REDUCE THE OBJECT MEMORY FOOTPRINT BY NOT STORING THE SAMPLING DATA THUS FETCHING THE PROPERTY MAY BE SLOWER THAN EXPECTED

FIT SELF\_X Y NONE SAMPLE WEIGHT NONE

FIT ESTIMATOR

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES USE

DTYPE NPFLOAT32 FOR MAXIMUM EFFICIENCY SPARSE MATRICES ARE ALSO SUPPORTED USE

SPARSE CSC MATRIX FOR MAXIMUM EFFICIENCY

SAMPLE WEIGHT ARRAYLIKE SHAPE NSAMPLES OR NONE SAMPLE WEIGHTS IF NONE THEN

SAMPLES ARE EQUALLY WEIGHTED

Y IGNORED NOT USED PRESENT FOR API CONSISTENCY BY CONVENTION

RETURNS

SELF OBJECT

FIT PREDICT SELF\_X Y NONE

PERFORMS FIT ON X AND RETURNS LABELS FOR X

RETURNS 1 FOR OUTLIERS AND 1 FOR INLIERS

PARAMETERS

X NDARRAY SHAPE NSAMPLES NFEATURES INPUT DATA

Y IGNORED NOT USED PRESENT FOR API CONSISTENCY BY CONVENTION

RETURNS

Y NDARRAY SHAPE NSAMPLES 1 FOR INLIERS 1 FOR OUTLIERS

GET PARAMS SELF DEEP TRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

1710 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELF

PREDICT IF A PARTICULAR SAMPLE IS AN OUTLIER OR NOT

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTERNALLY IT WILL BE CONVERTED TO DTYPE NPFLOAT32 AND IF A SPARSE MATRIX IS PROVIDED TO A SPARSE CSR MATRIX

RETURNS

IS INLIER ARRAY SHAPE NSAMPLES FOR EACH OBSERVATION TELLS WHETHER OR NOT 1 OR 0 IT SHOULD BE CONSIDERED AS AN INLIER ACCORDING TO THE FITTED MODEL

SCORES SAMPLES SELF

OPPOSITE OF THE ANOMALY SCORE DEFINED IN THE ORIGINAL PAPER

THE ANOMALY SCORE OF AN INPUT SAMPLE IS COMPUTED AS THE MEAN ANOMALY SCORE OF THE TREES IN THE FOREST

THE MEASURE OF NORMALITY OF AN OBSERVATION GIVEN A TREE IS THE DEPTH OF THE LEAF CONTAINING THIS OBSERVATION WHICH IS EQUIVALENT TO THE NUMBER OF SPLITTINGS REQUIRED TO ISOLATE THIS POINT IN CASE OF SEVERAL OBSERVATIONS

NLEFT IN THE LEAF THE AVERAGE PATH LENGTH OF A NLEFT SAMPLES ISOLATION TREE IS ADDED

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES

RETURNS

SCORES ARRAY SHAPE NSAMPLES THE ANOMALY SCORE OF THE INPUT SAMPLES THE LOWER THE MORE ABNORMAL

SETPARAMS SELF

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENT PARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

EXAMPLES USING SKLEARN ENSEMBLE ISOLATION FOREST

- COMPARING ANOMALY DETECTION ALGORITHMS FOR OUTLIER DETECTION ON TOY DATASETS
- ISOLATION FOREST EXAMPLE

612 SKLEARN ENSEMBLE ENSEMBLE METHODS 1711

SCIKITLEARN USER GUIDE RELEASE 0213  
6126SKLEARNENSEMBLE RANDOMTREEEMBEDDING  
CLASSSSKLEARNENSEMBLE RANDOMTREEEMBEDDING NESTIMATORS'WARN' MAXDEPTH5  
MINSAMPLESSPLIT2 MINSAMPLESLEAF1  
MINWEIGHTFRACTIONLEAF00  
MAXLEAFNODESNONE  
MINIMPURITYDECREASE00  
MINIMPURITYSPLITNONE SPARSEOUTPUTTRUE  
NJOBSNONE RANDOMSTATENONE VERBOSE0  
WARMSTARTFALSE  
AN ENSEMBLE OF TOTALLY RANDOM TREES  
AN UNSUPERVISED TRANSFORMATION OF A DATASET TO A HIGHDIMENSIONAL SPARSE REPRESENTATION A DATAPOINT IS CODED  
ACCORDING TO WHICH LEAF OF EACH TREE IT IS SORTED INTO USING A ONEHOT ENCODING OF THE LEAVES THIS LEADS TO A  
BINARY CODING WITH AS MANY ONES AS THERE ARE TREES IN THE FOREST  
THE DIMENSIONALITY OF THE RESULTING REPRESENTATION IS NOUT NESTIMATORS MAXLEAFNODES  
IFMAXLEAFNODES NONE THE NUMBER OF LEAF NODES IS AT MOST NESTIMATORS 2  
MAXDEPTH  
READ MORE IN THE USER GUIDE  
PARAMETERS  
NESTIMATORS INTEGER OPTIONAL DEFAULT10 NUMBER OF TREES IN THE FOREST  
CHANGED IN VERSION 020 THE DEFAULT VALUE OF NESTIMATORS WILL CHANGE FROM 10 IN  
VERSION 020 TO 100 IN VERSION 022  
MAXDEPTH INTEGER OPTIONAL DEFAULT5 THE MAXIMUM DEPTH OF EACH TREE IF NONE  
THEN NODES ARE EXPANDED UNTIL ALL LEAVES ARE PURE OR UNTIL ALL LEAVES CONTAIN LESS THAN  
MINSAMPLESSPLIT SAMPLES  
MINSAMPLESSPLIT INT FLOAT OPTIONAL DEFAULT2 THE MINIMUM NUMBER OF SAMPLES REQUIRED  
TO SPLIT AN INTERNAL NODE  
• IF INT THEN CONSIDER MINSAMPLESSPLIT AS THE MINIMUM NUMBER  
• IF FLOAT THEN MINSAMPLESSPLIT IS A FRACTION AND CEILMINSAMPLESSPLIT  
NSAMPLES IS THE MINIMUM NUMBER OF SAMPLES FOR EACH SPLIT  
CHANGED IN VERSION 018 ADDED FLOAT VALUES FOR FRACTIONS  
MINSAMPLESLEAF INT FLOAT OPTIONAL DEFAULT1 THE MINIMUM NUMBER OF SAMPLES REQUIRED  
TO BE AT A LEAF NODE A SPLIT POINT AT ANY DEPTH WILL ONLY BE CONSIDERED IF IT LEAVES AT LEAST  
MINSAMPLESLEAF TRAINING SAMPLES IN EACH OF THE LEFT AND RIGHT BRANCHES THIS MAY  
HAVE THE EFFECT OF SMOOTHING THE MODEL ESPECIALLY IN REGRESSION  
• IF INT THEN CONSIDER MINSAMPLESLEAF AS THE MINIMUM NUMBER  
• IF FLOAT THEN MINSAMPLESLEAF IS A FRACTION AND CEILMINSAMPLESLEAF  
NSAMPLES IS THE MINIMUM NUMBER OF SAMPLES FOR EACH NODE  
CHANGED IN VERSION 018 ADDED FLOAT VALUES FOR FRACTIONS  
MINWEIGHTFRACTIONLEAF FLOAT OPTIONAL DEFAULT0 THE MINIMUM WEIGHTED FRACTION OF THE  
SUM TOTAL OF WEIGHTS OF ALL THE INPUT SAMPLES REQUIRED TO BE AT A LEAF NODE SAMPLES HAVE  
EQUAL WEIGHT WHEN SAMPLEWEIGHT IS NOT PROVIDED  
MAXLEAFNODES INT OR NONE OPTIONAL DEFAULTNONE GROW TREES WITH MAXLEAFNODES  
IN BESTFIRST FASHION BEST NODES ARE DEFINED AS RELATIVE REDUCTION IN IMPURITY IF NONE THEN  
UNLIMITED NUMBER OF LEAF NODES  
1712 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

MINIMPURITYDECREASE FLOAT OPTIONAL DEFAULT0 A NODE WILL BE SPLIT IF THIS SPLIT INDUCES A DECREASE OF THE IMPURITY GREATER THAN OR EQUAL TO THIS VALUE

THE WEIGHTED IMPURITY DECREASE EQUATION IS THE FOLLOWING

NT N IMPURITY NTR NT RIGHTIMPURITY

NTL NT LEFTIMPURITY

WHERE N IS THE TOTAL NUMBER OF SAMPLES NT IS THE NUMBER OF SAMPLES AT THE CURRENT NODE

NTL IS THE NUMBER OF SAMPLES IN THE LEFT CHILD AND NTR IS THE NUMBER OF SAMPLES IN THE RIGHT CHILD

NNTNTR ANDNTL ALL REFER TO THE WEIGHTED SUM IF SAMPLEWEIGHT IS PASSED

NEW IN VERSION 019

MINIMPURITYSPLIT FLOAT DEFAULT1E7 THRESHOLD FOR EARLY STOPPING IN TREE GROWTH A NODE WILL SPLIT IF ITS IMPURITY IS ABOVE THE THRESHOLD OTHERWISE IT IS A LEAF

DEPRECATED SINCE VERSION 019 MINIMPURITYSPLIT HAS BEEN DEPRECATED IN FAVOR OF MINIMPURITYDECREASE IN 019 THE DEFAULT VALUE OF MINIMPURITYSPLIT WILL CHANGE FROM 1E7 TO 0 IN 023 AND IT WILL BE REMOVED IN 025 USE MINIMPURITYDECREASE INSTEAD

SPARSEOUTPUT BOOL OPTIONAL DEFAULTTRUE WHETHER OR NOT TO RETURN A SPARSE CSR MATRIX AS DEFAULT BEHAVIOR OR TO RETURN A DENSE ARRAY COMPATIBLE WITH DENSE PIPELINE OPERATORS

NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF JOBS TO RUN IN PARALLEL FOR BOTH FIT ANDPREDICT NONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT

1MEANS USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RANDOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NRANDOM

VERBOSE INT OPTIONAL DEFAULT0 CONTROLS THE VERBOSITY WHEN FITTING AND PREDICTING

WARMSTART BOOL OPTIONAL DEFAULTFALSE WHEN SET TO TRUE REUSE THE SOLUTION OF THE PREVIOUS CALL TO FIT AND ADD MORE ESTIMATORS TO THE ENSEMBLE OTHERWISE JUST FIT A WHOLE NEW FOREST SEE THE GLOSSARY

ATTRIBUTES

ESTIMATORS LIST OF DECISIONTREECLASSIFIER THE COLLECTION OF FITTED SUBESTIMATORS

REFERENCES

R6E47E53BACBD1 R6E47E53BACBD2

METHODS

APPLY SELF X APPLY TREES IN THE FOREST TO X RETURN LEAF INDICES

DECISIONPATH SELF X RETURN THE DECISION PATH IN THE FOREST

FITSELF X Y SAMPLEWEIGHT FIT ESTIMATOR

FITTRANSFORM SELF X Y SAMPLEWEIGHT FIT ESTIMATOR AND TRANSFORM DATASET

CONTINUED ON NEXT PAGE

612SKLEARNENSEMBLE ENSEMBLE METHODS 1713

SCIKITLEARN USER GUIDE RELEASE 0213

TABLE 677 – CONTINUED FROM PREVIOUS PAGE

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF X TRANSFORM DATASET

INIT SELFNESTIMATORS'WARN' MAXDEPTH5 MINSAMPLESSPLIT2 MINSAMPLESLEAF1

MINWEIGHTFRACTIONLEAF00 MAXLEAFNODESNONE MINIMPURITYDECREASE00

MINIMPURITYSPLITNONE SPARSEOUTPUTTRUE NJOBSNONE RANDOMSTATENONE

VERBOSE0 WARMSTARTFALSE

APPLYSELF X

APPLY TREES IN THE FOREST TO X RETURN LEAF INDICES

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTER

NALLY ITS DTYPE WILL BE CONVERTED TO DTYPENPFLOAT32 IF A SPARSE MATRIX IS PROVIDED

IT WILL BE CONVERTED INTO A SPARSE CSRMATRIX

RETURNS

XLEAVES ARRAYLIKE SHAPE NSAMPLES NESTIMATORS FOR EACH DATAPOINT X IN X AND FOR

EACH TREE IN THE FOREST RETURN THE INDEX OF THE LEAF X ENDS UP IN

DECISIONPATH SELF X

RETURN THE DECISION PATH IN THE FOREST

NEW IN VERSION 018

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTER

NALLY ITS DTYPE WILL BE CONVERTED TO DTYPENPFLOAT32 IF A SPARSE MATRIX IS PROVIDED

IT WILL BE CONVERTED INTO A SPARSE CSRMATRIX

RETURNS

INDICATOR SPARSE CSR ARRAY SHAPE NSAMPLES NNODES RETURN A NODE INDICATOR MATRIX

WHERE NON ZERO ELEMENTS INDICATES THAT THE SAMPLES GOES THROUGH THE NODES

NNODESPTR ARRAY OF SIZE NESTIMATORS 1 THE COLUMNS FROM INDICA

TORNNODESPTRINNNODESPTRI1 GIVES THE INDICATOR VALUE FOR THE ITH ESTIMATOR

FEATUREIMPORTANCES

RETURN THE FEATURE IMPORTANCES THE HIGHER THE MORE IMPORTANT THE FEATURE

RETURNS

FEATUREIMPORTANCES ARRAY SHAPE NFEATURES THE VALUES OF THIS ARRAY SUM TO 1 UNLESS

ALL TREES ARE SINGLE NODE TREES CONSISTING OF ONLY THE ROOT NODE IN WHICH CASE IT WILL BE AN

ARRAY OF ZEROS

FITSELFXYNONE SAMPLEWEIGHTNONE

FIT ESTIMATOR

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPENSAMPLES NFEATURES THE INPUT SAMPLES USE

DTYPENPFLOAT32 FOR MAXIMUM EFFICIENCY SPARSE MATRICES ARE ALSO SUPPORTED USE

SPARSECSCMATRIX FOR MAXIMUM EFFICIENCY

1714 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OR NONE SAMPLE WEIGHTS IF NONE THEN  
SAMPLES ARE EQUALLY WEIGHTED SPLITS THAT WOULD CREATE CHILD NODES WITH NET ZERO OR NEG  
ATIVE WEIGHT ARE IGNORED WHILE SEARCHING FOR A SPLIT IN EACH NODE IN THE CASE OF CLASSIFI  
CATION SPLITS ARE ALSO IGNORED IF THEY WOULD RESULT IN ANY SINGLE CLASS CARRYING A NEGATIVE  
WEIGHT IN EITHER CHILD NODE

RETURNS  
SELF OBJECT

FITTRANSFORM SELFXYNONE SAMPLEWEIGHTNONE  
FIT ESTIMATOR AND TRANSFORM DATASET

PARAMETERS  
XARRAYLIKE OR SPARSE MATRIX SHAPENSAMPLES NFEATURES INPUT DATA USED TO BUILD  
FORESTS USE DTYPENPFLOAT32 FOR MAXIMUM EFFICIENCY

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OR NONE SAMPLE WEIGHTS IF NONE THEN  
SAMPLES ARE EQUALLY WEIGHTED SPLITS THAT WOULD CREATE CHILD NODES WITH NET ZERO OR NEG  
ATIVE WEIGHT ARE IGNORED WHILE SEARCHING FOR A SPLIT IN EACH NODE IN THE CASE OF CLASSIFI  
CATION SPLITS ARE ALSO IGNORED IF THEY WOULD RESULT IN ANY SINGLE CLASS CARRYING A NEGATIVE  
WEIGHT IN EITHER CHILD NODE

RETURNS  
XTRANSFORMED SPARSE MATRIX SHAPENSAMPLES NOUT TRANSFORMED DATASET

GETPARAMS SELFDEEPTTRUE  
GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS  
DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBOBJECTS THAT ARE ESTIMATORS

RETURNS  
PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

SETPARAMS SELFPARAMS  
SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE  
PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT  
OF A NESTED OBJECT

RETURNS  
SELF

TRANSFORM SELF  
TRANSFORM DATASET

PARAMETERS  
XARRAYLIKE OR SPARSE MATRIX SHAPENSAMPLES NFEATURES INPUT DATA TO BE TRANSFORMED  
USEDTYPENPFLOAT32 FOR MAXIMUM EFFICIENCY SPARSE MATRICES ARE ALSO SUPPORTED  
USE SPARSECSRMATRIX FOR MAXIMUM EFFICIENCY

RETURNS  
XTRANSFORMED SPARSE MATRIX SHAPENSAMPLES NOUT TRANSFORMED DATASET

612SKLEARNENSEMBLE ENSEMBLE METHODS 1715

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES USING SKLEARNENSEMBLERANDOMTREESEMBEDDING

- HASHING FEATURE TRANSFORMATION USING TOTALLY RANDOM TREES
- FEATURE TRANSFORMATIONS WITH ENSEMBLES OF TREES
- MANIFOLD LEARNING ON HANDWRITTEN DIGITS LOCALLY LINEAR EMBEDDING ISOMAP

6127SKLEARNENSEMBLE VOTINGCLASSIFIER

CLASSSSKLEARNENSEMBLE VOTINGCLASSIFIER ESTIMATORS VOTING‘HARD’ WEIGHTSNONE

NJOBSNONE FLATTENTRANSFORMTRUE

SOFT V OTINGMAJORITY RULE CLASSIFIER FOR UNFITTED ESTIMATORS

NEW IN VERSION 017

READ MORE IN THE USER GUIDE

PARAMETERS

ESTIMATORS LIST OF STRING ESTIMATOR TUPLES INVOKING THE FIT METHOD ON THE VOTINGCLASSIFIER WILL FIT CLONES OF THOSE ORIGINAL ESTIMATORS THAT WILL BE STORED IN THE CLASS ATTRIBUTE SELFESTIMATORS AN ESTIMATOR CAN BE SET TO NONE ORDROP USINGSETPARAMS

VOTING STR ‘HARD’ ‘SOFT’ DEFAULT‘HARD’ IF ‘HARD’ USES PREDICTED CLASS LABELS FOR MAJORITY RULE VOTING ELSE IF ‘SOFT’ PREDICTS THE CLASS LABEL BASED ON THE ARGMAX OF THE SUMS OF THE PREDICTED PROBABILITIES WHICH IS RECOMMENDED FOR AN ENSEMBLE OF WELLCALIBRATED CLASSIFIERS

WEIGHTS ARRAYLIKE SHAPE NCLASSIFIERS OPTIONAL DEFAULT‘NONE’ SEQUENCE OF WEIGHTS

FLOAT ORINT TO WEIGHT THE OCCURRENCES OF PREDICTED CLASS LABELS HARD VOTING OR CLASS PROBABILITIES BEFORE AVERAGING SOFT VOTING USES UNIFORM WEIGHTS IF NONE

NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF JOBS TO RUN IN PARALLEL FOR FIT

NONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

FLATTENTRANSFORM BOOL OPTIONAL DEFAULTTRUE AFFECTS SHAPE OF TRANSFORM OUTPUT ONLY WHEN VOTING‘SOFT’ IF VOTING‘SOFT’ AND FLATTENTRANSFORMTRUE TRANSFORM METHOD RETURNS MATRIX WITH SHAPE NSAMPLES NCLASSIFIERS NCLASSES IF FLATTENTRANSFORMFALSE IT RETURNS NCLASSIFIERS NSAMPLES NCLASSES

ATTRIBUTES

ESTIMATORS LIST OF CLASSIFIERS THE COLLECTION OF FITTED SUBESTIMATORS AS DEFINED IN ESTIMATORS THAT ARE NOT NONE

NAMEDESTIMATORS BUNCH OBJECT A DICTIONARY WITH ATTRIBUTE ACCESS ATTRIBUTE TO ACCESS ANY FITTED SUBESTIMATORS BY NAME

NEW IN VERSION 020

CLASSES ARRAYLIKE SHAPE NPREDICTIONS THE CLASSES LABELS

SEE ALSO

VOTINGREGRESSOR PREDICTION VOTING REGRESSOR

1716 CHAPTER 6 API REFERENCE

```
SCIKITLEARN USER GUIDE RELEASE 0213
EXAMPLES
IMPORT NUMPY AS NP
FROM SKLEARNLINEARMODEL IMPORT LOGISTICREGRESSION
FROM SKLEARNNAIVEBAYES IMPORT GAUSSIANNB
FROM SKLEARNENSEMBLE IMPORT RANDOMFORESTCLASSIFIER VOTINGCLASSIFIER
CLF1 LOGISTICREGRESSIONSOLVERLBFGS MULTICLASSMULTINOMIAL
RANDOMSTATE1
CLF2 RANDOMFORESTCLASSIFIERNESTIMATOR550 RANDOMSTATE1
CLF3 GAUSSIANNB
X NPARRAY1 1 2 1 3 2 1 1 2 1 3 2
Y NPARRAY1 1 1 2 2 2
ECLF1 VOTINGCLASSIFIERESTIMATORS
LR CLF1 RF CLF2 GNB CLF3 VOTINGHARD
ECLF1 ECLF1FITX Y
PRINTECLF1PREDICTX
1 1 1 2 2 2
NPARRAYEQUALECLF1NAMEDESTIMATORSLSRPREDICTX
ECLF1NAMEDESTIMATORSLSRPREDICTX
TRUE
ECLF2 VOTINGCLASSIFIERESTIMATORS
LR CLF1 RF CLF2 GNB CLF3
VOTINGSOFT
ECLF2 ECLF2FITX Y
PRINTECLF2PREDICTX
1 1 1 2 2 2
ECLF3 VOTINGCLASSIFIERESTIMATORS
LR CLF1 RF CLF2 GNB CLF3
VOTINGSOFT WEIGHTS211
FLATTENTTRANSFORM TRUE
ECLF3 ECLF3FITX Y
PRINTECLF3PREDICTX
1 1 1 2 2 2
PRINTECLF3TRANSFORMXSHAPE
6 6
METHODS
FITSELF X Y SAMPLEWEIGHT FIT THE ESTIMATORS
FITTRANSFORM SELF X Y FIT TO DATA THEN TRANSFORM IT
GETPARAMS SELF DEEP GET THE PARAMETERS OF THE ENSEMBLE ESTIMATOR
PREDICT SELF X PREDICT CLASS LABELS FOR X
SCORE SELF X Y SAMPLEWEIGHT RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND
LABELS
SETPARAMS SELF PARAMS SETTING THE PARAMETERS FOR THE ENSEMBLE ESTIMATOR
TRANSFORM SELF X RETURN CLASS LABELS OR PROBABILITIES FOR X FOR EACH ESTI
MATOR
INIT SELFESTIMATORS VOTING'HARD' WEIGHTSNONE NJOBSNONE FLATTENTTRANSFORMTRUE
FITSELFXYSAMPLEWEIGHTNONE
FIT THE ESTIMATORS
PARAMETERS
612SKLEARNENSEMBLE ENSEMBLE METHODS 1717
```

SCIKITLEARN USER GUIDE RELEASE 0213

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES TRAINING VECTORS WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

YARRAYLIKE SHAPE NSAMPLES TARGET VALUES

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OR NONE SAMPLE WEIGHTS IF NONE THEN SAMPLES ARE EQUALLY WEIGHTED NOTE THAT THIS IS SUPPORTED ONLY IF ALL UNDERLYING ESTIMATORS SUPPORT SAMPLE WEIGHTS

RETURNS

SELF OBJECT

FITTRANSFORM SELFXYNONE FITPARAMS

FIT TO DATA THEN TRANSFORM IT

FITS TRANSFORMER TO X AND Y WITH OPTIONAL PARAMETERS FITPARAMS AND RETURNS A TRANSFORMED VERSION OF X

PARAMETERS

XNUMPY ARRAY OF SHAPE NSAMPLES NFEATURES TRAINING SET

YNUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES

RETURNS

XNEW NUMPY ARRAY OF SHAPE NSAMPLES NFEATURESNEW TRANSFORMED ARRAY

GETPARAMS SELFDEEPTTRUE

GET THE PARAMETERS OF THE ENSEMBLE ESTIMATOR

PARAMETERS

DEEP BOOL SETTING IT TO TRUE GETS THE VARIOUS ESTIMATORS AND THE PARAMETERS OF THE ESTIMATORS AS WELL

PREDICTSELF

PREDICT CLASS LABELS FOR X

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES

RETURNS

MAJ ARRAYLIKE SHAPE NSAMPLES PREDICTED CLASS LABELS

PREDICTPROBA

COMPUTE PROBABILITIES OF POSSIBLE OUTCOMES FOR SAMPLES IN X

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES

RETURNS

AVG ARRAYLIKE SHAPE NSAMPLES NCLASSES WEIGHTED AVERAGE PROBABILITY FOR EACH CLASS PER SAMPLE

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS

IN MULTILABEL CLASSIFICATION THIS IS THE SUBSET ACCURACY WHICH IS A HARSH METRIC SINCE YOU REQUIRE FOR EACH SAMPLE THAT EACH LABEL SET BE CORRECTLY PREDICTED

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES

1718 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE LABELS FOR X  
SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT MEAN ACCURACY OF SELF PREDICTX WRT Y

SETPARAMS SELFPARAMS

SETTING THE PARAMETERS FOR THE ENSEMBLE ESTIMATOR

VALID PARAMETER KEYS CAN BE LISTED WITH GETPARAMS

PARAMETERS

PARAMS KEYWORD ARGUMENTS SPECIFIC PARAMETERS USING EG  
SETPARAMSPARAMETERNAMEVALUE IN ADDITION TO SETTING THE PARAMETERS OF  
THE ENSEMBLE ESTIMATOR THE INDIVIDUAL ESTIMATORS OF THE ENSEMBLE ESTIMATOR CAN ALSO BE  
SET OR REPLACED BY SETTING THEM TO NONE

EXAMPLES

IN THIS EXAMPLE THE RANDOMFORESTCLASSIFIER IS REMOVED CLF1 LOGISTICREGRESSION CLF2 RANDOM  
FORESTCLASSIFIER ECLF VOTINGCLASSIFIERESTIMATORS'LR' CLF1 'RF' CLF2 ECLFSETPARAMSRFNONE

TRANSFORM SELF X

RETURN CLASS LABELS OR PROBABILITIES FOR X FOR EACH ESTIMATOR

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES TRAINING VECTORS WHERE  
NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

RETURNS

PROBABILITIESORLABELS

IFVOTINGSOFT AND FLATTENTTRANSFORMTRUE RETURNS ARRAYLIKE OF SHAPE  
NCLASSIFIERS NSAMPLES NCLASSES BEING CLASS PROBABILITIES CALCULATED BY EACH CLASSIFIER

IFVOTINGSOFT AND FLATTENTTRANSFORMFALSE ARRAYLIKE OF SHAPE  
NCLASSIFIERS NSAMPLES NCLASSES

IFVOTINGHARD ARRAYLIKE OF SHAPE NSAMPLES NCLASSIFIERS BEING CLASS LABELS  
PREDICTED BY EACH CLASSIFIER

EXAMPLES USING SKLEARNENSEMBLEVOTINGCLASSIFIER

- PLOT THE DECISION BOUNDARIES OF A VOTINGCLASSIFIER
- PLOT CLASS PROBABILITIES CALCULATED BY THE VOTINGCLASSIFIER

6128SKLEARNENSEMBLE VOTINGREGRESSOR

CLASSSSKLEARNENSEMBLE VOTINGREGRESSOR ESTIMATORS WEIGHTSNONE NJOBSNONE

PREDICTION VOTING REGRESSOR FOR UNFITTED ESTIMATORS

NEW IN VERSION 021

612SKLEARNENSEMBLE ENSEMBLE METHODS 1719

SCIKITLEARN USER GUIDE RELEASE 0213

A VOTING REGRESSOR IS AN ENSEMBLE METAESTIMATOR THAT FITS BASE REGRESSORS EACH ON THE WHOLE DATASET IT THEN AVERAGES THE INDIVIDUAL PREDICTIONS TO FORM A FINAL PREDICTION

READ MORE IN THE USER GUIDE

PARAMETERS

ESTIMATORS LIST OF STRING ESTIMATOR TUPLES INVOKING THE FIT METHOD ON THE

VOTINGREGRESSOR WILL FIT CLONES OF THOSE ORIGINAL ESTIMATORS THAT WILL BE STORED

IN THE CLASS ATTRIBUTE SELFESTIMATORS AN ESTIMATOR CAN BE SET TO NONE OR DROP

USING SETPARAMS

WEIGHTS ARRAYLIKE SHAPE NREGRESSORS OPTIONAL DEFAULT'NONE' SEQUENCE OF WEIGHTS

FLOAT OR INT TO WEIGHT THE OCCURRENCES OF PREDICTED VALUES BEFORE AVERAGING USES

UNIFORM WEIGHTS IF NONE

NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF JOBS TO RUN IN PARALLEL FOR FIT

NONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL

PROCESSORS SEE GLOSSARY FOR MORE DETAILS

ATTRIBUTES

ESTIMATORS LIST OF REGRESSORS THE COLLECTION OF FITTED SUBESTIMATORS AS DEFINED IN

ESTIMATORS THAT ARE NOT NONE

NAMEDESTIMATORS BUNCH OBJECT A DICTIONARY WITH ATTRIBUTE ACCESS ATTRIBUTE TO ACCESS ANY

FITTED SUBESTIMATORS BY NAME

SEE ALSO

VOTINGCLASSIFIER SOFT V OTINGMAJORITY RULE CLASSIFIER

EXAMPLES

```
import numpy as np
from sklearn.linear_model import LinearRegression
from sklearn.ensemble import RandomForestRegressor
from sklearn.ensemble import VotingRegressor

R1 = LinearRegression
R2 = RandomForestRegressor(n_estimators=10, random_state=1)
X = np.array([1, 2, 4, 3, 9, 4, 16, 5, 25, 6, 36])
Y = np.array([6, 12, 20, 30, 42])
ER = VotingRegressor([R1, R2])
print(ER.fit(X, Y).predict(X))
33 57 118 197 28 403
```

METHODS

FITSELF X Y SAMPLEWEIGHT FIT THE ESTIMATORS

FITTRANSFORM SELF X Y FIT TO DATA THEN TRANSFORM IT

GETPARAMS SELF DEEP GET THE PARAMETERS OF THE ENSEMBLE ESTIMATOR

PREDICT SELF X PREDICT REGRESSION TARGET FOR X

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

SETPARAMS SELF PARAMS SETTING THE PARAMETERS FOR THE ENSEMBLE ESTIMATOR

CONTINUED ON NEXT PAGE

1720 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

TABLE 679 – CONTINUED FROM PREVIOUS PAGE

TRANSFORM SELF X RETURN PREDICTIONS FOR X FOR EACH ESTIMATOR

INIT SELFESTIMATORS WEIGHTSNONE NJOBSNONE

FITSELFXYSAMPLEWEIGHTNONE

FIT THE ESTIMATORS

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES TRAINING VECTORS WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

YARRAYLIKE SHAPE NSAMPLES TARGET VALUES

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OR NONE SAMPLE WEIGHTS IF NONE THEN SAMPLES ARE EQUALLY WEIGHTED NOTE THAT THIS IS SUPPORTED ONLY IF ALL UNDERLYING ESTIMATORS SUPPORT SAMPLE WEIGHTS

RETURNS

SELF OBJECT

FITTRANSFORM SELFXYNONE FITPARAMS

FIT TO DATA THEN TRANSFORM IT

FITS TRANSFORMER TO X AND Y WITH OPTIONAL PARAMETERS FITPARAMS AND RETURNS A TRANSFORMED VERSION OF X

PARAMETERS

XNUMPY ARRAY OF SHAPE NSAMPLES NFEATURES TRAINING SET

YNUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES

RETURNS

XNEW NUMPY ARRAY OF SHAPE NSAMPLES NFEATURESNEW TRANSFORMED ARRAY

GETPARAMS SELFDEEPTTRUE

GET THE PARAMETERS OF THE ENSEMBLE ESTIMATOR

PARAMETERS

DEEP BOOL SETTING IT TO TRUE GETS THE VARIOUS ESTIMATORS AND THE PARAMETERS OF THE ESTIMATORS AS WELL

PREDICTSELFXY

PREDICT REGRESSION TARGET FOR X

THE PREDICTED REGRESSION TARGET OF AN INPUT SAMPLE IS COMPUTED AS THE MEAN PREDICTED REGRESSION TARGETS OF THE ESTIMATORS IN THE ENSEMBLE

PARAMETERS

XARRAYLIKE SPARSE MATRIX OF SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES

RETURNS

YARRAY OF SHAPE NSAMPLES THE PREDICTED VALUES

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{U}{V}$  WHERE U IS THE RESIDUAL SUM OF SQUARES  $Y_{TRUE} - Y_{PRED}$  2SUM AND V IS THE TOTAL SUM OF SQUARES  $Y_{TRUE} - Y_{TRUEMEAN}$  2SUM THE BEST POSSIBLE SCORE

612SKLEARNENSEMBLE ENSEMBLE METHODS 1721

SCIKITLEARN USER GUIDE RELEASE 0213

IS 10 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 00

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES FOR SOME ESTIMATORS THIS MAY BE A PRECOMPUTED KERNEL MATRIX INSTEAD SHAPE NSAMPLES NSAMPLESFITTED WHERE NSAMPLESFITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE VALUES FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT R2 OF SELF PREDICT X WRT Y

NOTES

THE R2 SCORE USED WHEN CALLING SCORE ON A REGRESSOR WILL USE MULTIOUTPUTUNIFORMAVERAGE FROM VERSION 023 TO KEEP CONSISTENT WITH METRICSR2SCORE THIS WILL INFLUENCE THE SCORE METHOD OF ALL THE MULTIOUTPUT REGRESSORS EXCEPT FOR MULTIOUTPUTMULTIOUTPUTREGRESSOR TO SPECIFY THE DEFAULT VALUE MANUALLY AND AVOID THE WARNING PLEASE EITHER CALL METRICSR2SCORE DIRECTLY OR MAKE A CUSTOM SCORER WITH METRICSMAKESCORER THE BUILTIN SCORER R2 USES MULTIOUTPUTUNIFORMAVERAGE

SETPARAMS SELFPARAMS

SETTING THE PARAMETERS FOR THE ENSEMBLE ESTIMATOR

VALID PARAMETER KEYS CAN BE LISTED WITH GETPARAMS

PARAMETERS

PARAMS KEYWORD ARGUMENTS SPECIFIC PARAMETERS USING EG

SETPARAMSPARAMETERNAMEVALUE IN ADDITION TO SETTING THE PARAMETERS OF THE ENSEMBLE ESTIMATOR THE INDIVIDUAL ESTIMATORS OF THE ENSEMBLE ESTIMATOR CAN ALSO BE SET OR REPLACED BY SETTING THEM TO NONE

EXAMPLES

IN THIS EXAMPLE THE RANDOMFORESTCLASSIFIER IS REMOVED CLF1 LOGISTICREGRESSION CLF2 RANDOM FORESTCLASSIFIER ECLF V OTINGCLASSIFIERESTIMATORS'LR' CLF1 'RF' CLF2 ECLFSETPARAMSRFNONE

TRANSFORM SELF X

RETURN PREDICTIONS FOR X FOR EACH ESTIMATOR

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES

RETURNS

PREDICTIONS ARRAYLIKE OF SHAPE NSAMPLES NCLASSIFIERS BEING VALUES PREDICTED BY EACH REGRESSOR

EXAMPLES USING SKLEARNENSEMBLEVOTINGREGRESSOR

- PLOT INDIVIDUAL AND VOTING REGRESSION PREDICTIONS

1722 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

6129SKLEARNENSEMBLE HISTGRADIENTBOOSTINGREGRESSOR  
CLASSSSKLEARNENSEMBLE HISTGRADIENTBOOSTINGREGRESSOR LOSS'LEASTSQUARES' LEARN  
INGRATE01 MAXITER100  
MAXLEAFNODES31  
MAXDEPTHNONE  
MINSAMPLESLEAF20  
L2REGULARIZATION00  
MAXBINS256 SCORINGNONE  
VALIDATIONFRACTION01  
NITERNOCHANGENONE  
TOL1E07 VERBOSE0 RAN  
DOMSTATENONE  
HISTOGRAMBASED GRADIENT BOOSTING REGRESSION TREE  
THIS ESTIMATOR IS MUCH FASTER THAN GRADIENTBOOSTINGREGRESSOR FOR BIG DATASETS NSAMPLES 10 000  
THE INPUT DATA XIS PREBINNED INTO INTEGERVERALUED BINS WHICH CONSIDERABLY REDUCES THE NUMBER OF SPLITTING  
POINTS TO CONSIDER AND ALLOWS THE ALGORITHM TO LEVERAGE INTEGERBASED DATA STRUCTURES FOR SMALL SAMPLE SIZES  
GRADIENTBOOSTINGREGRESSOR MIGHT BE PREFERRED SINCE BINNING MAY LEAD TO SPLIT POINTS THAT ARE TOO  
APPROXIMATE IN THIS SETTING  
THIS IMPLEMENTATION IS INSPIRED BY LIGHTGBM  
NOTE THIS ESTIMATOR IS STILL EXPERIMENTAL FOR NOW THE PREDICTIONS AND THE API MIGHT CHANGE WITHOUT ANY  
DEPRECATION CYCLE TO USE IT YOU NEED TO EXPLICITLY IMPORT ENABLEHISTGRADIENTBOOSTING  
EXPLICITLY REQUIRE THIS EXPERIMENTAL FEATURE  
FROM SKLEARNEXPERIMENTAL IMPORT ENABLEHISTGRADIENTBOOSTING NOQA  
NOW YOU CAN IMPORT NORMALLY FROM ENSEMBLE  
FROM SKLEARNENSEMBLE IMPORT HISTGRADIENTBOOSTINGCLASSIFIER  
PARAMETERS  
LOSS 'LEASTSQUARES' OPTIONAL DEFAULT'LEASTSQUARES' THE LOSS FUNCTION TO USE IN THE BOOST  
ING PROCESS NOTE THAT THE "LEAST SQUARES" LOSS ACTUALLY IMPLEMENTS AN "HALF LEAST SQUARES  
LOSS" TO SIMPLIFY THE COMPUTATION OF THE GRADIENT  
LEARNINGRATE FLOAT OPTIONAL DEFAULT01 THE LEARNING RATE ALSO KNOWN AS SHRINKAGE THIS  
IS USED AS A MULTIPLICATIVE FACTOR FOR THE LEAVES VALUES USE 1FOR NO SHRINKAGE  
MAXITER INT OPTIONAL DEFAULT100 THE MAXIMUM NUMBER OF ITERATIONS OF THE BOOSTING PRO  
CESS IE THE MAXIMUM NUMBER OF TREES  
MAXLEAFNODES INT OR NONE OPTIONAL DEFAULT31 THE MAXIMUM NUMBER OF LEAVES FOR EACH  
TREE MUST BE STRICTLY GREATER THAN 1 IF NONE THERE IS NO MAXIMUM LIMIT  
MAXDEPTH INT OR NONE OPTIONAL DEFAULTNONE THE MAXIMUM DEPTH OF EACH TREE THE  
DEPTH OF A TREE IS THE NUMBER OF NODES TO GO FROM THE ROOT TO THE DEEPEST LEAF MUST BE  
STRICTLY GREATER THAN 1 DEPTH ISN'T CONSTRAINED BY DEFAULT  
MINSAMPLESLEAF INT OPTIONAL DEFAULT20 THE MINIMUM NUMBER OF SAMPLES PER LEAF FOR  
SMALL DATASETS WITH LESS THAN A FEW HUNDRED SAMPLES IT IS RECOMMENDED TO LOWER THIS VALUE  
SINCE ONLY VERY SHALLOW TREES WOULD BE BUILT  
L2REGULARIZATION FLOAT OPTIONAL DEFAULT0 THE L2 REGULARIZATION PARAMETER USE 0FOR NO  
REGULARIZATION DEFAULT  
612SKLEARNENSEMBLE ENSEMBLE METHODS 1723

SCIKITLEARN USER GUIDE RELEASE 0213

MAXBINS INT OPTIONAL DEFAULT256 THE MAXIMUM NUMBER OF BINS TO USE BEFORE TRAINING  
EACH FEATURE OF THE INPUT ARRAY XIS BINNED INTO AT MOST MAXBINS BINS WHICH ALLOWS FOR A  
MUCH FASTER TRAINING STAGE FEATURES WITH A SMALL NUMBER OF UNIQUE VALUES MAY USE LESS THAN  
MAXBINS BINS MUST BE NO LARGER THAN 256  
SCORING STR OR CALLABLE OR NONE OPTIONAL DEFAULTNONE SCORING PARAMETER TO USE FOR EARLY  
STOPPING IT CAN BE A SINGLE STRING SEE THE SCORING PARAMETER DEFINING MODEL EVALUATION  
RULES OR A CALLABLE SEE DEFINING YOUR SCORING STRATEGY FROM METRIC FUNCTIONS IF NONE THE  
ESTIMATOR'S DEFAULT SCORER IS USED IF SCORINGLOSS EARLY STOPPING IS CHECKED WRT  
THE LOSS VALUE ONLY USED IF NITERNOCHANGE IS NOT NONE  
VALIDATIONFRACTION INT OR FLOAT OR NONE OPTIONAL DEFAULT01 PROPORTION OR ABSOLUTE SIZE  
OF TRAINING DATA TO SET ASIDE AS VALIDATION DATA FOR EARLY STOPPING IF NONE EARLY STOPPING IS  
DONE ON THE TRAINING DATA ONLY USED IF NITERNOCHANGE IS NOT NONE  
NITERNOCHANGE INT OR NONE OPTIONAL DEFAULTNONE USED TO DETERMINE WHEN TO "EARLY  
STOP" THE FITTING PROCESS IS STOPPED WHEN NONE OF THE LAST NITERNOCHANGE SCORES  
ARE BETTER THAN THE "NITERNOCHANGE 1" THTOLAST ONE UP TO SOME TOLERANCE IF NONE OR  
0 NO EARLYSTOPPING IS DONE  
TOLFLOAT OR NONE OPTIONAL DEFAULT1E7 THE ABSOLUTE TOLERANCE TO USE WHEN COMPARING  
SCORES DURING EARLY STOPPING THE HIGHER THE TOLERANCE THE MORE LIKELY WE ARE TO EARLY  
STOP HIGHER TOLERANCE MEANS THAT IT WILL BE HARDER FOR SUBSEQUENT ITERATIONS TO BE CONSIDERED  
AN IMPROVEMENT UPON THE REFERENCE SCORE  
VERBOSE INT OPTIONAL DEFAULT0 THE VERBOSITY LEVEL IF NOT ZERO PRINT SOME INFORMATION  
ABOUT THE FITTING PROCESS  
RANDOMSTATE INT NRANDOMRANDOMSTATEINSTANCE OR NONE OPTIONAL DEFAULTNONE  
PSEUDORANDOM NUMBER GENERATOR TO CONTROL THE SUBSAMPLING IN THE BINNING PROCESS AND  
THE TRAINVALIDATION DATA SPLIT IF EARLY STOPPING IS ENABLED SEE RANDOMSTATE  
ATTRIBUTES  
NITER INT THE NUMBER OF ITERATIONS AS SELECTED BY EARLY STOPPING IF NITERNOCHANGE IS NOT  
NONE OTHERWISE IT CORRESPONDS TO MAXITER  
NTREESPERITERATION INT THE NUMBER OF TREE THAT ARE BUILT AT EACH ITERATION FOR REGRESSORS  
THIS IS ALWAYS 1  
TRAINSCORE NDARRAY SHAPE MAXITER 1 THE SCORES AT EACH ITERATION ON THE TRAINING DATA  
THE FIRST ENTRY IS THE SCORE OF THE ENSEMBLE BEFORE THE FIRST ITERATION SCORES ARE COMPUTED  
ACCORDING TO THE SCORING PARAMETER IF SCORING IS NOT 'LOSS' SCORES ARE COMPUTED ON A  
SUBSET OF AT MOST 10 000 SAMPLES EMPTY IF NO EARLY STOPPING  
VALIDATIONSCORE NDARRAY SHAPE MAXITER 1 THE SCORES AT EACH ITERATION ON THE HELD  
OUT VALIDATION DATA THE FIRST ENTRY IS THE SCORE OF THE ENSEMBLE BEFORE THE FIRST ITERATION  
SCORES ARE COMPUTED ACCORDING TO THE SCORING PARAMETER EMPTY IF NO EARLY STOPPING OR  
IFVALIDATIONFRACTION IS NONE  
EXAMPLES  
TO USE THIS EXPERIMENTAL FEATURE WE NEED TO EXPLICITLY ASK FOR IT  
FROM SKLEARNEXPERIMENTAL IMPORT ENABLEHISTGRADIENTBOOSTING NOQA  
FROM SKLEARNENSEMBLE IMPORT HISTGRADIENTBOOSTINGREGRESSOR  
FROM SKLEARNDATASETS IMPORT LOADBOSTON  
X Y LOADBOSTONRETURNXY TRUE  
EST HISTGRADIENTBOOSTINGREGRESSORFITX Y  
1724 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

ESTSCOREX Y

098

METHODS

FITSELF X Y FIT THE GRADIENT BOOSTING MODEL

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PREDICT SELF X PREDICT VALUES FOR X

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELFLOSS'LEASTSQUARES' LEARNINGRATE01 MAXITER100 MAXLEAFNODES31

MAXDEPTHNONE MINSAMPLESLEAF20 L2REGULARIZATION00 MAXBINS256 SCOR

INGNONE VALIDATIONFRACTION01 NITERNOCHANGENONE TOL1E07 VERBOSE0

RANDOMSTATENONE

FITSELFXY

FIT THE GRADIENT BOOSTING MODEL

PARAMETERS

XARRAYLIKE SHAPENSAMPLES NFEATURES THE INPUT SAMPLES

YARRAYLIKE SHAPENSAMPLES TARGET VALUES

RETURNS

SELF OBJECT

GETPARAMS SELFDEEPTREE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELF

PREDICT VALUES FOR X

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES

RETURNS

YNDARRAY SHAPE NSAMPLES THE PREDICTED VALUES

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{U}{V}$  WHERE U IS THE RESIDUAL SUM OF SQUARES  $Y_{TRUE} - Y_{PRED}$

$2 \sum$  AND V IS THE TOTAL SUM OF SQUARES  $Y_{TRUE} - Y_{TRUEMEAN}$   $2 \sum$  THE BEST POSSIBLE SCORE

IS 10 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS

PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 00

612SKLEARNENSEMBLE ENSEMBLE METHODS 1725

SCIKITLEARN USER GUIDE RELEASE 0213

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES FOR SOME ESTIMATORS THIS MAY BE A PRECOMPUTED KERNEL MATRIX INSTEAD SHAPE NSAMPLES NSAMPLESFITTED WHERE NSAMPLESFITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE VALUES FOR X SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT R2 OF SELF PREDICT X WRT Y

NOTES

THE R2 SCORE USED WHEN CALLING SCORE ON A REGRESSOR WILL USE MULTIOUTPUTUNIFORMAVERAGE FROM VERSION 023 TO KEEP CONSISTENT WITH METRICSR2SCORE THIS WILL INFLUENCE THE SCORE METHOD OF ALL THE MULTIOUTPUT REGRESSORS EXCEPT FOR MULTIOUTPUTMULTIOUTPUTREGRESSOR TO SPECIFY THE DEFAULT VALUE MANUALLY AND AVOID THE WARNING PLEASE EITHER CALL METRICSR2SCORE DIRECTLY OR MAKE A CUSTOM SCORER WITH METRICSMAKESCORER THE BUILTIN SCORER R2 USES MULTIOUTPUTUNIFORMAVERAGE

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

61210SKLEARNENSEMBLE HISTGRADIENTBOOSTINGCLASSIFIER  
CLASSSSKLEARNENSEMBLE HISTGRADIENTBOOSTINGCLASSIFIER LOSS'AUTO' LEARN

INGRATE01 MAXITER100

MAXLEAFNODES31

MAXDEPTHNONE

MINSAMPLESLEAF20

L2REGULARIZATION00

MAXBINS256 SCORINGNONE

VALIDATIONFRACTION01

NITERNOCHANGENONE

TOL1E07 VERBOSE0 RAN

DOMSTATENONE

HISTOGRAMBASED GRADIENT BOOSTING CLASSIFICATION TREE

THIS ESTIMATOR IS MUCH FASTER THAN GRADIENTBOOSTINGCLASSIFIER FOR BIG DATASETS NSAMPLES 10000 THE INPUT DATA X IS PREBINNED INTO INTEGERVALUED BINS WHICH CONSIDERABLY REDUCES THE NUMBER OF SPLITTING POINTS TO CONSIDER AND ALLOWS THE ALGORITHM TO LEVERAGE INTEGERBASED DATA STRUCTURES FOR SMALL SAMPLE SIZES GRADIENTBOOSTINGCLASSIFIER MIGHT BE PREFERRED SINCE BINNING MAY LEAD TO SPLIT POINTS THAT ARE TOO APPROXIMATE IN THIS SETTING

THIS IMPLEMENTATION IS INSPIRED BY LIGHTGBM

1726 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

NOTE THIS ESTIMATOR IS STILL EXPERIMENTAL FOR NOW THE PREDICTIONS AND THE API MIGHT CHANGE WITHOUT ANY DEPRECATION CYCLE TO USE IT YOU NEED TO EXPLICITLY IMPORT `ENABLEHISTGRADIENTBOOSTING` EXPLICITLY REQUIRE THIS EXPERIMENTAL FEATURE

FROM `SKLEARNEXPERIMENTAL` IMPORT `ENABLEHISTGRADIENTBOOSTING` NOQA

NOW YOU CAN IMPORT NORMALLY FROM `ENSEMBLE`

FROM `SKLEARNENSEMBLE` IMPORT `HISTGRADIENTBOOSTINGCLASSIFIER`

PARAMETERS

LOSS 'AUTO' 'BINARYCROSSENTROPY' 'CATEGORICALCROSSENTROPY' OPTIONAL DEFAULT 'AUTO'

THE LOSS FUNCTION TO USE IN THE BOOSTING PROCESS 'BINARYCROSSENTROPY' ALSO KNOWN AS LOGISTIC LOSS IS USED FOR BINARY CLASSIFICATION AND GENERALIZES TO 'CATEGORICALCROSSENTROPY' FOR MULTICLASS CLASSIFICATION 'AUTO' WILL AUTOMATICALLY CHOOSE EITHER LOSS DEPENDING ON THE NATURE OF THE PROBLEM

LEARNINGRATE FLOAT OPTIONAL DEFAULT 0.1 THE LEARNING RATE ALSO KNOWN AS SHRINKAGE THIS IS USED AS A MULTIPLICATIVE FACTOR FOR THE LEAVES VALUES USE 1 FOR NO SHRINKAGE

MAXITER INT OPTIONAL DEFAULT 100 THE MAXIMUM NUMBER OF ITERATIONS OF THE BOOSTING PROCESS IE THE MAXIMUM NUMBER OF TREES FOR BINARY CLASSIFICATION FOR MULTICLASS CLASSIFICATION `N_CLASSES` TREES PER ITERATION ARE BUILT

MAXLEAFNODES INT OR NONE OPTIONAL DEFAULT 31 THE MAXIMUM NUMBER OF LEAVES FOR EACH TREE MUST BE STRICTLY GREATER THAN 1 IF NONE THERE IS NO MAXIMUM LIMIT

MAXDEPTH INT OR NONE OPTIONAL DEFAULT None THE MAXIMUM DEPTH OF EACH TREE THE DEPTH OF A TREE IS THE NUMBER OF NODES TO GO FROM THE ROOT TO THE DEEPEST LEAF MUST BE STRICTLY GREATER THAN 1 DEPTH ISN'T CONSTRAINED BY DEFAULT

MINSAMPLESLEAF INT OPTIONAL DEFAULT 20 THE MINIMUM NUMBER OF SAMPLES PER LEAF FOR SMALL DATASETS WITH LESS THAN A FEW HUNDRED SAMPLES IT IS RECOMMENDED TO LOWER THIS VALUE SINCE ONLY VERY SHALLOW TREES WOULD BE BUILT

L2REGULARIZATION FLOAT OPTIONAL DEFAULT 0 THE L2 REGULARIZATION PARAMETER USE 0 FOR NO REGULARIZATION

MAXBINS INT OPTIONAL DEFAULT 256 THE MAXIMUM NUMBER OF BINS TO USE BEFORE TRAINING EACH FEATURE OF THE INPUT ARRAY `X` IS BINNED INTO AT MOST `MAXBINS` BINS WHICH ALLOWS FOR A MUCH FASTER TRAINING STAGE FEATURES WITH A SMALL NUMBER OF UNIQUE VALUES MAY USE LESS THAN `MAXBINS` BINS MUST BE NO LARGER THAN 256

SCORING STR OR CALLABLE OR NONE OPTIONAL DEFAULT None SCORING PARAMETER TO USE FOR EARLY STOPPING IT CAN BE A SINGLE STRING SEE THE SCORING PARAMETER DEFINING MODEL EVALUATION RULES OR A CALLABLE SEE DEFINING YOUR SCORING STRATEGY FROM METRIC FUNCTIONS IF NONE THE ESTIMATOR'S DEFAULT SCORER IS USED IF `SCORINGLOSS` EARLY STOPPING IS CHECKED WRT THE LOSS VALUE ONLY USED IF `NITERNOCHANGE` IS NOT None

VALIDATIONFRACTION INT OR FLOAT OR NONE OPTIONAL DEFAULT 0.1 PROPORTION OR ABSOLUTE SIZE OF TRAINING DATA TO SET ASIDE AS VALIDATION DATA FOR EARLY STOPPING IF NONE EARLY STOPPING IS DONE ON THE TRAINING DATA

NITERNOCHANGE INT OR NONE OPTIONAL DEFAULT None USED TO DETERMINE WHEN TO "EARLY STOP" THE FITTING PROCESS IS STOPPED WHEN NONE OF THE LAST `NITERNOCHANGE` SCORES ARE BETTER THAN THE "NITERNOCHANGE - 1" TH TOLAST ONE UP TO SOME TOLERANCE IF NONE OR 0 NO EARLY STOPPING IS DONE

612 SKLEARNENSEMBLE ENSEMBLE METHODS 1727

SCIKITLEARN USER GUIDE RELEASE 0213

TOLFLOAT OR NONE OPTIONAL DEFAULT1E7 THE ABSOLUTE TOLERANCE TO USE WHEN COMPARING SCORES THE HIGHER THE TOLERANCE THE MORE LIKELY WE ARE TO EARLY STOP HIGHER TOLERANCE MEANS THAT IT WILL BE HARDER FOR SUBSEQUENT ITERATIONS TO BE CONSIDERED AN IMPROVEMENT UPON THE REFERENCE SCORE

VERBOSE INT OPTIONAL DEFAULT0 THE VERBOSITY LEVEL IF NOT ZERO PRINT SOME INFORMATION ABOUT THE FITTING PROCESS

RANDOMSTATE INT NPRANDOMRANDOMSTATEINSTANCE OR NONE OPTIONAL DEFAULTNONE PSEUDORANDOM NUMBER GENERATOR TO CONTROL THE SUBSAMPLING IN THE BINNING PROCESS AND THE TRAINVALIDATION DATA SPLIT IF EARLY STOPPING IS ENABLED SEE RANDOMSTATE

ATTRIBUTES

NITER INT THE NUMBER OF ESTIMATORS AS SELECTED BY EARLY STOPPING IF NITERNOCHANGE IS NOT NONE OTHERWISE IT CORRESPONDS TO MAXITER

NTREESPERITERATION INT THE NUMBER OF TREE THAT ARE BUILT AT EACH ITERATION THIS IS EQUAL TO 1 FOR BINARY CLASSIFICATION AND TO NCLASSES FOR MULTICLASS CLASSIFICATION

TRAINSORE NDARRAY SHAPE MAXITER 1 THE SCORES AT EACH ITERATION ON THE TRAINING DATA THE FIRST ENTRY IS THE SCORE OF THE ENSEMBLE BEFORE THE FIRST ITERATION SCORES ARE COMPUTED ACCORDING TO THE SCORING PARAMETER IF SCORING IS NOT 'LOSS' SCORES ARE COMPUTED ON A SUBSET OF AT MOST 10 000 SAMPLES EMPTY IF NO EARLY STOPPING

VALIDATIONSORE NDARRAY SHAPE MAXITER 1 THE SCORES AT EACH ITERATION ON THE HELD OUT VALIDATION DATA THE FIRST ENTRY IS THE SCORE OF THE ENSEMBLE BEFORE THE FIRST ITERATION SCORES ARE COMPUTED ACCORDING TO THE SCORING PARAMETER EMPTY IF NO EARLY STOPPING OR IFVALIDATIONFRACTION IS NONE

EXAMPLES

TO USE THIS EXPERIMENTAL FEATURE WE NEED TO EXPLICITLY ASK FOR IT

```
FROM SKLEARNEXPERIMENTAL IMPORT ENABLEHISTGRADIENTBOOSTING NOQA
FROM SKLEARNENSEMBLE IMPORT HISTGRADIENTBOOSTINGREGRESSOR
FROM SKLEARNDATASETS IMPORT LOADIRIS
X Y LOADIRISRETURNXY TRUE
CLF HISTGRADIENTBOOSTINGCLASSIFIERFITX Y
CLFSOREX Y
10
```

METHODS

DECISIONFUNCTION SELF X COMPUTE THE DECISION FUNCTION OF X

FITSELF X Y FIT THE GRADIENT BOOSTING MODEL

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PREDICT SELF X PREDICT CLASSES FOR X

PREDICTPROBA SELF X PREDICT CLASS PROBABILITIES FOR X

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

1728 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213  
INIT SELF LOSS'AUTO' LEARNINGRATE01 MAXITER100 MAXLEAFNODES31  
MAXDEPTHNONE MINSAMPLESLEAF20 L2REGULARIZATION00 MAXBINS256 SCOR  
INGNONE VALIDATIONFRACTION01 NITERNOCHANGENONE TOL1E07 VERBOSE0  
RANDOMSTATENONE  
DECISIONFUNCTION SELF  
COMPUTE THE DECISION FUNCTION OF X  
PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES  
RETURNS  
DECISION NDARRAY SHAPE NSAMPLES OR NSAMPLES NTREESPERITERATION THE RAW PRE  
DICTED VALUES IE THE SUM OF THE TREES LEAVES FOR EACH SAMPLE NTREESPERITERATION IS  
EQUAL TO THE NUMBER OF CLASSES IN MULTICLASS CLASSIFICATION  
FITSELFXY  
FIT THE GRADIENT BOOSTING MODEL  
PARAMETERS  
XARRAYLIKE SHAPENSAMPLES NFEATURES THE INPUT SAMPLES  
YARRAYLIKE SHAPENSAMPLES TARGET VALUES  
RETURNS  
SELF OBJECT  
GETPARAMS SELFDEEPTREE  
GET PARAMETERS FOR THIS ESTIMATOR  
PARAMETERS  
DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBOBJECTS THAT ARE ESTIMATORS  
RETURNS  
PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES  
PREDICTSELF  
PREDICT CLASSES FOR X  
PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES  
RETURNS  
YNDARRAY SHAPE NSAMPLES THE PREDICTED CLASSES  
PREDICTPROBA SELF  
PREDICT CLASS PROBABILITIES FOR X  
PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES  
RETURNS  
PNDARRAY SHAPE NSAMPLES NCLASSES THE CLASS PROBABILITIES OF THE INPUT SAMPLES  
612SKLEARNENSEMBLE ENSEMBLE METHODS 1729

SCIKITLEARN USER GUIDE RELEASE 0213

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS

IN MULTILABEL CLASSIFICATION THIS IS THE SUBSET ACCURACY WHICH IS A HARSH METRIC SINCE YOU REQUIRE FOR EACH SAMPLE THAT EACH LABEL SET BE CORRECTLY PREDICTED

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE LABELS FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT MEAN ACCURACY OF SELF-PREDICTX WRT Y

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

613SKLEARNEXCEPTIONS EXCEPTIONS AND WARNINGS

THESKLEARNEXCEPTIONS MODULE INCLUDES ALL CUSTOM WARNINGS AND ERROR CLASSES USED ACROSS SCIKITLEARN

EXCEPTIONSCHANGEDBEHAVIORWARNING WARNING CLASS USED TO NOTIFY THE USER OF ANY CHANGE IN THE BEHAVIOR

EXCEPTIONSCONVERGENCEWARNING CUSTOM WARNING TO CAPTURE CONVERGENCE PROBLEMS

EXCEPTIONSDATACONVERSIONWARNING WARNING USED TO NOTIFY IMPLICIT DATA CONVERSIONS HAPPENING IN THE CODE

EXCEPTIONSDATADIMENSIONALITYWARNING CUSTOM WARNING TO NOTIFY POTENTIAL ISSUES WITH DATA DIMENSIONALITY

EXCEPTIONSEFFICIENCYWARNING WARNING USED TO NOTIFY THE USER OF INEFFICIENT COMPUTATION

EXCEPTIONSFITFAILEDWARNING WARNING CLASS USED IF THERE IS AN ERROR WHILE FITTING THE ESTIMATOR

EXCEPTIONSNOTFITTEDERROR EXCEPTION CLASS TO RAISE IF ESTIMATOR IS USED BEFORE FITTING

EXCEPTIONSNONBLASDOTWARNING WARNING USED WHEN THE DOT OPERATION DOES NOT USE BLAS

EXCEPTIONSUNDEFINEDMETRICWARNING WARNING USED WHEN THE METRIC IS INVALID

6131SKLEARNEXCEPTIONS CHANGEDBEHAVIORWARNING

CLASSSSKLEARNEXCEPTIONS CHANGEDBEHAVIORWARNING

WARNING CLASS USED TO NOTIFY THE USER OF ANY CHANGE IN THE BEHAVIOR

CHANGED IN VERSION 018 MOVED FROM SKLEARNBASE

ATTRIBUTES

1730 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

ARGS

METHODS

WITHTRACEBACK EXCEPTIONWITHTRACEBACKTB - SET SELFTRACEBACK TO TB AND RETURN SELF

WITHTRACEBACK

EXCEPTIONWITHTRACEBACKTB - SET SELFTRACEBACK TO TB AND RETURN SELF

6132SKLEARNEXCEPTIONS CONVERGENCEWARNING

CLASSSSKLEARNEXCEPTIONS CONVERGENCEWARNING

CUSTOM WARNING TO CAPTURE CONVERGENCE PROBLEMS

CHANGED IN VERSION 018 MOVED FROM SKLEARNUTILS

ATTRIBUTES

ARGS

METHODS

WITHTRACEBACK EXCEPTIONWITHTRACEBACKTB - SET SELFTRACEBACK TO TB AND RETURN SELF

WITHTRACEBACK

EXCEPTIONWITHTRACEBACKTB - SET SELFTRACEBACK TO TB AND RETURN SELF

EXAMPLES USING SKLEARNEXCEPTIONS CONVERGENCEWARNING

- MULTICLASS SPARSE LOGISITIC REGRESSION ON NEWGROUPS20
- EARLY STOPPING OF STOCHASTIC GRADIENT DESCENT
- COMPARE STOCHASTIC LEARNING STRATEGIES FOR MLPCLASSIFIER
- FEATURE DISCRETIZATION

6133SKLEARNEXCEPTIONS DATACONVERSIONWARNING

CLASSSSKLEARNEXCEPTIONS DATACONVERSIONWARNING

WARNING USED TO NOTIFY IMPLICIT DATA CONVERSIONS HAPPENING IN THE CODE

THIS WARNING OCCURS WHEN SOME INPUT DATA NEEDS TO BE CONVERTED OR INTERPRETED IN A WAY THAT MAY NOT MATCH THE USER'S EXPECTATIONS

FOR EXAMPLE THIS WARNING MAY OCCUR WHEN THE USER

- PASSES AN INTEGER ARRAY TO A FUNCTION WHICH EXPECTS FLOAT INPUT AND WILL CONVERT THE INPUT

613SKLEARNEXCEPTIONS EXCEPTIONS AND WARNINGS 1731

SCIKITLEARN USER GUIDE RELEASE 0213

- REQUESTS A NONCOPYING OPERATION BUT A COPY IS REQUIRED TO MEET THE IMPLEMENTATION’S DATATYPE EXPECTATIONS

- PASSES AN INPUT WHOSE SHAPE CAN BE INTERPRETED AMBIGUOUSLY

CHANGED IN VERSION 018 MOVED FROM SKLEARNUTILSVALIDATION

ATTRIBUTES

ARGS

METHODS

WITHTRACEBACK EXCEPTIONWITHTRACEBACKTB – SET SELFTRACEBACK TO TB AND RETURN SELF

WITHTRACEBACK

EXCEPTIONWITHTRACEBACKTB – SET SELFTRACEBACK TO TB AND RETURN SELF

6134SKLEARNEXCEPTIONS DATADIMENSIONALITYWARNING

CLASSSSKLEARNEXCEPTIONS DATADIMENSIONALITYWARNING

CUSTOM WARNING TO NOTIFY POTENTIAL ISSUES WITH DATA DIMENSIONALITY

FOR EXAMPLE IN RANDOM PROJECTION THIS WARNING IS RAISED WHEN THE NUMBER OF COMPONENTS WHICH QUANTIFIES THE DIMENSIONALITY OF THE TARGET PROJECTION SPACE IS HIGHER THAN THE NUMBER OF FEATURES WHICH QUANTIFIES THE DIMENSIONALITY OF THE ORIGINAL SOURCE SPACE TO IMPLY THAT THE DIMENSIONALITY OF THE PROBLEM WILL NOT BE REDUCED

CHANGED IN VERSION 018 MOVED FROM SKLEARNUTILS

ATTRIBUTES

ARGS

METHODS

WITHTRACEBACK EXCEPTIONWITHTRACEBACKTB – SET SELFTRACEBACK TO TB AND RETURN SELF

WITHTRACEBACK

EXCEPTIONWITHTRACEBACKTB – SET SELFTRACEBACK TO TB AND RETURN SELF

6135SKLEARNEXCEPTIONS EFFICIENCYWARNING

CLASSSSKLEARNEXCEPTIONS EFFICIENCYWARNING

WARNING USED TO NOTIFY THE USER OF INEFFICIENT COMPUTATION

THIS WARNING NOTIFIES THE USER THAT THE EFFICIENCY MAY NOT BE OPTIMAL DUE TO SOME REASON WHICH MAY BE INCLUDED AS A PART OF THE WARNING MESSAGE THIS MAY BE SUBCLASSED INTO A MORE SPECIFIC WARNING CLASS

NEW IN VERSION 018

ATTRIBUTES

1732 CHAPTER 6 API REFERENCE

```
SCIKITLEARN USER GUIDE RELEASE 0213
ARGS
METHODS
WITHTRACEBACK EXCEPTIONWITHTRACEBACKTB - SET SELFTRACEBACK
TO TB AND RETURN SELF
WITHTRACEBACK
EXCEPTIONWITHTRACEBACKTB - SET SELFTRACEBACK TO TB AND RETURN SELF
6136SKLEARNEXCEPTIONS FITFAILEDWARNING
CLASSSSKLEARNEXCEPTIONS FITFAILEDWARNING
WARNING CLASS USED IF THERE IS AN ERROR WHILE FITTING THE ESTIMATOR
THIS WARNING IS USED IN META ESTIMATORS GRIDSEARCHCV AND RANDOMIZEDSEARCHCV AND THE CROSSVALIDATION
HELPER FUNCTION CROSSVALSCORE TO WARN WHEN THERE IS AN ERROR WHILE FITTING THE ESTIMATOR
ATTRIBUTES
ARGS
EXAMPLES
FROM SKLEARNMODELSELECTION IMPORT GRIDSEARCHCV
FROM SKLEARN SVM IMPORT LINEARSVC
FROM SKLEARNEXCEPTIONS IMPORT FITFAILEDWARNING
IMPORT WARNINGS
WARNINGSSIMPLEFILTERALWAYS FITFAILEDWARNING
GS GRIDSEARCHCVLINEARSVC C 1 2 ERRORSORE0 CV2
X Y 1 2 3 4 5 6 7 8 0 0 1 1
WITH WARNINGSCATCHWARNINGSRECORD TRUEASW
TRY
GSFITX Y THIS WILL RAISE A VALUEERROR SINCE C IS 0
EXCEPT VALUEERROR
PASS
PRINTREPRW1MESSAGE

FITFAILEDWARNINGESTIMATOR FIT FAILED THE SCORE ON THIS TRAI NTEST
PARTITION FOR THESE PARAMETERS WILL BE SET TO 0000000
DETAILS NVALUEERROR PENALTY TERM MUST BE POSITIVE GOT C2N
CHANGED IN VERSION 018 MOVED FROM SKLEARN CROSSVALIDATION
METHODS
WITHTRACEBACK EXCEPTIONWITHTRACEBACKTB - SET SELFTRACEBACK
TO TB AND RETURN SELF
613SKLEARNEXCEPTIONS EXCEPTIONS AND WARNINGS 1733
```

SCIKITLEARN USER GUIDE RELEASE 0213

WITHTRACEBACK

EXCEPTIONWITHTRACEBACKTB - SET SELFTRACEBACK TO TB AND RETURN SELF

6137SKLEARNEXCEPTIONS NOTFITTEDERROR

CLASSSKLEARNEXCEPTIONS NOTFITTEDERROR

EXCEPTION CLASS TO RAISE IF ESTIMATOR IS USED BEFORE FITTING

THIS CLASS INHERITS FROM BOTH VALUEERROR AND ATTRIBUTEERROR TO HELP WITH EXCEPTION HANDLING AND BACKWARD COMPATIBILITY

ATTRIBUTES

ARGS

EXAMPLES

```
FROM SKLEARN SVM IMPORT LINEAR SVC
FROM SKLEARN EXCEPTIONS IMPORT NOTFITTEDERROR
TRY
LINEAR SVC PREDICT 1 2 2 3 3 4
EXCEPT NOTFITTEDERROR ASE
PRINT REPRE
```

NOTFITTEDERROR THIS LINEAR SVC INSTANCE IS NOT FITTED YET

CHANGED IN VERSION 018 MOVED FROM SKLEARN UTILS VALIDATION

METHODS

WITHTRACEBACK EXCEPTIONWITHTRACEBACKTB - SET SELFTRACEBACK

TO TB AND RETURN SELF

WITHTRACEBACK

EXCEPTIONWITHTRACEBACKTB - SET SELFTRACEBACK TO TB AND RETURN SELF

6138SKLEARNEXCEPTIONS NONBLASDOTWARNING

CLASSSKLEARNEXCEPTIONS NONBLASDOTWARNING

WARNING USED WHEN THE DOT OPERATION DOES NOT USE BLAS

THIS WARNING IS USED TO NOTIFY THE USER THAT BLAS WAS NOT USED FOR DOT OPERATION AND HENCE THE EFFICIENCY MAY BE AFFECTED

CHANGED IN VERSION 018 MOVED FROM SKLEARN UTILS VALIDATION EXTENDS EFFICIENCYWARNING

ATTRIBUTES

ARGS

1734 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

METHODS

WITHTRACEBACK EXCEPTIONWITHTRACEBACKTB - SET SELFTRACEBACK  
TO TB AND RETURN SELF

WITHTRACEBACK

EXCEPTIONWITHTRACEBACKTB - SET SELFTRACEBACK TO TB AND RETURN SELF

6139SKLEARNEXCEPTIONS UNDEFINEDMETRICWARNING

CLASSSSKLEARNEXCEPTIONS UNDEFINEDMETRICWARNING

WARNING USED WHEN THE METRIC IS INVALID

CHANGED IN VERSION 018 MOVED FROM SKLEARNBASE

ATTRIBUTES

ARGS

METHODS

WITHTRACEBACK EXCEPTIONWITHTRACEBACKTB - SET SELFTRACEBACK  
TO TB AND RETURN SELF

WITHTRACEBACK

EXCEPTIONWITHTRACEBACKTB - SET SELFTRACEBACK TO TB AND RETURN SELF

614SKLEARNEXPERIMENTAL EXPERIMENTAL

THESKLEARNEXPERIMENTAL MODULE PROVIDES IMPORTABLE MODULES THAT ENABLE THE USE OF EXPERIMENTAL FEATURES  
OR ESTIMATORS

THE FEATURES AND ESTIMATORS THAT ARE EXPERIMENTAL AREN'T SUBJECT TO DEPRECATION CYCLES USE THEM AT YOUR OWN RISKS

EXPERIMENTALENABLEHISTGRADIENTBOOSTING ENABLES HISTOGRAMBASED GRADIENT BOOSTING ESTIMATORS

EXPERIMENTALENABLEITERATIVEIMPUTER ENABLES ITERATIVEIMPUTER

6141 SKLEARNEXPERIMENTALENABLEHISTGRADIENTBOOSTING

ENABLES HISTOGRAMBASED GRADIENT BOOSTING ESTIMATORS

THE API AND RESULTS OF THESE ESTIMATORS MIGHT CHANGE WITHOUT ANY DEPRECATION CYCLE

IMPORTING THIS FILE DYNAMICALLY SETS THE SKLEARNENSEMBLEHISTGRADIENTBOOSTINGCLASSIFIER AND  
SKLEARNENSEMBLEHISTGRADIENTBOOSTINGREGRESSOR AS ATTRIBUTES OF THE ENSEMBLE MODULE

EXPLICITLY REQUIRE THIS EXPERIMENTAL FEATURE

FROM SKLEARNEXPERIMENTAL IMPORT ENABLEHISTGRADIENTBOOSTING NOQA

NOW YOU CAN IMPORT NORMALLY FROM ENSEMBLE

614SKLEARNEXPERIMENTAL EXPERIMENTAL 1735

SCIKITLEARN USER GUIDE RELEASE 0213

FROM SKLEARNENSEMBLE IMPORT HISTGRADIENTBOOSTINGCLASSIFIER

FROM SKLEARNENSEMBLE IMPORT HISTGRADIENTBOOSTINGREGRESSOR

THE NOQA COMMENT COMMENT CAN BE REMOVED IT JUST TELLS LINTERS LIKE FLAKE8 TO IGNORE THE IMPORT WHICH APPEARS AS UNUSED

6142 SKLEARNEXPERIMENTALENABLEITERATIVEIMPUTER

ENABLES ITERATIVEIMPUTER

THE API AND RESULTS OF THIS ESTIMATOR MIGHT CHANGE WITHOUT ANY DEPRECATION CYCLE

IMPORTING THIS FILE DYNAMICALLY SETS SKLEARNIMPUTEITERATIVEIMPUTER AS AN ATTRIBUTE OF THE IMPUTE MODULE

EXPLICITLY REQUIRE THIS EXPERIMENTAL FEATURE

FROM SKLEARNEXPERIMENTAL IMPORT ENABLEITERATIVEIMPUTER NOQA

NOW YOU CAN IMPORT NORMALLY FROM IMPUTE

FROM SKLEARNIMPUTE IMPORT ITERATIVEIMPUTER

615SKLEARNFEATUREEXTRACTION FEATURE EXTRACTION

THESKLEARNFEATUREEXTRACTION MODULE DEALS WITH FEATURE EXTRACTION FROM RAW DATA IT CURRENTLY INCLUDES METHODS TO EXTRACT FEATURES FROM TEXT AND IMAGES

USER GUIDE SEE THE FEATURE EXTRACTION SECTION FOR FURTHER DETAILS

FEATUREEXTRACTIONDICTVECTORIZER DTYPE

TRANSFORMS LISTS OF FEATUREVALUE MAPPINGS TO VECTORS

FEATUREEXTRACTIONFEATUREHASHER IMPLEMENTS FEATURE HASHING AKA THE HASHING TRICK

6151SKLEARNFEATUREEXTRACTION DICTVECTORIZER

CLASSSSKLEARNFEATUREEXTRACTION DICTVECTORIZER DTYPECLASS 'NUMPYFLOAT64' SEPARATOR'' SPARSETRUE SORTTRUE

TRANSFORMS LISTS OF FEATUREVALUE MAPPINGS TO VECTORS

THIS TRANSFORMER TURNS LISTS OF MAPPINGS DICTLIKE OBJECTS OF FEATURE NAMES TO FEATURE VALUES INTO NUMPY ARRAYS OR SCIPYSPARSE MATRICES FOR USE WITH SCIKITLEARN ESTIMATORS

WHEN FEATURE VALUES ARE STRINGS THIS TRANSFORMER WILL DO A BINARY ONEHOT AKA ONEOFK CODING ONE BOOLEAN VALUED FEATURE IS CONSTRUCTED FOR EACH OF THE POSSIBLE STRING VALUES THAT THE FEATURE CAN TAKE ON FOR INSTANCE A FEATURE "F" THAT CAN TAKE ON THE VALUES "HAM" AND "SPAM" WILL BECOME TWO FEATURES IN THE OUTPUT ONE SIGNIFYING "FHAM" THE OTHER "FSPAM"

HOWEVER NOTE THAT THIS TRANSFORMER WILL ONLY DO A BINARY ONEHOT ENCODING WHEN FEATURE VALUES ARE OF TYPE STRING IF CATEGORICAL FEATURES ARE REPRESENTED AS NUMERIC VALUES SUCH AS INT THE DICTVECTORIZER CAN BE FOLLOWED BYSKLEARNPREPROCESSINGONEHOTENCODER TO COMPLETE BINARY ONEHOT ENCODING

FEATURES THAT DO NOT OCCUR IN A SAMPLE MAPPING WILL HAVE A ZERO VALUE IN THE RESULTING ARRAYMATRIX

READ MORE IN THE USER GUIDE

PARAMETERS

1736 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

DTTYPE CALLABLE OPTIONAL THE TYPE OF FEATURE VALUES PASSED TO NUMPY ARRAYS SCIPYSPARSE MATRICES CONSTRUCTORS AS THE DTTYPE ARGUMENT

SEPARATOR STRING OPTIONAL SEPARATOR STRING USED WHEN CONSTRUCTING NEW FEATURES FOR ONEHOT CODING

SPARSE BOOLEAN OPTIONAL WHETHER TRANSFORM SHOULD PRODUCE SCIPYSPARSE MATRICES TRUE BY DEFAULT

SORT BOOLEAN OPTIONAL WHETHER FEATURENAMES AND VOCABULARY SHOULD BE SORTED WHEN FITTING TRUE BY DEFAULT

ATTRIBUTES

VOCABULARY DICT A DICTIONARY MAPPING FEATURE NAMES TO FEATURE INDICES

FEATURENAMES LIST A LIST OF LENGTH NFEATRES CONTAINING THE FEATURE NAMES EG "FHAM" AND "FSPAM"

SEE ALSO

FEATUREHASHER PERFORMS VECTORIZATION USING ONLY A HASH FUNCTION

SKLEARNPREPROCESSINGORDINALENCODER HANDLES NOMINALCATEGORICAL FEATURES ENCODED AS COLUMNS OF ARBITRARY DATA TYPES

EXAMPLES

```
FROM SKLEARNFEATUREEXTRACTION IMPORT DICTVECTORIZER
V DICTVECTORIZERSPARSE FALSE
D FOO 1 BAR 2 FOO 3 BAZ 1
X VFITTRANSFORMD
X
ARRAY2 0 1
0 1 3
VINVERSETRANSFORMX BAR 20 FOO 10 BAZ 10 FOO
↪ 30
TRUE
VTRANSFORMFOO 4 UNSEENFEATURE 3
ARRAY0 0 4
```

METHODS

FITSELF X Y LEARN A LIST OF FEATURE NAME INDICES MAPPINGS

FITTRANSFORM SELF X Y LEARN A LIST OF FEATURE NAME INDICES MAPPINGS AND TRANSFORM X

GETFEATURENAMES SELF RETURNS A LIST OF FEATURE NAMES ORDERED BY THEIR INDICES

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

INVERSETRANSFORM SELF X DICTTYPE TRANSFORM ARRAY OR SPARSE MATRIX X BACK TO FEATURE MAPPINGS

RESTRICT SELF SUPPORT INDICES RESTRICT THE FEATURES TO THOSE IN SUPPORT USING FEATURE SELECTION

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

CONTINUED ON NEXT PAGE

615SKLEARNFEATUREEXTRACTION FEATURE EXTRACTION 1737

SCIKITLEARN USER GUIDE RELEASE 0213

TABLE 695 – CONTINUED FROM PREVIOUS PAGE

TRANSFORM SELF X TRANSFORM FEATUREVALUE DICTS TO ARRAY OR SPARSE MATRICES

INIT SELFDTYPECLASS ‘NUMPYFLOAT64’ SEPARATOR’’ SPARSETRUE SORTTRUE

FITSELFXYNONE

LEARN A LIST OF FEATURE NAME INDICES MAPPINGS

PARAMETERS

XMAPPING OR ITERABLE OVER MAPPINGS DICTS OR MAPPINGS FROM FEATURE NAMES ARBITRARY

PYTHON OBJECTS TO FEATURE VALUES STRINGS OR CONVERTIBLE TO DTYPE

YIGNORED

RETURNS

SELF

FITTRANSFORM SELFXYNONE

LEARN A LIST OF FEATURE NAME INDICES MAPPINGS AND TRANSFORM X

LIKE FITX FOLLOWED BY TRANSFORMX BUT DOES NOT REQUIRE MATERIALIZING X IN MEMORY

PARAMETERS

XMAPPING OR ITERABLE OVER MAPPINGS DICTS OR MAPPINGS FROM FEATURE NAMES ARBITRARY

PYTHON OBJECTS TO FEATURE VALUES STRINGS OR CONVERTIBLE TO DTYPE

YIGNORED

RETURNS

XAARRAY SPARSE MATRIX FEATURE VECTORS ALWAYS 2D

GETFEATURENAMES SELF

RETURNS A LIST OF FEATURE NAMES ORDERED BY THEIR INDICES

IF ONEOFK CODING IS APPLIED TO CATEGORICAL FEATURES THIS WILL INCLUDE THE CONSTRUCTED FEATURE NAMES BUT NOT THE ORIGINAL ONES

GETPARAMS SELFDEEPTTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

INVERSETRANSFORM SELFXDICTTYPECLASS ‘DICT’

TRANSFORM ARRAY OR SPARSE MATRIX X BACK TO FEATURE MAPPINGS

X MUST HAVE BEEN PRODUCED BY THIS DICTVECTORIZER’S TRANSFORM OR FITTRANSFORM METHOD IT MAY ONLY HAVE PASSED THROUGH TRANSFORMERS THAT PRESERVE THE NUMBER OF FEATURES AND THEIR ORDER

IN THE CASE OF ONEHOTONEOFK CODING THE CONSTRUCTED FEATURE NAMES AND VALUES ARE RETURNED RATHER THAN THE ORIGINAL ONES

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES SAMPLE MATRIX

1738 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

DICTTYPE CALLABLE OPTIONAL CONSTRUCTOR FOR FEATURE MAPPINGS MUST CONFORM TO THE COLLECTIONSMAPPING API

RETURNS

DLIST OF DICTTYPE OBJECTS LENGTH NSAMPLES FEATURE MAPPINGS FOR THE SAMPLES IN X

RESTRICT SELFSUPPORT INDICESFALSE

RESTRICT THE FEATURES TO THOSE IN SUPPORT USING FEATURE SELECTION

THIS FUNCTION MODIFIES THE ESTIMATOR INPLACE

PARAMETERS

SUPPORT ARRAYLIKE BOOLEAN MASK OR LIST OF INDICES AS RETURNED BY THE GETSUPPORT MEMBER OF FEATURE SELECTORS

INDICES BOOLEAN OPTIONAL WHETHER SUPPORT IS A LIST OF INDICES

RETURNS

SELF

EXAMPLES

```
FROM SKLEARNFEATUREEXTRACTION IMPORT DICTVECTORIZER
FROM SKLEARNFEATURESELECTION IMPORT SELECTKBEST CHI2
V DICTVECTORIZER
D FOO 1 BAR 2 FOO 3 BAZ 1
X VFITTRANSFORMD
SUPPORT SELECTKBESTCHI2 K2FITX 0 1
VGETFEATURENAMES
BAR BAZ FOO
VRESTRICTSUPPORTGETSUPPORT
```

DICTVECTORIZERDTYPE SEPARATOR SORTTRUE

SPARSETRUE

VGETFEATURENAMES

BAR FOO

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELF

TRANSFORM FEATUREVALUE DICTS TO ARRAY OR SPARSE MATRIX

NAMED FEATURES NOT ENCOUNTERED DURING FIT OR FITTRANSFORM WILL BE SILENTLY IGNORED

PARAMETERS

XMAPPING OR ITERABLE OVER MAPPINGS LENGTH NSAMPLES DICTS OR MAPPINGS FROM FEATURE NAMES ARBITRARY PYTHON OBJECTS TO FEATURE VALUES STRINGS OR CONVERTIBLE TO DTYPE

RETURNS

615SKLEARNFEATUREEXTRACTION FEATURE EXTRACTION 1739

SCIKITLEARN USER GUIDE RELEASE 0213

XAARRAY SPARSE MATRIX FEATURE VECTORS ALWAYS 2D

EXAMPLES USING SKLEARNFEATUREEXTRACTIONDICTVECTORIZER

- COLUMN TRANSFORMER WITH HETEROGENEOUS DATA SOURCES
- FEATUREHASHER AND DICTVECTORIZER COMPARISON

6152SKLEARNFEATUREEXTRACTION FEATUREHASHER

CLASSSSKLEARNFEATUREEXTRACTION FEATUREHASHER NFEATURES1048576 INPUTTYPE'DICT'

DTYPECLASS 'NUMPYFLOAT64' ALTER

NATESIGNTRUE

IMPLEMENTS FEATURE HASHING AKA THE HASHING TRICK

THIS CLASS TURNS SEQUENCES OF SYMBOLIC FEATURE NAMES STRINGS INTO SCIPYSPARSE MATRICES USING A HASH FUNCTION TO COMPUTE THE MATRIX COLUMN CORRESPONDING TO A NAME THE HASH FUNCTION EMPLOYED IS THE SIGNED 32BIT VERSION OF MURMURHASH3

FEATURE NAMES OF TYPE BYTE STRING ARE USED ASIS UNICODE STRINGS ARE CONVERTED TO UTF8 FIRST BUT NO UNICODE NORMALIZATION IS DONE FEATURE VALUES MUST BE FINITE NUMBERS

THIS CLASS IS A LOWMEMORY ALTERNATIVE TO DICTVECTORIZER AND COUNTVECTORIZER INTENDED FOR LARGESCALE ONLINE LEARNING AND SITUATIONS WHERE MEMORY IS TIGHT EG WHEN RUNNING PREDICTION CODE ON EMBEDDED DEVICES

READ MORE IN THE USER GUIDE

PARAMETERS

NFEATURES INTEGER OPTIONAL THE NUMBER OF FEATURES COLUMNS IN THE OUTPUT MATRICES SMALL NUMBERS OF FEATURES ARE LIKELY TO CAUSE HASH COLLISIONS BUT LARGE NUMBERS WILL CAUSE LARGER COEFFICIENT DIMENSIONS IN LINEAR LEARNERS

INPUTTYPE STRING OPTIONAL DEFAULT "DICT" EITHER "DICT" THE DEFAULT TO ACCEPT DICTIONARIES OVER FEATURENAME VALUE "PAIR" TO ACCEPT PAIRS OF FEATURENAME VALUE OR "STRING" TO ACCEPT SINGLE STRINGS FEATURENAME SHOULD BE A STRING WHILE VALUE SHOULD BE A NUMBER

IN THE CASE OF "STRING" A VALUE OF 1 IS IMPLIED THE FEATURENAME IS HASHED TO FIND THE APPROPRIATE COLUMN FOR THE FEATURE THE VALUE'S SIGN MIGHT BE FLIPPED IN THE OUTPUT BUT SEE NONNEGATIVE BELOW

DTYPE NUMPY TYPE OPTIONAL DEFAULT NPFLOAT64 THE TYPE OF FEATURE VALUES PASSED TO SCIPYSPARSE MATRIX CONSTRUCTORS AS THE DTYPE ARGUMENT DO NOT SET THIS TO BOOL NPBOOLEAN OR ANY UNSIGNED INTEGER TYPE

ALTERNATESIGN BOOLEAN OPTIONAL DEFAULT TRUE WHEN TRUE AN ALTERNATING SIGN IS ADDED TO THE FEATURES AS TO APPROXIMATELY CONSERVE THE INNER PRODUCT IN THE HASHED SPACE EVEN FOR SMALL NFEATURES THIS APPROACH IS SIMILAR TO SPARSE RANDOM PROJECTION

SEE ALSO

DICTVECTORIZER VECTORIZES STRINGVALUED FEATURES USING A HASH TABLE

SKLEARNPREPROCESSINGONEHOTENCODER HANDLES NOMINALCATEGORICAL FEATURES

1740 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES

FROM SKLEARNFEATUREEXTRACTION IMPORT FEATUREHASHER

H FEATUREHASHERNFEATURES10

D DOG 1 CAT2 ELEPHANT4DOG 2 RUN 5

F HTRANSFORMD

FTOARRAY

ARRAY 0 0 4 1 0 0 0 0 2

0 0 0 2 5 0 0 0 0

METHODS

FITSELF X Y NOOP

FITTRANSFORM SELF X Y FIT TO DATA THEN TRANSFORM IT

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF RAWX TRANSFORM A SEQUENCE OF INSTANCES TO A SCIPYSPARSE MATR

INIT SELFNFEATURES1048576 INPUTTYPE'DICT' DTYPECLASS 'NUMPYFLOAT64' ALTER

NATESIGNTRUE

FITSELFXNONE YNONE

NOOP

THIS METHOD DOESN'T DO ANYTHING IT EXISTS PURELY FOR COMPATIBILITY WITH THE SCIKITLEARN TRANSFORMER API

PARAMETERS

XARRAYLIKE

RETURNS

SELF FEATUREHASHER

FITTRANSFORM SELFXYNONE FITPARAMS

FIT TO DATA THEN TRANSFORM IT

FITS TRANSFORMER TO X AND Y WITH OPTIONAL PARAMETERS FITPARAMS AND RETURNS A TRANSFORMED VERSION OF X

PARAMETERS

XNUMPY ARRAY OF SHAPE NSAMPLES NFEATURES TRAINING SET

YNUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES

RETURNS

XNEW NUMPY ARRAY OF SHAPE NSAMPLES NFEATURESNEW TRANSFORMED ARRAY

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

615SKLEARNFEATUREEXTRACTION FEATURE EXTRACTION 1741

SCIKITLEARN USER GUIDE RELEASE 0213

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELFRAWX

TRANSFORM A SEQUENCE OF INSTANCES TO A SCIPYSPARSE MATRIX

PARAMETERS

RAWX ITERABLE OVER ITERABLE OVER RAW FEATURES LENGTH NSAMPLES SAMPLES EACH SAMPLE MUST BE ITERABLE AN EG A LIST OR TUPLE CONTAININGGENERATING FEATURE NAMES AND OPTION ALLY VALUES SEE THE INPUTTYPE CONSTRUCTOR ARGUMENT WHICH WILL BE HASHED RAWX NEED NOT SUPPORT THE LEN FUNCTION SO IT CAN BE THE RESULT OF A GENERATOR NSAMPLES IS DETERMINED ON THE FLY

RETURNS

XSCIPYSPARSE MATRIX SHAPE NSAMPLES SELFNFEATURES FEATURE MATRIX FOR USE WITH ESTIMATORS OR FURTHER TRANSFORMERS

EXAMPLES USING SKLEARNFEATUREEXTRACTIONFEATUREHASHER

- FEATUREHASHER AND DICTVECTORIZER COMPARISON

6153 FROM IMAGES

THESKLEARNFEATUREEXTRACTIONIMAGE SUBMODULE GATHERS UTILITIES TO EXTRACT FEATURES FROM IMAGES

FEATUREEXTRACTIONIMAGE

EXTRACTPATCHES2D RESHAPE A 2D IMAGE INTO A COLLECTION OF PATCHES

FEATUREEXTRACTIONIMAGE

GRIDTOGRAPH NX NYGRAPH OF THE PIXELTOPIXEL CONNECTIONS

FEATUREEXTRACTIONIMAGE

IMGTOGRAPH IMG GRAPH OF THE PIXELTOPIXEL GRADIENT CONNECTIONS

FEATUREEXTRACTIONIMAGE

RECONSTRUCTFROMPATCHES2D RECONSTRUCT THE IMAGE FROM ALL OF ITS PATCHES

FEATUREEXTRACTIONIMAGE

PATCHEXTRACTOR EXTRACTS PATCHES FROM A COLLECTION OF IMAGES

SKLEARNFEATUREEXTRACTIONIMAGE EXTRACTPATCHES2D

SKLEARNFEATUREEXTRACTIONIMAGE EXTRACTPATCHES2D IMAGE PATCHSIZE

MAXPATCHESNONE RAN

DOMSTATENONE

RESHAPE A 2D IMAGE INTO A COLLECTION OF PATCHES

THE RESULTING PATCHES ARE ALLOCATED IN A DEDICATED ARRAY

1742 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

READ MORE IN THE USER GUIDE

PARAMETERS

IMAGE ARRAY SHAPE IMAGEHEIGHT IMAGEWIDTH OR IMAGEHEIGHT IMAGEWIDTH

NCHANNELS THE ORIGINAL IMAGE DATA FOR COLOR IMAGES THE LAST DIMENSION SPECIFIES THE CHANNEL A RGB IMAGE WOULD HAVE NCHANNELS3

PATCHSIZE TUPLE OF INTS PATCHHEIGHT PATCHWIDTH THE DIMENSIONS OF ONE PATCH

MAXPATCHES INTEGER OR FLOAT OPTIONAL DEFAULT IS NONE THE MAXIMUM NUMBER OF PATCHES TO EXTRACT IF MAXPATCHES IS A FLOAT BETWEEN 0 AND 1 IT IS TAKEN TO BE A PROPORTION OF THE TOTAL NUMBER OF PATCHES

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE PSEUDO NUMBER

GENERATOR STATE USED FOR RANDOM SAMPLING TO USE IF MAXPATCHES IS NOT NONE IF INT

RANDOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE

RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NPRANDOM

RETURNS

PATCHES ARRAY SHAPE NPATCHES PATCHHEIGHT PATCHWIDTH OR NPATCHES PATCHHEIGHT

PATCHWIDTH NCHANNELS THE COLLECTION OF PATCHES EXTRACTED FROM THE IMAGE WHERE

NPATCHES IS EITHERMAXPATCHES OR THE TOTAL NUMBER OF PATCHES THAT CAN BE EXTRACTED

EXAMPLES

FROM SKLEARNDATASETS IMPORT LOADSAMPLEIMAGE

FROM SKLEARNFEATUREEXTRACTION IMPORT IMAGE

USE THE ARRAY DATA FROM THE FIRST IMAGE IN THIS DATASET

ONEIMAGE LOADSAMPLEIMAGECHINAJPG

PRINTIMAGE SHAPE FORMATONEIMAGESHAPE

IMAGE SHAPE 427 640 3

PATCHES IMAGEEXTRACTPATCHES2DONEIMAGE 2 2

PRINTPATCHES SHAPE FORMATPATCHESSHAPE

PATCHES SHAPE 272214 2 2 3

HERE ARE JUST TWO OF THESE PATCHES

PRINTPATCHES1

174 201 231

174 201 231

173 200 230

173 200 230

PRINTPATCHES800

187 214 243

188 215 244

187 214 243

188 215 244

EXAMPLES USING SKLEARNFEATUREEXTRACTIONIMAGEEXTRACTPATCHES2D

•ONLINE LEARNING OF A DICTIONARY OF PARTS OF FACES

•IMAGE DENOISING USING DICTIONARY LEARNING

615SKLEARNFEATUREEXTRACTION FEATURE EXTRACTION 1743

SCIKITLEARN USER GUIDE RELEASE 0213  
SKLEARNFEATUREEXTRACTIONIMAGE GRIDTOGRAPH  
SKLEARNFEATUREEXTRACTIONIMAGE GRIDTOGRAPH NX NY NZ1  
MASKNONE RETURNASCLASS  
'SCIPYSPARSECOOCOOMATRIX'  
DTYPECLASS 'INT'  
GRAPH OF THE PIXELTOPIXEL CONNECTIONS  
EDGES EXIST IF 2 VOXELS ARE CONNECTED  
PARAMETERS  
NX INT DIMENSION IN X AXIS  
NY INT DIMENSION IN Y AXIS  
NZ INT OPTIONAL DEFAULT 1 DIMENSION IN Z AXIS  
MASK NDARRAY OF BOOLEANS OPTIONAL AN OPTIONAL MASK OF THE IMAGE TO CONSIDER ONLY PART OF  
THE PIXELS  
RETURNAS NPNDARRAY OR A SPARSE MATRIX CLASS OPTIONAL THE CLASS TO USE TO BUILD THE RETURNED  
ADJACENCY MATRIX  
DTYPE DTYPE OPTIONAL DEFAULT INT THE DATA OF THE RETURNED SPARSE MATRIX BY DEFAULT IT IS INT  
NOTES  
FOR SCIKITLEARN VERSIONS 0141 AND PRIOR RETURNASNPNDARRAY WAS HANDLED BY RETURNING A DENSE NPMATRIX  
INSTANCE GOING FORWARD NPNDARRAY RETURNS AN NPNDARRAY AS EXPECTED  
FOR COMPATIBILITY USER CODE RELYING ON THIS METHOD SHOULD WRAP ITS CALLS IN NPASARRAY TO AVOID TYPE ISSUES  
SKLEARNFEATUREEXTRACTIONIMAGE IMGTOGRAPH  
SKLEARNFEATUREEXTRACTIONIMAGE IMGTOGRAPH IMGMASKNONE RETURNASCLASS  
'SCIPYSPARSECOOCOOMATRIX'  
DTYPENONE  
GRAPH OF THE PIXELTOPIXEL GRADIENT CONNECTIONS  
EDGES ARE WEIGHTED WITH THE GRADIENT VALUES  
READ MORE IN THE USER GUIDE  
PARAMETERS  
IMG NDARRAY 2D OR 3D 2D OR 3D IMAGE  
MASK NDARRAY OF BOOLEANS OPTIONAL AN OPTIONAL MASK OF THE IMAGE TO CONSIDER ONLY PART OF  
THE PIXELS  
RETURNAS NPNDARRAY OR A SPARSE MATRIX CLASS OPTIONAL THE CLASS TO USE TO BUILD THE RETURNED  
ADJACENCY MATRIX  
DTYPE NONE OR DTYPE OPTIONAL THE DATA OF THE RETURNED SPARSE MATRIX BY DEFAULT IT IS THE  
DTYPE OF IMG  
1744 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

NOTES

FOR SCIKITLEARN VERSIONS 0141 AND PRIOR RETURNASNPNDARRAY WAS HANDLED BY RETURNING A DENSE NPMATRIX  
INSTANCE GOING FORWARD NPNDARRAY RETURNS AN NPNDARRAY AS EXPECTED  
FOR COMPATIBILITY USER CODE RELYING ON THIS METHOD SHOULD WRAP ITS CALLS IN NPASARRAY TO AVOID TYPE ISSUES  
SKLEARNFEATUREEXTRACTIONIMAGE RECONSTRUCTFROMPATCHES2D  
SKLEARNFEATUREEXTRACTIONIMAGE RECONSTRUCTFROMPATCHES2D PATCHES IM

AGESIZE

RECONSTRUCT THE IMAGE FROM ALL OF ITS PATCHES

PATCHES ARE ASSUMED TO OVERLAP AND THE IMAGE IS CONSTRUCTED BY FILLING IN THE PATCHES FROM LEFT TO RIGHT TOP TO  
BOTTOM AVERAGING THE OVERLAPPING REGIONS

READ MORE IN THE USER GUIDE

PARAMETERS

PATCHES ARRAY SHAPE NPATCHES PATCHHEIGHT PATCHWIDTH OR NPATCHES PATCHHEIGHT  
PATCHWIDTH NCHANNELS THE COMPLETE SET OF PATCHES IF THE PATCHES CONTAIN COLOUR  
INFORMATION CHANNELS ARE INDEXED ALONG THE LAST DIMENSION RGB PATCHES WOULD HAVE  
NCHANNELS3

IMAGESIZE TUPLE OF INTS IMAGEHEIGHT IMAGEWIDTH OR IMAGEHEIGHT IMAGEWIDTH  
NCHANNELS THE SIZE OF THE IMAGE THAT WILL BE RECONSTRUCTED

RETURNS

IMAGE ARRAY SHAPE IMAGESIZE THE RECONSTRUCTED IMAGE

EXAMPLES USING SKLEARNFEATUREEXTRACTIONIMAGERECONSTRUCTFROMPATCHES2D

•IMAGE DENOISING USING DICTIONARY LEARNING

SKLEARNFEATUREEXTRACTIONIMAGE PATCHEXTRACTOR

CLASSSSKLEARNFEATUREEXTRACTIONIMAGE PATCHEXTRACTOR PATCHSIZENONE

MAXPATCHESNONE RAN

DOMSTATENONE

EXTRACTS PATCHES FROM A COLLECTION OF IMAGES

READ MORE IN THE USER GUIDE

PARAMETERS

PATCHSIZE TUPLE OF INTS PATCHHEIGHT PATCHWIDTH THE DIMENSIONS OF ONE PATCH

MAXPATCHES INTEGER OR FLOAT OPTIONAL DEFAULT IS NONE THE MAXIMUM NUMBER OF PATCHES PER

IMAGE TO EXTRACT IF MAXPATCHES IS A FLOAT IN 0 1 IT IS TAKEN TO MEAN A PROPORTION OF THE

TOTAL NUMBER OF PATCHES

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RAN

DOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE

RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS

THE RANDOMSTATE INSTANCE USED BY NPRANDOM

615SKLEARNFEATUREEXTRACTION FEATURE EXTRACTION 1745

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES

FROM SKLEARNDATASETS IMPORT LOADSAMPLEIMAGES

FROM SKLEARNFEATUREEXTRACTION IMPORT IMAGE

USE THE ARRAY DATA FROM THE SECOND IMAGE IN THIS DATASET

X LOADSAMPLEIMAGESIMAGES1

PRINTIMAGE SHAPE FORMATXSHAPE

IMAGE SHAPE 427 640 3

PE IMAGEPATCHEXTRACTORPATCHSIZE2 2

PEFIT PEFITX

PETRANS PETRANSFORMX

PRINTPATCHES SHAPE FORMATPETRANSSHAPE

PATCHES SHAPE 545706 2 2

METHODS

FITSELF X Y DO NOTHING AND RETURN THE ESTIMATOR UNCHANGED

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF X TRANSFORMS THE IMAGE SAMPLES IN X INTO A MATRIX OF PATCH DATA

INIT SELFPATCHSIZENONE MAXPATCHESNONE RANDOMSTATENONE

FITSELFXYNONE

DO NOTHING AND RETURN THE ESTIMATOR UNCHANGED

THIS METHOD IS JUST THERE TO IMPLEMENT THE USUAL API AND HENCE WORK IN PIPELINES

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA

GETPARAMS SELFDEEPTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE

PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT

OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELFX

TRANSFORMS THE IMAGE SAMPLES IN X INTO A MATRIX OF PATCH DATA

1746 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

PARAMETERS

XARRAY SHAPE NSAMPLES IMAGEHEIGHT IMAGEWIDTH OR NSAMPLES IMAGEHEIGHT  
IMAGEWIDTH NCHANNELS ARRAY OF IMAGES FROM WHICH TO EXTRACT PATCHES FOR  
COLOR IMAGES THE LAST DIMENSION SPECIFIES THE CHANNEL A RGB IMAGE WOULD HAVE  
NCHANNELS3

RETURNS

PATCHES ARRAY SHAPE NPATCHES PATCHHEIGHT PATCHWIDTH OR NPATCHES  
PATCHHEIGHT PATCHWIDTH NCHANNELS THE COLLECTION OF PATCHES EXTRACTED FROM THE IM  
AGES WHERE NPATCHES IS EITHERNSAMPLES MAXPATCHES OR THE TOTAL NUMBER  
OF PATCHES THAT CAN BE EXTRACTED

6154 FROM TEXT

THESKLEARNFEATUREEXTRACTIONTEXT SUBMODULE GATHERS UTILITIES TO BUILD FEATURE VECTORS FROM TEXT DOC  
UMENTS

FEATUREEXTRACTIONTEXT

COUNTVECTORIZER CONVERT A COLLECTION OF TEXT DOCUMENTS TO A MATRIX OF TOKEN  
COUNTS

FEATUREEXTRACTIONTEXT

HASHINGVECTORIZER CONVERT A COLLECTION OF TEXT DOCUMENTS TO A MATRIX OF TOKEN  
OCCURRENCES

FEATUREEXTRACTIONTEXT

TFIDFTRANSFORMER TRANSFORM A COUNT MATRIX TO A NORMALIZED TF OR TFIDF REPRE  
SENTATION

FEATUREEXTRACTIONTEXT

TFIDFVECTORIZER CONVERT A COLLECTION OF RAW DOCUMENTS TO A MATRIX OF TF  
IDF FEATURES

SKLEARNFEATUREEXTRACTIONTEXT COUNTVECTORIZER

CLASSSSKLEARNFEATUREEXTRACTIONTEXT COUNTVECTORIZER INPUT'CONTENT' ENCODING'UTF  
8' DECODEERROR'STRICT'

STRIPACCENTSNONE LOW

ERCASETRUE PREPROCES

SORNONE TOKENIZERNONE

STOPWORDSNONE TO

KENPATTERN'UBWWB'

NGRAMRANGE1 1 ANA

LYZER'WORD' MAXDF10

MINDF1 MAXFEATURESNONE

VOCABULARYNONE BINARYFALSE

DTYPECLASS 'NUMPYINT64'

CONVERT A COLLECTION OF TEXT DOCUMENTS TO A MATRIX OF TOKEN COUNTS

THIS IMPLEMENTATION PRODUCES A SPARSE REPRESENTATION OF THE COUNTS USING SCIPYSPARSECSRMATRIX  
IF YOU DO NOT PROVIDE AN APRIORI DICTIONARY AND YOU DO NOT USE AN ANALYZER THAT DOES SOME KIND OF FEATURE  
SELECTION THEN THE NUMBER OF FEATURES WILL BE EQUAL TO THE VOCABULARY SIZE FOUND BY ANALYZING THE DATA  
READ MORE IN THE USER GUIDE

PARAMETERS

INPUT STRING 'FILENAME' 'FILE' 'CONTENT' IF 'FILENAME' THE SEQUENCE PASSED AS AN ARGUMENT TO  
FIT IS EXPECTED TO BE A LIST OF FILENAMES THAT NEED READING TO FETCH THE RAW CONTENT TO ANALYZE

615SKLEARNFEATUREEXTRACTION FEATURE EXTRACTION 1747

SCIKITLEARN USER GUIDE RELEASE 0213

IF 'FILE' THE SEQUENCE ITEMS MUST HAVE A 'READ' METHOD FILELIKE OBJECT THAT IS CALLED TO FETCH THE BYTES IN MEMORY OTHERWISE THE INPUT IS EXPECTED TO BE THE SEQUENCE STRINGS OR BYTES ITEMS ARE EXPECTED TO BE ANALYZED DIRECTLY ENCODING STRING 'UTF8' BY DEFAULT IF BYTES OR FILES ARE GIVEN TO ANALYZE THIS ENCODING IS USED TO DECODE

DECODEERROR 'STRICT' 'IGNORE' 'REPLACE' INSTRUCTION ON WHAT TO DO IF A BYTE SEQUENCE IS GIVEN TO ANALYZE THAT CONTAINS CHARACTERS NOT OF THE GIVEN ENCODING BY DEFAULT IT IS 'STRICT' MEANING THAT A UNICODEDECODEERROR WILL BE RAISED OTHER VALUES ARE 'IGNORE' AND 'REPLACE'

STRIPACCENTS 'ASCII' 'UNICODE' NONE REMOVE ACCENTS AND PERFORM OTHER CHARACTER NOR MALIZATION DURING THE PREPROCESSING STEP 'ASCII' IS A FAST METHOD THAT ONLY WORKS ON CHAR ACTERS THAT HAVE AN DIRECT ASCII MAPPING 'UNICODE' IS A SLIGHTLY SLOWER METHOD THAT WORKS ON ANY CHARACTERS NONE DEFAULT DOES NOTHING

BOTH 'ASCII' AND 'UNICODE' USE NFKD NORMALIZATION FROM UNICODEDATANORMALIZE LOWERCASE BOOLEAN TRUE BY DEFAULT CONVERT ALL CHARACTERS TO LOWERCASE BEFORE TOKENIZING PREPROCESSOR CALLABLE OR NONE DEFAULT OVERRIDE THE PREPROCESSING STRING TRANSFORMATION STAGE WHILE PRESERVING THE TOKENIZING AND NGRAMS GENERATION STEPS

TOKENIZER CALLABLE OR NONE DEFAULT OVERRIDE THE STRING TOKENIZATION STEP WHILE PRESERVING THE PREPROCESSING AND NGRAMS GENERATION STEPS ONLY APPLIES IF ANALYZER WORD

STOPWORDS STRING 'ENGLISH' LIST OR NONE DEFAULT IF 'ENGLISH' A BUILTIN STOP WORD LIST FOR ENGLISH IS USED THERE ARE SEVERAL KNOWN ISSUES WITH 'ENGLISH' AND YOU SHOULD CONSIDER AN ALTERNATIVE SEE USING STOP WORDS

IF A LIST THAT LIST IS ASSUMED TO CONTAIN STOP WORDS ALL OF WHICH WILL BE REMOVED FROM THE RESULTING TOKENS ONLY APPLIES IF ANALYZER WORD

IF NONE NO STOP WORDS WILL BE USED MAXDF CAN BE SET TO A VALUE IN THE RANGE 07 10 TO AUTOMATICALLY DETECT AND FILTER STOP WORDS BASED ON INTRA CORPUS DOCUMENT FREQUENCY OF TERMS

TOKENPATTERN STRING REGULAR EXPRESSION DENOTING WHAT CONSTITUTES A "TOKEN" ONLY USED IF ANALYZER WORD THE DEFAULT REGEXP SELECT TOKENS OF 2 OR MORE ALPHANUMERIC

CHARACTERS PUNCTUATION IS COMPLETELY IGNORED AND ALWAYS TREATED AS A TOKEN SEPARATOR NGRAMRANGE TUPLE MINN MAXN THE LOWER AND UPPER BOUNDARY OF THE RANGE OF NVALUES

FOR DIFFERENT NGRAMS TO BE EXTRACTED ALL VALUES OF N SUCH THAT MINN N MAXN WILL BE USED

ANALYZER STRING 'WORD' 'CHAR' 'CHARWB' OR CALLABLE WHETHER THE FEATURE SHOULD BE MADE OF WORD OR CHARACTER NGRAMS OPTION 'CHARWB' CREATES CHARACTER NGRAMS ONLY FROM TEXT INSIDE WORD BOUNDARIES NGRAMS AT THE EDGES OF WORDS ARE PADDED WITH SPACE

IF A CALLABLE IS PASSED IT IS USED TO EXTRACT THE SEQUENCE OF FEATURES OUT OF THE RAW UNPROCESSED INPUT

CHANGED IN VERSION 021

SINCE V021 IF INPUT ISFILENAME ORFILE THE DATA IS FIRST READ FROM THE FILE AND THEN PASSED TO THE GIVEN CALLABLE ANALYZER

MAXDF FLOAT IN RANGE 00 10 OR INT DEFAULT10 WHEN BUILDING THE VOCABULARY IGNORE TERMS THAT HAVE A DOCUMENT FREQUENCY STRICTLY HIGHER THAN THE GIVEN THRESHOLD CORPUSSPECIFIC STOP

1748 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

WORDS IF FLOAT THE PARAMETER REPRESENTS A PROPORTION OF DOCUMENTS INTEGER ABSOLUTE COUNTS

THIS PARAMETER IS IGNORED IF VOCABULARY IS NOT NONE

MINDF FLOAT IN RANGE 00 10 OR INT DEFAULT1 WHEN BUILDING THE VOCABULARY IGNORE TERMS

THAT HAVE A DOCUMENT FREQUENCY STRICTLY LOWER THAN THE GIVEN THRESHOLD THIS VALUE IS ALSO

CALLED CUTOFF IN THE LITERATURE IF FLOAT THE PARAMETER REPRESENTS A PROPORTION OF DOCUMENTS

INTEGER ABSOLUTE COUNTS THIS PARAMETER IS IGNORED IF VOCABULARY IS NOT NONE

MAXFEATURES INT OR NONE DEFAULTNONE IF NOT NONE BUILD A VOCABULARY THAT ONLY CONSIDER

THE TOP MAXFEATURES ORDERED BY TERM FREQUENCY ACROSS THE CORPUS

THIS PARAMETER IS IGNORED IF VOCABULARY IS NOT NONE

VOCABULARY MAPPING OR ITERABLE OPTIONAL EITHER A MAPPING EG A DICT WHERE KEYS ARE TERMS

AND VALUES ARE INDICES IN THE FEATURE MATRIX OR AN ITERABLE OVER TERMS IF NOT GIVEN A VOCABU

LARY IS DETERMINED FROM THE INPUT DOCUMENTS INDICES IN THE MAPPING SHOULD NOT BE REPEATED

AND SHOULD NOT HAVE ANY GAP BETWEEN 0 AND THE LARGEST INDEX

BINARY BOOLEAN DEFAULTFALSE IF TRUE ALL NON ZERO COUNTS ARE SET TO 1 THIS IS USEFUL FOR

DISCRETE PROBABILISTIC MODELS THAT MODEL BINARY EVENTS RATHER THAN INTEGER COUNTS

DTYPE TYPE OPTIONAL TYPE OF THE MATRIX RETURNED BY FITTRANSFORM OR TRANSFORM

ATTRIBUTES

VOCABULARY DICT A MAPPING OF TERMS TO FEATURE INDICES

STOPWORDS SET TERMS THAT WERE IGNORED BECAUSE THEY EITHER

- OCCURRED IN TOO MANY DOCUMENTS MAXDF

- OCCURRED IN TOO FEW DOCUMENTS MINDF

- WERE CUT OFF BY FEATURE SELECTION MAXFEATURES

THIS IS ONLY AVAILABLE IF NO VOCABULARY WAS GIVEN

SEE ALSO

HASHINGVECTORIZER TFIDFVECTORIZER

NOTES

THESSTOPWORDS ATTRIBUTE CAN GET LARGE AND INCREASE THE MODEL SIZE WHEN PICKLING THIS ATTRIBUTE IS PROVIDED

ONLY FOR INTROSPECTION AND CAN BE SAFELY REMOVED USING DELATTR OR SET TO NONE BEFORE PICKLING

EXAMPLES

FROM SKLEARNFEATUREEXTRACTIONTEXT IMPORT COUNTVECTORIZER

CORPUS

THIS IS THE FIRST DOCUMENT

THIS DOCUMENT IS THE SECOND DOCUMENT

AND THIS IS THE THIRD ONE

IS THIS THE FIRST DOCUMENT

VECTORIZER COUNTVECTORIZER

X VECTORIZERFITTRANSFORMCORPUS

PRINTVECTORIZERGETFEATURENAMES

AND DOCUMENT FIRST IS ONE SECOND THE THIRD THIS

615SKLEARNFEATUREEXTRACTION FEATURE EXTRACTION 1749

SCIKITLEARN USER GUIDE RELEASE 0213

PRINTXTOARRAY

0 1 1 1 0 0 1 0 1

0 2 0 1 0 1 1 0 1

1 0 0 1 1 0 1 1 1

0 1 1 1 0 0 1 0 1

METHODS

BUILDANALYZER SELF RETURN A CALLABLE THAT HANDLES PREPROCESSING AND TOK

ENIZATION

BUILDPREPROCESSOR SELF RETURN A FUNCTION TO PREPROCESS THE TEXT BEFORE TOKENIZA

TION

BUILDTOKENIZER SELF RETURN A FUNCTION THAT SPLITS A STRING INTO A SEQUENCE OF

TOKENS

DECODE SELF DOC DECODE THE INPUT INTO A STRING OF UNICODE SYMBOLS

FITSELF RAWDOCUMENTS Y LEARN A VOCABULARY DICTIONARY OF ALL TOKENS IN THE RAW

DOCUMENTS

FITTRANSFORM SELF RAWDOCUMENTS Y LEARN THE VOCABULARY DICTIONARY AND RETURN TERM

DOCUMENT MATRIX

GETFEATURENAMES SELF ARRAY MAPPING FROM FEATURE INTEGER INDICES TO FEATURE

NAME

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

GETSTOPWORDS SELF BUILD OR FETCH THE EFFECTIVE STOP WORDS LIST

INVERSETRANSFORM SELF X RETURN TERMS PER DOCUMENT WITH NONZERO ENTRIES IN X

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF RAWDOCUMENTS TRANSFORM DOCUMENTS TO DOCUMENTTERM MATRIX

INIT SELFINPUT'CONTENT' ENCODING'UTF8' DECODEERROR'STRICT' STRIPACCENTSNONE

LOWERCASETRUE PREPROCESSORNONE TOKENIZERNONE STOPWORDSNONE TO

KENPATTERN'UBWWB' NGRAMRANGE1 1ANALYZER'WORD' MAXDF10

MINDF1 MAXFEATURESNONE VOCABULARYNONE BINARYFALSE DTYPECLASS

'NUMPYINT64'

BUILDANALYZER SELF

RETURN A CALLABLE THAT HANDLES PREPROCESSING AND TOKENIZATION

BUILDPREPROCESSOR SELF

RETURN A FUNCTION TO PREPROCESS THE TEXT BEFORE TOKENIZATION

BUILDTOKENIZER SELF

RETURN A FUNCTION THAT SPLITS A STRING INTO A SEQUENCE OF TOKENS

DECODESELFDOC

DECODE THE INPUT INTO A STRING OF UNICODE SYMBOLS

THE DECODING STRATEGY DEPENDS ON THE VECTORIZER PARAMETERS

PARAMETERS

DOC STRING THE STRING TO DECODE

FITSELFRAWDOCUMENTS YNONE

LEARN A VOCABULARY DICTIONARY OF ALL TOKENS IN THE RAW DOCUMENTS

PARAMETERS

1750 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

RAWDOCUMENTS ITERABLE AN ITERABLE WHICH YIELDS EITHER STR UNICODE OR FILE OBJECTS

RETURNS

SELF

FITTRANSFORM SELFRAWDOCUMENTS YNONE

LEARN THE VOCABULARY DICTIONARY AND RETURN TERMDOCUMENT MATRIX

THIS IS EQUIVALENT TO FIT FOLLOWED BY TRANSFORM BUT MORE EFFICIENTLY IMPLEMENTED

PARAMETERS

RAWDOCUMENTS ITERABLE AN ITERABLE WHICH YIELDS EITHER STR UNICODE OR FILE OBJECTS

RETURNS

XARRAY NSAMPLES NFEATURES DOCUMENTTERM MATRIX

GETFEATURENAMES SELF

ARRAY MAPPING FROM FEATURE INTEGER INDICES TO FEATURE NAME

GETPARAMS SELFDEEPTTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

GETSTOPWORDS SELF

BUILD OR FETCH THE EFFECTIVE STOP WORDS LIST

INVERSETRANSFORM SELF

RETURN TERMS PER DOCUMENT WITH NONZERO ENTRIES IN X

PARAMETERS

XARRAY SPARSE MATRIX SHAPE NSAMPLES NFEATURES

RETURNS

XINV LIST OF ARRAYS LEN NSAMPLES LIST OF ARRAYS OF TERMS

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE

PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT

OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELFRAWDOCUMENTS

TRANSFORM DOCUMENTS TO DOCUMENTTERM MATRIX

EXTRACT TOKEN COUNTS OUT OF RAW TEXT DOCUMENTS USING THE VOCABULARY FITTED WITH FIT OR THE ONE PROVIDED TO

THE CONSTRUCTOR

PARAMETERS

RAWDOCUMENTS ITERABLE AN ITERABLE WHICH YIELDS EITHER STR UNICODE OR FILE OBJECTS

615SKLEARNFEATUREEXTRACTION FEATURE EXTRACTION 1751

SCIKITLEARN USER GUIDE RELEASE 0213

RETURNS

XSPARSE MATRIX NSAMPLES NFEATURES DOCUMENTTERM MATRIX

EXAMPLES USING SKLEARNFEATUREEXTRACTIONTEXTCOUNTVECTORIZER

- TOPIC EXTRACTION WITH NONNEGATIVE MATRIX FACTORIZATION AND LATENT DIRICHLET ALLOCATION
- SAMPLE PIPELINE FOR TEXT FEATURE EXTRACTION AND EVALUATION

SKLEARNFEATUREEXTRACTIONTEXT HASHINGVECTORIZER

CLASSSSKLEARNFEATUREEXTRACTIONTEXT HASHINGVECTORIZER INPUT'CONTENT'

ENCODING'UTF8' DE

CODEERROR'STRICT'

STRIPACCENTSNONE LOW

ERCASETTRUE PREPROCES

SORNONE TOKENIZERNONE

STOPWORDSNONE TO

KENPATTERN'UBWWB'

NGRAMRANGE1

1 ANALYZER'WORD'

NFEATURES1048576 BI

NARYFALSE NORM'L2'

ALTERNATESIGNTRUE

DTYPECLASS

'NUMPYFLOAT64'

CONVERT A COLLECTION OF TEXT DOCUMENTS TO A MATRIX OF TOKEN OCCURRENCES

IT TURNS A COLLECTION OF TEXT DOCUMENTS INTO A SCIPYSPARSE MATRIX HOLDING TOKEN OCCURRENCE COUNTS OR BINARY OCCURRENCE INFORMATION POSSIBLY NORMALIZED AS TOKEN FREQUENCIES IF NORM'L1' OR PROJECTED ON THE EUCLIDEAN UNIT SPHERE IF NORM'L2'

THIS TEXT VECTORIZER IMPLEMENTATION USES THE HASHING TRICK TO FIND THE TOKEN STRING NAME TO FEATURE INTEGER INDEX MAPPING

THIS STRATEGY HAS SEVERAL ADVANTAGES

- IT IS VERY LOW MEMORY SCALABLE TO LARGE DATASETS AS THERE IS NO NEED TO STORE A VOCABULARY DICTIONARY IN MEMORY
- IT IS FAST TO PICKLE AND UNPICKLE AS IT HOLDS NO STATE BESIDES THE CONSTRUCTOR PARAMETERS
- IT CAN BE USED IN A STREAMING PARTIAL FIT OR PARALLEL PIPELINE AS THERE IS NO STATE COMPUTED DURING FIT

THERE ARE ALSO A COUPLE OF CONS VS USING A COUNTVECTORIZER WITH AN INMEMORY VOCABULARY

- THERE IS NO WAY TO COMPUTE THE INVERSE TRANSFORM FROM FEATURE INDICES TO STRING FEATURE NAMES WHICH CAN BE A PROBLEM WHEN TRYING TO INTROSPECT WHICH FEATURES ARE MOST IMPORTANT TO A MODEL
- THERE CAN BE COLLISIONS DISTINCT TOKENS CAN BE MAPPED TO THE SAME FEATURE INDEX HOWEVER IN PRACTICE THIS IS RARELY AN ISSUE IF NFEATURES IS LARGE ENOUGH EG 2<sup>18</sup> FOR TEXT CLASSIFICATION PROBLEMS
- NO IDF WEIGHTING AS THIS WOULD RENDER THE TRANSFORMER STATEFUL

THE HASH FUNCTION EMPLOYED IS THE SIGNED 32BIT VERSION OF MURMURHASH3

READ MORE IN THE USER GUIDE

PARAMETERS

1752 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

INPUT STRING 'FILENAME' 'FILE' 'CONTENT' IF 'FILENAME' THE SEQUENCE PASSED AS AN ARGUMENT TO FIT IS EXPECTED TO BE A LIST OF FILENAMES THAT NEED READING TO FETCH THE RAW CONTENT TO ANALYZE IF 'FILE' THE SEQUENCE ITEMS MUST HAVE A 'READ' METHOD FILELIKE OBJECT THAT IS CALLED TO FETCH THE BYTES IN MEMORY OTHERWISE THE INPUT IS EXPECTED TO BE THE SEQUENCE STRINGS OR BYTES ITEMS ARE EXPECTED TO BE ANALYZED DIRECTLY ENCODING STRING DEFAULT 'UTF8' IF BYTES OR FILES ARE GIVEN TO ANALYZE THIS ENCODING IS USED TO DECODE DECODEERROR 'STRICT' 'IGNORE' 'REPLACE' INSTRUCTION ON WHAT TO DO IF A BYTE SEQUENCE IS GIVEN TO ANALYZE THAT CONTAINS CHARACTERS NOT OF THE GIVEN ENCODING BY DEFAULT IT IS 'STRICT' MEANING THAT A UNICODEDECODEERROR WILL BE RAISED OTHER VALUES ARE 'IGNORE' AND 'REPLACE'

STRIPACCENTS 'ASCII' 'UNICODE' NONE REMOVE ACCENTS AND PERFORM OTHER CHARACTER NOR MALIZATION DURING THE PREPROCESSING STEP 'ASCII' IS A FAST METHOD THAT ONLY WORKS ON CHAR ACTERS THAT HAVE AN DIRECT ASCII MAPPING 'UNICODE' IS A SLIGHTLY SLOWER METHOD THAT WORKS ON ANY CHARACTERS NONE DEFAULT DOES NOTHING BOTH 'ASCII' AND 'UNICODE' USE NFKD NORMALIZATION FROM UNICODATANORMALIZE LOWERCASE BOOLEAN DEFAULT TRUE CONVERT ALL CHARACTERS TO LOWERCASE BEFORE TOKENIZING PREPROCESSOR CALLABLE OR NONE DEFAULT OVERRIDE THE PREPROCESSING STRING TRANSFORMATION STAGE WHILE PRESERVING THE TOKENIZING AND NGRAMS GENERATION STEPS TOKENIZER CALLABLE OR NONE DEFAULT OVERRIDE THE STRING TOKENIZATION STEP WHILE PRESERVING THE PREPROCESSING AND NGRAMS GENERATION STEPS ONLY APPLIES IF ANALYZER WORD STOPWORDS STRING 'ENGLISH' LIST OR NONE DEFAULT IF 'ENGLISH' A BUILTIN STOP WORD LIST FOR ENGLISH IS USED THERE ARE SEVERAL KNOWN ISSUES WITH 'ENGLISH' AND YOU SHOULD CONSIDER AN ALTERNATIVE SEE USING STOP WORDS IF A LIST THAT LIST IS ASSUMED TO CONTAIN STOP WORDS ALL OF WHICH WILL BE REMOVED FROM THE RESULTING TOKENS ONLY APPLIES IF ANALYZER WORD TOKENPATTERN STRING REGULAR EXPRESSION DENOTING WHAT CONSTITUTES A "TOKEN" ONLY USED IF ANALYZER WORD THE DEFAULT REGEXP SELECTS TOKENS OF 2 OR MORE ALPHANUMERIC CHARACTERS PUNCTUATION IS COMPLETELY IGNORED AND ALWAYS TREATED AS A TOKEN SEPARATOR NGRAMRANGE TUPLE MINN MAXN DEFAULT 1 1 THE LOWER AND UPPER BOUNDARY OF THE RANGE OF NVALUES FOR DIFFERENT NGRAMS TO BE EXTRACTED ALL VALUES OF N SUCH THAT MINN N MAXN WILL BE USED ANALYZER STRING 'WORD' 'CHAR' 'CHARWB' OR CALLABLE WHETHER THE FEATURE SHOULD BE MADE OF WORD OR CHARACTER NGRAMS OPTION 'CHARWB' CREATES CHARACTER NGRAMS ONLY FROM TEXT INSIDE WORD BOUNDARIES NGRAMS AT THE EDGES OF WORDS ARE PADDED WITH SPACE IF A CALLABLE IS PASSED IT IS USED TO EXTRACT THE SEQUENCE OF FEATURES OUT OF THE RAW UNPROCESSED INPUT

CHANGED IN VERSION 021

SINCE V021 IF INPUT ISFILENAME ORFILE THE DATA IS FIRST READ FROM THE FILE AND THEN PASSED TO THE GIVEN CALLABLE ANALYZER NFEATURES INTEGER DEFAULT 20 THE NUMBER OF FEATURES COLUMNS IN THE OUTPUT MATRI CES SMALL NUMBERS OF FEATURES ARE LIKELY TO CAUSE HASH COLLISIONS BUT LARGE NUMBERS WILL CAUSE LARGER COEFFICIENT DIMENSIONS IN LINEAR LEARNERS

615SKLEARNFEATUREEXTRACTION FEATURE EXTRACTION 1753

SCIKITLEARN USER GUIDE RELEASE 0213

BINARY BOOLEAN DEFAULTFALSE IF TRUE ALL NON ZERO COUNTS ARE SET TO 1 THIS IS USEFUL FOR DISCRETE PROBABILISTIC MODELS THAT MODEL BINARY EVENTS RATHER THAN INTEGER COUNTS

NORM ‘L1’ ‘L2’ OR NONE OPTIONAL NORM USED TO NORMALIZE TERM VECTORS NONE FOR NO NORMALIZATION

ALTERNATESIGN BOOLEAN OPTIONAL DEFAULT TRUE WHEN TRUE AN ALTERNATING SIGN IS ADDED TO THE FEATURES AS TO APPROXIMATELY CONSERVE THE INNER PRODUCT IN THE HASHED SPACE EVEN FOR SMALL NFEATURES THIS APPROACH IS SIMILAR TO SPARSE RANDOM PROJECTION

NEW IN VERSION 019

DTYPE TYPE OPTIONAL TYPE OF THE MATRIX RETURNED BY FITTRANSFORM OR TRANSFORM

SEE ALSO

COUNTVECTORIZER TFIDFVECTORIZER

EXAMPLES

FROM SKLEARNFEATUREEXTRACTIONTEXT IMPORT HASHINGVECTORIZER

CORPUS

THIS IS THE FIRST DOCUMENT

THIS DOCUMENT IS THE SECOND DOCUMENT

AND THIS IS THE THIRD ONE

IS THIS THE FIRST DOCUMENT

VECTORIZER HASHINGVECTORIZERNFEATURES2 4

X VECTORIZERFITTRANSFORMCORPUS

PRINTXSHAPE

4 16

METHODS

BUILDANALYZER SELF RETURN A CALLABLE THAT HANDLES PREPROCESSING AND TOKENIZATION

BUILDPREPROCESSOR SELF RETURN A FUNCTION TO PREPROCESS THE TEXT BEFORE TOKENIZATION

BUILDTOKENIZER SELF RETURN A FUNCTION THAT SPLITS A STRING INTO A SEQUENCE OF TOKENS

DECODE SELF DOC DECODE THE INPUT INTO A STRING OF UNICODE SYMBOLS

FITSELF X Y DOES NOTHING THIS TRANSFORMER IS STATELESS

FITTRANSFORM SELF X Y TRANSFORM A SEQUENCE OF DOCUMENTS TO A DOCUMENTTERM MATRIX

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

GETSTOPWORDS SELF BUILD OR FETCH THE EFFECTIVE STOP WORDS LIST

PARTIALFIT SELF X Y DOES NOTHING THIS TRANSFORMER IS STATELESS

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF X TRANSFORM A SEQUENCE OF DOCUMENTS TO A DOCUMENTTERM MATRIX

1754 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

INIT SELFINPUT'CONTENT' ENCODING'UTF8' DECODEERROR'STRICT' STRIPACCENTSNONE

LOWERCASETRUE PREPROCESSORNONE TOKENIZERNONE STOPWORDSNONE TO

KENPATTERN'UBWWB' NGRAMRANGE1 1ANALYZER'WORD' NFEATURES1048576

BINARYFALSE NORM'L2' ALTERNATESIGNTRUE DTYPECLASS 'NUMPYFLOAT64'

BUILDANALYZER SELF

RETURN A CALLABLE THAT HANDLES PREPROCESSING AND TOKENIZATION

BUILDPREPROCESSOR SELF

RETURN A FUNCTION TO PREPROCESS THE TEXT BEFORE TOKENIZATION

BUILDTOKENIZER SELF

RETURN A FUNCTION THAT SPLITS A STRING INTO A SEQUENCE OF TOKENS

DECODESELFDOC

DECODE THE INPUT INTO A STRING OF UNICODE SYMBOLS

THE DECODING STRATEGY DEPENDS ON THE VECTORIZER PARAMETERS

PARAMETERS

DOC STRING THE STRING TO DECODE

FITSELFXYNONE

DOES NOTHING THIS TRANSFORMER IS STATELESS

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA

FITTRANSFORM SELFXYNONE

TRANSFORM A SEQUENCE OF DOCUMENTS TO A DOCUMENTTERM MATRIX

PARAMETERS

XITERABLE OVER RAW TEXT DOCUMENTS LENGTH NSAMPLES SAMPLES EACH SAMPLE MUST BE A

TEXT DOCUMENT EITHER BYTES OR UNICODE STRINGS FILE NAME OR FILE OBJECT DEPENDING ON THE

CONSTRUCTOR ARGUMENT WHICH WILL BE TOKENIZED AND HASHED

YANY IGNORED THIS PARAMETER EXISTS ONLY FOR COMPATIBILITY WITH SKLEARNPIPELINEPIPELINE

RETURNS

XSCIPYSPARSE MATRIX SHAPE NSAMPLES SELFNFEATURES DOCUMENTTERM MATRIX

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

GETSTOPWORDS SELF

BUILD OR FETCH THE EFFECTIVE STOP WORDS LIST

PARTIALFIT SELFXYNONE

DOES NOTHING THIS TRANSFORMER IS STATELESS

THIS METHOD IS JUST THERE TO MARK THE FACT THAT THIS TRANSFORMER CAN WORK IN A STREAMING SETUP

PARAMETERS

615SKLEARNFEATUREEXTRACTION FEATURE EXTRACTION 1755

SCIKITLEARN USER GUIDE RELEASE 0213

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELF

TRANSFORM A SEQUENCE OF DOCUMENTS TO A DOCUMENTTERM MATRIX

PARAMETERS

XITERABLE OVER RAW TEXT DOCUMENTS LENGTH NSAMPLES SAMPLES EACH SAMPLE MUST BE A TEXT DOCUMENT EITHER BYTES OR UNICODE STRINGS FILE NAME OR FILE OBJECT DEPENDING ON THE CONSTRUCTOR ARGUMENT WHICH WILL BE TOKENIZED AND HASHED

RETURNS

XSCIPYSPARSE MATRIX SHAPE NSAMPLES SELFNFEATURES DOCUMENTTERM MATRIX

EXAMPLES USING SKLEARNFEATUREEXTRACTIONTEXTHASHINGVECTORIZER

- OUTOFCORE CLASSIFICATION OF TEXT DOCUMENTS
- CLUSTERING TEXT DOCUMENTS USING KMEANS
- CLASSIFICATION OF TEXT DOCUMENTS USING SPARSE FEATURES

SKLEARNFEATUREEXTRACTIONTEXT TFIDFTRANSFORMER

CLASSSSKLEARNFEATUREEXTRACTIONTEXT TFIDFTRANSFORMER NORM'L2' USEIDFTRUE

SMOOTHIDFTRUE SUBLIN

EARTFFFALSE

TRANSFORM A COUNT MATRIX TO A NORMALIZED TF OR TFIDF REPRESENTATION

TF MEANS TERMFREQUENCY WHILE TFIDF MEANS TERMFREQUENCY TIMES INVERSE DOCUMENTFREQUENCY THIS IS A COMMON TERM WEIGHTING SCHEME IN INFORMATION RETRIEVAL THAT HAS ALSO FOUND GOOD USE IN DOCUMENT CLASSIFICATION

THE GOAL OF USING TFIDF INSTEAD OF THE RAW FREQUENCIES OF OCCURRENCE OF A TOKEN IN A GIVEN DOCUMENT IS TO SCALE DOWN THE IMPACT OF TOKENS THAT OCCUR VERY FREQUENTLY IN A GIVEN CORPUS AND THAT ARE HENCE EMPIRICALLY LESS INFORMATIVE THAN FEATURES THAT OCCUR IN A SMALL FRACTION OF THE TRAINING CORPUS

THE FORMULA THAT IS USED TO COMPUTE THE TFIDF FOR A TERM T OF A DOCUMENT D IN A DOCUMENT SET IS  $TFIDF_{DT} = TF_{DT} \cdot IDF_{DT}$  AND THE IDF IS COMPUTED AS  $IDF_{DT} = \log \frac{N}{DF_{DT} + 1}$  IF SMOOTHIDFFALSE WHERE N IS THE TOTAL NUMBER OF DOCUMENTS IN THE DOCUMENT SET AND DF<sub>DT</sub> IS THE DOCUMENT FREQUENCY OF T THE DOCUMENT FREQUENCY IS THE NUMBER OF DOCUMENTS IN THE DOCUMENT SET THAT CONTAIN THE TERM T THE EFFECT OF ADDING "1" TO THE IDF IN THE EQUATION ABOVE IS THAT TERMS WITH ZERO IDF IE TERMS THAT OCCUR IN ALL DOCUMENTS IN A TRAINING SET WILL NOT BE ENTIRELY IGNORED NOTE THAT THE IDF FORMULA ABOVE DIFFERS FROM THE STANDARD TEXTBOOK NOTATION THAT DEFINES THE IDF AS  $IDF_{DT} = \log \frac{N}{DF_{DT}}$

IFSMOOTHIDFTRUE THE DEFAULT THE CONSTANT "1" IS ADDED TO THE NUMERATOR AND DENOMINATOR OF THE IDF AS IF AN EXTRA DOCUMENT WAS SEEN CONTAINING EVERY TERM IN THE COLLECTION EXACTLY ONCE WHICH PREVENTS ZERO DIVISIONS  $IDF_{DT} = \log \frac{N + 1}{DF_{DT} + 1}$

1756 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

FURTHERMORE THE FORMULAS USED TO COMPUTE TF AND IDF DEPEND ON PARAMETER SETTINGS THAT CORRESPOND TO THE SMART NOTATION USED IN IR AS FOLLOWS

TF IS “N” NATURAL BY DEFAULT “L” LOGARITHMIC WHEN SUBLINEARTFTRUE IDF IS “T” WHEN USEIDF IS GIVEN “N” NONE OTHERWISE NORMALIZATION IS “C” COSINE WHEN NORML2 “N” NONE WHEN NORMNONE

READ MORE IN THE USER GUIDE

PARAMETERS

NORM ‘L1’ ‘L2’ OR NONE OPTIONAL DEFAULT‘L2’ EACH OUTPUT ROW WILL HAVE UNIT NORM EITHER ‘L2’ SUM OF SQUARES OF VECTOR ELEMENTS IS 1 THE COSINE SIMILARITY BETWEEN TWO VECTORS IS THEIR DOT PRODUCT WHEN L2 NORM HAS BEEN APPLIED ‘L1’ SUM OF ABSOLUTE VALUES OF VECTOR ELEMENTS IS 1 SEE PREPROCESSINGNORMALIZE

USEIDF BOOLEAN DEFAULTTRUE ENABLE INVERSEDOCUMENTFREQUENCY REWEIGHTING

SMOOTHIDF BOOLEAN DEFAULTTRUE SMOOTH IDF WEIGHTS BY ADDING ONE TO DOCUMENT FREQUENCIES AS IF AN EXTRA DOCUMENT WAS SEEN CONTAINING EVERY TERM IN THE COLLECTION EXACTLY ONCE PREVENTS ZERO DIVISIONS

SUBLINEARTF BOOLEAN DEFAULTFALSE APPLY SUBLINEAR TF SCALING IE REPLACE TF WITH 1 LOGTF

ATTRIBUTES

IDF ARRAY SHAPE NFEATURES THE INVERSE DOCUMENT FREQUENCY IDF VECTOR ONLY DEFINED IF USEIDF IS TRUE

REFERENCES

R1B90AC3CA370YATES2011 R1B90AC3CA370MRS2008

METHODS

FITSELF X Y LEARN THE IDF VECTOR GLOBAL TERM WEIGHTS

FITTRANSFORM SELF X Y FIT TO DATA THEN TRANSFORM IT

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF X COPY TRANSFORM A COUNT MATRIX TO A TF OR TFIDF REPRESENTATION

INIT SELFNORM‘L2’ USEIDFTRUE SMOOTHIDFTRUE SUBLINEARTFFALSE

FITSELFXYNONE

LEARN THE IDF VECTOR GLOBAL TERM WEIGHTS

PARAMETERS

XSPARSE MATRIX NSAMPLES NFEATURES A MATRIX OF TERMTOKEN COUNTS

FITTRANSFORM SELFXYNONE FITPARAMS

FIT TO DATA THEN TRANSFORM IT

FITS TRANSFORMER TO X AND Y WITH OPTIONAL PARAMETERS FITPARAMS AND RETURNS A TRANSFORMED VERSION OF X

PARAMETERS

XNUMPY ARRAY OF SHAPE NSAMPLES NFEATURES TRAINING SET

615SKLEARNFEATUREEXTRACTION FEATURE EXTRACTION 1757

SCIKITLEARN USER GUIDE RELEASE 0213

Y: NUMPY ARRAY OF SHAPE (NSAMPLES, TARGET VALUES)

RETURNS

X: NEW NUMPY ARRAY OF SHAPE (NSAMPLES, NFEATURES)

NEW TRANSFORMED ARRAY

GETPARAMS: SELF

DEEPT: TRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP: BOOLEAN, OPTIONAL. IF TRUE, WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS: MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

SETPARAMS: SELF

PARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES. THE LATTER HAVE PARAMETERS OF THE FORM COMPONENT.PARAMETER, SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM: SELF

XCOPY: TRUE

TRANSFORM A COUNT MATRIX TO A TF OR TFIDF REPRESENTATION

PARAMETERS

X: SPARSE MATRIX OF SHAPE (NSAMPLES, NFEATURES)

A MATRIX OF TERM-TOKEN COUNTS

COPY: BOOLEAN, DEFAULT TRUE. WHETHER TO COPY X AND OPERATE ON THE COPY OR PERFORM INPLACE OPERATIONS

RETURNS

VECTORS: SPARSE MATRIX OF SHAPE (NSAMPLES, NFEATURES)

EXAMPLES USING SKLEARN.FEATURE.EXTRACTION.TEXT.TFIDF.TRANSFORMER

- SAMPLE PIPELINE FOR TEXT FEATURE EXTRACTION AND EVALUATION
- CLUSTERING TEXT DOCUMENTS USING KMEANS

1758 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213  
SKLEARNFEATUREEXTRACTIONTEXT TFIDFVECTORIZER  
CLASSSSKLEARNFEATUREEXTRACTIONTEXT TFIDFVECTORIZER INPUT'CONTENT' ENCODING'UTF  
8' DECODEERROR'STRICT'  
STRIPACCENTSNONE LOW  
ERCASETRUE PREPROCES  
SORNONE TOKENIZERNONE ANA  
LYZER'WORD' STOPWORDSNONE  
TOKENPATTERN'UBWWB'  
NGRAMRANGE1 1  
MAXDF10 MINDF1  
MAXFEATURESNONE VOCAB  
ULARYNONE BINARYFALSE  
DTYPECLASS 'NUMPYFLOAT64'  
NORM'L2' USEIDFTRUE  
SMOOTHIDFTRUE SUBLIN  
EARTFFALSE  
CONVERT A COLLECTION OF RAW DOCUMENTS TO A MATRIX OF TFIDF FEATURES  
EQUIVALENT TO COUNTVECTORIZER FOLLOWED BY TFIDFTRANSFORMER  
READ MORE IN THE USER GUIDE  
PARAMETERS  
INPUT STRING 'FILENAME' 'FILE' 'CONTENT' IF 'FILENAME' THE SEQUENCE PASSED AS AN ARGUMENT TO  
FIT IS EXPECTED TO BE A LIST OF FILENAMES THAT NEED READING TO FETCH THE RAW CONTENT TO ANALYZE  
IF 'FILE' THE SEQUENCE ITEMS MUST HAVE A 'READ' METHOD FILELIKE OBJECT THAT IS CALLED TO FETCH  
THE BYTES IN MEMORY  
OTHERWISE THE INPUT IS EXPECTED TO BE THE SEQUENCE STRINGS OR BYTES ITEMS ARE EXPECTED TO BE  
ANALYZED DIRECTLY  
ENCODING STRING 'UTF8' BY DEFAULT IF BYTES OR FILES ARE GIVEN TO ANALYZE THIS ENCODING IS USED  
TO DECODE  
DECODEERROR 'STRICT' 'IGNORE' 'REPLACE' DEFAULT'STRICT' INSTRUCTION ON WHAT TO DO IF A  
BYTE SEQUENCE IS GIVEN TO ANALYZE THAT CONTAINS CHARACTERS NOT OF THE GIVEN ENCODING BY  
DEFAULT IT IS 'STRICT' MEANING THAT A UNICODEDECODEERROR WILL BE RAISED OTHER VALUES ARE  
'IGNORE' AND 'REPLACE'  
STRIPACCENTS 'ASCII' 'UNICODE' NONE DEFAULTNONE REMOVE ACCENTS AND PERFORM OTHER  
CHARACTER NORMALIZATION DURING THE PREPROCESSING STEP 'ASCII' IS A FAST METHOD THAT ONLY  
WORKS ON CHARACTERS THAT HAVE AN DIRECT ASCII MAPPING 'UNICODE' IS A SLIGHTLY SLOWER  
METHOD THAT WORKS ON ANY CHARACTERS NONE DEFAULT DOES NOTHING  
BOTH 'ASCII' AND 'UNICODE' USE NFKD NORMALIZATION FROM UNICODEDATANORMALIZE  
LOWERCASE BOOLEAN DEFAULTTRUE CONVERT ALL CHARACTERS TO LOWERCASE BEFORE TOKENIZING  
PREPROCESSOR CALLABLE OR NONE DEFAULTNONE OVERRIDE THE PREPROCESSING STRING TRANSFORMA  
TION STAGE WHILE PRESERVING THE TOKENIZING AND NGRAMS GENERATION STEPS  
TOKENIZER CALLABLE OR NONE DEFAULTNONE OVERRIDE THE STRING TOKENIZATION STEP WHILE PRE  
SERVING THE PREPROCESSING AND NGRAMS GENERATION STEPS ONLY APPLIES IF ANALYZER  
WORD  
ANALYZER STRING 'WORD' 'CHAR' 'CHARWB' OR CALLABLE WHETHER THE FEATURE SHOULD BE MADE  
OF WORD OR CHARACTER NGRAMS OPTION 'CHARWB' CREATES CHARACTER NGRAMS ONLY FROM TEXT  
INSIDE WORD BOUNDARIES NGRAMS AT THE EDGES OF WORDS ARE PADDED WITH SPACE  
615SKLEARNFEATUREEXTRACTION FEATURE EXTRACTION 1759

SCIKITLEARN USER GUIDE RELEASE 0213

IF A CALLABLE IS PASSED IT IS USED TO EXTRACT THE SEQUENCE OF FEATURES OUT OF THE RAW UNPROCESSED INPUT

CHANGED IN VERSION 021

SINCE V021 IF INPUT ISFILENAME ORFILE THE DATA IS FIRST READ FROM THE FILE AND THEN PASSED TO THE GIVEN CALLABLE ANALYZER

STOPWORDS STRING ‘ENGLISH’ LIST OR NONE DEFAULTNONE IF A STRING IT IS PASSED TO CHECKSTOPLIST AND THE APPROPRIATE STOP LIST IS RETURNED ‘ENGLISH’ IS CURRENTLY THE ONLY SUPPORTED STRING VALUE THERE ARE SEVERAL KNOWN ISSUES WITH ‘ENGLISH’ AND YOU SHOULD CONSIDER AN ALTERNATIVE SEE USING STOP WORDS

IF A LIST THAT LIST IS ASSUMED TO CONTAIN STOP WORDS ALL OF WHICH WILL BE REMOVED FROM THE RESULTING TOKENS ONLY APPLIES IF ANALYZER WORD

IF NONE NO STOP WORDS WILL BE USED MAXDF CAN BE SET TO A VALUE IN THE RANGE 07 10 TO AUTOMATICALLY DETECT AND FILTER STOP WORDS BASED ON INTRA CORPUS DOCUMENT FREQUENCY OF TERMS

TOKENPATTERN STRING REGULAR EXPRESSION DENOTING WHAT CONSTITUTES A “TOKEN” ONLY USED IF ANALYZER WORD THE DEFAULT REGEXP SELECTS TOKENS OF 2 OR MORE ALPHANUMERIC CHARACTERS PUNCTUATION IS COMPLETELY IGNORED AND ALWAYS TREATED AS A TOKEN SEPARATOR

NGRAMRANGE TUPLE MINN MAXN DEFAULT1 1 THE LOWER AND UPPER BOUNDARY OF THE RANGE OF NVALUES FOR DIFFERENT NGRAMS TO BE EXTRACTED ALL VALUES OF N SUCH THAT MINN N MAXN WILL BE USED

MAXDF FLOAT IN RANGE 00 10 OR INT DEFAULT10 WHEN BUILDING THE VOCABULARY IGNORE TERMS THAT HAVE A DOCUMENT FREQUENCY STRICTLY HIGHER THAN THE GIVEN THRESHOLD CORPUS SPECIFIC STOP WORDS IF FLOAT THE PARAMETER REPRESENTS A PROPORTION OF DOCUMENTS INTEGER ABSOLUTE COUNTS THIS PARAMETER IS IGNORED IF VOCABULARY IS NOT NONE

MINDF FLOAT IN RANGE 00 10 OR INT DEFAULT1 WHEN BUILDING THE VOCABULARY IGNORE TERMS THAT HAVE A DOCUMENT FREQUENCY STRICTLY LOWER THAN THE GIVEN THRESHOLD THIS VALUE IS ALSO CALLED CUTOFF IN THE LITERATURE IF FLOAT THE PARAMETER REPRESENTS A PROPORTION OF DOCUMENTS INTEGER ABSOLUTE COUNTS THIS PARAMETER IS IGNORED IF VOCABULARY IS NOT NONE

MAXFEATURES INT OR NONE DEFAULTNONE IF NOT NONE BUILD A VOCABULARY THAT ONLY CONSIDER THE TOP MAXFEATURES ORDERED BY TERM FREQUENCY ACROSS THE CORPUS

THIS PARAMETER IS IGNORED IF VOCABULARY IS NOT NONE

VOCABULARY MAPPING OR ITERABLE OPTIONAL DEFAULTNONE EITHER A MAPPING EG A DICT WHERE KEYS ARE TERMS AND VALUES ARE INDICES IN THE FEATURE MATRIX OR AN ITERABLE OVER TERMS IF NOT GIVEN A VOCABULARY IS DETERMINED FROM THE INPUT DOCUMENTS

BINARY BOOLEAN DEFAULTFALSE IF TRUE ALL NONZERO TERM COUNTS ARE SET TO 1 THIS DOES NOT MEAN OUTPUTS WILL HAVE ONLY 01 VALUES ONLY THAT THE TF TERM IN TFIDF IS BINARY SET IDF AND NORMALIZATION TO FALSE TO GET 01 OUTPUTS

DTYPE TYPE OPTIONAL DEFAULTFLOAT64 TYPE OF THE MATRIX RETURNED BY FITTRANSFORM OR TRANSFORM

NORM ‘L1’ ‘L2’ OR NONE OPTIONAL DEFAULT‘L2’ EACH OUTPUT ROW WILL HAVE UNIT NORM EITHER ‘L2’ SUM OF SQUARES OF VECTOR ELEMENTS IS 1 THE COSINE SIMILARITY BETWEEN TWO VECTORS IS THEIR DOT PRODUCT WHEN L2 NORM HAS BEEN APPLIED ‘L1’ SUM OF ABSOLUTE VALUES OF VECTOR ELEMENTS IS 1 SEE PREPROCESSINGNORMALIZE

USEIDF BOOLEAN DEFAULTTRUE ENABLE INVERSEDOCUMENTFREQUENCY REWEIGHTING

1760 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

SMOOTHIDF BOOLEAN DEFAULTTRUE SMOOTH IDF WEIGHTS BY ADDING ONE TO DOCUMENT FREQUENCIES AS IF AN EXTRA DOCUMENT WAS SEEN CONTAINING EVERY TERM IN THE COLLECTION EXACTLY ONCE PREVENTS ZERO DIVISIONS

SUBLINEARTF BOOLEAN DEFAULTFALSE APPLY SUBLINEAR TF SCALING IE REPLACE TF WITH 1 LOGTF

ATTRIBUTES

VOCABULARY DICT A MAPPING OF TERMS TO FEATURE INDICES

IDF ARRAY SHAPE NFEATURES THE INVERSE DOCUMENT FREQUENCY IDF VECTOR ONLY DEFINED IF USEIDF IS TRUE

STOPWORDS SET TERMS THAT WERE IGNORED BECAUSE THEY EITHER

- OCCURRED IN TOO MANY DOCUMENTS MAXDF
- OCCURRED IN TOO FEW DOCUMENTS MINDF
- WERE CUT OFF BY FEATURE SELECTION MAXFEATURES

THIS IS ONLY AVAILABLE IF NO VOCABULARY WAS GIVEN

SEE ALSO

COUNTVECTORIZER TRANSFORMS TEXT INTO A SPARSE MATRIX OF NGRAM COUNTS

TFIDFTRANSFORMER PERFORMS THE TFIDF TRANSFORMATION FROM A PROVIDED MATRIX OF COUNTS

NOTES

THESTOPWORDS ATTRIBUTE CAN GET LARGE AND INCREASE THE MODEL SIZE WHEN PICKLING THIS ATTRIBUTE IS PROVIDED ONLY FOR INTROSPECTION AND CAN BE SAFELY REMOVED USING DELATTR OR SET TO NONE BEFORE PICKLING

EXAMPLES

```
FROM SKLEARNFEATUREEXTRACTIONTEXT IMPORT TFIDFVECTORIZER
CORPUS
THIS IS THE FIRST DOCUMENT
THIS DOCUMENT IS THE SECOND DOCUMENT
AND THIS IS THE THIRD ONE
IS THIS THE FIRST DOCUMENT
```

VECTORIZER TFIDFVECTORIZER

```
X VECTORIZERFITTRANSFORMCORPUS
PRINTVECTORIZERGETFEATURENAMES
AND DOCUMENT FIRST IS ONE SECOND THE THIRD THIS
PRINTXSHAPE
4 9
```

METHODS

BUILDANALYZER SELF RETURN A CALLABLE THAT HANDLES PREPROCESSING AND TOKENIZATION

CONTINUED ON NEXT PAGE

615SKLEARNFEATUREEXTRACTION FEATURE EXTRACTION 1761

SCIKITLEARN USER GUIDE RELEASE 0213

TABLE 6103 – CONTINUED FROM PREVIOUS PAGE

BUILDPREPROCESSOR SELF RETURN A FUNCTION TO PREPROCESS THE TEXT BEFORE TOKENIZATION

BUILDTOKENIZER SELF RETURN A FUNCTION THAT SPLITS A STRING INTO A SEQUENCE OF TOKENS

DECODE SELF DOC DECODE THE INPUT INTO A STRING OF UNICODE SYMBOLS

FITSELF RAWDOCUMENTS Y LEARN VOCABULARY AND IDF FROM TRAINING SET

FITTRANSFORM SELF RAWDOCUMENTS Y LEARN VOCABULARY AND IDF RETURN TERMDOCUMENT MATRIX

GETFEATURENAMES SELF ARRAY MAPPING FROM FEATURE INTEGER INDICES TO FEATURE NAME

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

GETSTOPWORDS SELF BUILD OR FETCH THE EFFECTIVE STOP WORDS LIST

INVERSETRANSFORM SELF X RETURN TERMS PER DOCUMENT WITH NONZERO ENTRIES IN X

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF RAWDOCUMENTS COPY TRANSFORM DOCUMENTS TO DOCUMENTTERM MATRIX

INIT SELFINPUT‘CONTENT’ ENCODING‘UTF8’ DECODEERROR‘STRICT’ STRIPACCENTSNONE LOW ERCASETTRUE PREPROCESSORNONE TOKENIZERNONE ANALYZER‘WORD’ STOPWORDSNONE TOKENPATTERN‘UBWWB’ NGRAMRANGE1 1MAXDF10 MINDF1 MAXFEATURESNONE VOCABULARYNONE BINARYFALSE DTYPECLASS ‘NUMPYFLOAT64’ NORM‘L2’ USEIDFTRUE SMOOTHIDFTRUE SUBLINEARTFFALSE

BUILDANALYZER SELF

RETURN A CALLABLE THAT HANDLES PREPROCESSING AND TOKENIZATION

BUILDPREPROCESSOR SELF

RETURN A FUNCTION TO PREPROCESS THE TEXT BEFORE TOKENIZATION

BUILDTOKENIZER SELF

RETURN A FUNCTION THAT SPLITS A STRING INTO A SEQUENCE OF TOKENS

DECODESELFDOC

DECODE THE INPUT INTO A STRING OF UNICODE SYMBOLS

THE DECODING STRATEGY DEPENDS ON THE VECTORIZER PARAMETERS

PARAMETERS

DOC STRING THE STRING TO DECODE

FITSELFRAWDOCUMENTS YNONE

LEARN VOCABULARY AND IDF FROM TRAINING SET

PARAMETERS

RAWDOCUMENTS ITERABLE AN ITERABLE WHICH YIELDS EITHER STR UNICODE OR FILE OBJECTS

RETURNS

SELF TFIDFVECTORIZER

FITTRANSFORM SELFRAWDOCUMENTS YNONE

LEARN VOCABULARY AND IDF RETURN TERMDOCUMENT MATRIX

THIS IS EQUIVALENT TO FIT FOLLOWED BY TRANSFORM BUT MORE EFFICIENTLY IMPLEMENTED

PARAMETERS

RAWDOCUMENTS ITERABLE AN ITERABLE WHICH YIELDS EITHER STR UNICODE OR FILE OBJECTS

RETURNS

XSPARSE MATRIX NSAMPLES NFEATURES TFIDFWEIGHTED DOCUMENTTERM MATRIX

1762 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

GETFEATURENAMES SELF

ARRAY MAPPING FROM FEATURE INTEGER INDICES TO FEATURE NAME

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

GETSTOPWORDS SELF

BUILD OR FETCH THE EFFECTIVE STOP WORDS LIST

INVERSETRANSFORM SELF

RETURN TERMS PER DOCUMENT WITH NONZERO ENTRIES IN X

PARAMETERS

XARRAY SPARSE MATRIX SHAPE NSAMPLES NFEATURES

RETURNS

XINV LIST OF ARRAYS LEN NSAMPLES LIST OF ARRAYS OF TERMS

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE

PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT

OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELFRAWDOCUMENTS COPYTRUE

TRANSFORM DOCUMENTS TO DOCUMENTTERM MATRIX

USES THE VOCABULARY AND DOCUMENT FREQUENCIES DF LEARNED BY FIT OR FITTRANSFORM

PARAMETERS

RAWDOCUMENTS ITERABLE AN ITERABLE WHICH YIELDS EITHER STR UNICODE OR FILE OBJECTS

COPY BOOLEAN DEFAULT TRUE WHETHER TO COPY X AND OPERATE ON THE COPY OR PERFORM INPLACE

OPERATIONS

RETURNS

XSPARSE MATRIX NSAMPLES NFEATURES TFIDFWEIGHTED DOCUMENTTERM MATRIX

EXAMPLES USING SKLEARNFEATUREEXTRACTIONTEXTTFIDFVECTORIZER

- TOPIC EXTRACTION WITH NONNEGATIVE MATRIX FACTORIZATION AND LATENT DIRICHLET ALLOCATION
- BICLUSTERING DOCUMENTS WITH THE SPECTRAL COCLUSTERING ALGORITHM
- COLUMN TRANSFORMER WITH HETEROGENEOUS DATA SOURCES
- CLUSTERING TEXT DOCUMENTS USING KMEANS

615SKLEARNFEATUREEXTRACTION FEATURE EXTRACTION 1763

SCIKITLEARN USER GUIDE RELEASE 0213

- CLASSIFICATION OF TEXT DOCUMENTS USING SPARSE FEATURES

6161SKLEARNFEATURESELECTION FEATURE SELECTION

THESKLEARNFEATURESELECTION MODULE IMPLEMENTS FEATURE SELECTION ALGORITHMS IT CURRENTLY INCLUDES UNIVARIATE FILTER SELECTION METHODS AND THE RECURSIVE FEATURE ELIMINATION ALGORITHM

USER GUIDE SEE THE FEATURE SELECTION SECTION FOR FURTHER DETAILS

FEATURESELECTION

GENERICUNIVARIATESELECT UNIVARIATE FEATURE SELECTOR WITH CONFIGURABLE STRATEGY

FEATURESELECTIONSELECTPERCENTILE SELECT FEATURES ACCORDING TO A PERCENTILE OF THE HIGHEST SCORES

FEATURESELECTIONSELECTKBEST SCOREFUNC

KSELECT FEATURES ACCORDING TO THE K HIGHEST SCORES

FEATURESELECTIONSELECTFPR SCOREFUNC AL

PHAFILTER SELECT THE PVALUES BELOW ALPHA BASED ON A FPR TEST

FEATURESELECTIONSELECTFDR SCOREFUNC AL

PHAFILTER SELECT THE PVALUES FOR AN ESTIMATED FALSE DISCOVERY RATE

FEATURESELECTION

SELECTFROMMODEL ESTIMATORMETATransformer FOR SELECTING FEATURES BASED ON IMPORTANCE WEIGHTS

FEATURESELECTIONSELECTFWE SCOREFUNC AL

PHAFILTER SELECT THE PVALUES CORRESPONDING TO FAMILYWISE ERROR RATE

FEATURESELECTIONRFE ESTIMATOR FEATURE RANKING WITH RECURSIVE FEATURE ELIMINATION

FEATURESELECTIONRFECV ESTIMATOR STEP FEATURE RANKING WITH RECURSIVE FEATURE ELIMINATION AND CROSSVALIDATED SELECTION OF THE BEST NUMBER OF FEATURES

FEATURESELECTION

VARIANCETHRESHOLD THRESHOLDFEATURE SELECTOR THAT REMOVES ALL LOWVARIANCE FEATURES

6161SKLEARNFEATURESELECTION GENERICUNIVARIATESELECT

CLASSSSKLEARNFEATURESELECTION GENERICUNIVARIATESELECT SCOREFUNCFUNCTION

FCLASSIF MODE'PERCENTILE'

PARAM1E05

UNIVARIATE FEATURE SELECTOR WITH CONFIGURABLE STRATEGY

READ MORE IN THE USER GUIDE

PARAMETERS

SCOREFUNC CALLABLE FUNCTION TAKING TWO ARRAYS X AND Y AND RETURNING A PAIR OF ARRAYS SCORES PVALUES FOR MODES 'PERCENTILE' OR 'KBEST' IT CAN RETURN A SINGLE ARRAY SCORES

MODE 'PERCENTILE' 'KBEST' 'FPR' 'FDR' 'FWE' FEATURE SELECTION MODE

PARAM FLOAT OR INT DEPENDING ON THE FEATURE SELECTION MODE PARAMETER OF THE CORRESPONDING MODE

ATTRIBUTES

SCORES ARRAYLIKE SHAPENFEATURES SCORES OF FEATURES

PVALUES ARRAYLIKE SHAPENFEATURES PVALUES OF FEATURE SCORES NONE IF SCOREFUNC

RETURNED SCORES ONLY

SEE ALSO

1764 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

FCLASSIF ANOV A FVALUE BETWEEN LABELFEATURE FOR CLASSIFICATION TASKS

MUTUALINFOCLASSIF MUTUAL INFORMATION FOR A DISCRETE TARGET

CHI2 CHISQUARED STATS OF NONNEGATIVE FEATURES FOR CLASSIFICATION TASKS

FREGRESSION FVALUE BETWEEN LABELFEATURE FOR REGRESSION TASKS

MUTUALINFOREGRESSION MUTUAL INFORMATION FOR A CONTINUOUS TARGET

SELECTPERCENTILE SELECT FEATURES BASED ON PERCENTILE OF THE HIGHEST SCORES

SELECTKBEST SELECT FEATURES BASED ON THE K HIGHEST SCORES

SELECTFPR SELECT FEATURES BASED ON A FALSE POSITIVE RATE TEST

SELECTFDR SELECT FEATURES BASED ON AN ESTIMATED FALSE DISCOVERY RATE

SELECTFWE SELECT FEATURES BASED ON FAMILYWISE ERROR RATE

EXAMPLES

FROM SKLEARNDATASETS IMPORT LOADBREASTCANCER

FROM SKLEARNFEATURESELECTION IMPORT GENERICUNIVARIATESELECT CHI2

X Y LOADBREASTCANCERRETURNXY TRUE

XSHAPE

569 30

TRANSFORMER GENERICUNIVARIATESELECTCHI2 KBEST PARAM20

XNEW TRANSFORMERFITTRANSFORMX Y

XNEWSHAPE

569 20

METHODS

FITSELF X Y RUN SCORE FUNCTION ON X Y AND GET THE APPROPRIATE

FEATURES

FITTRANSFORM SELF X Y FIT TO DATA THEN TRANSFORM IT

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

GETSUPPORT SELF INDICES GET A MASK OR INTEGER INDEX OF THE FEATURES SELECTED

INVERSETRANSFORM SELF X REVERSE THE TRANSFORMATION OPERATION

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF X REDUCE X TO THE SELECTED FEATURES

INIT SELFSCOREFUNCFUNCTION FCLASSIF AT 0X7EFE30BB2158 MODE'PERCENTILE' PARAM1E

05

FITSELFXY

RUN SCORE FUNCTION ON X Y AND GET THE APPROPRIATE FEATURES

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES THE TRAINING INPUT SAMPLES

YARRAYLIKE SHAPE NSAMPLES THE TARGET VALUES CLASS LABELS IN CLASSIFICATION REAL NUM

BERS IN REGRESSION

RETURNS

616SKLEARNFEATURESELECTION FEATURE SELECTION 1765

SCIKITLEARN USER GUIDE RELEASE 0213

SELF OBJECT

FITTRANSFORM SELFXYNONE FITPARAMS

FIT TO DATA THEN TRANSFORM IT

FITS TRANSFORMER TO X AND Y WITH OPTIONAL PARAMETERS FITPARAMS AND RETURNS A TRANSFORMED VERSION OF X

PARAMETERS

XNUMPY ARRAY OF SHAPE NSAMPLES NFEATURES TRAINING SET

YNUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES

RETURNS

XNEW NUMPY ARRAY OF SHAPE NSAMPLES NFEATURESNEW TRANSFORMED ARRAY

GETPARAMS SELFDEEPTTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

GETSUPPORT SELFINDICESFALSE

GET A MASK OR INTEGER INDEX OF THE FEATURES SELECTED

PARAMETERS

INDICES BOOLEAN DEFAULT FALSE IF TRUE THE RETURN VALUE WILL BE AN ARRAY OF INTEGERS RATHER

THAN A BOOLEAN MASK

RETURNS

SUPPORT ARRAY AN INDEX THAT SELECTS THE RETAINED FEATURES FROM A FEATURE VECTOR IF

INDICES IS FALSE THIS IS A BOOLEAN ARRAY OF SHAPE INPUT FEATURES IN WHICH AN ELE

MENT IS TRUE IFF ITS CORRESPONDING FEATURE IS SELECTED FOR RETENTION IF INDICES IS TRUE

THIS IS AN INTEGER ARRAY OF SHAPE OUTPUT FEATURES WHOSE VALUES ARE INDICES INTO THE INPUT

FEATURE VECTOR

INVERSETRANSFORM SELF

REVERSE THE TRANSFORMATION OPERATION

PARAMETERS

XARRAY OF SHAPE NSAMPLES NSELECTEDFEATURES THE INPUT SAMPLES

RETURNS

XR ARRAY OF SHAPE NSAMPLES NORIGINALFEATURES XWITH COLUMNS OF ZEROS INSERTED

WHERE FEATURES WOULD HAVE BEEN REMOVED BY TRANSFORM

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE

PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT

OF A NESTED OBJECT

RETURNS

SELF

1766 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

TRANSFORM SELF X

REDUCE X TO THE SELECTED FEATURES

PARAMETERS

X ARRAY OF SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES

RETURNS

XR ARRAY OF SHAPE NSAMPLES NSELECTEDFEATURES THE INPUT SAMPLES WITH ONLY THE SE  
LECTED FEATURES

6162SKLEARNFEATURESELECTION SELECTPERCENTILE

CLASSSKLEARNFEATURESELECTION SELECTPERCENTILE SCOREFUNCFUNCTION FCLASSIF PER  
CENTILE10

SELECT FEATURES ACCORDING TO A PERCENTILE OF THE HIGHEST SCORES

READ MORE IN THE USER GUIDE

PARAMETERS

SCOREFUNC CALLABLE FUNCTION TAKING TWO ARRAYS X AND Y AND RETURNING A PAIR OF ARRAYS SCORES  
PVALUES OR A SINGLE ARRAY WITH SCORES DEFAULT IS FCLASSIF SEE BELOW “SEE ALSO” THE  
DEFAULT FUNCTION ONLY WORKS WITH CLASSIFICATION TASKS

PERCENTILE INT OPTIONAL DEFAULT10 PERCENT OF FEATURES TO KEEP

ATTRIBUTES

SCORES ARRAYLIKE SHAPENFEATURES SCORES OF FEATURES

PVALUES ARRAYLIKE SHAPENFEATURES PVALUES OF FEATURE SCORES NONE IF SCOREFUNC  
RETURNED ONLY SCORES

SEE ALSO

FCLASSIF ANOV A FVALUE BETWEEN LABELFEATURE FOR CLASSIFICATION TASKS

MUTUALINFOCLASSIF MUTUAL INFORMATION FOR A DISCRETE TARGET

CHI2 CHISQUARED STATS OF NONNEGATIVE FEATURES FOR CLASSIFICATION TASKS

FREGRESSION FVALUE BETWEEN LABELFEATURE FOR REGRESSION TASKS

MUTUALINFOREGRESSION MUTUAL INFORMATION FOR A CONTINUOUS TARGET

SELECTKBEST SELECT FEATURES BASED ON THE K HIGHEST SCORES

SELECTFPR SELECT FEATURES BASED ON A FALSE POSITIVE RATE TEST

SELECTFDR SELECT FEATURES BASED ON AN ESTIMATED FALSE DISCOVERY RATE

SELECTFWE SELECT FEATURES BASED ON FAMILYWISE ERROR RATE

GENERICUNIVARIATESELECT UNIVARIATE FEATURE SELECTOR WITH CONFIGURABLE MODE

NOTES

TIES BETWEEN FEATURES WITH EQUAL SCORES WILL BE BROKEN IN AN UNSPECIFIED WAY

616SKLEARNFEATURESELECTION FEATURE SELECTION 1767

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES

FROM SKLEARNDATASETS IMPORT LOADDIGITS

FROM SKLEARNFEATURESELECTION IMPORT SELECTPERCENTILE CHI2

X Y LOADDIGITSRETURNXY TRUE

XSHAPE

1797 64

XNEW SELECTPERCENTILECHI2 PERCENTILE10FITTRANSFORMX Y

XNEWSHAPE

1797 7

METHODS

FITSELF X Y RUN SCORE FUNCTION ON X Y AND GET THE APPROPRIATE FEATURES

FITTRANSFORM SELF X Y FIT TO DATA THEN TRANSFORM IT

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

GETSUPPORT SELF INDICES GET A MASK OR INTEGER INDEX OF THE FEATURES SELECTED

INVERSETRANSFORM SELF X REVERSE THE TRANSFORMATION OPERATION

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF X REDUCE X TO THE SELECTED FEATURES

INIT SELFSCOREFUNCFUNCTION FCLASSIF AT 0X7F3C10E93840 PERCENTILE10

FITSELFXY

RUN SCORE FUNCTION ON X Y AND GET THE APPROPRIATE FEATURES

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES THE TRAINING INPUT SAMPLES

YARRAYLIKE SHAPE NSAMPLES THE TARGET VALUES CLASS LABELS IN CLASSIFICATION REAL NUMBERS IN REGRESSION

RETURNS

SELF OBJECT

FITTRANSFORM SELFXYNONE FITPARAMS

FIT TO DATA THEN TRANSFORM IT

FITS TRANSFORMER TO X AND Y WITH OPTIONAL PARAMETERS FITPARAMS AND RETURNS A TRANSFORMED VERSION OF X

PARAMETERS

XNUMPY ARRAY OF SHAPE NSAMPLES NFEATURES TRAINING SET

YNUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES

RETURNS

XNEW NUMPY ARRAY OF SHAPE NSAMPLES NFEATURESNEW TRANSFORMED ARRAY

GETPARAMS SELFDEEPTTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

1768 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

GETSUPPORT SELFINDICESFALSE

GET A MASK OR INTEGER INDEX OF THE FEATURES SELECTED

PARAMETERS

INDICES BOOLEAN DEFAULT FALSE IF TRUE THE RETURN VALUE WILL BE AN ARRAY OF INTEGERS RATHER THAN A BOOLEAN MASK

RETURNS

SUPPORT ARRAY AN INDEX THAT SELECTS THE RETAINED FEATURES FROM A FEATURE VECTOR IF INDICES IS FALSE THIS IS A BOOLEAN ARRAY OF SHAPE INPUT FEATURES IN WHICH AN ELEMENT IS TRUE IFF ITS CORRESPONDING FEATURE IS SELECTED FOR RETENTION IF INDICES IS TRUE THIS IS AN INTEGER ARRAY OF SHAPE OUTPUT FEATURES WHOSE VALUES ARE INDICES INTO THE INPUT FEATURE VECTOR

INVERSETRANSFORM SELF

REVERSE THE TRANSFORMATION OPERATION

PARAMETERS

XARRAY OF SHAPE NSAMPLES NSELECTEDFEATURES THE INPUT SAMPLES

RETURNS

XR ARRAY OF SHAPE NSAMPLES NORIGINALFEATURES XWITH COLUMNS OF ZEROS INSERTED WHERE FEATURES WOULD HAVE BEEN REMOVED BY TRANSFORM

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELF

REDUCE X TO THE SELECTED FEATURES

PARAMETERS

XARRAY OF SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES

RETURNS

XR ARRAY OF SHAPE NSAMPLES NSELECTEDFEATURES THE INPUT SAMPLES WITH ONLY THE SELECTED FEATURES

EXAMPLES USING SKLEARNFEATURESELECTIONSELECTPERCENTILE

- FEATURE AGGLOMERATION VS UNIVARIATE SELECTION
- UNIVARIATE FEATURE SELECTION
- SVMANOVA SVM WITH UNIVARIATE FEATURE SELECTION

616SKLEARNFEATURESELECTION FEATURE SELECTION 1769

SCIKITLEARN USER GUIDE RELEASE 0213  
6163SKLEARNFEATURESELECTION SELECTKBEST  
CLASSSSKLEARNFEATURESELECTION SELECTKBEST SCOREFUNCFUNCTION FCLASSIF K10  
SELECT FEATURES ACCORDING TO THE K HIGHEST SCORES  
READ MORE IN THE USER GUIDE  
PARAMETERS  
SCOREFUNC CALLABLE FUNCTION TAKING TWO ARRAYS X AND Y AND RETURNING A PAIR OF ARRAYS SCORES  
PVALUES OR A SINGLE ARRAY WITH SCORES DEFAULT IS FCLASSIF SEE BELOW “SEE ALSO” THE  
DEFAULT FUNCTION ONLY WORKS WITH CLASSIFICATION TASKS  
KINT OR “ALL” OPTIONAL DEFAULT10 NUMBER OF TOP FEATURES TO SELECT THE “ALL” OPTION BYPASSES  
SELECTION FOR USE IN A PARAMETER SEARCH  
ATTRIBUTES  
SCORES ARRAYLIKE SHAPENFEATURES SCORES OF FEATURES  
PVALUES ARRAYLIKE SHAPENFEATURES PVALUES OF FEATURE SCORES NONE IF SCOREFUNC  
RETURNED ONLY SCORES  
SEE ALSO  
FCLASSIF ANOV A FVALUE BETWEEN LABELFEATURE FOR CLASSIFICATION TASKS  
MUTUALINFOCLASSIF MUTUAL INFORMATION FOR A DISCRETE TARGET  
CHI2 CHISQUARED STATS OF NONNEGATIVE FEATURES FOR CLASSIFICATION TASKS  
FREGRESSION FVALUE BETWEEN LABELFEATURE FOR REGRESSION TASKS  
MUTUALINFOREGRESSION MUTUAL INFORMATION FOR A CONTINUOUS TARGET  
SELECTPERCENTILE SELECT FEATURES BASED ON PERCENTILE OF THE HIGHEST SCORES  
SELECTFPR SELECT FEATURES BASED ON A FALSE POSITIVE RATE TEST  
SELECTFDR SELECT FEATURES BASED ON AN ESTIMATED FALSE DISCOVERY RATE  
SELECTFWE SELECT FEATURES BASED ON FAMILYWISE ERROR RATE  
GENERICUNIVARIATESELECT UNIVARIATE FEATURE SELECTOR WITH CONFIGURABLE MODE  
NOTES  
TIES BETWEEN FEATURES WITH EQUAL SCORES WILL BE BROKEN IN AN UNSPECIFIED WAY  
EXAMPLES  
FROM SKLEARNDATASETS IMPORT LOADDIGITS  
FROM SKLEARNFEATURESELECTION IMPORT SELECTKBEST CHI2  
X Y LOADDIGITSRETURNXY TRUE  
XSHAPE  
1797 64  
XNEW SELECTKBESTCHI2 K20FITTRANSFORMX Y  
XNEWSHAPE  
1797 20  
1770 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

METHODS

FITSELF X Y RUN SCORE FUNCTION ON X Y AND GET THE APPROPRIATE FEATURES

FITTRANSFORM SELF X Y FIT TO DATA THEN TRANSFORM IT

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

GETSUPPORT SELF INDICES GET A MASK OR INTEGER INDEX OF THE FEATURES SELECTED

INVERSETRANSFORM SELF X REVERSE THE TRANSFORMATION OPERATION

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF X REDUCE X TO THE SELECTED FEATURES

INIT SELFSCOREFUNCFUNCTION FCLASSIF AT 0X7F3C10E93840 K10

FITSELFXY

RUN SCORE FUNCTION ON X Y AND GET THE APPROPRIATE FEATURES PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES THE TRAINING INPUT SAMPLES

YARRAYLIKE SHAPE NSAMPLES THE TARGET VALUES CLASS LABELS IN CLASSIFICATION REAL NUMBERS IN REGRESSION

RETURNS

SELF OBJECT

FITTRANSFORM SELFXYNONE FITPARAMS

FIT TO DATA THEN TRANSFORM IT

FITS TRANSFORMER TO X AND Y WITH OPTIONAL PARAMETERS FITPARAMS AND RETURNS A TRANSFORMED VERSION OF X PARAMETERS

XNUMPY ARRAY OF SHAPE NSAMPLES NFEATURES TRAINING SET

YNUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES

RETURNS

XNEW NUMPY ARRAY OF SHAPE NSAMPLES NFEATURESNEW TRANSFORMED ARRAY

GETPARAMS SELFDEEPTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

GETSUPPORT SELFINDICESFALSE

GET A MASK OR INTEGER INDEX OF THE FEATURES SELECTED

PARAMETERS

INDICES BOOLEAN DEFAULT FALSE IF TRUE THE RETURN VALUE WILL BE AN ARRAY OF INTEGERS RATHER THAN A BOOLEAN MASK

RETURNS

616SKLEARNFEATURESELECTION FEATURE SELECTION 1771

SCIKITLEARN USER GUIDE RELEASE 0213

SUPPORT ARRAY AN INDEX THAT SELECTS THE RETAINED FEATURES FROM A FEATURE VECTOR IF INDICES IS FALSE THIS IS A BOOLEAN ARRAY OF SHAPE INPUT FEATURES IN WHICH AN ELEMENT IS TRUE IFF ITS CORRESPONDING FEATURE IS SELECTED FOR RETENTION IF INDICES IS TRUE THIS IS AN INTEGER ARRAY OF SHAPE OUTPUT FEATURES WHOSE VALUES ARE INDICES INTO THE INPUT FEATURE VECTOR

INVERSETRANSFORM SELF

REVERSE THE TRANSFORMATION OPERATION

PARAMETERS

XARRAY OF SHAPE NSAMPLES NSELECTEDFEATURES THE INPUT SAMPLES

RETURNS

XR ARRAY OF SHAPE NSAMPLES NORIGINALFEATURES XWITH COLUMNS OF ZEROS INSERTED WHERE FEATURES WOULD HAVE BEEN REMOVED BY TRANSFORM

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELF

REDUCE X TO THE SELECTED FEATURES

PARAMETERS

XARRAY OF SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES

RETURNS

XR ARRAY OF SHAPE NSAMPLES NSELECTEDFEATURES THE INPUT SAMPLES WITH ONLY THE SELECTED FEATURES

EXAMPLES USING SKLEARNFEATURESELECTIONSELECTKBEST

- CONCATENATING MULTIPLE FEATURE EXTRACTION METHODS
- SELECTING DIMENSIONALITY REDUCTION WITH PIPELINE AND GRIDSEARCHCV
- PIPELINE ANOVA SVM
- CLASSIFICATION OF TEXT DOCUMENTS USING SPARSE FEATURES

6164SKLEARNFEATURESELECTION SELECTFPR

CLASSSSKLEARNFEATURESELECTION SELECTFPR SCOREFUNCFUNCTION FCLASSIF ALPHA005

FILTER SELECT THE PVALUES BELOW ALPHA BASED ON A FPR TEST

FPR TEST STANDS FOR FALSE POSITIVE RATE TEST IT CONTROLS THE TOTAL AMOUNT OF FALSE DETECTIONS

READ MORE IN THE USER GUIDE

PARAMETERS

1772 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

SCOREFUNC CALLABLE FUNCTION TAKING TWO ARRAYS X AND Y AND RETURNING A PAIR OF ARRAYS SCORES  
PVALUES DEFAULT IS FCLASSIF SEE BELOW “SEE ALSO” THE DEFAULT FUNCTION ONLY WORKS WITH  
CLASSIFICATION TASKS

ALPHA FLOAT OPTIONAL THE HIGHEST PVALUE FOR FEATURES TO BE KEPT

ATTRIBUTES

SCORES ARRAYLIKE SHAPENFEATURES SCORES OF FEATURES

PVALUES ARRAYLIKE SHAPENFEATURES PVALUES OF FEATURE SCORES

SEE ALSO

FCLASSIF ANOV A FVALUE BETWEEN LABELFEATURE FOR CLASSIFICATION TASKS

CHI2 CHISQUARED STATS OF NONNEGATIVE FEATURES FOR CLASSIFICATION TASKS

MUTUALINFOCLASSIF

FREGRESSION FVALUE BETWEEN LABELFEATURE FOR REGRESSION TASKS

MUTUALINFOREGRESSION MUTUAL INFORMATION BETWEEN FEATURES AND THE TARGET

SELECTPERCENTILE SELECT FEATURES BASED ON PERCENTILE OF THE HIGHEST SCORES

SELECTKBEST SELECT FEATURES BASED ON THE K HIGHEST SCORES

SELECTFDR SELECT FEATURES BASED ON AN ESTIMATED FALSE DISCOVERY RATE

SELECTFWE SELECT FEATURES BASED ON FAMILYWISE ERROR RATE

GENERICUNIVARIATESELECT UNIVARIATE FEATURE SELECTOR WITH CONFIGURABLE MODE

EXAMPLES

FROM SKLEARNDATASETS IMPORT LOADBREASTCANCER

FROM SKLEARNFEATURESELECTION IMPORT SELECTFPR CHI2

X Y LOADBREASTCANCERRETURNXY TRUE

XSHAPE

569 30

XNEW SELECTFPRCHI2 ALPHA001FITTRANSFORMX Y

XNEWSHAPE

569 16

METHODS

FITSELF X Y RUN SCORE FUNCTION ON X Y AND GET THE APPROPRIATE  
FEATURES

FITTRANSFORM SELF X Y FIT TO DATA THEN TRANSFORM IT

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

GETSUPPORT SELF INDICES GET A MASK OR INTEGER INDEX OF THE FEATURES SELECTED

INVERSETRANSFORM SELF X REVERSE THE TRANSFORMATION OPERATION

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF X REDUCE X TO THE SELECTED FEATURES

INIT SELFSCOREFUNCFUNCTION FCLASSIF AT 0X7EFE30BB2158 ALPHA005

616SKLEARNFEATURESELECTION FEATURE SELECTION 1773

SCIKITLEARN USER GUIDE RELEASE 0213

**FITSELFXY**  
RUN SCORE FUNCTION ON X Y AND GET THE APPROPRIATE FEATURES  
PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES THE TRAINING INPUT SAMPLES  
YARRAYLIKE SHAPE NSAMPLES THE TARGET VALUES CLASS LABELS IN CLASSIFICATION REAL NUM  
BERS IN REGRESSION  
RETURNS  
SELF OBJECT

**FITTRANSFORM SELFXYNONE FITPARAMS**  
FIT TO DATA THEN TRANSFORM IT  
FITS TRANSFORMER TO X AND Y WITH OPTIONAL PARAMETERS FITPARAMS AND RETURNS A TRANSFORMED VERSION OF X  
PARAMETERS  
XNUMPY ARRAY OF SHAPE NSAMPLES NFEATURES TRAINING SET  
YNUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES  
RETURNS  
XNEW NUMPY ARRAY OF SHAPE NSAMPLES NFEATURESNEW TRANSFORMED ARRAY

**GETPARAMS SELFDEEPTTRUE**  
GET PARAMETERS FOR THIS ESTIMATOR  
PARAMETERS  
DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBOBJECTS THAT ARE ESTIMATORS  
RETURNS  
PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

**GETSUPPORT SELFINDICESFALSE**  
GET A MASK OR INTEGER INDEX OF THE FEATURES SELECTED  
PARAMETERS  
INDICES BOOLEAN DEFAULT FALSE IF TRUE THE RETURN VALUE WILL BE AN ARRAY OF INTEGERS RATHER  
THAN A BOOLEAN MASK  
RETURNS  
SUPPORT ARRAY AN INDEX THAT SELECTS THE RETAINED FEATURES FROM A FEATURE VECTOR IF  
INDICES IS FALSE THIS IS A BOOLEAN ARRAY OF SHAPE INPUT FEATURES IN WHICH AN ELE  
MENT IS TRUE IFF ITS CORRESPONDING FEATURE IS SELECTED FOR RETENTION IF INDICES IS TRUE  
THIS IS AN INTEGER ARRAY OF SHAPE OUTPUT FEATURES WHOSE VALUES ARE INDICES INTO THE INPUT  
FEATURE VECTOR

**INVERSETRANSFORM SELFXY**  
REVERSE THE TRANSFORMATION OPERATION  
PARAMETERS  
XARRAY OF SHAPE NSAMPLES NSELECTEDFEATURES THE INPUT SAMPLES  
RETURNS

1774 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

XR ARRAY OF SHAPE NSAMPLES NORIGINALFEATURES XWITH COLUMNS OF ZEROS INSERTED  
WHERE FEATURES WOULD HAVE BEEN REMOVED BY TRANSFORM

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE  
PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT  
OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELF

REDUCE X TO THE SELECTED FEATURES

PARAMETERS

XARRAY OF SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES

RETURNS

XR ARRAY OF SHAPE NSAMPLES NSELECTEDFEATURES THE INPUT SAMPLES WITH ONLY THE SE  
LECTED FEATURES

6165SKLEARNFEATURESELECTION SELECTFDR

CLASSSSKLEARNFEATURESELECTION SELECTFDR SCOREFUNCFUNCTION FCLASSIF ALPHA005

FILTER SELECT THE PVALUES FOR AN ESTIMATED FALSE DISCOVERY RATE

THIS USES THE BENJAMINIHOCHBERG PROCEDURE ALPHA IS AN UPPER BOUND ON THE EXPECTED FALSE DISCOVERY RATE

READ MORE IN THE USER GUIDE

PARAMETERS

SCOREFUNC CALLABLE FUNCTION TAKING TWO ARRAYS X AND Y AND RETURNING A PAIR OF ARRAYS SCORES

PVALUES DEFAULT IS FCLASSIF SEE BELOW “SEE ALSO” THE DEFAULT FUNCTION ONLY WORKS WITH  
CLASSIFICATION TASKS

ALPHA FLOAT OPTIONAL THE HIGHEST UNCORRECTED PVALUE FOR FEATURES TO KEEP

ATTRIBUTES

SCORES ARRAYLIKE SHAPENFEATURES SCORES OF FEATURES

PVALUES ARRAYLIKE SHAPENFEATURES PVALUES OF FEATURE SCORES

SEE ALSO

FCLASSIF ANOV A FVALUE BETWEEN LABELFEATURE FOR CLASSIFICATION TASKS

MUTUALINFOCLASSIF MUTUAL INFORMATION FOR A DISCRETE TARGET

CHI2 CHISQUARED STATS OF NONNEGATIVE FEATURES FOR CLASSIFICATION TASKS

FREGRESSION FVALUE BETWEEN LABELFEATURE FOR REGRESSION TASKS

MUTUALINFOREGRESSION MUTUAL INFORMATION FOR A CONTNUOUS TARGET

SELECTPERCENTILE SELECT FEATURES BASED ON PERCENTILE OF THE HIGHEST SCORES

SELECTKBEST SELECT FEATURES BASED ON THE K HIGHEST SCORES

SELECTFPR SELECT FEATURES BASED ON A FALSE POSITIVE RATE TEST

6165SKLEARNFEATURESELECTION FEATURE SELECTION 1775

SCIKITLEARN USER GUIDE RELEASE 0213

SELECTFWE SELECT FEATURES BASED ON FAMILYWISE ERROR RATE

GENERICUNIVARIATESELECT UNIVARIATE FEATURE SELECTOR WITH CONFIGURABLE MODE

REFERENCES

HTTPSENWIKIPEDIAORGWIKIFALSEDISCOVERYRATE

EXAMPLES

FROM SKLEARNDATASETS IMPORT LOADBREASTCANCER

FROM SKLEARNFEATURESELECTION IMPORT SELECTFDR CHI2

X Y LOADBREASTCANCERRETURNXY TRUE

XSHAPE

569 30

XNEW SELECTFDRCHI2 ALPHA001FITTRANSFORMX Y

XNEWSHAPE

569 16

METHODS

FITSELF X Y RUN SCORE FUNCTION ON X Y AND GET THE APPROPRIATE

FEATURES

FITTRANSFORM SELF X Y FIT TO DATA THEN TRANSFORM IT

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

GETSUPPORT SELF INDICES GET A MASK OR INTEGER INDEX OF THE FEATURES SELECTED

INVERSETRANSFORM SELF X REVERSE THE TRANSFORMATION OPERATION

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF X REDUCE X TO THE SELECTED FEATURES

INIT SELFSCOREFUNCFUNCTION FCLASSIF AT 0X7EFE30BB2158 ALPHA005

FITSELFXY

RUN SCORE FUNCTION ON X Y AND GET THE APPROPRIATE FEATURES

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES THE TRAINING INPUT SAMPLES

YARRAYLIKE SHAPE NSAMPLES THE TARGET VALUES CLASS LABELS IN CLASSIFICATION REAL NUM

BERS IN REGRESSION

RETURNS

SELF OBJECT

FITTRANSFORM SELFXYNONE FITPARAMS

FIT TO DATA THEN TRANSFORM IT

FITS TRANSFORMER TO X AND Y WITH OPTIONAL PARAMETERS FITPARAMS AND RETURNS A TRANSFORMED VERSION OF X

PARAMETERS

XNUMPY ARRAY OF SHAPE NSAMPLES NFEATURES TRAINING SET

1776 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

YNUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES

RETURNS

XNEW NUMPY ARRAY OF SHAPE NSAMPLES NFEATURESNEW TRANSFORMED ARRAY

GETPARAMS SELFDEEPTTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

GETSUPPORT SELFINDICESFALSE

GET A MASK OR INTEGER INDEX OF THE FEATURES SELECTED

PARAMETERS

INDICES BOOLEAN DEFAULT FALSE IF TRUE THE RETURN VALUE WILL BE AN ARRAY OF INTEGERS RATHER

THAN A BOOLEAN MASK

RETURNS

SUPPORT ARRAY AN INDEX THAT SELECTS THE RETAINED FEATURES FROM A FEATURE VECTOR IF

INDICES IS FALSE THIS IS A BOOLEAN ARRAY OF SHAPE INPUT FEATURES IN WHICH AN ELE

MENT IS TRUE IFF ITS CORRESPONDING FEATURE IS SELECTED FOR RETENTION IF INDICES IS TRUE

THIS IS AN INTEGER ARRAY OF SHAPE OUTPUT FEATURES WHOSE VALUES ARE INDICES INTO THE INPUT

FEATURE VECTOR

INVERSETRANSFORM SELF

REVERSE THE TRANSFORMATION OPERATION

PARAMETERS

XARRAY OF SHAPE NSAMPLES NSELECTEDFEATURES THE INPUT SAMPLES

RETURNS

XR ARRAY OF SHAPE NSAMPLES NORIGINALFEATURES XWITH COLUMNS OF ZEROS INSERTED

WHERE FEATURES WOULD HAVE BEEN REMOVED BY TRANSFORM

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE

PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT

OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELF

REDUCE X TO THE SELECTED FEATURES

PARAMETERS

XARRAY OF SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES

RETURNS

616SKLEARNFEATURESELECTION FEATURE SELECTION 1777

SCIKITLEARN USER GUIDE RELEASE 0213

XR ARRAY OF SHAPE NSAMPLES NSELECTEDFEATURES THE INPUT SAMPLES WITH ONLY THE SE  
LECTED FEATURES

6166SKLEARNFEATURESELECTION SELECTFROMMODEL

CLASSSKLEARNFEATURESELECTION SELECTFROMMODEL ESTIMATOR THRESHOLDNONE PREFITFALSE

NORMORDER1 MAXFEATURESNONE

METATransformer FOR SELECTING FEATURES BASED ON IMPORTANCE WEIGHTS

NEW IN VERSION 017

PARAMETERS

ESTIMATOR OBJECT THE BASE ESTIMATOR FROM WHICH THE TRANSFORMER IS BUILT THIS CAN BE BOTH A  
FITTED IFPREFIT IS SET TO TRUE OR A NONFITTED ESTIMATOR THE ESTIMATOR MUST HAVE EITHER A  
FEATUREIMPORTANCES ORCOEF ATTRIBUTE AFTER FITTING

THRESHOLD STRING FLOAT OPTIONAL DEFAULT NONE THE THRESHOLD VALUE TO USE FOR FEATURE SELECTION  
FEATURES WHOSE IMPORTANCE IS GREATER OR EQUAL ARE KEPT WHILE THE OTHERS ARE DISCARDED IF  
“MEDIAN” RESP “MEAN” THEN THE THRESHOLD VALUE IS THE MEDIAN RESP THE MEAN OF THE  
FEATURE IMPORTANCES A SCALING FACTOR EG “125MEAN” MAY ALSO BE USED IF NONE AND  
IF THE ESTIMATOR HAS A PARAMETER PENALTY SET TO L1 EITHER EXPLICITLY OR IMPLICITLY EG LASSO  
THE THRESHOLD USED IS 1E5 OTHERWISE “MEAN” IS USED BY DEFAULT

PREFIT BOOL DEFAULT FALSE WHETHER A PREFIT MODEL IS EXPECTED TO BE PASSED INTO THE CONSTRUCTOR  
DIRECTLY OR NOT IF TRUE TRANSFORM MUST BE CALLED DIRECTLY AND SELECTFROMMODEL CAN  
NOT BE USED WITH CROSSVALSCORE GRIDSEARCHCV AND SIMILAR UTILITIES THAT CLONE  
THE ESTIMATOR OTHERWISE TRAIN THE MODEL USING FIT AND THENTRANSFORM TO DO FEATURE  
SELECTION

NORMORDER NONZERO INT INF INF DEFAULT 1 ORDER OF THE NORM USED TO FILTER THE VECTORS OF  
COEFFICIENTS BELOW THRESHOLD IN THE CASE WHERE THE COEF ATTRIBUTE OF THE ESTIMATOR IS  
OF DIMENSION 2

MAXFEATURES INT OR NONE OPTIONAL THE MAXIMUM NUMBER OF FEATURES SELECTED SCORING ABOVE  
THRESHOLD TO DISABLE THRESHOLD AND ONLY SELECT BASED ON MAXFEATURES SET

THRESHOLDNPNINF

NEW IN VERSION 020

ATTRIBUTES

ESTIMATOR AN ESTIMATOR THE BASE ESTIMATOR FROM WHICH THE TRANSFORMER IS BUILT THIS IS STORED  
ONLY WHEN A NONFITTED ESTIMATOR IS PASSED TO THE SELECTFROMMODEL IE WHEN PREFIT IS  
FALSE

THRESHOLD FLOAT THE THRESHOLD VALUE USED FOR FEATURE SELECTION

METHODS

FITSELF X Y FIT THE SELECTFROMMODEL METATransformer

FITTRANSFORM SELF X Y FIT TO DATA THEN TRANSFORM IT

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

GETSUPPORT SELF INDICES GET A MASK OR INTEGER INDEX OF THE FEATURES SELECTED

INVERSETRANSFORM SELF X REVERSE THE TRANSFORMATION OPERATION

PARTIALFIT SELF X Y FIT THE SELECTFROMMODEL METATransformer ONLY ONCE

CONTINUED ON NEXT PAGE

1778 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

TABLE 6110 – CONTINUED FROM PREVIOUS PAGE

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF X REDUCE X TO THE SELECTED FEATURES

INIT SELFESTIMATOR THRESHOLDNONE PREFITFALSE NORMORDER1 MAXFEATURESNONE

FITSELFXYNONE FITPARAMS

FIT THE SELECTFROMMODEL METATRANSFORMER

PARAMETERS

XARRAYLIKE OF SHAPE NSAMPLES NFEATURES THE TRAINING INPUT SAMPLES

YARRAYLIKE SHAPE NSAMPLES THE TARGET VALUES INTEGERS THAT CORRESPOND TO CLASSES IN CLASSIFICATION REAL NUMBERS IN REGRESSION

FITPARAMS OTHER ESTIMATOR SPECIFIC PARAMETERS

RETURNS

SELF OBJECT

FITTRANSFORM SELFXYNONE FITPARAMS

FIT TO DATA THEN TRANSFORM IT

FITS TRANSFORMER TO X AND Y WITH OPTIONAL PARAMETERS FITPARAMS AND RETURNS A TRANSFORMED VERSION OF X

PARAMETERS

XNUMPY ARRAY OF SHAPE NSAMPLES NFEATURES TRAINING SET

YNUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES

RETURNS

XNEW NUMPY ARRAY OF SHAPE NSAMPLES NFEATURESNEW TRANSFORMED ARRAY

GETPARAMS SELFDEEPTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

GETSUPPORT SELFINDICESFALSE

GET A MASK OR INTEGER INDEX OF THE FEATURES SELECTED

PARAMETERS

INDICES BOOLEAN DEFAULT FALSE IF TRUE THE RETURN VALUE WILL BE AN ARRAY OF INTEGERS RATHER THAN A BOOLEAN MASK

RETURNS

SUPPORT ARRAY AN INDEX THAT SELECTS THE RETAINED FEATURES FROM A FEATURE VECTOR IF INDICES IS FALSE THIS IS A BOOLEAN ARRAY OF SHAPE INPUT FEATURES IN WHICH AN ELEMENT IS TRUE IFF ITS CORRESPONDING FEATURE IS SELECTED FOR RETENTION IF INDICES IS TRUE THIS IS AN INTEGER ARRAY OF SHAPE OUTPUT FEATURES WHOSE VALUES ARE INDICES INTO THE INPUT FEATURE VECTOR

616SKLEARNFEATURESELECTION FEATURE SELECTION 1779

SCIKITLEARN USER GUIDE RELEASE 0213

INVERSETRANSFORM SELF  
REVERSE THE TRANSFORMATION OPERATION

PARAMETERS

XARRAY OF SHAPE NSAMPLES NSELECTEDFEATURES THE INPUT SAMPLES

RETURNS

XR ARRAY OF SHAPE NSAMPLES NORIGINALFEATURES XWITH COLUMNS OF ZEROS INSERTED  
WHERE FEATURES WOULD HAVE BEEN REMOVED BY TRANSFORM

PARTIALFIT SELFXYNONE FITPARAMS

FIT THE SELECTFROMMODEL METATransformer ONLY ONCE

PARAMETERS

XARRAYLIKE OF SHAPE NSAMPLES NFEATURES THE TRAINING INPUT SAMPLES

YARRAYLIKE SHAPE NSAMPLES THE TARGET VALUES INTEGERS THAT CORRESPOND TO CLASSES IN  
CLASSIFICATION REAL NUMBERS IN REGRESSION

FITPARAMS OTHER ESTIMATOR SPECIFIC PARAMETERS

RETURNS

SELF OBJECT

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE  
PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT  
OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELF  
REDUCE X TO THE SELECTED FEATURES

PARAMETERS

XARRAY OF SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES

RETURNS

XR ARRAY OF SHAPE NSAMPLES NSELECTEDFEATURES THE INPUT SAMPLES WITH ONLY THE SE  
LECTED FEATURES

EXAMPLES USING SKLEARNFEATURESELECTIONSELECTFROMMODEL

- FEATURE SELECTION USING SELECTFROMMODEL AND LASSOCV
- CLASSIFICATION OF TEXT DOCUMENTS USING SPARSE FEATURES

6167SKLEARNFEATURESELECTION SELECTFWE

CLASSSSKLEARNFEATURESELECTION SELECTFWE SCOREFUNCFUNCTION FCLASSIF ALPHA005

FILTER SELECT THE PVALUES CORRESPONDING TO FAMILYWISE ERROR RATE

READ MORE IN THE USER GUIDE

1780 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

PARAMETERS

SCOREFUNC CALLABLE FUNCTION TAKING TWO ARRAYS X AND Y AND RETURNING A PAIR OF ARRAYS SCORES  
PVALUES DEFAULT IS FCLASSIF SEE BELOW “SEE ALSO” THE DEFAULT FUNCTION ONLY WORKS WITH  
CLASSIFICATION TASKS

ALPHA FLOAT OPTIONAL THE HIGHEST UNCORRECTED PVALUE FOR FEATURES TO KEEP

ATTRIBUTES

SCORES ARRAYLIKE SHAPENFEATURES SCORES OF FEATURES

PVALUES ARRAYLIKE SHAPENFEATURES PVALUES OF FEATURE SCORES

SEE ALSO

FCLASSIF ANOV A FVALUE BETWEEN LABELFEATURE FOR CLASSIFICATION TASKS

CHI2 CHISQUARED STATS OF NONNEGATIVE FEATURES FOR CLASSIFICATION TASKS

FREGRESSION FVALUE BETWEEN LABELFEATURE FOR REGRESSION TASKS

SELECTPERCENTILE SELECT FEATURES BASED ON PERCENTILE OF THE HIGHEST SCORES

SELECTKBEST SELECT FEATURES BASED ON THE K HIGHEST SCORES

SELECTFPR SELECT FEATURES BASED ON A FALSE POSITIVE RATE TEST

SELECTFDR SELECT FEATURES BASED ON AN ESTIMATED FALSE DISCOVERY RATE

GENERICUNIVARIATESELECT UNIVARIATE FEATURE SELECTOR WITH CONFIGURABLE MODE

EXAMPLES

FROM SKLEARNDATASETS IMPORT LOADBREASTCANCER

FROM SKLEARNFEATURESELECTION IMPORT SELECTFWE CHI2

X Y LOADBREASTCANCERRETURNXY TRUE

XSHAPE

569 30

XNEW SELECTFWECHI2 ALPHA001FITTRANSFORMX Y

XNEWSHAPE

569 15

METHODS

FITSELF X Y RUN SCORE FUNCTION ON X Y AND GET THE APPROPRIATE  
FEATURES

FITTRANSFORM SELF X Y FIT TO DATA THEN TRANSFORM IT

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

GETSUPPORT SELF INDICES GET A MASK OR INTEGER INDEX OF THE FEATURES SELECTED

INVERSETRANSFORM SELF X REVERSE THE TRANSFORMATION OPERATION

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF X REDUCE X TO THE SELECTED FEATURES

INIT SELFSCOREFUNCFUNCTION FCLASSIF AT 0X7EFE30BB2158 ALPHA005

FITSELFXY

616SKLEARNFEATURESELECTION FEATURE SELECTION 1781

SCIKITLEARN USER GUIDE RELEASE 0213

RUN SCORE FUNCTION ON X Y AND GET THE APPROPRIATE FEATURES

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES THE TRAINING INPUT SAMPLES

YARRAYLIKE SHAPE NSAMPLES THE TARGET VALUES CLASS LABELS IN CLASSIFICATION REAL NUMBERS IN REGRESSION

RETURNS

SELF OBJECT

FITTRANSFORM SELFXYNONE FITPARAMS

FIT TO DATA THEN TRANSFORM IT

FITS TRANSFORMER TO X AND Y WITH OPTIONAL PARAMETERS FITPARAMS AND RETURNS A TRANSFORMED VERSION OF X

PARAMETERS

XNUMPY ARRAY OF SHAPE NSAMPLES NFEATURES TRAINING SET

YNUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES

RETURNS

XNEW NUMPY ARRAY OF SHAPE NSAMPLES NFEATURESNEW TRANSFORMED ARRAY

GETPARAMS SELFDEEPTTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

GETSUPPORT SELFINDICESFALSE

GET A MASK OR INTEGER INDEX OF THE FEATURES SELECTED

PARAMETERS

INDICES BOOLEAN DEFAULT FALSE IF TRUE THE RETURN VALUE WILL BE AN ARRAY OF INTEGERS RATHER THAN A BOOLEAN MASK

RETURNS

SUPPORT ARRAY AN INDEX THAT SELECTS THE RETAINED FEATURES FROM A FEATURE VECTOR IF INDICES IS FALSE THIS IS A BOOLEAN ARRAY OF SHAPE INPUT FEATURES IN WHICH AN ELEMENT IS TRUE IFF ITS CORRESPONDING FEATURE IS SELECTED FOR RETENTION IF INDICES IS TRUE THIS IS AN INTEGER ARRAY OF SHAPE OUTPUT FEATURES WHOSE VALUES ARE INDICES INTO THE INPUT FEATURE VECTOR

INVERSETRANSFORM SELFXY

REVERSE THE TRANSFORMATION OPERATION

PARAMETERS

XARRAY OF SHAPE NSAMPLES NSELECTEDFEATURES THE INPUT SAMPLES

RETURNS

XR ARRAY OF SHAPE NSAMPLES NORIGINALFEATURES XWITH COLUMNS OF ZEROS INSERTED WHERE FEATURES WOULD HAVE BEEN REMOVED BY TRANSFORM

1782 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELF X

REDUCE X TO THE SELECTED FEATURES

PARAMETERS

X ARRAY OF SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES

RETURNS

XR ARRAY OF SHAPE NSAMPLES NSELECTEDFEATURES THE INPUT SAMPLES WITH ONLY THE SELECTED FEATURES

6168SKLEARNFEATURESELECTION RFE

CLASSSKLEARNFEATURESELECTION RFEESTIMATOR NFEATURESTOSELECTNONE STEP1 VERBOSE0

FEATURE RANKING WITH RECURSIVE FEATURE ELIMINATION

GIVEN AN EXTERNAL ESTIMATOR THAT ASSIGNS WEIGHTS TO FEATURES EG THE COEFFICIENTS OF A LINEAR MODEL THE GOAL OF RECURSIVE FEATURE ELIMINATION RFE IS TO SELECT FEATURES BY RECURSIVELY CONSIDERING SMALLER AND SMALLER SETS OF FEATURES FIRST THE ESTIMATOR IS TRAINED ON THE INITIAL SET OF FEATURES AND THE IMPORTANCE OF EACH FEATURE IS OBTAINED EITHER THROUGH A COEF ATTRIBUTE OR THROUGH A FEATUREIMPORTANCES ATTRIBUTE THEN THE LEAST IMPORTANT FEATURES ARE PRUNED FROM CURRENT SET OF FEATURES THAT PROCEDURE IS RECURSIVELY REPEATED ON THE PRUNED SET UNTIL THE DESIRED NUMBER OF FEATURES TO SELECT IS EVENTUALLY REACHED

READ MORE IN THE USER GUIDE

PARAMETERS

ESTIMATOR OBJECT A SUPERVISED LEARNING ESTIMATOR WITH A FIT METHOD THAT PROVIDES INFORMATION ABOUT FEATURE IMPORTANCE EITHER THROUGH A COEF ATTRIBUTE OR THROUGH A FEATUREIMPORTANCES ATTRIBUTE

NFEATURESTOSELECT INT OR NONE DEFAULTNONE THE NUMBER OF FEATURES TO SELECT IF NONE HALF OF THE FEATURES ARE SELECTED

STEP INT OR FLOAT OPTIONAL DEFAULT1 IF GREATER THAN OR EQUAL TO 1 THEN STEP CORRESPONDS TO THE INTEGER NUMBER OF FEATURES TO REMOVE AT EACH ITERATION IF WITHIN 00 10 THEN STEP CORRESPONDS TO THE PERCENTAGE ROUNDED DOWN OF FEATURES TO REMOVE AT EACH ITERATION

VERBOSE INT DEFAULT0 CONTROLS VERBOSITY OF OUTPUT

ATTRIBUTES

NFEATURES INT THE NUMBER OF SELECTED FEATURES

SUPPORT ARRAY OF SHAPE NFEATURES THE MASK OF SELECTED FEATURES

RANKING ARRAY OF SHAPE NFEATURES THE FEATURE RANKING SUCH THAT RANKINGI CORRESPONDS TO THE RANKING POSITION OF THE ITH FEATURE SELECTED IE ESTIMATED BEST FEATURES ARE ASSIGNED RANK 1

ESTIMATOR OBJECT THE EXTERNAL ESTIMATOR FIT ON THE REDUCED DATASET

616SKLEARNFEATURESELECTION FEATURE SELECTION 1783

SCIKITLEARN USER GUIDE RELEASE 0213

SEE ALSO

RFE CV RECURSIVE FEATURE ELIMINATION WITH BUILT IN CROSSVALIDATED SELECTION OF THE BEST NUMBER OF FEATURES

REFERENCES

RE310F679C81E1

EXAMPLES

THE FOLLOWING EXAMPLE SHOWS HOW TO RETRIEVE THE 5 RIGHT INFORMATIVE FEATURES IN THE FRIEDMAN 1 DATASET

```
FROM SKLEARN DATASETS IMPORT MAKEFRIEDMAN1
FROM SKLEARN FEATURE SELECTION IMPORT RFE
FROM SKLEARN SVM IMPORT SVR
X Y MAKEFRIEDMAN1 NSAMPLES50 NFEATURES10 RANDOMSTATE0
ESTIMATOR SVRKERNELLINEAR
SELECTOR RFEESTIMATOR 5 STEP1
SELECTOR SELECTORFITX Y
SELECTORSUPPORT
ARRAY TRUE TRUE TRUE TRUE TRUE FALSE FALSE FALSE FALSE
FALSE
SELECTORRANKING
ARRAY1 1 1 1 1 6 4 3 2 5
```

METHODS

DECISIONFUNCTION SELF X COMPUTE THE DECISION FUNCTION OF X

FITSELF X Y FIT THE RFE MODEL AND THEN THE UNDERLYING ESTIMATOR ON THE SELECTED FEATURES

FITTRANSFORM SELF X Y FIT TO DATA THEN TRANSFORM IT

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

GETSUPPORT SELF INDICES GET A MASK OR INTEGER INDEX OF THE FEATURES SELECTED

INVERSETRANSFORM SELF X REVERSE THE TRANSFORMATION OPERATION

PREDICT SELF X REDUCE X TO THE SELECTED FEATURES AND THEN PREDICT USING THE UNDERLYING ESTIMATOR

PREDICTLOGPROBA SELF X PREDICT CLASS LOGPROBABILITIES FOR X

PREDICTPROBA SELF X PREDICT CLASS PROBABILITIES FOR X

SCORE SELF X Y REDUCE X TO THE SELECTED FEATURES AND THEN RETURN THE SCORE OF THE UNDERLYING ESTIMATOR

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF X REDUCE X TO THE SELECTED FEATURES

INIT SELF ESTIMATOR NFEATURES TO SELECT NONE STEP1 VERBOSE0

DECISIONFUNCTION SELF X

COMPUTE THE DECISION FUNCTION OF X

PARAMETERS

X ARRAY LIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTERNALLY IT WILL BE CONVERTED TO DTYPENP FLOAT32 AND IF A SPARSE MATRIX IS PROVIDED TO

1784 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

A SPARSECSRMATRIX

RETURNS

SCORE ARRAY SHAPE NSAMPLES NCLASSES OR NSAMPLES THE DECISION FUNCTION OF THE INPUT SAMPLES THE ORDER OF THE CLASSES CORRESPONDS TO THAT IN THE ATTRIBUTE CLASSES REGRESSION AND BINARY CLASSIFICATION PRODUCE AN ARRAY OF SHAPE NSAMPLES

FITSELFXY

FIT THE RFE MODEL AND THEN THE UNDERLYING ESTIMATOR ON THE SELECTED FEATURES

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE TRAINING INPUT SAMPLES

YARRAYLIKE SHAPE NSAMPLES THE TARGET VALUES

FITTRANSFORM SELFXYNONE FITPARAMS

FIT TO DATA THEN TRANSFORM IT

FITS TRANSFORMER TO X AND Y WITH OPTIONAL PARAMETERS FITPARAMS AND RETURNS A TRANSFORMED VERSION OF X

PARAMETERS

XNUMPY ARRAY OF SHAPE NSAMPLES NFEATURES TRAINING SET

YNUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES

RETURNS

XNEW NUMPY ARRAY OF SHAPE NSAMPLES NFEATURESNEW TRANSFORMED ARRAY

GETPARAMS SELFDEEPTTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

GETSUPPORT SELFINDICESFALSE

GET A MASK OR INTEGER INDEX OF THE FEATURES SELECTED

PARAMETERS

INDICES BOOLEAN DEFAULT FALSE IF TRUE THE RETURN VALUE WILL BE AN ARRAY OF INTEGERS RATHER THAN A BOOLEAN MASK

RETURNS

SUPPORT ARRAY AN INDEX THAT SELECTS THE RETAINED FEATURES FROM A FEATURE VECTOR IF INDICES IS FALSE THIS IS A BOOLEAN ARRAY OF SHAPE INPUT FEATURES IN WHICH AN ELEMENT IS TRUE IFF ITS CORRESPONDING FEATURE IS SELECTED FOR RETENTION IF INDICES IS TRUE THIS IS AN INTEGER ARRAY OF SHAPE OUTPUT FEATURES WHOSE VALUES ARE INDICES INTO THE INPUT FEATURE VECTOR

INVERSETRANSFORM SELFXY

REVERSE THE TRANSFORMATION OPERATION

PARAMETERS

616SKLEARNFEATURESELECTION FEATURE SELECTION 1785

SCIKITLEARN USER GUIDE RELEASE 0213

XARRAY OF SHAPE NSAMPLES NSELECTEDFEATURES THE INPUT SAMPLES

RETURNS

XR ARRAY OF SHAPE NSAMPLES NORIGINALFEATURES XWITH COLUMNS OF ZEROS INSERTED WHERE FEATURES WOULD HAVE BEEN REMOVED BY TRANSFORM

PREDICTSELF X

REDUCE X TO THE SELECTED FEATURES AND THEN PREDICT USING THE UNDERLYING ESTIMATOR

PARAMETERS

XARRAY OF SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES

RETURNS

YARRAY OF SHAPE NSAMPLES THE PREDICTED TARGET VALUES

PREDICTLOGPROBA SELF X

PREDICT CLASS LOGPROBABILITIES FOR X

PARAMETERS

XARRAY OF SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES

RETURNS

PARRAY OF SHAPE NSAMPLES NCLASSES THE CLASS LOGPROBABILITIES OF THE INPUT SAMPLES THE ORDER OF THE CLASSES CORRESPONDS TO THAT IN THE ATTRIBUTE CLASSES

PREDICTPROBA SELF X

PREDICT CLASS PROBABILITIES FOR X

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTER NALLY IT WILL BE CONVERTED TO DTYPENPFLOAT32 AND IF A SPARSE MATRIX IS PROVIDED TO A SPARSECSRMATRIX

RETURNS

PARRAY OF SHAPE NSAMPLES NCLASSES THE CLASS PROBABILITIES OF THE INPUT SAMPLES THE ORDER OF THE CLASSES CORRESPONDS TO THAT IN THE ATTRIBUTE CLASSES

SCORESELF X

REDUCE X TO THE SELECTED FEATURES AND THEN RETURN THE SCORE OF THE UNDERLYING ESTIMATOR

PARAMETERS

XARRAY OF SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES

YARRAY OF SHAPE NSAMPLES THE TARGET VALUES

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

1786 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

TRANSFORM SELF

REDUCE X TO THE SELECTED FEATURES

PARAMETERS

XARRAY OF SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES

RETURNS

XR ARRAY OF SHAPE NSAMPLES NSELECTEDFEATURES THE INPUT SAMPLES WITH ONLY THE SE  
LECTED FEATURES

EXAMPLES USING SKLEARNFEATURESELECTIONRFE

- RECURSIVE FEATURE ELIMINATION

6169SKLEARNFEATURESELECTION RFECV

CLASSSSKLEARNFEATURESELECTION RFECVESTIMATOR STEP1 MINFEATURESTOSELECT1 CV'WARN'

SCORINGNONE VERBOSE0 NJOBSNONE

FEATURE RANKING WITH RECURSIVE FEATURE ELIMINATION AND CROSSVALIDATED SELECTION OF THE BEST NUMBER OF FEATURES

SEE GLOSSARY ENTRY FOR CROSSVALIDATION ESTIMATOR

READ MORE IN THE USER GUIDE

PARAMETERS

ESTIMATOR OBJECT A SUPERVISED LEARNING ESTIMATOR WITH A FIT METHOD THAT PROVIDES IN  
FORMATION ABOUT FEATURE IMPORTANCE EITHER THROUGH A COEF ATTRIBUTE OR THROUGH A  
FEATUREIMPORTANCES ATTRIBUTE

STEP INT OR FLOAT OPTIONAL DEFAULT1 IF GREATER THAN OR EQUAL TO 1 THEN STEP CORRESPONDS  
TO THE INTEGER NUMBER OF FEATURES TO REMOVE AT EACH ITERATION IF WITHIN 00 10 THEN  
STEP CORRESPONDS TO THE PERCENTAGE ROUNDED DOWN OF FEATURES TO REMOVE AT EACH ITERA  
TION NOTE THAT THE LAST ITERATION MAY REMOVE FEWER THAN STEP FEATURES IN ORDER TO REACH  
MINFEATURESTOSELECT

MINFEATURESTOSELECT INT DEFAULT1 THE MINIMUM NUMBER OF FEATURES TO BE SELECTED  
THIS NUMBER OF FEATURES WILL ALWAYS BE SCORED EVEN IF THE DIFFERENCE BETWEEN THE ORIGI  
NAL FEATURE COUNT AND MINFEATURESTOSELECT ISN'T DIVISIBLE BY STEP

CVINT CROSSVALIDATION GENERATOR OR AN ITERABLE OPTIONAL DETERMINES THE CROSSVALIDATION SPLIT  
TING STRATEGY POSSIBLE INPUTS FOR CV ARE

- NONE TO USE THE DEFAULT 3FOLD CROSSVALIDATION
- INTEGER TO SPECIFY THE NUMBER OF FOLDS
- CV SPLITTER
- AN ITERABLE YIELDING TRAIN TEST SPLITS AS ARRAYS OF INDICES

FOR INTEGERNONE INPUTS IF YIS BINARY OR MULTICLASS SKLEARNMODELSELECTION  
STRATIFIEDKFOLD IS USED IF THE ESTIMATOR IS A CLASSIFIER OR IF YIS NEITHER BINARY NOR  
MULTICLASS SKLEARNMODELSELECTIONKFOLD IS USED

REFER USER GUIDE FOR THE VARIOUS CROSSVALIDATION STRATEGIES THAT CAN BE USED HERE

CHANGED IN VERSION 020 CVDEFAULT VALUE OF NONE WILL CHANGE FROM 3FOLD TO 5FOLD IN  
V022

616SKLEARNFEATURESELECTION FEATURE SELECTION 1787

SCIKITLEARN USER GUIDE RELEASE 0213

SCORING STRING CALLABLE OR NONE OPTIONAL DEFAULTNONE A STRING SEE MODEL EVALUATION DOCUMENTATION OR A SCORER CALLABLE OBJECT FUNCTION WITH SIGNATURE SCORERESTIMATOR

X Y

VERBOSE INT DEFAULT0 CONTROLS VERBOSITY OF OUTPUT

NJOBS INT OR NONE OPTIONAL DEFAULTNONE NUMBER OF CORES TO RUN IN PARALLEL WHILE FITTING

ACROSS FOLDS NONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1

MEANS USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

ATTRIBUTES

NFEATURES INT THE NUMBER OF SELECTED FEATURES WITH CROSSVALIDATION

SUPPORT ARRAY OF SHAPE NFEATURES THE MASK OF SELECTED FEATURES

RANKING ARRAY OF SHAPE NFEATURES THE FEATURE RANKING SUCH THAT RANKINGI CORRESPONDS TO THE RANKING POSITION OF THE ITH FEATURE SELECTED IE ESTIMATED BEST FEATURES ARE ASSIGNED RANK 1

GRIDSCORES ARRAY OF SHAPE NSUBSETSOFFEATURES THE CROSSVALIDATION SCORES SUCH THAT GRIDSCORESI CORRESPONDS TO THE CV SCORE OF THE ITH SUBSET OF FEATURES

ESTIMATOR OBJECT THE EXTERNAL ESTIMATOR FIT ON THE REDUCED DATASET

SEE ALSO

RFE RECURSIVE FEATURE ELIMINATION

NOTES

THE SIZE OF GRIDSCORES IS EQUAL TO CEILNFEATURES MINFEATURESTOSELECT

STEP 1 WHERE STEP IS THE NUMBER OF FEATURES REMOVED AT EACH ITERATION

REFERENCES

R6F4D61CEB4111

EXAMPLES

THE FOLLOWING EXAMPLE SHOWS HOW TO RETRIEVE THE APRIORI NOT KNOWN 5 INFORMATIVE FEATURES IN THE FRIEDMAN 1 DATASET

```
FROM SKLEARNDATASETS IMPORT MAKEFRIEDMAN1
FROM SKLEARNFEATURESELECTION IMPORT RFECV
FROM SKLEARN SVM IMPORT SVR
X Y MAKEFRIEDMAN1NSAMPLES50 NFEATURES10 RANDOMSTATE0
ESTIMATOR SVRKERNELLINEAR
SELECTOR RFECVESTIMATOR STEP1 CV5
SELECTOR SELECTORFITX Y
SELECTORSUPPORT
ARRAY TRUE TRUE TRUE TRUE TRUE FALSE FALSE FALSE FALSE
FALSE
SELECTORRANKING
ARRAY1 1 1 1 1 6 4 3 2 5
```

1788 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

METHODS

DECISIONFUNCTION SELF X COMPUTE THE DECISION FUNCTION OF X

FITSELF X Y GROUPS FIT THE RFE MODEL AND AUTOMATICALLY TUNE THE NUMBER OF SELECTED FEATURES

FITTRANSFORM SELF X Y FIT TO DATA THEN TRANSFORM IT

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

GETSUPPORT SELF INDICES GET A MASK OR INTEGER INDEX OF THE FEATURES SELECTED

INVERSETRANSFORM SELF X REVERSE THE TRANSFORMATION OPERATION

PREDICT SELF X REDUCE X TO THE SELECTED FEATURES AND THEN PREDICT USING THE UNDERLYING ESTIMATOR

PREDICTLOGPROBA SELF X PREDICT CLASS LOGPROBABILITIES FOR X

PREDICTPROBA SELF X PREDICT CLASS PROBABILITIES FOR X

SCORE SELF X Y REDUCE X TO THE SELECTED FEATURES AND THEN RETURN THE SCORE OF THE UNDERLYING ESTIMATOR

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF X REDUCE X TO THE SELECTED FEATURES

INIT SELFESTIMATOR STEP1 MINFEATURESTOSELECT1 CV'WARN' SCORINGNONE VERBOSE0 NJOBSNONE

DECISIONFUNCTION SELF X

COMPUTE THE DECISION FUNCTION OF X

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTERNALLY IT WILL BE CONVERTED TO DTYPENPFLOAT32 AND IF A SPARSE MATRIX IS PROVIDED TO A SPARSECSRMATRIX

RETURNS

SCORE ARRAY SHAPE NSAMPLES NCLASSES OR NSAMPLES THE DECISION FUNCTION OF THE INPUT SAMPLES THE ORDER OF THE CLASSES CORRESPONDS TO THAT IN THE ATTRIBUTE CLASSES REGRESSION AND BINARY CLASSIFICATION PRODUCE AN ARRAY OF SHAPE NSAMPLES

FITSELFXYGROUPSNONE

FIT THE RFE MODEL AND AUTOMATICALLY TUNE THE NUMBER OF SELECTED FEATURES

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES TRAINING VECTOR WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE TOTAL NUMBER OF FEATURES

YARRAYLIKE SHAPE NSAMPLES TARGET VALUES INTEGERS FOR CLASSIFICATION REAL NUMBERS FOR REGRESSION

GROUPS ARRAYLIKE SHAPE NSAMPLES OPTIONAL GROUP LABELS FOR THE SAMPLES USED WHILE SPLITTING THE DATASET INTO TRAINTEST SET ONLY USED IN CONJUNCTION WITH A "GROUP" CVINSTANCE

EGGROUPOKFOLD

FITTRANSFORM SELFXYNONE FITPARAMS

FIT TO DATA THEN TRANSFORM IT

FITS TRANSFORMER TO X AND Y WITH OPTIONAL PARAMETERS FITPARAMS AND RETURNS A TRANSFORMED VERSION OF X

PARAMETERS

616SKLEARNFEATURESELECTION FEATURE SELECTION 1789

SCIKITLEARN USER GUIDE RELEASE 0213

XNUMPY ARRAY OF SHAPE NSAMPLES NFEATURES TRAINING SET

YNUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES

RETURNS

XNEW NUMPY ARRAY OF SHAPE NSAMPLES NFEATURESNEW TRANSFORMED ARRAY

GETPARAMS SELFDEEPTTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

GETSUPPORT SELFINDICESFALSE

GET A MASK OR INTEGER INDEX OF THE FEATURES SELECTED

PARAMETERS

INDICES BOOLEAN DEFAULT FALSE IF TRUE THE RETURN VALUE WILL BE AN ARRAY OF INTEGERS RATHER

THAN A BOOLEAN MASK

RETURNS

SUPPORT ARRAY AN INDEX THAT SELECTS THE RETAINED FEATURES FROM A FEATURE VECTOR IF

INDICES IS FALSE THIS IS A BOOLEAN ARRAY OF SHAPE INPUT FEATURES IN WHICH AN ELE

MENT IS TRUE IFF ITS CORRESPONDING FEATURE IS SELECTED FOR RETENTION IF INDICES IS TRUE

THIS IS AN INTEGER ARRAY OF SHAPE OUTPUT FEATURES WHOSE VALUES ARE INDICES INTO THE INPUT

FEATURE VECTOR

INVERSETRANSFORM SELF

REVERSE THE TRANSFORMATION OPERATION

PARAMETERS

XARRAY OF SHAPE NSAMPLES NSELECTEDFEATURES THE INPUT SAMPLES

RETURNS

XR ARRAY OF SHAPE NSAMPLES NORIGINALFEATURES XWITH COLUMNS OF ZEROS INSERTED

WHERE FEATURES WOULD HAVE BEEN REMOVED BY TRANSFORM

PREDICTSELF

REDUCE X TO THE SELECTED FEATURES AND THEN PREDICT USING THE UNDERLYING ESTIMATOR

PARAMETERS

XARRAY OF SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES

RETURNS

YARRAY OF SHAPE NSAMPLES THE PREDICTED TARGET VALUES

PREDICTLOGPROBA SELF

PREDICT CLASS LOGPROBABILITIES FOR X

PARAMETERS

XARRAY OF SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES

1790 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

RETURNS

PARRAY OF SHAPE NSAMPLES NCLASSES THE CLASS LOGPROBABILITIES OF THE INPUT SAMPLES  
THE ORDER OF THE CLASSES CORRESPONDS TO THAT IN THE ATTRIBUTE CLASSES

PREDICTPROBA SELF X

PREDICT CLASS PROBABILITIES FOR X

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTER  
NALLY IT WILL BE CONVERTED TO DTYPENPFLOAT32 AND IF A SPARSE MATRIX IS PROVIDED TO  
A SPARSECSRMATRIX

RETURNS

PARRAY OF SHAPE NSAMPLES NCLASSES THE CLASS PROBABILITIES OF THE INPUT SAMPLES THE  
ORDER OF THE CLASSES CORRESPONDS TO THAT IN THE ATTRIBUTE CLASSES

SCORESELF X

REDUCE X TO THE SELECTED FEATURES AND THEN RETURN THE SCORE OF THE UNDERLYING ESTIMATOR

PARAMETERS

XARRAY OF SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES

YARRAY OF SHAPE NSAMPLES THE TARGET VALUES

SETPARAMS SELF PARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE  
PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT  
OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELF X

REDUCE X TO THE SELECTED FEATURES

PARAMETERS

XARRAY OF SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES

RETURNS

XR ARRAY OF SHAPE NSAMPLES NSELECTEDFEATURES THE INPUT SAMPLES WITH ONLY THE SE  
LECTED FEATURES

EXAMPLES USING SKLEARNFEATURESELECTIONRFE CV

- RECURSIVE FEATURE ELIMINATION WITH CROSSVALIDATION

61610 SKLEARNFEATURESELECTION VARIANCETHRESHOLD

CLASS SKLEARNFEATURESELECTION VARIANCETHRESHOLD THRESHOLD 00

FEATURE SELECTOR THAT REMOVES ALL LOWVARIANCE FEATURES

616 SKLEARNFEATURESELECTION FEATURE SELECTION 1791

SCIKITLEARN USER GUIDE RELEASE 0213

THIS FEATURE SELECTION ALGORITHM LOOKS ONLY AT THE FEATURES X NOT THE DESIRED OUTPUTS Y AND CAN THUS BE USED FOR UNSUPERVISED LEARNING

READ MORE IN THE USER GUIDE

PARAMETERS

THRESHOLD FLOAT OPTIONAL FEATURES WITH A TRAININGSET VARIANCE LOWER THAN THIS THRESHOLD WILL BE REMOVED THE DEFAULT IS TO KEEP ALL FEATURES WITH NONZERO VARIANCE IE REMOVE THE FEATURES THAT HAVE THE SAME VALUE IN ALL SAMPLES

ATTRIBUTES

VARIANCES ARRAY SHAPE NFEATURES VARIANCES OF INDIVIDUAL FEATURES

EXAMPLES

THE FOLLOWING DATASET HAS INTEGER FEATURES TWO OF WHICH ARE THE SAME IN EVERY SAMPLE THESE ARE REMOVED WITH THE DEFAULT SETTING FOR THRESHOLD

X 0 2 0 3 0 1 4 3 0 1 1 3

SELECTOR VARIANCETHRESHOLD

SELECTORFITTRANSFORMX

ARRAY2 0

1 4

1 1

METHODS

FITSELF X Y LEARN EMPIRICAL VARIANCES FROM X

FITTRANSFORM SELF X Y FIT TO DATA THEN TRANSFORM IT

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

GETSUPPORT SELF INDICES GET A MASK OR INTEGER INDEX OF THE FEATURES SELECTED

INVERSETRANSFORM SELF X REVERSE THE TRANSFORMATION OPERATION

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF X REDUCE X TO THE SELECTED FEATURES

INIT SELFTHRESHOLD00

FITSELFXYNONE

LEARN EMPIRICAL VARIANCES FROM X

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES SAMPLE VECTORS FROM WHICH TO COMPUTE VARIANCES

YANY IGNORED THIS PARAMETER EXISTS ONLY FOR COMPATIBILITY WITH SKLEARNPIPELINEPIPELINE

RETURNS

SELF

FITTRANSFORM SELFXYNONE FITPARAMS

FIT TO DATA THEN TRANSFORM IT

FITS TRANSFORMER TO X AND Y WITH OPTIONAL PARAMETERS FITPARAMS AND RETURNS A TRANSFORMED VERSION OF X

1792 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

PARAMETERS

XNUMPY ARRAY OF SHAPE NSAMPLES NFEATURES TRAINING SET

YNUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES

RETURNS

XNEW NUMPY ARRAY OF SHAPE NSAMPLES NFEATURESNEW TRANSFORMED ARRAY

GETPARAMS SELFDEEPTTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

GETSUPPORT SELFINDICESFALSE

GET A MASK OR INTEGER INDEX OF THE FEATURES SELECTED

PARAMETERS

INDICES BOOLEAN DEFAULT FALSE IF TRUE THE RETURN VALUE WILL BE AN ARRAY OF INTEGERS RATHER

THAN A BOOLEAN MASK

RETURNS

SUPPORT ARRAY AN INDEX THAT SELECTS THE RETAINED FEATURES FROM A FEATURE VECTOR IF

INDICES IS FALSE THIS IS A BOOLEAN ARRAY OF SHAPE INPUT FEATURES IN WHICH AN ELE

MENT IS TRUE IFF ITS CORRESPONDING FEATURE IS SELECTED FOR RETENTION IF INDICES IS TRUE

THIS IS AN INTEGER ARRAY OF SHAPE OUTPUT FEATURES WHOSE VALUES ARE INDICES INTO THE INPUT

FEATURE VECTOR

INVERSETRANSFORM SELF

REVERSE THE TRANSFORMATION OPERATION

PARAMETERS

XARRAY OF SHAPE NSAMPLES NSELECTEDFEATURES THE INPUT SAMPLES

RETURNS

XR ARRAY OF SHAPE NSAMPLES NORIGINALFEATURES XWITH COLUMNS OF ZEROS INSERTED

WHERE FEATURES WOULD HAVE BEEN REMOVED BY TRANSFORM

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE

PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT

OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELF

REDUCE X TO THE SELECTED FEATURES

PARAMETERS

XARRAY OF SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES

616SKLEARNFEATURESELECTION FEATURE SELECTION 1793

SCIKITLEARN USER GUIDE RELEASE 0213

RETURNS

XR ARRAY OF SHAPE NSAMPLES NSELECTEDFEATURES THE INPUT SAMPLES WITH ONLY THE SE  
LECTED FEATURES

FEATURESELECTIONCHI2 X Y COMPUTE CHISQUARED STATS BETWEEN EACH NONNEGATIVE FEA  
TURE AND CLASS

FEATURESELECTIONFCLASSIF X Y COMPUTE THE ANOV A FVALUE FOR THE PROVIDED SAMPLE

FEATURESELECTIONFREGRESSION X Y CEN

TERUNIVARIATE LINEAR REGRESSION TESTS

FEATURESELECTION

MUTUALINFOCLASSIF X YESTIMATE MUTUAL INFORMATION FOR A DISCRETE TARGET VARIABLE

FEATURESELECTION

MUTUALINFOREGRESSION X YESTIMATE MUTUAL INFORMATION FOR A CONTINUOUS TARGET VARI  
ABLE

61611SKLEARNFEATURESELECTION CHI2

SKLEARNFEATURESELECTION CHI2XY

COMPUTE CHISQUARED STATS BETWEEN EACH NONNEGATIVE FEATURE AND CLASS

THIS SCORE CAN BE USED TO SELECT THE NFEATURES FEATURES WITH THE HIGHEST VALUES FOR THE TEST CHISQUARED STATISTIC  
FROM X WHICH MUST CONTAIN ONLY NONNEGATIVE FEATURES SUCH AS BOOLEANS OR FREQUENCIES EG TERM COUNTS IN  
DOCUMENT CLASSIFICATION RELATIVE TO THE CLASSES

RECALL THAT THE CHISQUARE TEST MEASURES DEPENDENCE BETWEEN STOCHASTIC VARIABLES SO USING THIS FUNCTION “WEEDS  
OUT” THE FEATURES THAT ARE THE MOST LIKELY TO BE INDEPENDENT OF CLASS AND THEREFORE IRRELEVANT FOR CLASSIFICATION  
READ MORE IN THE USER GUIDE

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURESIN SAMPLE VECTORS

YARRAYLIKE SHAPE NSAMPLES TARGET VECTOR CLASS LABELS

RETURNS

CHI2 ARRAY SHAPE NFEATURES CHI2 STATISTICS OF EACH FEATURE

PVAL ARRAY SHAPE NFEATURES PVALUES OF EACH FEATURE

SEE ALSO

FCLASSIF ANOV A FVALUE BETWEEN LABELFEATURE FOR CLASSIFICATION TASKS

FREGRESSION FVALUE BETWEEN LABELFEATURE FOR REGRESSION TASKS

NOTES

COMPLEXITY OF THIS ALGORITHM IS ONCLASSES NFEATURES

EXAMPLES USING SKLEARNFEATURESELECTIONCHI2

- SELECTING DIMENSIONALITY REDUCTION WITH PIPELINE AND GRIDSEARCHCV
- SVMANOVA SVM WITH UNIVARIATE FEATURE SELECTION
- CLASSIFICATION OF TEXT DOCUMENTS USING SPARSE FEATURES

1794 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213  
61612SKLEARNFEATURESELECTION FCLASSIF  
SKLEARNFEATURESELECTION FCLASSIF XY  
COMPUTE THE ANOV A FVALUE FOR THE PROVIDED SAMPLE  
READ MORE IN THE USER GUIDE  
PARAMETERS  
XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE SET OF REGRESSORS THAT WILL  
BE TESTED SEQUENTIALLY  
YARRAY OF SHAPENSAMPLES THE DATA MATRIX  
RETURNS  
FARRAY SHAPE NFEATURES THE SET OF F VALUES  
PVAL ARRAY SHAPE NFEATURES THE SET OF PVALUES  
SEE ALSO  
CHI2 CHISQUARED STATS OF NONNEGATIVE FEATURES FOR CLASSIFICATION TASKS  
FREGRESSION FVALUE BETWEEN LABELFEATURE FOR REGRESSION TASKS  
EXAMPLES USING SKLEARNFEATURESELECTIONFCLASSIF  
•UNIVARIATE FEATURE SELECTION  
61613SKLEARNFEATURESELECTION FREGRESSION  
SKLEARNFEATURESELECTION FREGRESSION XYCENTERTRUE  
UNIVARIATE LINEAR REGRESSION TESTS  
LINEAR MODEL FOR TESTING THE INDIVIDUAL EFFECT OF EACH OF MANY REGRESSORS THIS IS A SCORING FUNCTION TO BE USED  
IN A FEATURE SELECTION PROCEDURE NOT A FREE STANDING FEATURE SELECTION PROCEDURE  
THIS IS DONE IN 2 STEPS  
1 THE CORRELATION BETWEEN EACH REGRESSOR AND THE TARGET IS COMPUTED THAT IS  $X_i$   $MEANX_i$   $Y$   
 $MEANY$   $STDX_i$   $STDY$   
2 IT IS CONVERTED TO AN F SCORE THEN TO A PVALUE  
FOR MORE ON USAGE SEE THE USER GUIDE  
PARAMETERS  
XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE SET OF REGRESSORS THAT WILL  
BE TESTED SEQUENTIALLY  
YARRAY OF SHAPENSAMPLES THE DATA MATRIX  
CENTER TRUE BOOL IF TRUE X AND Y WILL BE CENTERED  
RETURNS  
FARRAY SHAPENFEATURES F VALUES OF FEATURES  
PVAL ARRAY SHAPENFEATURES PVALUES OF FSCORES  
SEE ALSO  
616SKLEARNFEATURESELECTION FEATURE SELECTION 1795

SCIKITLEARN USER GUIDE RELEASE 0213

MUTUALINFOREGRESSION MUTUAL INFORMATION FOR A CONTINUOUS TARGET

FCLASSIF ANOV A FVALUE BETWEEN LABELFEATURE FOR CLASSIFICATION TASKS

CHI2 CHISQUARED STATS OF NONNEGATIVE FEATURES FOR CLASSIFICATION TASKS

SELECTKBEST SELECT FEATURES BASED ON THE K HIGHEST SCORES

SELECTFPR SELECT FEATURES BASED ON A FALSE POSITIVE RATE TEST

SELECTFDR SELECT FEATURES BASED ON AN ESTIMATED FALSE DISCOVERY RATE

SELECTFWE SELECT FEATURES BASED ON FAMILYWISE ERROR RATE

SELECTPERCENTILE SELECT FEATURES BASED ON PERCENTILE OF THE HIGHEST SCORES

EXAMPLES USING SKLEARNFEATURESELECTIONFREGRESSION

•FEATURE AGGLOMERATION VS UNIVARIATE SELECTION

•COMPARISON OF FTEST AND MUTUAL INFORMATION

•PIPELINE ANOVA SVM

61614SKLEARNFEATURESELECTION MUTUALINFOCLASSIF

SKLEARNFEATURESELECTION MUTUALINFOCLASSIF X Y DISCRETEFEATURES'AUTO'

NNEIGHBORS3 COPYTRUE RAN

DOMSTATENONE

ESTIMATE MUTUAL INFORMATION FOR A DISCRETE TARGET VARIABLE

MUTUAL INFORMATION MI 1BETWEEN TWO RANDOM VARIABLES IS A NONNEGATIVE VALUE WHICH MEASURES THE DEPENDENCY

BETWEEN THE VARIABLES IT IS EQUAL TO ZERO IF AND ONLY IF TWO RANDOM VARIABLES ARE INDEPENDENT AND HIGHER

VALUES MEAN HIGHER DEPENDENCY

THE FUNCTION RELIES ON NONPARAMETRIC METHODS BASED ON ENTROPY ESTIMATION FROM KNEAREST NEIGHBORS DISTANCES

AS DESCRIBED IN 2AND3 BOTH METHODS ARE BASED ON THE IDEA ORIGINALLY PROPOSED IN 4

IT CAN BE USED FOR UNIVARIATE FEATURES SELECTION READ MORE IN THE USER GUIDE

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES FEATURE MATRIX

YARRAYLIKE SHAPE NSAMPLES TARGET VECTOR

DISCRETEFEATURES 'AUTO' BOOL ARRAYLIKE DEFAULT 'AUTO' IF BOOL THEN DETERMINES WHETHER

TO CONSIDER ALL FEATURES DISCRETE OR CONTINUOUS IF ARRAY THEN IT SHOULD BE EITHER A BOOLEAN

MASK WITH SHAPE NFEATURES OR ARRAY WITH INDICES OF DISCRETE FEATURES IF 'AUTO' IT IS

ASSIGNED TO FALSE FOR DENSE XAND TO TRUE FOR SPARSE X

NNEIGHBORS INT DEFAULT 3 NUMBER OF NEIGHBORS TO USE FOR MI ESTIMATION FOR CONTINUOUS VARIABLES

SEE 2AND3 HIGHER VALUES REDUCE VARIANCE OF THE ESTIMATION BUT COULD INTRODUCE

A BIAS

COPY BOOL DEFAULT TRUE WHETHER TO MAKE A COPY OF THE GIVEN DATA IF SET TO FALSE THE INITIAL

DATA WILL BE OVERWRITTEN

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULT NONE THE SEED OF THE

PSEUDO RANDOM NUMBER GENERATOR FOR ADDING SMALL NOISE TO CONTINUOUS VARIABLES IN ORDER

TO REMOVE REPEATED VALUES IF INT RANDOMSTATE IS THE SEED USED BY THE RANDOM NUMBER

1796 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE  
THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NPRANDOM  
RETURNS

MINDARRAY SHAPE NFEATURES ESTIMATED MUTUAL INFORMATION BETWEEN EACH FEATURE AND THE  
TARGET

NOTES

1 THE TERM “DISCRETE FEATURES” IS USED INSTEAD OF NAMING THEM “CATEGORICAL” BECAUSE IT DESCRIBES THE ESSENCE  
MORE ACCURATELY FOR EXAMPLE PIXEL INTENSITIES OF AN IMAGE ARE DISCRETE FEATURES BUT HARDLY CATEGORICAL  
AND YOU WILL GET BETTER RESULTS IF MARK THEM AS SUCH ALSO NOTE THAT TREATING A CONTINUOUS VARIABLE AS  
DISCRETE AND VICE VERSA WILL USUALLY GIVE INCORRECT RESULTS SO BE ATTENTIVE ABOUT THAT

2 TRUE MUTUAL INFORMATION CAN’T BE NEGATIVE IF ITS ESTIMATE TURNS OUT TO BE NEGATIVE IT IS REPLACED BY ZERO  
REFERENCES

1234

61615SKLEARNFEATURESELECTION MUTUALINFOREGRESSION

SKLEARNFEATURESELECTION MUTUALINFOREGRESSION XY DISCRETEFEATURES’AUTO’

NNEIGHBORS3 COPYTRUE RAN

DOMSTATENONE

ESTIMATE MUTUAL INFORMATION FOR A CONTINUOUS TARGET VARIABLE

MUTUAL INFORMATION MI 1BETWEEN TWO RANDOM VARIABLES IS A NONNEGATIVE VALUE WHICH MEASURES THE DEPENDENCY  
BETWEEN THE VARIABLES IT IS EQUAL TO ZERO IF AND ONLY IF TWO RANDOM VARIABLES ARE INDEPENDENT AND HIGHER  
VALUES MEAN HIGHER DEPENDENCY

THE FUNCTION RELIES ON NONPARAMETRIC METHODS BASED ON ENTROPY ESTIMATION FROM KNEAREST NEIGHBORS DISTANCES  
AS DESCRIBED IN 2AND3 BOTH METHODS ARE BASED ON THE IDEA ORIGINALLY PROPOSED IN 4  
IT CAN BE USED FOR UNIVARIATE FEATURES SELECTION READ MORE IN THE USER GUIDE

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES FEATURE MATRIX

YARRAYLIKE SHAPE NSAMPLES TARGET VECTOR

DISCRETEFEATURES ‘AUTO’ BOOL ARRAYLIKE DEFAULT ‘AUTO’ IF BOOL THEN DETERMINES WHETHER  
TO CONSIDER ALL FEATURES DISCRETE OR CONTINUOUS IF ARRAY THEN IT SHOULD BE EITHER A BOOLEAN  
MASK WITH SHAPE NFEATURES OR ARRAY WITH INDICES OF DISCRETE FEATURES IF ‘AUTO’ IT IS  
ASSIGNED TO FALSE FOR DENSE XAND TO TRUE FOR SPARSE X

NNEIGHBORS INT DEFAULT 3 NUMBER OF NEIGHBORS TO USE FOR MI ESTIMATION FOR CONTINUOUS VARIABLES  
SEE 2AND3 HIGHER VALUES REDUCE VARIANCE OF THE ESTIMATION BUT COULD INTRODUCE  
A BIAS

COPY BOOL DEFAULT TRUE WHETHER TO MAKE A COPY OF THE GIVEN DATA IF SET TO FALSE THE INITIAL  
DATA WILL BE OVERWRITTEN

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULT NONE THE SEED OF THE  
PSEUDO RANDOM NUMBER GENERATOR FOR ADDING SMALL NOISE TO CONTINUOUS VARIABLES IN ORDER  
TO REMOVE REPEATED VALUES IF INT RANDOMSTATE IS THE SEED USED BY THE RANDOM NUMBER  
616SKLEARNFEATURESELECTION FEATURE SELECTION 1797

SCIKITLEARN USER GUIDE RELEASE 0213

GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE  
THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NPRANDOM  
RETURNS

MINDARRAY SHAPE NFEATURES ESTIMATED MUTUAL INFORMATION BETWEEN EACH FEATURE AND THE  
TARGET

NOTES

1 THE TERM “DISCRETE FEATURES” IS USED INSTEAD OF NAMING THEM “CATEGORICAL” BECAUSE IT DESCRIBES THE ESSENCE  
MORE ACCURATELY FOR EXAMPLE PIXEL INTENSITIES OF AN IMAGE ARE DISCRETE FEATURES BUT HARDLY CATEGORICAL  
AND YOU WILL GET BETTER RESULTS IF MARK THEM AS SUCH ALSO NOTE THAT TREATING A CONTINUOUS VARIABLE AS  
DISCRETE AND VICE VERSA WILL USUALLY GIVE INCORRECT RESULTS SO BE ATTENTIVE ABOUT THAT

2 TRUE MUTUAL INFORMATION CAN’T BE NEGATIVE IF ITS ESTIMATE TURNS OUT TO BE NEGATIVE IT IS REPLACED BY ZERO  
REFERENCES

1234

EXAMPLES USING SKLEARNFEATURESELECTIONMUTUALINFOREGRESSION

•COMPARISON OF FTEST AND MUTUAL INFORMATION

617SKLEARNGAUSSIANPROCESS GAUSSIAN PROCESSES

THESKLEARNGAUSSIANPROCESS MODULE IMPLEMENTS GAUSSIAN PROCESS BASED REGRESSION AND CLASSIFICATION

USER GUIDE SEE THE GAUSSIAN PROCESSES SECTION FOR FURTHER DETAILS

GAUSSIANPROCESSGAUSSIANPROCESSCLASSIFIER GAUSSIAN PROCESS CLASSIFICATION GPC BASED ON LAPLACE  
APPROXIMATION

GAUSSIANPROCESSGAUSSIANPROCESSREGRESSOR GAUSSIAN PROCESS REGRESSION GPR

6171SKLEARNGAUSSIANPROCESS GAUSSIANPROCESSCLASSIFIER

CLASSSSKLEARNGAUSSIANPROCESS GAUSSIANPROCESSCLASSIFIER KERNELNONE OPTI

MIZER’FMINLBFGSB’

NRESTARTSOPTIMIZER0

MAXITERPREDICT100

WARMSTARTFALSE

COPYXTRAINTRUE

RANDOMSTATENONE

MULTICLASS’ONEVSREST’

NJOBSNONE

GAUSSIAN PROCESS CLASSIFICATION GPC BASED ON LAPLACE APPROXIMATION

THE IMPLEMENTATION IS BASED ON ALGORITHM 31 32 AND 51 OF GAUSSIAN PROCESSES FOR MACHINE LEARNING

GPML BY RASMUSSEN AND WILLIAMS

1798 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

INTERNALLY THE LAPLACE APPROXIMATION IS USED FOR APPROXIMATING THE NONGAUSSIAN POSTERIOR BY A GAUSSIAN CURRENTLY THE IMPLEMENTATION IS RESTRICTED TO USING THE LOGISTIC LINK FUNCTION FOR MULTICLASS CLASSIFICATION SEVERAL BINARY ONEVERSUS REST CLASSIFIERS ARE FITTED NOTE THAT THIS CLASS THUS DOES NOT IMPLEMENT A TRUE MULTI CLASS LAPLACE APPROXIMATION

PARAMETERS

KERNEL KERNEL OBJECT THE KERNEL SPECIFYING THE COVARIANCE FUNCTION OF THE GP IF NONE IS PASSED THE KERNEL “10 RBF10” IS USED AS DEFAULT NOTE THAT THE KERNEL’S HYPERPARAMETERS ARE OPTIMIZED DURING FITTING

OPTIMIZER STRING OR CALLABLE OPTIONAL DEFAULT “FMINLBFGSB” CAN EITHER BE ONE OF THE INTERNALLY SUPPORTED OPTIMIZERS FOR OPTIMIZING THE KERNEL’S PARAMETERS SPECIFIED BY A STRING OR AN EXTERNALLY DEFINED OPTIMIZER PASSED AS A CALLABLE IF A CALLABLE IS PASSED IT MUST HAVE THE SIGNATURE

DEFOPTIMIZEROBJFUNC INITIALTHETA BOUNDS

OBJFUNC IS THE OBJECTIVE FUNCTION TO BE MAXIMIZED WHICH TAKES THE HYPERPARAMETERS THETA AS PARAMETER AND AN OPTIONAL FLAG EVALGRADIENT WHICH DETERMINES IF THE GRADIENT IS RETURNED ADDITIONALLY TO THE FUNCTION VALUE

INITIALTHETA THE INITIAL VALUE FOR THETA WHICH CAN BE USED BY LOCAL OPTIMIZERS

BOUNDS THE BOUNDS ON THE VALUES OF THETA

RETURNED ARE THE BEST FOUND HYPERPARAMETERS THETA AND THE CORRESPONDING VALUE OF THE TARGET FUNCTION

RETURNTHETAOPT FUNCMIN

PER DEFAULT THE ‘FMINLBFGSB’ ALGORITHM FROM SCIPYOPTIMIZE IS USED IF NONE IS PASSED THE KERNEL’S PARAMETERS ARE KEPT FIXED AVAILABLE INTERNAL OPTIMIZERS ARE

FMINLBFGSB

NRESTARTSOPTIMIZER INT OPTIONAL DEFAULT 0 THE NUMBER OF RESTARTS OF THE OPTIMIZER FOR FINDING THE KERNEL’S PARAMETERS WHICH MAXIMIZE THE LOGMARGINAL LIKELIHOOD THE FIRST RUN OF THE OPTIMIZER IS PERFORMED FROM THE KERNEL’S INITIAL PARAMETERS THE REMAINING ONES IF ANY FROM THETAS SAMPLED LOGUNIFORM RANDOMLY FROM THE SPACE OF ALLOWED THETAVALUES IF GREATER THAN 0 ALL BOUNDS MUST BE FINITE NOTE THAT NRESTARTSOPTIMIZER0 IMPLIES THAT ONE RUN IS PERFORMED

MAXITERPREDICT INT OPTIONAL DEFAULT 100 THE MAXIMUM NUMBER OF ITERATIONS IN NEWTON’S METHOD FOR APPROXIMATING THE POSTERIOR DURING PREDICT SMALLER VALUES WILL REDUCE COMPUTATION TIME AT THE COST OF WORSE RESULTS

WARMSTART BOOL OPTIONAL DEFAULT FALSE IF WARMSTARTS ARE ENABLED THE SOLUTION OF THE LAST NEWTON ITERATION ON THE LAPLACE APPROXIMATION OF THE POSTERIOR MODE IS USED AS INITIALIZATION FOR THE NEXT CALL OF POSTERIORMODE THIS CAN SPEED UP CONVERGENCE WHEN POSTERIORMODE IS CALLED SEVERAL TIMES ON SIMILAR PROBLEMS AS IN HYPERPARAMETER OPTIMIZATION

SEETHE GLOSSARY

COPYXTRAIN BOOL OPTIONAL DEFAULT TRUE IF TRUE A PERSISTENT COPY OF THE TRAINING DATA IS STORED IN THE OBJECT OTHERWISE JUST A REFERENCE TO THE TRAINING DATA IS STORED WHICH MIGHT CAUSE PREDICTIONS TO CHANGE IF THE DATA IS MODIFIED EXTERNALLY

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULT NONE THE GENERATOR USED TO INITIALIZE THE CENTERS IF INT RANDOMSTATE IS THE SEED USED BY THE RANDOM NUMBER

617SKLEARNGAUSSIANPROCESS GAUSSIAN PROCESSES 1799

SCIKITLEARN USER GUIDE RELEASE 0213

GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE  
THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NPRANDOM

MULTICLASS STRING DEFAULT SPECIFIES HOW MULTICLASS CLASSIFICATION PROBLEMS ARE HANDLED  
SUPPORTED ARE “ONEVSREST” AND “ONEVSONE” IN “ONEVSREST” ONE BINARY GAUSSIAN  
PROCESS CLASSIFIER IS FITTED FOR EACH CLASS WHICH IS TRAINED TO SEPARATE THIS CLASS FROM THE REST  
IN “ONEVSONE” ONE BINARY GAUSSIAN PROCESS CLASSIFIER IS FITTED FOR EACH PAIR OF CLASSES  
WHICH IS TRAINED TO SEPARATE THESE TWO CLASSES THE PREDICTIONS OF THESE BINARY PREDICTORS ARE  
COMBINED INTO MULTICLASS PREDICTIONS NOTE THAT “ONEVSONE” DOES NOT SUPPORT PREDICTING  
PROBABILITY ESTIMATES

NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF JOBS TO USE FOR THE COMPUTATION  
NONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL  
PROCESSORS SEE GLOSSARY FOR MORE DETAILS

ATTRIBUTES

KERNEL KERNEL OBJECT THE KERNEL USED FOR PREDICTION IN CASE OF BINARY CLASSIFICATION THE  
STRUCTURE OF THE KERNEL IS THE SAME AS THE ONE PASSED AS PARAMETER BUT WITH OPTIMIZED HY  
PERPARAMETERS IN CASE OF MULTICLASS CLASSIFICATION A COMPOUNDKERNEL IS RETURNED WHICH  
CONSISTS OF THE DIFFERENT KERNELS USED IN THE ONEVERSUSREST CLASSIFIERS

LOGMARGINALLIKELIHOODVALUE FLOAT THE LOGMARGINALLIKELIHOOD OF SELFKERNEL

THETA

CLASSES ARRAYLIKE SHAPE NCLASSES UNIQUE CLASS LABELS  
NCLASSES INT THE NUMBER OF CLASSES IN THE TRAINING DATA

EXAMPLES

```
FROM SKLEARNDATASETS IMPORT LOADIRIS
FROM SKLEARNGAUSSIANPROCESS IMPORT GAUSSIANPROCESSCLASSIFIER
FROM SKLEARNGAUSSIANPROCESSKERNELS IMPORT RBF
X Y LOADIRISRETURNXY TRUE
KERNEL 10 RBF10
GPC GAUSSIANPROCESSCLASSIFIERKERNELKERNEL
RANDOMSTATE0FITX Y
GPCSCOREX Y
09866
GPCPREDICTPROBAX2
ARRAY083548752 003228706 013222543
079064206 006525643 014410151
NEW IN VERSION 018
```

METHODS

FITSELF X Y FIT GAUSSIAN PROCESS CLASSIFICATION MODEL

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

LOGMARGINALLIKELIHOOD SELF THETA RETURNS LOGMARGINAL LIKELIHOOD OF THETA FOR TRAINING  
DATA

PREDICT SELF X PERFORM CLASSIFICATION ON AN ARRAY OF TEST VECTORS X

PREDICTPROBA SELF X RETURN PROBABILITY ESTIMATES FOR THE TEST VECTOR X

CONTINUED ON NEXT PAGE

1800 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

TABLE 6117 – CONTINUED FROM PREVIOUS PAGE

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELF KERNELNONE OPTIMIZER'FMINLBFGSB' NRESTARTSOPTIMIZERO

MAXITERPREDICT100 WARMSTARTFALSE COPYXTRAINTRUE RANDOMSTATENONE

MULTICLASS'ONEVSREST' NJOBSNONE

FITSELFXY

FIT GAUSSIAN PROCESS CLASSIFICATION MODEL

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA

YARRAYLIKE SHAPE NSAMPLES TARGET VALUES MUST BE BINARY

RETURNS

SELF RETURNS AN INSTANCE OF SELF

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

LOGMARGINALLIKELIHOOD SELFTHETANONE EVALGRADIENTFALSE

RETURNS LOGMARGINAL LIKELIHOOD OF THETA FOR TRAINING DATA

IN THE CASE OF MULTICLASS CLASSIFICATION THE MEAN LOGMARGINAL LIKELIHOOD OF THE ONEVERSUSREST CLASSIFIERS ARE RETURNED

PARAMETERS

THETA ARRAYLIKE SHAPE NKERNELPARAMS OR NONE KERNEL HYPERPARAMETERS FOR WHICH THE LOGMARGINAL LIKELIHOOD IS EVALUATED IN THE CASE OF MULTICLASS CLASSIFICATION THETA MAY BE THE HYPERPARAMETERS OF THE COMPOUND KERNEL OR OF AN INDIVIDUAL KERNEL IN THE LATTER CASE ALL INDIVIDUAL KERNEL GET ASSIGNED THE SAME THETA VALUES IF NONE THE PRECOMPUTED LOGMARGINALLIKELIHOOD OF SELFKERNELTHETA IS RETURNED

EVALGRADIENT BOOL DEFAULT FALSE IF TRUE THE GRADIENT OF THE LOGMARGINAL LIKELIHOOD WITH RESPECT TO THE KERNEL HYPERPARAMETERS AT POSITION THETA IS RETURNED ADDITIONALLY NOTE THAT GRADIENT COMPUTATION IS NOT SUPPORTED FOR NONBINARY CLASSIFICATION IF TRUE THETA MUST NOT BE NONE

RETURNS

LOGLIKELIHOOD FLOAT LOGMARGINAL LIKELIHOOD OF THETA FOR TRAINING DATA

LOGLIKELIHOODGRADIENT ARRAY SHAPE NKERNELPARAMS OPTIONAL GRADIENT OF THE LOG MARGINAL LIKELIHOOD WITH RESPECT TO THE KERNEL HYPERPARAMETERS AT POSITION THETA ONLY RETURNED WHEN EVALGRADIENT IS TRUE

PREDICTSELFXY

PERFORM CLASSIFICATION ON AN ARRAY OF TEST VECTORS X

6175KLEARNGAUSSIANPROCESS GAUSSIAN PROCESSES 1801

SCIKITLEARN USER GUIDE RELEASE 0213

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES

RETURNS

CARRAY SHAPE NSAMPLES PREDICTED TARGET VALUES FOR X VALUES ARE FROM CLASSES

PREDICTPROBA SELF

RETURN PROBABILITY ESTIMATES FOR THE TEST VECTOR X

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES

RETURNS

CARRAYLIKE SHAPE NSAMPLES NCLASSES RETURNS THE PROBABILITY OF THE SAMPLES FOR EACH CLASS IN THE MODEL THE COLUMNS CORRESPOND TO THE CLASSES IN SORTED ORDER AS THEY APPEAR IN THE ATTRIBUTE CLASSES

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS

IN MULTILABEL CLASSIFICATION THIS IS THE SUBSET ACCURACY WHICH IS A HARSH METRIC SINCE YOU REQUIRE FOR EACH SAMPLE THAT EACH LABEL SET BE CORRECTLY PREDICTED

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE LABELS FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT MEAN ACCURACY OF SELF PREDICT X WRT Y

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

EXAMPLES USING SKLEARN GAUSSIAN PROCESS GAUSSIAN PROCESS CLASSIFIER

- PLOT CLASSIFICATION PROBABILITY
- CLASSIFIER COMPARISON
- ILLUSTRATION OF GAUSSIAN PROCESS CLASSIFICATION GPC ON THE XOR DATASET
- GAUSSIAN PROCESS CLASSIFICATION GPC ON IRIS DATASET
- ISOPROBABILITY LINES FOR GAUSSIAN PROCESSES CLASSIFICATION GPC
- PROBABILISTIC PREDICTIONS WITH GAUSSIAN PROCESS CLASSIFICATION GPC

1802 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

61725KLEARNGAUSSIANPROCESS GAUSSIANPROCESSREGRESSOR

CLASSSSKLEARNGAUSSIANPROCESS GAUSSIANPROCESSREGRESSOR KERNELNONE ALPHA1E10

OPTIMIZER'FMINLBFGSB'

NRESTARTSOPTIMIZER0

NORMALIZEYFALSE

COPYXTRAINTRUE RAN

DOMSTATENONE

GAUSSIAN PROCESS REGRESSION GPR

THE IMPLEMENTATION IS BASED ON ALGORITHM 21 OF GAUSSIAN PROCESSES FOR MACHINE LEARNING GPML BY RAS

MUSSEN AND WILLIAMS

IN ADDITION TO STANDARD SCIKITLEARN ESTIMATOR API GAUSSIANPROCESSREGRESSOR

- ALLOWS PREDICTION WITHOUT PRIOR FITTING BASED ON THE GP PRIOR
- PROVIDES AN ADDITIONAL METHOD SAMPLEYX WHICH EVALUATES SAMPLES DRAWN FROM THE GPR PRIOR OR POS

TERIOR AT GIVEN INPUTS

- EXPOSES A METHOD LOGMARGINALLIKELIHOODTHETA WHICH CAN BE USED EXTERNALLY FOR OTHER WAYS OF SELECTING

HYPERPARAMETERS EG VIA MARKOV CHAIN MONTE CARLO

READ MORE IN THE USER GUIDE

NEW IN VERSION 018

PARAMETERS

KERNEL KERNEL OBJECT THE KERNEL SPECIFYING THE COVARIANCE FUNCTION OF THE GP IF NONE IS

PASSED THE KERNEL "10 RBF10" IS USED AS DEFAULT NOTE THAT THE KERNEL'S HYPERPA

RAMETERS ARE OPTIMIZED DURING FITTING

ALPHA FLOAT OR ARRAYLIKE OPTIONAL DEFAULT 1E10 VALUE ADDED TO THE DIAGONAL OF THE KERNEL

MATRIX DURING FITTING LARGER VALUES CORRESPOND TO INCREASED NOISE LEVEL IN THE OBSERVATIONS

THIS CAN ALSO PREVENT A POTENTIAL NUMERICAL ISSUE DURING FITTING BY ENSURING THAT THE CALCU

LATED VALUES FORM A POSITIVE DEFINITE MATRIX IF AN ARRAY IS PASSED IT MUST HAVE THE SAME

NUMBER OF ENTRIES AS THE DATA USED FOR FITTING AND IS USED AS DATAPOINTDEPENDENT NOISE LEVEL

NOTE THAT THIS IS EQUIVALENT TO ADDING A WHITEKERNEL WITH CALPHA ALLOWING TO SPECIFY THE

NOISE LEVEL DIRECTLY AS A PARAMETER IS MAINLY FOR CONVENIENCE AND FOR CONSISTENCY WITH RIDGE

OPTIMIZER STRING OR CALLABLE OPTIONAL DEFAULT "FMINLBFGSB" CAN EITHER BE ONE OF THE

INTERNALLY SUPPORTED OPTIMIZERS FOR OPTIMIZING THE KERNEL'S PARAMETERS SPECIFIED BY A STRING

OR AN EXTERNALLY DEFINED OPTIMIZER PASSED AS A CALLABLE IF A CALLABLE IS PASSED IT MUST HAVE

THE SIGNATURE

DEFOPTIMIZEROBJFUNC INITIALTHETA BOUNDS

OBJFUNC IS THE OBJECTIVE FUNCTION TO BE MINIMIZED WHICH

TAKES THE HYPERPARAMETERS THETA AS PARAMETER AND AN

OPTIONAL FLAG EVALGRADIENT WHICH DETERMINES IF THE

GRADIENT IS RETURNED ADDITIONALLY TO THE FUNCTION VALUE

INITIALTHETA THE INITIAL VALUE FOR THETA WHICH CAN BE

USED BY LOCAL OPTIMIZERS

BOUNDS THE BOUNDS ON THE VALUES OF THETA

RETURNED ARE THE BEST FOUND HYPERPARAMETERS THETA AND

THE CORRESPONDING VALUE OF THE TARGET FUNCTION

RETURNTHETAOPT FUNCMIN

PER DEFAULT THE 'FMINLBFGSB' ALGORITHM FROM SCIPYOPTIMIZE IS USED IF NONE IS PASSED

THE KERNEL'S PARAMETERS ARE KEPT FIXED AVAILABLE INTERNAL OPTIMIZERS ARE

6175KLEARNGAUSSIANPROCESS GAUSSIAN PROCESSES 1803

SCIKITLEARN USER GUIDE RELEASE 0213

FMINLBFGSB

NRESTARTSOPTIMIZER INT OPTIONAL DEFAULT 0 THE NUMBER OF RESTARTS OF THE OPTIMIZER FOR FINDING THE KERNEL'S PARAMETERS WHICH MAXIMIZE THE LOGMARGINAL LIKELIHOOD THE FIRST RUN OF THE OPTIMIZER IS PERFORMED FROM THE KERNEL'S INITIAL PARAMETERS THE REMAINING ONES IF ANY FROM THETAS SAMPLED LOGUNIFORM RANDOMLY FROM THE SPACE OF ALLOWED THETAVALUES IF GREATER THAN 0 ALL BOUNDS MUST BE FINITE NOTE THAT NRESTARTSOPTIMIZER 0 IMPLIES THAT ONE RUN IS PERFORMED

NORMALIZEY BOOLEAN OPTIONAL DEFAULT FALSE WHETHER THE TARGET VALUES Y ARE NORMALIZED IE THE MEAN OF THE OBSERVED TARGET VALUES BECOME ZERO THIS PARAMETER SHOULD BE SET TO TRUE IF THE TARGET VALUES' MEAN IS EXPECTED TO DIFFER CONSIDERABLE FROM ZERO WHEN ENABLED THE NORMALIZATION EFFECTIVELY MODIFIES THE GP'S PRIOR BASED ON THE DATA WHICH CONTRADICTS THE LIKELIHOOD PRINCIPLE NORMALIZATION IS THUS DISABLED PER DEFAULT

COPYXTRAIN BOOL OPTIONAL DEFAULT TRUE IF TRUE A PERSISTENT COPY OF THE TRAINING DATA IS STORED IN THE OBJECT OTHERWISE JUST A REFERENCE TO THE TRAINING DATA IS STORED WHICH MIGHT CAUSE PREDICTIONS TO CHANGE IF THE DATA IS MODIFIED EXTERNALLY

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULT NONE THE GENERATOR USED TO INITIALIZE THE CENTERS IF INT RANDOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NPROBANDOM

ATTRIBUTES

XTRAIN ARRAYLIKE SHAPE NSAMPLES NFEATURES FEATURE VALUES IN TRAINING DATA ALSO REQUIRED FOR PREDICTION

YTRAIN ARRAYLIKE SHAPE NSAMPLES NOUTPUTDIMS TARGET VALUES IN TRAINING DATA ALSO REQUIRED FOR PREDICTION

KERNEL KERNEL OBJECT THE KERNEL USED FOR PREDICTION THE STRUCTURE OF THE KERNEL IS THE SAME AS THE ONE PASSED AS PARAMETER BUT WITH OPTIMIZED HYPERPARAMETERS

LARRAYLIKE SHAPE NSAMPLES NSAMPLES LOWERTRIANGULAR CHOLESKY DECOMPOSITION OF THE KERNEL IN XTRAIN

ALPHA ARRAYLIKE SHAPE NSAMPLES DUAL COEFFICIENTS OF TRAINING DATA POINTS IN KERNEL SPACE

LOGMARGINALLIKELIHOODVALUE FLOAT THE LOGMARGINALLIKELIHOOD OF SELFKERNEL

THETA

EXAMPLES

```
FROM SKLEARN DATASETS IMPORT MAKEFRIEDMAN2
FROM SKLEARN GAUSSIANPROCESS IMPORT GAUSSIANPROCESSREGRESSOR
FROM SKLEARN GAUSSIANPROCESS KERNELS IMPORT DOTPRODUCT WHITEKERNEL
X Y MAKEFRIEDMAN2 NSAMPLES500 NOISE0 RANDOMSTATE0
KERNEL DOTPRODUCT WHITEKERNEL
GPR GAUSSIANPROCESSREGRESSOR KERNEL KERNEL
RANDOMSTATE0 FIT X Y
GPR SCORE X Y
03680
GPR PREDICT X2 RETURN STD TRUE
ARRAY6530 5921 ARRAY3166 3166
1804 CHAPTER 6 API REFERENCE
```

SCIKITLEARN USER GUIDE RELEASE 0213

METHODS

FITSELF X Y FIT GAUSSIAN PROCESS REGRESSION MODEL  
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR  
LOGMARGINALLIKELIHOOD SELF THETA RETURNS LOGMARGINAL LIKELIHOOD OF THETA FOR TRAINING DATA  
PREDICT SELF X RETURNSTD RETURNCOV PREDICT USING THE GAUSSIAN PROCESS REGRESSION MODEL  
SAMPLEY SELF X NSAMPLES RANDOMSTATE DRAW SAMPLES FROM GAUSSIAN PROCESS AND EVALUATE AT X  
SCORE SELF X Y SAMPLEWEIGHT RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION  
SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR  
INIT SELFKERNELNONE ALPHA1E10 OPTIMIZER'FMINLBFGSB' NRESTARTSOPTIMIZER0 NORMALIZEYFALSE COPYXTRAINTRUE RANDOMSTATENONE  
FITSELFXY  
FIT GAUSSIAN PROCESS REGRESSION MODEL  
PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA  
YARRAYLIKE SHAPE NSAMPLES NOUTPUTDIMS TARGET VALUES  
RETURNS  
SELF RETURNS AN INSTANCE OF SELF  
GETPARAMS SELFDEEPTURE  
GET PARAMETERS FOR THIS ESTIMATOR  
PARAMETERS  
DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS  
RETURNS  
PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES  
LOGMARGINALLIKELIHOOD SELFTHETANONE EVALGRADIENTFALSE  
RETURNS LOGMARGINAL LIKELIHOOD OF THETA FOR TRAINING DATA  
PARAMETERS  
THETA ARRAYLIKE SHAPE NKERNELPARAMS OR NONE KERNEL HYPERPARAMETERS FOR WHICH THE LOGMARGINAL LIKELIHOOD IS EVALUATED IF NONE THE PRECOMPUTED LOGMARGINALLIKELIHOOD OF SELFKERNELTHETA IS RETURNED  
EVALGRADIENT BOOL DEFAULT FALSE IF TRUE THE GRADIENT OF THE LOGMARGINAL LIKELIHOOD WITH RESPECT TO THE KERNEL HYPERPARAMETERS AT POSITION THETA IS RETURNED ADDITIONALLY IF TRUE THETA MUST NOT BE NONE  
RETURNS  
LOGLIKELIHOOD FLOAT LOGMARGINAL LIKELIHOOD OF THETA FOR TRAINING DATA  
LOGLIKELIHOODGRADIENT ARRAY SHAPE NKERNELPARAMS OPTIONAL GRADIENT OF THE LOG MARGINAL LIKELIHOOD WITH RESPECT TO THE KERNEL HYPERPARAMETERS AT POSITION THETA ONLY RETURNED WHEN EVALGRADIENT IS TRUE  
617SKLEARNGAUSSIANPROCESS GAUSSIAN PROCESSES 1805

SCIKITLEARN USER GUIDE RELEASE 0213

PREDICTSELF  
XRETURNSTD FALSE RETURNCOV FALSE

PREDICT USING THE GAUSSIAN PROCESS REGRESSION MODEL

WE CAN ALSO PREDICT BASED ON AN UNFITTED MODEL BY USING THE GP PRIOR IN ADDITION TO THE MEAN OF THE PREDICTIVE DISTRIBUTION ALSO ITS STANDARD DEVIATION RETURNSTD TRUE OR COVARIANCE RETURNCOV TRUE NOTE THAT AT MOST ONE OF THE TWO CAN BE REQUESTED

PARAMETERS

X ARRAYLIKE SHAPE NSAMPLES NFEATURES QUERY POINTS WHERE THE GP IS EVALUATED

RETURNSTD BOOL DEFAULT FALSE IF TRUE THE STANDARD DEVIATION OF THE PREDICTIVE DISTRIBUTION AT THE QUERY POINTS IS RETURNED ALONG WITH THE MEAN

RETURNCOV BOOL DEFAULT FALSE IF TRUE THE COVARIANCE OF THE JOINT PREDICTIVE DISTRIBUTION AT THE QUERY POINTS IS RETURNED ALONG WITH THE MEAN

RETURNS

YMEAN ARRAY SHAPE NSAMPLES NOUTPUTDIMS MEAN OF PREDICTIVE DISTRIBUTION AT QUERY POINTS

YSTD ARRAY SHAPE NSAMPLES OPTIONAL STANDARD DEVIATION OF PREDICTIVE DISTRIBUTION AT QUERY POINTS ONLY RETURNED WHEN RETURNSTD IS TRUE

YCOV ARRAY SHAPE NSAMPLES NSAMPLES OPTIONAL COVARIANCE OF JOINT PREDICTIVE DISTRIBUTION AT QUERY POINTS ONLY RETURNED WHEN RETURNCOV IS TRUE

SAMPLE SELF X NSAMPLES 1 RANDOMSTATE 0

DRAW SAMPLES FROM GAUSSIAN PROCESS AND EVALUATE AT X

PARAMETERS

X ARRAYLIKE SHAPE NSAMPLES X NFEATURES QUERY POINTS WHERE THE GP SAMPLES ARE EVALUATED

NSAMPLES INT DEFAULT 1 THE NUMBER OF SAMPLES DRAWN FROM THE GAUSSIAN PROCESS

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULT 0 IF INT RANDOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NRANDOM

RETURNS

Y SAMPLES ARRAY SHAPE NSAMPLES X NOUTPUTDIMS NSAMPLES VALUES OF NSAMPLES SAMPLES DRAWN FROM GAUSSIAN PROCESS AND EVALUATED AT QUERY POINTS

SCORE SELF X Y SAMPLEWEIGHT NONE

RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{U}{V}$  WHERE U IS THE RESIDUAL SUM OF SQUARES  $Y_{TRUE} - Y_{PRED}$  2SUM AND V IS THE TOTAL SUM OF SQUARES  $Y_{TRUE} - Y_{TRUEMEAN}$  2SUM THE BEST POSSIBLE SCORE IS 10 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 00

PARAMETERS

X ARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES FOR SOME ESTIMATORS THIS MAY BE A PRECOMPUTED KERNEL MATRIX INSTEAD SHAPE NSAMPLES NSAMPLESFITTED WHERE NSAMPLESFITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR

Y ARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE VALUES FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

1806 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

RETURNS

SCORE FLOAT R2 OF SELF PREDICT X WRT Y

NOTES

THE R2 SCORE USED WHEN CALLING SCORE ON A REGRESSOR WILL USE MULTIOUTPUT UNIFORM AVERAGE FROM VERSION 023 TO KEEP CONSISTENT WITH METRICS R2 SCORE THIS WILL INFLUENCE THE SCORE METHOD OF ALL THE MULTIOUTPUT REGRESSORS EXCEPT FOR MULTIOUTPUT MULTIOUTPUT REGRESSOR TO SPECIFY THE DEFAULT VALUE MANUALLY AND AVOID THE WARNING PLEASE EITHER CALL METRICS R2 SCORE DIRECTLY OR MAKE A CUSTOM SCORER WITH METRICS MAKE SCORER THE BUILT IN SCORER R2 USES

MULTIOUTPUT UNIFORM AVERAGE

SET PARAMS SELF PARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENT PARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT

OF A NESTED OBJECT

RETURNS

SELF

EXAMPLES USING SKLEARN GAUSSIAN PROCESS GAUSSIAN PROCESS REGRESSOR

- COMPARISON OF KERNEL RIDGE AND GAUSSIAN PROCESS REGRESSION
- ILLUSTRATION OF PRIOR AND POSTERIOR GAUSSIAN PROCESS FOR DIFFERENT KERNELS
- GAUSSIAN PROCESS REGRESSION GPR WITH NOISE LEVEL ESTIMATION
- GAUSSIAN PROCESSES REGRESSION BASIC INTRODUCTORY EXAMPLE
- GAUSSIAN PROCESS REGRESSION GPR ON MAUNA LOA CO2 DATA

KERNELS

GAUSSIAN PROCESS KERNELS

COMPOUND KERNEL KERNELS KERNEL WHICH IS COMPOSED OF A SET OF OTHER KERNELS

GAUSSIAN PROCESS KERNELS

CONSTANT KERNEL CONSTANT KERNEL

GAUSSIAN PROCESS KERNELS

DOT PRODUCT DOT PRODUCT KERNEL

GAUSSIAN PROCESS KERNELS

EXP SINES SQUARED EXP SINES SQUARED KERNEL

GAUSSIAN PROCESS KERNELS

EXPONENTIATION EXPONENTIATE KERNEL BY GIVEN EXPONENT

GAUSSIAN PROCESS KERNELS HYPERPARAMETER A KERNEL HYPERPARAMETER'S SPECIFICATION IN FORM OF A NAMED TUPLE

GAUSSIAN PROCESS KERNELS KERNEL BASE CLASS FOR ALL KERNELS

GAUSSIAN PROCESS KERNELS MATERN MATERN KERNEL

GAUSSIAN PROCESS KERNELS

PAIRWISE KERNEL WRAPPER FOR KERNELS IN SKLEARN METRIC PAIRWISE

GAUSSIAN PROCESS KERNELS PRODUCT K1 K2 PRODUCT KERNEL K1 K2 OF TWO KERNELS K1 AND K2

CONTINUED ON NEXT PAGE

617 SKLEARN GAUSSIAN PROCESS GAUSSIAN PROCESSES 1807

SCIKITLEARN USER GUIDE RELEASE 0213  
TABLE 6119 – CONTINUED FROM PREVIOUS PAGE  
GAUSSIANPROCESSKERNELSRBF LENGTHSCALE  
RADIALBASIS FUNCTION KERNEL AKA SQUAREDEXPONENTIAL KERNEL  
NEL  
GAUSSIANPROCESSKERNELS  
RATIONALQUADRATIC RATIONAL QUADRATIC KERNEL  
GAUSSIANPROCESSKERNELSSUM K1 K2 SUMKERNEL K1 K2 OF TWO KERNELS K1 AND K2  
GAUSSIANPROCESSKERNELS  
WHITEKERNEL WHITE KERNEL  
6173SKLEARNGAUSSIANPROCESSKERNELS COMPOUNDKERNEL  
CLASSSSKLEARNGAUSSIANPROCESSKERNELS COMPOUNDKERNEL KERNELS  
KERNEL WHICH IS COMPOSED OF A SET OF OTHER KERNELS  
NEW IN VERSION 018  
PARAMETERS  
KERNELS LIST OF KERNEL OBJECTS THE OTHER KERNELS  
ATTRIBUTES  
BOUNDS RETURNS THE LOGTRANSFORMED BOUNDS ON THE THETA  
HYPERPARAMETERS RETURNS A LIST OF ALL HYPERPARAMETER SPECIFICATIONS  
NDIMS RETURNS THE NUMBER OF NONFIXED HYPERPARAMETERS OF THE KERNEL  
THETA RETURNS THE FLATTENED LOGTRANSFORMED NONFIXED HYPERPARAMETERS  
METHODS  
CALL SELF X Y EVALGRADIENT RETURN THE KERNEL KX Y AND OPTIONALLY ITS GRADIENT  
CLONEWITHTHETA SELF THETA RETURNS A CLONE OF SELF WITH GIVEN HYPERPARAMETERS  
THETA  
DIAG SELF X RETURNS THE DIAGONAL OF THE KERNEL KX X  
GETPARAMS SELF DEEP GET PARAMETERS OF THIS KERNEL  
ISSTATIONARY SELF RETURNS WHETHER THE KERNEL IS STATIONARY  
SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS KERNEL  
INIT SELFKERNELS  
CALL SELFXYNONE EVALGRADIENTFALSE  
RETURN THE KERNEL KX Y AND OPTIONALLY ITS GRADIENT  
NOTE THAT THIS COMPOUND KERNEL RETURNS THE RESULTS OF ALL SIMPLE KERNEL STACKED ALONG AN ADDITIONAL AXIS  
PARAMETERS  
XARRAY SHAPE NSAMPLESX NFEATURES LEFT ARGUMENT OF THE RETURNED KERNEL KX Y  
YARRAY SHAPE NSAMPLESY NFEATURES OPTIONAL DEFAULTNONE RIGHT ARGUMENT OF THE  
RETURNED KERNEL KX Y IF NONE KX X IF EVALUATED INSTEAD  
EVALGRADIENT BOOL OPTIONAL DEFAULTFALSE DETERMINES WHETHER THE GRADIENT WITH RESPECT TO THE KERNEL HYPERPARAMETER IS DETERMINED  
RETURNS  
1808 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

KARRAY SHAPE NSAMPLESX NSAMPLESY NKERNELS KERNEL KX Y

KGRADIENT ARRAY SHAPE NSAMPLESX NSAMPLESX NDIMS NKERNELS THE GRADIENT OF THE KERNEL KX X WITH RESPECT TO THE HYPERPARAMETER OF THE KERNEL ONLY RETURNED WHEN EVALGRADIENT IS TRUE

BOUNDS

RETURNS THE LOGTRANSFORMED BOUNDS ON THE THETA

RETURNS

BOUNDS ARRAY SHAPE NDIMS 2 THE LOGTRANSFORMED BOUNDS ON THE KERNEL’S HYPERPARAMETERS THETA

CLONEWITHTHETA SELFTHETA

RETURNS A CLONE OF SELF WITH GIVEN HYPERPARAMETERS THETA

PARAMETERS

THETA ARRAY SHAPE NDIMS THE HYPERPARAMETERS

DIAGSELF

RETURNS THE DIAGONAL OF THE KERNEL KX X

THE RESULT OF THIS METHOD IS IDENTICAL TO NPDIAGSELF HOWEVER IT CAN BE EVALUATED MORE EFFICIENTLY SINCE ONLY THE DIAGONAL IS EVALUATED

PARAMETERS

XARRAY SHAPE NSAMPLESX NFEATURES LEFT ARGUMENT OF THE RETURNED KERNEL KX Y

RETURNS

KDIAG ARRAY SHAPE NSAMPLESX NKERNELS DIAGONAL OF KERNEL KX X

GETPARAMS SELFDEEPTURE

GET PARAMETERS OF THIS KERNEL

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

HYPERPARAMETERS

RETURNS A LIST OF ALL HYPERPARAMETER SPECIFICATIONS

ISSTATIONARY SELF

RETURNS WHETHER THE KERNEL IS STATIONARY

NDIMS

RETURNS THE NUMBER OF NONFIXED HYPERPARAMETERS OF THE KERNEL

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS KERNEL

THE METHOD WORKS ON SIMPLE KERNELS AS WELL AS ON NESTED KERNELS THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

617SKLEARNGAUSSIANPROCESS GAUSSIAN PROCESSES 1809

SCIKITLEARN USER GUIDE RELEASE 0213

THETA

RETURNS THE FLATTENED LOGTRANSFORMED NONFIXED HYPERPARAMETERS

NOTE THAT THETA ARE TYPICALLY THE LOGTRANSFORMED VALUES OF THE KERNEL'S HYPERPARAMETERS AS THIS REPRESENTATION OF THE SEARCH SPACE IS MORE AMENABLE FOR HYPERPARAMETER SEARCH AS HYPERPARAMETERS LIKE LENGTHSCALES NATURALLY LIVE ON A LOGSCALE

RETURNS

THETA ARRAY SHAPE NDIMS THE NONFIXED LOGTRANSFORMED HYPERPARAMETERS OF THE KERNEL

6174SKLEARNGAUSSIANPROCESSKERNELS CONSTANTKERNEL

CLASSSKLEARNGAUSSIANPROCESSKERNELS CONSTANTKERNEL CONSTANTVALUE10

CONSTANTVALUEBOUNDS1E

051000000

CONSTANT KERNEL

CAN BE USED AS PART OF A PRODUCTKERNEL WHERE IT SCALES THE MAGNITUDE OF THE OTHER FACTOR KERNEL OR AS PART OF A SUMKERNEL WHERE IT MODIFIES THE MEAN OF THE GAUSSIAN PROCESS

KX1 X2 CONSTANTVALUE FOR ALL X1 X2

NEW IN VERSION 018

PARAMETERS

CONSTANTVALUE FLOAT DEFAULT 10 THE CONSTANT VALUE WHICH DEFINES THE COVARIANCE KX1

X2 CONSTANTVALUE

CONSTANTVALUEBOUNDS PAIR OF FLOATS 0 DEFAULT 1E5 1E5 THE LOWER AND UPPER BOUND

ON CONSTANTVALUE

ATTRIBUTES

BOUNDS RETURNS THE LOGTRANSFORMED BOUNDS ON THE THETA

HYPERPARAMETERCONSTANTVALUE

HYPERPARAMETERS RETURNS A LIST OF ALL HYPERPARAMETER SPECIFICATIONS

NDIMS RETURNS THE NUMBER OF NONFIXED HYPERPARAMETERS OF THE KERNEL

THETA RETURNS THE FLATTENED LOGTRANSFORMED NONFIXED HYPERPARAMETERS

METHODS

CALL SELF X Y EVALGRADIENT RETURN THE KERNEL KX Y AND OPTIONALLY ITS GRADIENT

CLONEWITHTHETA SELF THETA RETURNS A CLONE OF SELF WITH GIVEN HYPERPARAMETERS

THETA

DIAG SELF X RETURNS THE DIAGONAL OF THE KERNEL KX X

GETPARAMS SELF DEEP GET PARAMETERS OF THIS KERNEL

ISSTATIONARY SELF RETURNS WHETHER THE KERNEL IS STATIONARY

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS KERNEL

INIT SELFCONSTANTVALUE10 CONSTANTVALUEBOUNDS1E05 1000000

CALL SELFXYNONE EVALGRADIENTFALSE

RETURN THE KERNEL KX Y AND OPTIONALLY ITS GRADIENT

1810 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

PARAMETERS

XARRAY SHAPE NSAMPLESX NFEATURES LEFT ARGUMENT OF THE RETURNED KERNEL KX Y  
YARRAY SHAPE NSAMPLESY NFEATURES OPTIONAL DEFAULTNONE RIGHT ARGUMENT OF THE  
RETURNED KERNEL KX Y IF NONE KX X IF EVALUATED INSTEAD  
EVALGRADIENT BOOL OPTIONAL DEFAULTFALSE DETERMINES WHETHER THE GRADIENT WITH RE  
SPECT TO THE KERNEL HYPERPARAMETER IS DETERMINED ONLY SUPPORTED WHEN Y IS NONE  
RETURNS

KARRAY SHAPE NSAMPLESX NSAMPLESY KERNEL KX Y  
KGRADIENT ARRAY OPT SHAPE NSAMPLESX NSAMPLESX NDIMS THE GRADIENT OF THE  
KERNEL KX X WITH RESPECT TO THE HYPERPARAMETER OF THE KERNEL ONLY RETURNED WHEN  
EVALGRADIENT IS TRUE

BOUNDS

RETURNS THE LOGTRANSFORMED BOUNDS ON THE THETA  
RETURNS  
BOUNDS ARRAY SHAPE NDIMS 2 THE LOGTRANSFORMED BOUNDS ON THE KERNEL’S HYPERPARAM  
ETERS THETA

CLONEWITHTHETA SELFTHETA  
RETURNS A CLONE OF SELF WITH GIVEN HYPERPARAMETERS THETA

PARAMETERS

THETA ARRAY SHAPE NDIMS THE HYPERPARAMETERS

DIAGSELF

RETURNS THE DIAGONAL OF THE KERNEL KX X  
THE RESULT OF THIS METHOD IS IDENTICAL TO NPDIAGSELF HOWEVER IT CAN BE EVALUATED MORE EFFICIENTLY  
SINCE ONLY THE DIAGONAL IS EVALUATED

PARAMETERS

XARRAY SHAPE NSAMPLESX NFEATURES LEFT ARGUMENT OF THE RETURNED KERNEL KX Y  
RETURNS

KDIAG ARRAY SHAPE NSAMPLESX DIAGONAL OF KERNEL KX X

GETPARAMS SELFDEEPT

GET PARAMETERS OF THIS KERNEL

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

HYPERPARAMETERS

RETURNS A LIST OF ALL HYPERPARAMETER SPECIFICATIONS

ISSTATIONARY SELF

RETURNS WHETHER THE KERNEL IS STATIONARY

NDIMS

RETURNS THE NUMBER OF NONFIXED HYPERPARAMETERS OF THE KERNEL

SCIKITLEARN USER GUIDE RELEASE 0213

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS KERNEL

THE METHOD WORKS ON SIMPLE KERNELS AS WELL AS ON NESTED KERNELS THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

THETA

RETURNS THE FLATTENED LOGTRANSFORMED NONFIXED HYPERPARAMETERS

NOTE THAT THETA ARE TYPICALLY THE LOGTRANSFORMED VALUES OF THE KERNEL'S HYPERPARAMETERS AS THIS REPRESENTATION OF THE SEARCH SPACE IS MORE AMENABLE FOR HYPERPARAMETER SEARCH AS HYPERPARAMETERS LIKE LENGTHSCALES NATURALLY LIVE ON A LOGSCALE

RETURNS

THETA ARRAY SHAPE NDIMS THE NONFIXED LOGTRANSFORMED HYPERPARAMETERS OF THE KERNEL

EXAMPLES USING SKLEARNGAUSSIANPROCESSKERNELSCONSTANTKERNEL

- ILLUSTRATION OF PRIOR AND POSTERIOR GAUSSIAN PROCESS FOR DIFFERENT KERNELS
- ISOPROBABILITY LINES FOR GAUSSIAN PROCESSES CLASSIFICATION GPC
- GAUSSIAN PROCESSES REGRESSION BASIC INTRODUCTORY EXAMPLE

6175SKLEARNGAUSSIANPROCESSKERNELS DOTPRODUCT

CLASSSSKLEARNGAUSSIANPROCESSKERNELS DOTPRODUCT SIGMA010 SIGMA0BOUNDS1E

051000000

DOTPRODUCT KERNEL

THE DOTPRODUCT KERNEL IS NONSTATIONARY AND CAN BE OBTAINED FROM LINEAR REGRESSION BY PUTTING NO 1 PRIORS ON THE COEFFICIENTS OF  $XD D^{-1} D$  AND A PRIOR OF NO SIGMA02 ON THE BIAS THE DOTPRODUCT KERNEL IS INVARIANT TO A ROTATION OF THE COORDINATES ABOUT THE ORIGIN BUT NOT TRANSLATIONS IT IS PARAMETERIZED BY A PARAMETER SIGMA02 FOR SIGMA02 0 THE KERNEL IS CALLED THE HOMOGENEOUS LINEAR KERNEL OTHERWISE IT IS INHOMOGENEOUS THE KERNEL IS GIVEN BY

$$K(X_i, X_j) = \text{SIGMA0}^{-2} X_i^T C \text{DOT} X_j$$

THE DOTPRODUCT KERNEL IS COMMONLY COMBINED WITH EXPONENTIATION

NEW IN VERSION 018

PARAMETERS

SIGMA0 FLOAT 0 DEFAULT 10 PARAMETER CONTROLLING THE INHOMOGENITY OF THE KERNEL IF SIGMA00 THE KERNEL IS HOMOGENOUS

SIGMA0BOUNDS PAIR OF FLOATS 0 DEFAULT 1E5 1E5 THE LOWER AND UPPER BOUND ON L

ATTRIBUTES

BOUNDS RETURNS THE LOGTRANSFORMED BOUNDS ON THE THETA

HYPERPARAMETERSSIGMA0

HYPERPARAMETERS RETURNS A LIST OF ALL HYPERPARAMETER SPECIFICATIONS

NDIMS RETURNS THE NUMBER OF NONFIXED HYPERPARAMETERS OF THE KERNEL

1812 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

THETA RETURNS THE FLATTENED LOGTRANSFORMED NONFIXED HYPERPARAMETERS

METHODS

CALL SELF X Y EVALGRADIENT RETURN THE KERNEL KX Y AND OPTIONALLY ITS GRADIENT

CLONEWITHTHETA SELF THETA RETURNS A CLONE OF SELF WITH GIVEN HYPERPARAMETERS

THETA

DIAG SELF X RETURNS THE DIAGONAL OF THE KERNEL KX X

GETPARAMS SELF DEEP GET PARAMETERS OF THIS KERNEL

ISSTATIONARY SELF RETURNS WHETHER THE KERNEL IS STATIONARY

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS KERNEL

INIT SELFSIGMA010 SIGMA0BOUNDS1E05 1000000

CALL SELFXYNONE EVALGRADIENTFALSE

RETURN THE KERNEL KX Y AND OPTIONALLY ITS GRADIENT

PARAMETERS

XARRAY SHAPE NSAMPLESX NFEATURES LEFT ARGUMENT OF THE RETURNED KERNEL KX Y

YARRAY SHAPE NSAMPLESY NFEATURES OPTIONAL DEFAULTNONE RIGHT ARGUMENT OF THE

RETURNED KERNEL KX Y IF NONE KX X IF EVALUATED INSTEAD

EVALGRADIENT BOOL OPTIONAL DEFAULTFALSE DETERMINES WHETHER THE GRADIENT WITH RE

SPECT TO THE KERNEL HYPERPARAMETER IS DETERMINED ONLY SUPPORTED WHEN Y IS NONE

RETURNS

KARRAY SHAPE NSAMPLESX NSAMPLESY KERNEL KX Y

KGRADIENT ARRAY OPT SHAPE NSAMPLESX NSAMPLESX NDIMS THE GRADIENT OF THE

KERNEL KX X WITH RESPECT TO THE HYPERPARAMETER OF THE KERNEL ONLY RETURNED WHEN

EVALGRADIENT IS TRUE

BOUNDS

RETURNS THE LOGTRANSFORMED BOUNDS ON THE THETA

RETURNS

BOUNDS ARRAY SHAPE NDIMS 2 THE LOGTRANSFORMED BOUNDS ON THE KERNEL’S HYPERPARAM

ETERS THETA

CLONEWITHTHETA SELFTHETA

RETURNS A CLONE OF SELF WITH GIVEN HYPERPARAMETERS THETA

PARAMETERS

THETA ARRAY SHAPE NDIMS THE HYPERPARAMETERS

DIAGSELF

RETURNS THE DIAGONAL OF THE KERNEL KX X

THE RESULT OF THIS METHOD IS IDENTICAL TO NPDIAGSELF HOWEVER IT CAN BE EVALUATED MORE EFFICIENTLY

SINCE ONLY THE DIAGONAL IS EVALUATED

PARAMETERS

XARRAY SHAPE NSAMPLESX NFEATURES LEFT ARGUMENT OF THE RETURNED KERNEL KX Y

RETURNS

617SKLEARNGAUSSIANPROCESS GAUSSIAN PROCESSES 1813

SCIKITLEARN USER GUIDE RELEASE 0213

KDIAG ARRAY SHAPE NSAMPLESX DIAGONAL OF KERNEL KX X

GETPARAMS SELFDEEPTTRUE

GET PARAMETERS OF THIS KERNEL

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

HYPERPARAMETERS

RETURNS A LIST OF ALL HYPERPARAMETER SPECIFICATIONS

ISSTATIONARY SELF

RETURNS WHETHER THE KERNEL IS STATIONARY

NDIMS

RETURNS THE NUMBER OF NONFIXED HYPERPARAMETERS OF THE KERNEL

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS KERNEL

THE METHOD WORKS ON SIMPLE KERNELS AS WELL AS ON NESTED KERNELS THE LATTER HAVE PARAMETERS OF THE FORM

COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

THETA

RETURNS THE FLATTENED LOGTRANSFORMED NONFIXED HYPERPARAMETERS

NOTE THAT THETA ARE TYPICALLY THE LOGTRANSFORMED VALUES OF THE KERNEL'S HYPERPARAMETERS AS THIS REPRESENTA

TION OF THE SEARCH SPACE IS MORE AMENABLE FOR HYPERPARAMETER SEARCH AS HYPERPARAMETERS LIKE LENGTHSCALES

NATURALLY LIVE ON A LOGSCALE

RETURNS

THETA ARRAY SHAPE NDIMS THE NONFIXED LOGTRANSFORMED HYPERPARAMETERS OF THE KERNEL

EXAMPLES USING SKLEARNGAUSSIANPROCESSKERNELSDOTPRODUCT

- ILLUSTRATION OF GAUSSIAN PROCESS CLASSIFICATION GPC ON THE XOR DATASET
- ILLUSTRATION OF PRIOR AND POSTERIOR GAUSSIAN PROCESS FOR DIFFERENT KERNELS
- ISOPROBABILITY LINES FOR GAUSSIAN PROCESSES CLASSIFICATION GPC

6176SKLEARNGAUSSIANPROCESSKERNELS EXPSINESQUARED

CLASSSSKLEARNGAUSSIANPROCESSKERNELS EXPSINESQUARED LENGTHSCALE10

PERIODICITY10

LENGTHSCALEBOUNDS1E

05 1000000

PERIODICITYBOUNDS1E05

1000000

EXPSINESQUARED KERNEL

1814 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

THE EXPSINESQUARED KERNEL ALLOWS MODELING PERIODIC FUNCTIONS IT IS PARAMETERIZED BY A LENGTHSCALE PARAMETER LENGTHSCALE0 AND A PERIODICITY PARAMETER PERIODICITY0 ONLY THE ISOTROPIC VARIANT WHERE L IS A SCALAR IS SUPPORTED AT THE MOMENT THE KERNEL GIVEN BY

$K(x_i, x_j) = \exp(-\frac{1}{2} \frac{\|x_i - x_j\|^2}{L^2}) \cdot \sin(\frac{2\pi}{P} \frac{\|x_i - x_j\|}{L})$

NEW IN VERSION 018

PARAMETERS

LENGTHSCALE FLOAT 0 DEFAULT 10 THE LENGTH SCALE OF THE KERNEL

PERIODICITY FLOAT 0 DEFAULT 10 THE PERIODICITY OF THE KERNEL

LENGTHSCALEBOUNDS PAIR OF FLOATS 0 DEFAULT 1E5 1E5 THE LOWER AND UPPER BOUND ON LENGTHSCALE

PERIODICITYBOUNDS PAIR OF FLOATS 0 DEFAULT 1E5 1E5 THE LOWER AND UPPER BOUND ON PERIODICITY

ATTRIBUTES

BOUNDS RETURNS THE LOGTRANSFORMED BOUNDS ON THE THETA

HYPERPARAMETERLENGTHSCALE

HYPERPARAMETERPERIODICITY

HYPERPARAMETERS RETURNS A LIST OF ALL HYPERPARAMETER SPECIFICATIONS

NDIMS RETURNS THE NUMBER OF NONFIXED HYPERPARAMETERS OF THE KERNEL

THETA RETURNS THE FLATTENED LOGTRANSFORMED NONFIXED HYPERPARAMETERS

METHODS

CALL SELF X Y EVALGRADIENT RETURN THE KERNEL KX Y AND OPTIONALLY ITS GRADIENT

CLONEWITHTHETA SELF THETA RETURNS A CLONE OF SELF WITH GIVEN HYPERPARAMETERS THETA

DIAG SELF X RETURNS THE DIAGONAL OF THE KERNEL KX X

GETPARAMS SELF DEEP GET PARAMETERS OF THIS KERNEL

ISSTATIONARY SELF RETURNS WHETHER THE KERNEL IS STATIONARY

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS KERNEL

INIT SELFLENGTHSCALE10 PERIODICITY10 LENGTHSCALEBOUNDS1E05 1000000 PERIODICITYBOUNDS1E05 1000000

CALL SELFXYNONE EVALGRADIENTFALSE

RETURN THE KERNEL KX Y AND OPTIONALLY ITS GRADIENT

PARAMETERS

XARRAY SHAPE NSAMPLESX NFEATURES LEFT ARGUMENT OF THE RETURNED KERNEL KX Y

YARRAY SHAPE NSAMPLESY NFEATURES OPTIONAL DEFAULTNONE RIGHT ARGUMENT OF THE RETURNED KERNEL KX Y IF NONE KX X IF EVALUATED INSTEAD

EVALGRADIENT BOOL OPTIONAL DEFAULTFALSE DETERMINES WHETHER THE GRADIENT WITH RESPECT TO THE KERNEL HYPERPARAMETER IS DETERMINED ONLY SUPPORTED WHEN Y IS NONE

RETURNS

617SKLEARNGAUSSIANPROCESS GAUSSIAN PROCESSES 1815

SCIKITLEARN USER GUIDE RELEASE 0213

KARRAY SHAPE NSAMPLESX NSAMPLESY KERNEL KX Y

KGRADIENT ARRAY OPT SHAPE NSAMPLESX NSAMPLESX NDIMS THE GRADIENT OF THE  
KERNEL KX X WITH RESPECT TO THE HYPERPARAMETER OF THE KERNEL ONLY RETURNED WHEN  
EVALGRADIENT IS TRUE

BOUNDS

RETURNS THE LOGTRANSFORMED BOUNDS ON THE THETA

RETURNS

BOUNDS ARRAY SHAPE NDIMS 2 THE LOGTRANSFORMED BOUNDS ON THE KERNEL’S HYPERPARAM  
ETERS THETA

CLONEWITHTHETA SELFTHETA

RETURNS A CLONE OF SELF WITH GIVEN HYPERPARAMETERS THETA

PARAMETERS

THETA ARRAY SHAPE NDIMS THE HYPERPARAMETERS

DIAGSELF

RETURNS THE DIAGONAL OF THE KERNEL KX X

THE RESULT OF THIS METHOD IS IDENTICAL TO NPDIAGSELF HOWEVER IT CAN BE EVALUATED MORE EFFICIENTLY  
SINCE ONLY THE DIAGONAL IS EVALUATED

PARAMETERS

XARRAY SHAPE NSAMPLESX NFEATURES LEFT ARGUMENT OF THE RETURNED KERNEL KX Y

RETURNS

KDIAG ARRAY SHAPE NSAMPLESX DIAGONAL OF KERNEL KX X

GETPARAMS SELFDEEPTURE

GET PARAMETERS OF THIS KERNEL

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

HYPERPARAMETERS

RETURNS A LIST OF ALL HYPERPARAMETER SPECIFICATIONS

ISSTATIONARY SELF

RETURNS WHETHER THE KERNEL IS STATIONARY

NDIMS

RETURNS THE NUMBER OF NONFIXED HYPERPARAMETERS OF THE KERNEL

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS KERNEL

THE METHOD WORKS ON SIMPLE KERNELS AS WELL AS ON NESTED KERNELS THE LATTER HAVE PARAMETERS OF THE FORM  
COMPONENTPARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

1816 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

THETA

RETURNS THE FLATTENED LOGTRANSFORMED NONFIXED HYPERPARAMETERS

NOTE THAT THETA ARE TYPICALLY THE LOGTRANSFORMED VALUES OF THE KERNEL'S HYPERPARAMETERS AS THIS REPRESENTATION OF THE SEARCH SPACE IS MORE AMENABLE FOR HYPERPARAMETER SEARCH AS HYPERPARAMETERS LIKE LENGTHSCALES NATURALLY LIVE ON A LOGSCALE

RETURNS

THETA ARRAY SHAPE NDIMS THE NONFIXED LOGTRANSFORMED HYPERPARAMETERS OF THE KERNEL

EXAMPLES USING SKLEARNGAUSSIANPROCESSKERNELSEXPSINESQUARED

- COMPARISON OF KERNEL RIDGE AND GAUSSIAN PROCESS REGRESSION
- ILLUSTRATION OF PRIOR AND POSTERIOR GAUSSIAN PROCESS FOR DIFFERENT KERNELS
- GAUSSIAN PROCESS REGRESSION GPR ON MAUNA LOA CO2 DATA

6177SKLEARNGAUSSIANPROCESSKERNELS EXPONENTIATION

CLASSSSKLEARNGAUSSIANPROCESSKERNELS EXPONENTIATION KERNEL EXPONENT

EXPONENTIATE KERNEL BY GIVEN EXPONENT

THE RESULTING KERNEL IS DEFINED AS  $K_{EXP} X Y = K_X Y^{EXPONENT}$

NEW IN VERSION 018

PARAMETERS

KERNEL KERNEL OBJECT THE BASE KERNEL

EXPONENT FLOAT THE EXPONENT FOR THE BASE KERNEL

ATTRIBUTES

BOUNDS RETURNS THE LOGTRANSFORMED BOUNDS ON THE THETA

HYPERPARAMETERS RETURNS A LIST OF ALL HYPERPARAMETER

NDIMS RETURNS THE NUMBER OF NONFIXED HYPERPARAMETERS OF THE KERNEL

THETA RETURNS THE FLATTENED LOGTRANSFORMED NONFIXED HYPERPARAMETERS

METHODS

CALL SELF X Y EVALGRADIENT RETURN THE KERNEL  $K_X Y$  AND OPTIONALLY ITS GRADIENT

CLONEWITHTHETA SELF THETA RETURNS A CLONE OF SELF WITH GIVEN HYPERPARAMETERS

THETA

DIAG SELF X RETURNS THE DIAGONAL OF THE KERNEL  $K_X X$

GETPARAMS SELF DEEP GET PARAMETERS OF THIS KERNEL

ISSTATIONARY SELF RETURNS WHETHER THE KERNEL IS STATIONARY

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS KERNEL

INIT SELFKERNEL EXPONENT

CALL SELFXYNONE EVALGRADIENTFALSE

RETURN THE KERNEL  $K_X Y$  AND OPTIONALLY ITS GRADIENT

617SKLEARNGAUSSIANPROCESS GAUSSIAN PROCESSES 1817

SCIKITLEARN USER GUIDE RELEASE 0213

PARAMETERS

XARRAY SHAPE NSAMPLESX NFEATURES LEFT ARGUMENT OF THE RETURNED KERNEL KX Y  
YARRAY SHAPE NSAMPLESY NFEATURES OPTIONAL DEFAULTNONE RIGHT ARGUMENT OF THE  
RETURNED KERNEL KX Y IF NONE KX X IF EVALUATED INSTEAD  
EVALGRADIENT BOOL OPTIONAL DEFAULTFALSE DETERMINES WHETHER THE GRADIENT WITH RE  
SPECT TO THE KERNEL HYPERPARAMETER IS DETERMINED

RETURNS

KARRAY SHAPE NSAMPLESX NSAMPLESY KERNEL KX Y  
KGRADIENT ARRAY OPT SHAPE NSAMPLESX NSAMPLESX NDIMS THE GRADIENT OF THE  
KERNEL KX X WITH RESPECT TO THE HYPERPARAMETER OF THE KERNEL ONLY RETURNED WHEN  
EVALGRADIENT IS TRUE

BOUNDS

RETURNS THE LOGTRANSFORMED BOUNDS ON THE THETA  
RETURNS  
BOUNDS ARRAY SHAPE NDIMS 2 THE LOGTRANSFORMED BOUNDS ON THE KERNEL’S HYPERPARAM  
ETERS THETA

CLONEWITHTHETA SELFTHETA

RETURNS A CLONE OF SELF WITH GIVEN HYPERPARAMETERS THETA

PARAMETERS

THETA ARRAY SHAPE NDIMS THE HYPERPARAMETERS

DIAGSELF

RETURNS THE DIAGONAL OF THE KERNEL KX X  
THE RESULT OF THIS METHOD IS IDENTICAL TO NPDIAGSELF HOWEVER IT CAN BE EVALUATED MORE EFFICIENTLY  
SINCE ONLY THE DIAGONAL IS EVALUATED

PARAMETERS

XARRAY SHAPE NSAMPLESX NFEATURES LEFT ARGUMENT OF THE RETURNED KERNEL KX Y

RETURNS

KDIAG ARRAY SHAPE NSAMPLESX DIAGONAL OF KERNEL KX X

GETPARAMS SELFDEEPT

GET PARAMETERS OF THIS KERNEL

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

HYPERPARAMETERS

RETURNS A LIST OF ALL HYPERPARAMETER

ISSTATIONARY SELF

RETURNS WHETHER THE KERNEL IS STATIONARY

NDIMS

RETURNS THE NUMBER OF NONFIXED HYPERPARAMETERS OF THE KERNEL



SCIKITLEARN USER GUIDE RELEASE 0213

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS KERNEL

THE METHOD WORKS ON SIMPLE KERNELS AS WELL AS ON NESTED KERNELS THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

THETA

RETURNS THE FLATTENED LOGTRANSFORMED NONFIXED HYPERPARAMETERS

NOTE THAT THETA ARE TYPICALLY THE LOGTRANSFORMED VALUES OF THE KERNEL’S HYPERPARAMETERS AS THIS REPRESENTATION OF THE SEARCH SPACE IS MORE AMENABLE FOR HYPERPARAMETER SEARCH AS HYPERPARAMETERS LIKE LENGTHSCALES NATURALLY LIVE ON A LOGSCALE

RETURNS

THETA ARRAY SHAPE NDIMS THE NONFIXED LOGTRANSFORMED HYPERPARAMETERS OF THE KERNEL

6178SKLEARNGAUSSIANPROCESSKERNELS HYPERPARAMETER

CLASSSSKLEARNGAUSSIANPROCESSKERNELS HYPERPARAMETER

A KERNEL HYPERPARAMETER’S SPECIFICATION IN FORM OF A NAMEDTUPLE

NEW IN VERSION 018

ATTRIBUTES

NAME STRING ALIAS FOR FIELD NUMBER 0

VALUETYPE STRING ALIAS FOR FIELD NUMBER 1

BOUNDS PAIR OF FLOATS 0 OR “FIXED” ALIAS FOR FIELD NUMBER 2

NELEMENTS INT DEFAULT1 ALIAS FOR FIELD NUMBER 3

FIXED BOOL DEFAULT NONE ALIAS FOR FIELD NUMBER 4

METHODS

COUNT

INDEX RAISES VALUEERROR IF THE VALUE IS NOT PRESENT

INIT SELFARGS KWARGS

INITIALIZE SELF SEE HELPTYPESELF FOR ACCURATE SIGNATURE

CALL ARGS KWARGS

CALL SELF AS A FUNCTION

BOUNDS

ALIAS FOR FIELD NUMBER 2

COUNT

FIXED

ALIAS FOR FIELD NUMBER 4

INDEX

617SKLEARNGAUSSIANPROCESS GAUSSIAN PROCESSES 1819

SCIKITLEARN USER GUIDE RELEASE 0213  
RAISES VALUEERROR IF THE VALUE IS NOT PRESENT  
NELEMENTS  
ALIAS FOR FIELD NUMBER 3  
NAME  
ALIAS FOR FIELD NUMBER 0  
VALUETYPE  
ALIAS FOR FIELD NUMBER 1  
6179SKLEARNGAUSSIANPROCESSKERNELS KERNEL  
CLASSSKLEARNGAUSSIANPROCESSKERNELS KERNEL  
BASE CLASS FOR ALL KERNELS  
NEW IN VERSION 018  
ATTRIBUTES  
BOUNDS RETURNS THE LOGTRANSFORMED BOUNDS ON THE THETA  
HYPERPARAMETERS RETURNS A LIST OF ALL HYPERPARAMETER SPECIFICATIONS  
NDIMS RETURNS THE NUMBER OF NONFIXED HYPERPARAMETERS OF THE KERNEL  
THETA RETURNS THE FLATTENED LOGTRANSFORMED NONFIXED HYPERPARAMETERS  
METHODS  
CALL SELF X Y EVALGRADIENT EVALUATE THE KERNEL  
CLONewithTHETA SELF THETA RETURNS A CLONE OF SELF WITH GIVEN HYPERPARAMETERS  
THETA  
DIAG SELF X RETURNS THE DIAGONAL OF THE KERNEL KX X  
GETPARAMS SELF DEEP GET PARAMETERS OF THIS KERNEL  
ISSTATIONARY SELF RETURNS WHETHER THE KERNEL IS STATIONARY  
SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS KERNEL  
INIT SELFARGS KWARGS  
INITIALIZE SELF SEE HELPTYPESELF FOR ACCURATE SIGNATURE  
CALL SELFXYNONE EVALGRADIENTFALSE  
EVALUATE THE KERNEL  
BOUNDS  
RETURNS THE LOGTRANSFORMED BOUNDS ON THE THETA  
RETURNS  
BOUNDS ARRAY SHAPE NDIMS 2 THE LOGTRANSFORMED BOUNDS ON THE KERNEL’S HYPERPARAM  
ETERS THETA  
CLONewithTHETA SELFTHETA  
RETURNS A CLONE OF SELF WITH GIVEN HYPERPARAMETERS THETA  
PARAMETERS  
THETA ARRAY SHAPE NDIMS THE HYPERPARAMETERS  
1820 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

DIAGSELF  
RETURNS THE DIAGONAL OF THE KERNEL KX X  
THE RESULT OF THIS METHOD IS IDENTICAL TO NPDIAGSELF HOWEVER IT CAN BE EVALUATED MORE EFFICIENTLY  
SINCE ONLY THE DIAGONAL IS EVALUATED  
PARAMETERS  
XARRAY SHAPE NSAMPLESX NFEATURES LEFT ARGUMENT OF THE RETURNED KERNEL KX Y  
RETURNS  
KDIAG ARRAY SHAPE NSAMPLESX DIAGONAL OF KERNEL KX X  
GETPARAMS SELFDEEPTURE  
GET PARAMETERS OF THIS KERNEL  
PARAMETERS  
DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBOBJECTS THAT ARE ESTIMATORS  
RETURNS  
PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES  
HYPERPARAMETERS  
RETURNS A LIST OF ALL HYPERPARAMETER SPECIFICATIONS  
ISSTATIONARY SELF  
RETURNS WHETHER THE KERNEL IS STATIONARY  
NDIMS  
RETURNS THE NUMBER OF NONFIXED HYPERPARAMETERS OF THE KERNEL  
SETPARAMS SELFPARAMS  
SET THE PARAMETERS OF THIS KERNEL  
THE METHOD WORKS ON SIMPLE KERNELS AS WELL AS ON NESTED KERNELS THE LATTER HAVE PARAMETERS OF THE FORM  
COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT  
RETURNS  
SELF  
THETA  
RETURNS THE FLATTENED LOGTRANSFORMED NONFIXED HYPERPARAMETERS  
NOTE THAT THETA ARE TYPICALLY THE LOGTRANSFORMED VALUES OF THE KERNEL'S HYPERPARAMETERS AS THIS REPRESENTA  
TION OF THE SEARCH SPACE IS MORE AMENABLE FOR HYPERPARAMETER SEARCH AS HYPERPARAMETERS LIKE LENGTHSCALES  
NATURALLY LIVE ON A LOGSCALE  
RETURNS  
THETA ARRAY SHAPE NDIMS THE NONFIXED LOGTRANSFORMED HYPERPARAMETERS OF THE KERNEL  
61710SKLEARNGAUSSIANPROCESSKERNELS MATERN  
CLASSSKLEARNGAUSSIANPROCESSKERNELS MATERNLENGTHSCALE10 LENGTHSCALEBOUNDS1E  
051000000 NU15  
MATERN KERNEL  
THE CLASS OF MATERN KERNELS IS A GENERALIZATION OF THE RBF AND THE ABSOLUTE EXPONENTIAL KERNEL PARAMETERIZED BY  
AN ADDITIONAL PARAMETER NU THE SMALLER NU THE LESS SMOOTH THE APPROXIMATED FUNCTION IS FOR NUINF THE KERNEL  
617SKLEARNGAUSSIANPROCESS GAUSSIAN PROCESSES 1821

SCIKITLEARN USER GUIDE RELEASE 0213

BECOMES EQUIVALENT TO THE RBF KERNEL AND FOR NU05 TO THE ABSOLUTE EXPONENTIAL KERNEL IMPORTANT INTERMEDIATE VALUES ARE NU15 ONCE DIFFERENTIABLE FUNCTIONS AND NU25 TWICE DIFFERENTIABLE FUNCTIONS

SEE RASMUSSEN AND WILLIAMS 2006 PP84 FOR DETAILS REGARDING THE DIFFERENT VARIANTS OF THE MATERN KERNEL

NEW IN VERSION 018

PARAMETERS

LENGTHSCALE FLOAT OR ARRAY WITH SHAPE NFEATURES DEFAULT 10 THE LENGTH SCALE OF THE KERNEL

IF A FLOAT AN ISOTROPIC KERNEL IS USED IF AN ARRAY AN ANISOTROPIC KERNEL IS USED WHERE EACH DIMENSION OF L DEFINES THE LENGTHSCALE OF THE RESPECTIVE FEATURE DIMENSION

LENGTHSCALEBOUNDS PAIR OF FLOATS 0 DEFAULT 1E5 1E5 THE LOWER AND UPPER BOUND ON LENGTHSCALE

NUFLOAT DEFAULT 15 THE PARAMETER NU CONTROLLING THE SMOOTHNESS OF THE LEARNED FUNCTION

THE SMALLER NU THE LESS SMOOTH THE APPROXIMATED FUNCTION IS FOR NUINF THE KERNEL BECOMES EQUIVALENT TO THE RBF KERNEL AND FOR NU05 TO THE ABSOLUTE EXPONENTIAL KERNEL IM

PORTANT INTERMEDIATE VALUES ARE NU15 ONCE DIFFERENTIABLE FUNCTIONS AND NU25 TWICE DIFFERENTIABLE FUNCTIONS NOTE THAT VALUES OF NU NOT IN 05 15 25 INF INCUR A CONSID

ERABLY HIGHER COMPUTATIONAL COST APPR 10 TIMES HIGHER SINCE THEY REQUIRE TO EVALUATE THE MODIFIED BESSEL FUNCTION FURTHERMORE IN CONTRAST TO L NU IS KEPT FIXED TO ITS INITIAL VALUE AND NOT OPTIMIZED

ATTRIBUTES

ANISOTROPIC

BOUNDS RETURNS THE LOGTRANSFORMED BOUNDS ON THE THETA

HYPERPARAMETERLENGTHSCALE

HYPERPARAMETERS RETURNS A LIST OF ALL HYPERPARAMETER SPECIFICATIONS

NDIMS RETURNS THE NUMBER OF NONFIXED HYPERPARAMETERS OF THE KERNEL

THETA RETURNS THE FLATTENED LOGTRANSFORMED NONFIXED HYPERPARAMETERS

METHODS

CALL SELF X Y EVALGRADIENT RETURN THE KERNEL KX Y AND OPTIONALLY ITS GRADIENT

CLONewithTHETA SELF THETA RETURNS A CLONE OF SELF WITH GIVEN HYPERPARAMETERS

THETA

DIAG SELF X RETURNS THE DIAGONAL OF THE KERNEL KX X

GETPARAMS SELF DEEP GET PARAMETERS OF THIS KERNEL

ISSTATIONARY SELF RETURNS WHETHER THE KERNEL IS STATIONARY

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS KERNEL

INIT SELFLENGTHSCALE10 LENGTHSCALEBOUNDS1E05 1000000 NU15

CALL SELFXYNONE EVALGRADIENTFALSE

RETURN THE KERNEL KX Y AND OPTIONALLY ITS GRADIENT

PARAMETERS

XARRAY SHAPE NSAMPLESX NFEATURES LEFT ARGUMENT OF THE RETURNED KERNEL KX Y

YARRAY SHAPE NSAMPLESY NFEATURES OPTIONAL DEFAULTNONE RIGHT ARGUMENT OF THE RETURNED KERNEL KX Y IF NONE KX X IF EVALUATED INSTEAD

1822 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

EVALGRADIENT BOOL OPTIONAL DEFAULTFALSE DETERMINES WHETHER THE GRADIENT WITH RESPECT TO THE KERNEL HYPERPARAMETER IS DETERMINED ONLY SUPPORTED WHEN Y IS NONE

RETURNS

KARRAY SHAPE NSAMPLESX NSAMPLESY KERNEL KX Y

KGRADIENT ARRAY OPT SHAPE NSAMPLESX NSAMPLESX NDIMS THE GRADIENT OF THE KERNEL KX X WITH RESPECT TO THE HYPERPARAMETER OF THE KERNEL ONLY RETURNED WHEN EVALGRADIENT IS TRUE

BOUNDS

RETURNS THE LOGTRANSFORMED BOUNDS ON THE THETA

RETURNS

BOUNDS ARRAY SHAPE NDIMS 2 THE LOGTRANSFORMED BOUNDS ON THE KERNEL’S HYPERPARAMETERS THETA

CLONEWITHTHETA SELFTHETA

RETURNS A CLONE OF SELF WITH GIVEN HYPERPARAMETERS THETA

PARAMETERS

THETA ARRAY SHAPE NDIMS THE HYPERPARAMETERS

DIAGSELF

RETURNS THE DIAGONAL OF THE KERNEL KX X

THE RESULT OF THIS METHOD IS IDENTICAL TO NPDIASELF HOWEVER IT CAN BE EVALUATED MORE EFFICIENTLY SINCE ONLY THE DIAGONAL IS EVALUATED

PARAMETERS

XARRAY SHAPE NSAMPLESX NFEATURES LEFT ARGUMENT OF THE RETURNED KERNEL KX Y

RETURNS

KDIAG ARRAY SHAPE NSAMPLESX DIAGONAL OF KERNEL KX X

GETPARAMS SELFDEEPTURE

GET PARAMETERS OF THIS KERNEL

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

HYPERPARAMETERS

RETURNS A LIST OF ALL HYPERPARAMETER SPECIFICATIONS

ISSTATIONARY SELF

RETURNS WHETHER THE KERNEL IS STATIONARY

NDIMS

RETURNS THE NUMBER OF NONFIXED HYPERPARAMETERS OF THE KERNEL

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS KERNEL

THE METHOD WORKS ON SIMPLE KERNELS AS WELL AS ON NESTED KERNELS THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

617SKLEARNGAUSSIANPROCESS GAUSSIAN PROCESSES 1823

SCIKITLEARN USER GUIDE RELEASE 0213

RETURNS  
SELF  
THETA  
RETURNS THE FLATTENED LOGTRANSFORMED NONFIXED HYPERPARAMETERS  
NOTE THAT THETA ARE TYPICALLY THE LOGTRANSFORMED VALUES OF THE KERNEL’S HYPERPARAMETERS AS THIS REPRESENTATION OF THE SEARCH SPACE IS MORE AMENABLE FOR HYPERPARAMETER SEARCH AS HYPERPARAMETERS LIKE LENGTHSCALES NATURALLY LIVE ON A LOGSCALE

RETURNS  
THETA ARRAY SHAPE NDIMS THE NONFIXED LOGTRANSFORMED HYPERPARAMETERS OF THE KERNEL

EXAMPLES USING SKLEARNGAUSSIANPROCESSKERNELSMATERN

- ILLUSTRATION OF PRIOR AND POSTERIOR GAUSSIAN PROCESS FOR DIFFERENT KERNELS

61711SKLEARNGAUSSIANPROCESSKERNELS PAIRWISEKERNEL  
CLASSSSKLEARNGAUSSIANPROCESSKERNELS PAIRWISEKERNEL GAMMA10  
GAMMABOUNDS1E05  
1000000 METRIC’LINEAR’ PAIR  
WISEKERNELSKWARGSNONE

WRAPPER FOR KERNELS IN SKLEARNMETRICSPAIRWISE  
A THIN WRAPPER AROUND THE FUNCTIONALITY OF THE KERNELS IN SKLEARNMETRICSPAIRWISE  
NOTE EVALUATION OF EVALGRADIENT IS NOT ANALYTIC BUT NUMERIC AND ALL KERNELS SUPPORT ONLY ISOTROPIC DISTANCES THE PARAMETER GAMMA IS CONSIDERED TO BE A HYPERPARAMETER AND MAY BE OPTIMIZED THE OTHER KERNEL PARAMETERS ARE SET DIRECTLY AT INITIALIZATION AND ARE KEPT FIXED

NEW IN VERSION 018  
PARAMETERS  
GAMMA FLOAT 0 DEFAULT 10 PARAMETER GAMMA OF THE PAIRWISE KERNEL SPECIFIED BY METRIC  
GAMMABOUNDS PAIR OF FLOATS 0 DEFAULT 1E5 1E5 THE LOWER AND UPPER BOUND ON  
GAMMA  
METRIC STRING OR CALLABLE DEFAULT “LINEAR” THE METRIC TO USE WHEN CALCULATING KERNEL BETWEEN INSTANCES IN A FEATURE ARRAY IF METRIC IS A STRING IT MUST BE ONE OF THE METRICS IN PAIRWISEPAIRWISEKERNELFUNCTIONS IF METRIC IS “PRECOMPUTED” X IS ASSUMED TO BE A KERNEL MATRIX ALTERNATIVELY IF METRIC IS A CALLABLE FUNCTION IT IS CALLED ON EACH PAIR OF INSTANCES ROWS AND THE RESULTING VALUE RECORDED THE CALLABLE SHOULD TAKE TWO ARRAYS FROM X AS INPUT AND RETURN A VALUE INDICATING THE DISTANCE BETWEEN THEM  
PAIRWISEKERNELSKWARGS DICT DEFAULT NONE ALL ENTRIES OF THIS DICT IF ANY ARE PASSED AS KEYWORD ARGUMENTS TO THE PAIRWISE KERNEL FUNCTION

ATTRIBUTES  
BOUNDS RETURNS THE LOGTRANSFORMED BOUNDS ON THE THETA  
HYPERPARAMETERGAMMA  
HYPERPARAMETERS RETURNS A LIST OF ALL HYPERPARAMETER SPECIFICATIONS  
NDIMS RETURNS THE NUMBER OF NONFIXED HYPERPARAMETERS OF THE KERNEL

1824 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

THETA RETURNS THE FLATTENED LOGTRANSFORMED NONFIXED HYPERPARAMETERS

METHODS

CALL SELF X Y EVALGRADIENT RETURN THE KERNEL KX Y AND OPTIONALLY ITS GRADIENT

CLONewithTHETA SELF THETA RETURNS A CLONE OF SELF WITH GIVEN HYPERPARAMETERS

THETA

DIAG SELF X RETURNS THE DIAGONAL OF THE KERNEL KX X

GETPARAMS SELF DEEP GET PARAMETERS OF THIS KERNEL

ISSTATIONARY SELF RETURNS WHETHER THE KERNEL IS STATIONARY

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS KERNEL

INIT SELF GAMMA10 GAMMABOUNDS1E05 1000000 METRIC'LINEAR' PAIR

WISEKERNELSKWARGSNONE

CALL SELF X Y NONE EVALGRADIENT FALSE

RETURN THE KERNEL KX Y AND OPTIONALLY ITS GRADIENT

PARAMETERS

XARRAY SHAPE NSAMPLESX NFEATURES LEFT ARGUMENT OF THE RETURNED KERNEL KX Y

YARRAY SHAPE NSAMPLESY NFEATURES OPTIONAL DEFAULT NONE RIGHT ARGUMENT OF THE

RETURNED KERNEL KX Y IF NONE KX X IF EVALUATED INSTEAD

EVALGRADIENT BOOL OPTIONAL DEFAULT FALSE DETERMINES WHETHER THE GRADIENT WITH RE

SPECT TO THE KERNEL HYPERPARAMETER IS DETERMINED ONLY SUPPORTED WHEN Y IS NONE

RETURNS

KARRAY SHAPE NSAMPLESX NSAMPLESY KERNEL KX Y

KGRADIENT ARRAY OPT SHAPE NSAMPLESX NSAMPLESX NDIMS THE GRADIENT OF THE

KERNEL KX X WITH RESPECT TO THE HYPERPARAMETER OF THE KERNEL ONLY RETURNED WHEN

EVALGRADIENT IS TRUE

BOUNDS

RETURNS THE LOGTRANSFORMED BOUNDS ON THE THETA

RETURNS

BOUNDS ARRAY SHAPE NDIMS 2 THE LOGTRANSFORMED BOUNDS ON THE KERNEL'S HYPERPARAM

ETERS THETA

CLONewithTHETA SELF THETA

RETURNS A CLONE OF SELF WITH GIVEN HYPERPARAMETERS THETA

PARAMETERS

THETA ARRAY SHAPE NDIMS THE HYPERPARAMETERS

DIAGSELF X

RETURNS THE DIAGONAL OF THE KERNEL KX X

THE RESULT OF THIS METHOD IS IDENTICAL TO NPDIAGSELF X HOWEVER IT CAN BE EVALUATED MORE EFFICIENTLY

SINCE ONLY THE DIAGONAL IS EVALUATED

PARAMETERS

XARRAY SHAPE NSAMPLESX NFEATURES LEFT ARGUMENT OF THE RETURNED KERNEL KX Y

617SKLEARN GAUSSIAN PROCESS GAUSSIAN PROCESSES 1825

SCIKITLEARN USER GUIDE RELEASE 0213

RETURNS

KDIAG ARRAY SHAPE NSAMPLESX DIAGONAL OF KERNEL KX X

GETPARAMS SELFDEEPTTRUE

GET PARAMETERS OF THIS KERNEL

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

HYPERPARAMETERS

RETURNS A LIST OF ALL HYPERPARAMETER SPECIFICATIONS

ISSTATIONARY SELF

RETURNS WHETHER THE KERNEL IS STATIONARY

NDIMS

RETURNS THE NUMBER OF NONFIXED HYPERPARAMETERS OF THE KERNEL

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS KERNEL

THE METHOD WORKS ON SIMPLE KERNELS AS WELL AS ON NESTED KERNELS THE LATTER HAVE PARAMETERS OF THE FORM

COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

THETA

RETURNS THE FLATTENED LOGTRANSFORMED NONFIXED HYPERPARAMETERS

NOTE THAT THETA ARE TYPICALLY THE LOGTRANSFORMED VALUES OF THE KERNEL'S HYPERPARAMETERS AS THIS REPRESENTA

TION OF THE SEARCH SPACE IS MORE AMENABLE FOR HYPERPARAMETER SEARCH AS HYPERPARAMETERS LIKE LENGTHSCALES

NATURALLY LIVE ON A LOGSCALE

RETURNS

THETA ARRAY SHAPE NDIMS THE NONFIXED LOGTRANSFORMED HYPERPARAMETERS OF THE KERNEL

61712SKLEARNGAUSSIANPROCESSKERNELS PRODUCT

CLASSSSKLEARNGAUSSIANPROCESSKERNELS PRODUCTK1K2

PRODUCTKERNEL K1 K2 OF TWO KERNELS K1 AND K2

THE RESULTING KERNEL IS DEFINED AS KPRODX Y K1X Y K2X Y

NEW IN VERSION 018

PARAMETERS

K1KERNEL OBJECT THE FIRST BASEKERNEL OF THE PRODUCTKERNEL

K2KERNEL OBJECT THE SECOND BASEKERNEL OF THE PRODUCTKERNEL

ATTRIBUTES

BOUNDS RETURNS THE LOGTRANSFORMED BOUNDS ON THE THETA

HYPERPARAMETERS RETURNS A LIST OF ALL HYPERPARAMETER

1826 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

NDIMS RETURNS THE NUMBER OF NONFIXED HYPERPARAMETERS OF THE KERNEL

THETA RETURNS THE FLATTENED LOGTRANSFORMED NONFIXED HYPERPARAMETERS

METHODS

CALL SELF X Y EVALGRADIENT RETURN THE KERNEL KX Y AND OPTIONALLY ITS GRADIENT

CLONEWITHTHETA SELF THETA RETURNS A CLONE OF SELF WITH GIVEN HYPERPARAMETERS

THETA

DIAG SELF X RETURNS THE DIAGONAL OF THE KERNEL KX X

GETPARAMS SELF DEEP GET PARAMETERS OF THIS KERNEL

ISSTATIONARY SELF RETURNS WHETHER THE KERNEL IS STATIONARY

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS KERNEL

INIT SELF K1 K2

CALL SELF X Y NONE EVALGRADIENT FALSE

RETURN THE KERNEL KX Y AND OPTIONALLY ITS GRADIENT

PARAMETERS

XARRAY SHAPE NSAMPLESX NFEATURES LEFT ARGUMENT OF THE RETURNED KERNEL KX Y

YARRAY SHAPE NSAMPLESY NFEATURES OPTIONAL DEFAULT NONE RIGHT ARGUMENT OF THE

RETURNED KERNEL KX Y IF NONE KX X IF EVALUATED INSTEAD

EVALGRADIENT BOOL OPTIONAL DEFAULT FALSE DETERMINES WHETHER THE GRADIENT WITH RE

SPECT TO THE KERNEL HYPERPARAMETER IS DETERMINED

RETURNS

KARRAY SHAPE NSAMPLESX NSAMPLESY KERNEL KX Y

KGRADIENT ARRAY OPT SHAPE NSAMPLESX NSAMPLESX NDIMS THE GRADIENT OF THE

KERNEL KX X WITH RESPECT TO THE HYPERPARAMETER OF THE KERNEL ONLY RETURNED WHEN

EVALGRADIENT IS TRUE

BOUNDS

RETURNS THE LOGTRANSFORMED BOUNDS ON THE THETA

RETURNS

BOUNDS ARRAY SHAPE NDIMS 2 THE LOGTRANSFORMED BOUNDS ON THE KERNEL'S HYPERPARAM

ETERS THETA

CLONEWITHTHETA SELF THETA

RETURNS A CLONE OF SELF WITH GIVEN HYPERPARAMETERS THETA

PARAMETERS

THETA ARRAY SHAPE NDIMS THE HYPERPARAMETERS

DIAGSELF X

RETURNS THE DIAGONAL OF THE KERNEL KX X

THE RESULT OF THIS METHOD IS IDENTICAL TO NPDIAGSELF X HOWEVER IT CAN BE EVALUATED MORE EFFICIENTLY

SINCE ONLY THE DIAGONAL IS EVALUATED

PARAMETERS

XARRAY SHAPE NSAMPLESX NFEATURES LEFT ARGUMENT OF THE RETURNED KERNEL KX Y

6175 SKLEARN GAUSSIAN PROCESS GAUSSIAN PROCESSES 1827

SCIKITLEARN USER GUIDE RELEASE 0213

RETURNS

KDIAG ARRAY SHAPE NSAMPLESX DIAGONAL OF KERNEL KX X

GETPARAMS SELFDEEPTURE

GET PARAMETERS OF THIS KERNEL

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

HYPERPARAMETERS

RETURNS A LIST OF ALL HYPERPARAMETER

ISSTATIONARY SELF

RETURNS WHETHER THE KERNEL IS STATIONARY

NDIMS

RETURNS THE NUMBER OF NONFIXED HYPERPARAMETERS OF THE KERNEL

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS KERNEL

THE METHOD WORKS ON SIMPLE KERNELS AS WELL AS ON NESTED KERNELS THE LATTER HAVE PARAMETERS OF THE FORM

COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

THETA

RETURNS THE FLATTENED LOGTRANSFORMED NONFIXED HYPERPARAMETERS

NOTE THAT THETA ARE TYPICALLY THE LOGTRANSFORMED VALUES OF THE KERNEL'S HYPERPARAMETERS AS THIS REPRESENTA

TION OF THE SEARCH SPACE IS MORE AMENABLE FOR HYPERPARAMETER SEARCH AS HYPERPARAMETERS LIKE LENGTHSCALES

NATURALLY LIVE ON A LOGSCALE

RETURNS

THETA ARRAY SHAPE NDIMS THE NONFIXED LOGTRANSFORMED HYPERPARAMETERS OF THE KERNEL

61713SKLEARNGAUSSIANPROCESSKERNELS RBF

CLASSSSKLEARNGAUSSIANPROCESSKERNELS RBFLENGTHSCALE10 LENGTHSCALEBOUNDS1E05

1000000

RADIALBASIS FUNCTION KERNEL AKA SQUAREDEXPONENTIAL KERNEL

THE RBF KERNEL IS A STATIONARY KERNEL IT IS ALSO KNOWN AS THE "SQUARED EXPONENTIAL" KERNEL IT IS PARAMETERIZED

BY A LENGTHSCALE PARAMETER LENGTHSCALE0 WHICH CAN EITHER BE A SCALAR ISOTROPIC VARIANT OF THE KERNEL OR A

VECTOR WITH THE SAME NUMBER OF DIMENSIONS AS THE INPUTS X ANISOTROPIC VARIANT OF THE KERNEL THE KERNEL IS

GIVEN BY

KXI XJ EXP1 2 DXI LENGTHSCALE XJ LENGTHSCALE2

THIS KERNEL IS INFINITELY DIFFERENTIABLE WHICH IMPLIES THAT GPS WITH THIS KERNEL AS COVARIANCE FUNCTION HAVE MEAN

SQUARE DERIVATIVES OF ALL ORDERS AND ARE THUS VERY SMOOTH

NEW IN VERSION 018

1828 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

PARAMETERS

LENGTHSCALE FLOAT OR ARRAY WITH SHAPE NFEATURES DEFAULT 10 THE LENGTH SCALE OF THE KERNEL  
IF A FLOAT AN ISOTROPIC KERNEL IS USED IF AN ARRAY AN ANISOTROPIC KERNEL IS USED WHERE EACH  
DIMENSION OF L DEFINES THE LENGTHSCALE OF THE RESPECTIVE FEATURE DIMENSION  
LENGTHSCALEBOUNDS PAIR OF FLOATS 0 DEFAULT 1E5 1E5 THE LOWER AND UPPER BOUND ON  
LENGTHSCALE

ATTRIBUTES

ANISOTROPIC

BOUNDS RETURNS THE LOGTRANSFORMED BOUNDS ON THE THETA

HYPERPARAMETERLENGTHSCALE

HYPERPARAMETERS RETURNS A LIST OF ALL HYPERPARAMETER SPECIFICATIONS

NDIMS RETURNS THE NUMBER OF NONFIXED HYPERPARAMETERS OF THE KERNEL

THETA RETURNS THE FLATTENED LOGTRANSFORMED NONFIXED HYPERPARAMETERS

METHODS

CALL SELF X Y EVALGRADIENT RETURN THE KERNEL KX Y AND OPTIONALLY ITS GRADIENT  
CLONewithTHETA SELF THETA RETURNS A CLONE OF SELF WITH GIVEN HYPERPARAMETERS

THETA

DIAG SELF X RETURNS THE DIAGONAL OF THE KERNEL KX X

GETPARAMS SELF DEEP GET PARAMETERS OF THIS KERNEL

ISSTATIONARY SELF RETURNS WHETHER THE KERNEL IS STATIONARY

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS KERNEL

INIT SELFLENGTHSCALE10 LENGTHSCALEBOUNDS1E05 1000000

CALL SELFXYNONE EVALGRADIENTFALSE

RETURN THE KERNEL KX Y AND OPTIONALLY ITS GRADIENT

PARAMETERS

XARRAY SHAPE NSAMPLESX NFEATURES LEFT ARGUMENT OF THE RETURNED KERNEL KX Y

YARRAY SHAPE NSAMPLESY NFEATURES OPTIONAL DEFAULTNONE RIGHT ARGUMENT OF THE  
RETURNED KERNEL KX Y IF NONE KX X IF EVALUATED INSTEAD

EVALGRADIENT BOOL OPTIONAL DEFAULTFALSE DETERMINES WHETHER THE GRADIENT WITH RE  
SPECT TO THE KERNEL HYPERPARAMETER IS DETERMINED ONLY SUPPORTED WHEN Y IS NONE

RETURNS

KARRAY SHAPE NSAMPLESX NSAMPLESY KERNEL KX Y

KGRADIENT ARRAY OPT SHAPE NSAMPLESX NSAMPLESX NDIMS THE GRADIENT OF THE  
KERNEL KX X WITH RESPECT TO THE HYPERPARAMETER OF THE KERNEL ONLY RETURNED WHEN  
EVALGRADIENT IS TRUE

BOUNDS

RETURNS THE LOGTRANSFORMED BOUNDS ON THE THETA

RETURNS

617SKLEARNGAUSSIANPROCESS GAUSSIAN PROCESSES 1829

SCIKITLEARN USER GUIDE RELEASE 0213

BOUNDS ARRAY SHAPE NDIMS 2 THE LOGTRANSFORMED BOUNDS ON THE KERNEL’S HYPERPARAMETERS THETA

CLONEWITHTHETA SELFTHETA

RETURNS A CLONE OF SELF WITH GIVEN HYPERPARAMETERS THETA

PARAMETERS

THETA ARRAY SHAPE NDIMS THE HYPERPARAMETERS

DIAGSELF

RETURNS THE DIAGONAL OF THE KERNEL KX X

THE RESULT OF THIS METHOD IS IDENTICAL TO NPDIAGSELF HOWEVER IT CAN BE EVALUATED MORE EFFICIENTLY SINCE ONLY THE DIAGONAL IS EVALUATED

PARAMETERS

XARRAY SHAPE NSAMPLESX NFEATURES LEFT ARGUMENT OF THE RETURNED KERNEL KX Y

RETURNS

KDIAG ARRAY SHAPE NSAMPLESX DIAGONAL OF KERNEL KX X

GETPARAMS SELFDEEPTURE

GET PARAMETERS OF THIS KERNEL

PARAMETERS

DEEPTURE BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

HYPERPARAMETERS

RETURNS A LIST OF ALL HYPERPARAMETER SPECIFICATIONS

ISSTATIONARY SELF

RETURNS WHETHER THE KERNEL IS STATIONARY

NDIMS

RETURNS THE NUMBER OF NONFIXED HYPERPARAMETERS OF THE KERNEL

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS KERNEL

THE METHOD WORKS ON SIMPLE KERNELS AS WELL AS ON NESTED KERNELS THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

THETA

RETURNS THE FLATTENED LOGTRANSFORMED NONFIXED HYPERPARAMETERS

NOTE THAT THETA ARE TYPICALLY THE LOGTRANSFORMED VALUES OF THE KERNEL’S HYPERPARAMETERS AS THIS REPRESENTATION OF THE SEARCH SPACE IS MORE AMENABLE FOR HYPERPARAMETER SEARCH AS HYPERPARAMETERS LIKE LENGTHSCALES NATURALLY LIVE ON A LOGSCALE

RETURNS

THETA ARRAY SHAPE NDIMS THE NONFIXED LOGTRANSFORMED HYPERPARAMETERS OF THE KERNEL

1830 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES USING SKLEARNGAUSSIANPROCESSKERNELSRBF

- PLOT CLASSIFICATION PROBABILITY
- CLASSIFIER COMPARISON
- ILLUSTRATION OF GAUSSIAN PROCESS CLASSIFICATION GPC ON THE XOR DATASET
- GAUSSIAN PROCESS CLASSIFICATION GPC ON IRIS DATASET
- ILLUSTRATION OF PRIOR AND POSTERIOR GAUSSIAN PROCESS FOR DIFFERENT KERNELS
- PROBABILISTIC PREDICTIONS WITH GAUSSIAN PROCESS CLASSIFICATION GPC
- GAUSSIAN PROCESS REGRESSION GPR WITH NOISELEVEL ESTIMATION
- GAUSSIAN PROCESSES REGRESSION BASIC INTRODUCTORY EXAMPLE
- GAUSSIAN PROCESS REGRESSION GPR ON MAUNA LOA CO2 DATA

61714SKLEARNGAUSSIANPROCESSKERNELS RATIONALQUADRATIC

CLASSSSKLEARNGAUSSIANPROCESSKERNELS RATIONALQUADRATIC LENGTHSCALE10

ALPHA10

LENGTHSCALEBOUNDS1E

05 1000000

ALPHABOUNDS1E05

1000000

RATIONAL QUADRATIC KERNEL

THE RATIONALQUADRATIC KERNEL CAN BE SEEN AS A SCALE MIXTURE AN INFINITE SUM OF RBF KERNELS WITH DIFFERENT CHARACTERISTIC LENGTHSCALES IT IS PARAMETERIZED BY A LENGTHSCALE PARAMETER LENGTHSCALE0 AND A SCALE MIXTURE PARAMETER ALPHA0 ONLY THE ISOTROPIC VARIANT WHERE LENGTHSCALE IS A SCALAR IS SUPPORTED AT THE MOMENT THE KERNEL GIVEN BY

$K(x_i, x_j) = \frac{1}{D(x_i, x_j)^2} \frac{2\alpha}{\alpha + \sqrt{1 + \alpha^2 L^2}}$

NEW IN VERSION 018

PARAMETERS

LENGTHSCALE FLOAT 0 DEFAULT 10 THE LENGTH SCALE OF THE KERNEL

ALPHA FLOAT 0 DEFAULT 10 SCALE MIXTURE PARAMETER

LENGTHSCALEBOUNDS PAIR OF FLOATS 0 DEFAULT 1E5 1E5 THE LOWER AND UPPER BOUND ON LENGTHSCALE

ALPHABOUNDS PAIR OF FLOATS 0 DEFAULT 1E5 1E5 THE LOWER AND UPPER BOUND ON ALPHA

ATTRIBUTES

BOUNDS RETURNS THE LOGTRANSFORMED BOUNDS ON THE THETA

HYPERPARAMETERALPHA

HYPERPARAMETERLENGTHSCALE

HYPERPARAMETERS RETURNS A LIST OF ALL HYPERPARAMETER SPECIFICATIONS

NDIMS RETURNS THE NUMBER OF NONFIXED HYPERPARAMETERS OF THE KERNEL

THETA RETURNS THE FLATTENED LOGTRANSFORMED NONFIXED HYPERPARAMETERS

6175SKLEARNGAUSSIANPROCESS GAUSSIAN PROCESSES 1831

SCIKITLEARN USER GUIDE RELEASE 0213

METHODS

CALL SELF X Y EVALGRADIENT RETURN THE KERNEL KX Y AND OPTIONALLY ITS GRADIENT

CLONEWITHTHETA SELF THETA RETURNS A CLONE OF SELF WITH GIVEN HYPERPARAMETERS THETA

DIAG SELF X RETURNS THE DIAGONAL OF THE KERNEL KX X

GETPARAMS SELF DEEP GET PARAMETERS OF THIS KERNEL

ISSTATIONARY SELF RETURNS WHETHER THE KERNEL IS STATIONARY

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS KERNEL

INIT SELF LENGTHSCALE10 ALPHA10 LENGTHSCALEBOUNDS1E05 1000000

ALPHABOUNDS1E05 1000000

CALL SELFXYNONE EVALGRADIENTFALSE

RETURN THE KERNEL KX Y AND OPTIONALLY ITS GRADIENT

PARAMETERS

XARRAY SHAPE NSAMPLESX NFEATURES LEFT ARGUMENT OF THE RETURNED KERNEL KX Y

YARRAY SHAPE NSAMPLESY NFEATURES OPTIONAL DEFAULTNONE RIGHT ARGUMENT OF THE

RETURNED KERNEL KX Y IF NONE KX X IF EVALUATED INSTEAD

EVALGRADIENT BOOL OPTIONAL DEFAULTFALSE DETERMINES WHETHER THE GRADIENT WITH RE

SPECT TO THE KERNEL HYPERPARAMETER IS DETERMINED ONLY SUPPORTED WHEN Y IS NONE

RETURNS

KARRAY SHAPE NSAMPLESX NSAMPLESY KERNEL KX Y

KGRADIENT ARRAY OPT SHAPE NSAMPLESX NSAMPLESX NDIMS THE GRADIENT OF THE

KERNEL KX X WITH RESPECT TO THE HYPERPARAMETER OF THE KERNEL ONLY RETURNED WHEN

EVALGRADIENT IS TRUE

BOUNDS

RETURNS THE LOGTRANSFORMED BOUNDS ON THE THETA

RETURNS

BOUNDS ARRAY SHAPE NDIMS 2 THE LOGTRANSFORMED BOUNDS ON THE KERNEL’S HYPERPARAM

ETERS THETA

CLONEWITHTHETA SELFTHETA

RETURNS A CLONE OF SELF WITH GIVEN HYPERPARAMETERS THETA

PARAMETERS

THETA ARRAY SHAPE NDIMS THE HYPERPARAMETERS

DIAGSELF

RETURNS THE DIAGONAL OF THE KERNEL KX X

THE RESULT OF THIS METHOD IS IDENTICAL TO NPDIAGSELF HOWEVER IT CAN BE EVALUATED MORE EFFICIENTLY

SINCE ONLY THE DIAGONAL IS EVALUATED

PARAMETERS

XARRAY SHAPE NSAMPLESX NFEATURES LEFT ARGUMENT OF THE RETURNED KERNEL KX Y

RETURNS

KDIAG ARRAY SHAPE NSAMPLESX DIAGONAL OF KERNEL KX X

1832 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

GETPARAMS SELFDEEPTURE

GET PARAMETERS OF THIS KERNEL

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

HYPERPARAMETERS

RETURNS A LIST OF ALL HYPERPARAMETER SPECIFICATIONS

ISSTATIONARY SELF

RETURNS WHETHER THE KERNEL IS STATIONARY

NDIMS

RETURNS THE NUMBER OF NONFIXED HYPERPARAMETERS OF THE KERNEL

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS KERNEL

THE METHOD WORKS ON SIMPLE KERNELS AS WELL AS ON NESTED KERNELS THE LATTER HAVE PARAMETERS OF THE FORM

COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

THETA

RETURNS THE FLATTENED LOGTRANSFORMED NONFIXED HYPERPARAMETERS

NOTE THAT THETA ARE TYPICALLY THE LOGTRANSFORMED VALUES OF THE KERNEL'S HYPERPARAMETERS AS THIS REPRESENTA

TION OF THE SEARCH SPACE IS MORE AMENABLE FOR HYPERPARAMETER SEARCH AS HYPERPARAMETERS LIKE LENGTHSCALES

NATURALLY LIVE ON A LOGSCALE

RETURNS

THETA ARRAY SHAPE NDIMS THE NONFIXED LOGTRANSFORMED HYPERPARAMETERS OF THE KERNEL

EXAMPLES USING SKLEARNGAUSSIANPROCESSKERNELSRATIONALQUADRATIC

- ILLUSTRATION OF PRIOR AND POSTERIOR GAUSSIAN PROCESS FOR DIFFERENT KERNELS
- GAUSSIAN PROCESS REGRESSION GPR ON MAUNA LOA CO2 DATA

61715SKLEARNGAUSSIANPROCESSKERNELS SUM

CLASSSKLEARNGAUSSIANPROCESSKERNELS SUMK1K2

SUMKERNEL K1 K2 OF TWO KERNELS K1 AND K2

THE RESULTING KERNEL IS DEFINED AS KSUMX Y K1X Y K2X Y

NEW IN VERSION 018

PARAMETERS

K1KERNEL OBJECT THE FIRST BASEKERNEL OF THE SUMKERNEL

K2KERNEL OBJECT THE SECOND BASEKERNEL OF THE SUMKERNEL

6175SKLEARNGAUSSIANPROCESS GAUSSIAN PROCESSES 1833

SCIKITLEARN USER GUIDE RELEASE 0213

ATTRIBUTES

BOUNDS RETURNS THE LOGTRANSFORMED BOUNDS ON THE THETA

HYPERPARAMETERS RETURNS A LIST OF ALL HYPERPARAMETER

NDIMS RETURNS THE NUMBER OF NONFIXED HYPERPARAMETERS OF THE KERNEL

THETA RETURNS THE FLATTENED LOGTRANSFORMED NONFIXED HYPERPARAMETERS

METHODS

CALL SELF X Y EVALGRADIENT RETURN THE KERNEL KX Y AND OPTIONALLY ITS GRADIENT

CLONEWITHTHETA SELF THETA RETURNS A CLONE OF SELF WITH GIVEN HYPERPARAMETERS

THETA

DIAG SELF X RETURNS THE DIAGONAL OF THE KERNEL KX X

GETPARAMS SELF DEEP GET PARAMETERS OF THIS KERNEL

ISSTATIONARY SELF RETURNS WHETHER THE KERNEL IS STATIONARY

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS KERNEL

INIT SELF K1 K2

CALL SELF X Y NONE EVALGRADIENT FALSE

RETURN THE KERNEL KX Y AND OPTIONALLY ITS GRADIENT

PARAMETERS

XARRAY SHAPE NSAMPLESX NFEATURES LEFT ARGUMENT OF THE RETURNED KERNEL KX Y

YARRAY SHAPE NSAMPLESY NFEATURES OPTIONAL DEFAULT NONE RIGHT ARGUMENT OF THE

RETURNED KERNEL KX Y IF NONE KX X IF EVALUATED INSTEAD

EVALGRADIENT BOOL OPTIONAL DEFAULT FALSE DETERMINES WHETHER THE GRADIENT WITH RE

SPECT TO THE KERNEL HYPERPARAMETER IS DETERMINED

RETURNS

KARRAY SHAPE NSAMPLESX NSAMPLESY KERNEL KX Y

KGRADIENT ARRAY OPT SHAPE NSAMPLESX NSAMPLESX NDIMS THE GRADIENT OF THE

KERNEL KX X WITH RESPECT TO THE HYPERPARAMETER OF THE KERNEL ONLY RETURNED WHEN

EVALGRADIENT IS TRUE

BOUNDS

RETURNS THE LOGTRANSFORMED BOUNDS ON THE THETA

RETURNS

BOUNDS ARRAY SHAPE NDIMS 2 THE LOGTRANSFORMED BOUNDS ON THE KERNEL'S HYPERPARAM

ETERS THETA

CLONEWITHTHETA SELF THETA

RETURNS A CLONE OF SELF WITH GIVEN HYPERPARAMETERS THETA

PARAMETERS

THETA ARRAY SHAPE NDIMS THE HYPERPARAMETERS

DIAG SELF X

RETURNS THE DIAGONAL OF THE KERNEL KX X

1834 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

THE RESULT OF THIS METHOD IS IDENTICAL TO NPDIAGSELF<sub>X</sub> HOWEVER IT CAN BE EVALUATED MORE EFFICIENTLY SINCE ONLY THE DIAGONAL IS EVALUATED

PARAMETERS

XARRAY SHAPE NSAMPLESX NFEATURES LEFT ARGUMENT OF THE RETURNED KERNEL KX Y

RETURNS

KDIAG ARRAY SHAPE NSAMPLESX DIAGONAL OF KERNEL KX X

GETPARAMS SELFDEEPT<sub>TRUE</sub>

GET PARAMETERS OF THIS KERNEL

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

HYPERPARAMETERS

RETURNS A LIST OF ALL HYPERPARAMETER

ISSTATIONARY SELF

RETURNS WHETHER THE KERNEL IS STATIONARY

NDIMS

RETURNS THE NUMBER OF NONFIXED HYPERPARAMETERS OF THE KERNEL

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS KERNEL

THE METHOD WORKS ON SIMPLE KERNELS AS WELL AS ON NESTED KERNELS THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

THETA

RETURNS THE FLATTENED LOGTRANSFORMED NONFIXED HYPERPARAMETERS

NOTE THAT THETA ARE TYPICALLY THE LOGTRANSFORMED VALUES OF THE KERNEL'S HYPERPARAMETERS AS THIS REPRESENTATION OF THE SEARCH SPACE IS MORE AMENABLE FOR HYPERPARAMETER SEARCH AS HYPERPARAMETERS LIKE LENGTHSCALES NATURALLY LIVE ON A LOGSCALE

RETURNS

THETA ARRAY SHAPE NDIMS THE NONFIXED LOGTRANSFORMED HYPERPARAMETERS OF THE KERNEL

61716SKLEARNGAUSSIANPROCESSKERNELS WHITEKERNEL

CLASSSSKLEARNGAUSSIANPROCESSKERNELS WHITEKERNEL NOISELEVEL10

NOISELEVELBOUNDS1E05

1000000

WHITE KERNEL

THE MAIN USECASE OF THIS KERNEL IS AS PART OF A SUMKERNEL WHERE IT EXPLAINS THE NOISECOMPONENT OF THE SIGNAL TUNING ITS PARAMETER CORRESPONDS TO ESTIMATING THE NOISELEVEL

KX1 X2 NOISELEVEL IF X1 X2 ELSE 0

617SKLEARNGAUSSIANPROCESS GAUSSIAN PROCESSES 1835

SCIKITLEARN USER GUIDE RELEASE 0213

NEW IN VERSION 018

PARAMETERS

NOISELEVEL FLOAT DEFAULT 10 PARAMETER CONTROLLING THE NOISE LEVEL

NOISELEVELBOUNDS PAIR OF FLOATS 0 DEFAULT 1E5 1E5 THE LOWER AND UPPER BOUND ON NOISELEVEL

ATTRIBUTES

BOUNDS RETURNS THE LOGTRANSFORMED BOUNDS ON THE THETA

HYPERPARAMETERNOISELEVEL

HYPERPARAMETERS RETURNS A LIST OF ALL HYPERPARAMETER SPECIFICATIONS

NDIMS RETURNS THE NUMBER OF NONFIXED HYPERPARAMETERS OF THE KERNEL

THETA RETURNS THE FLATTENED LOGTRANSFORMED NONFIXED HYPERPARAMETERS

METHODS

CALL SELF X Y EVALGRADIENT RETURN THE KERNEL KX Y AND OPTIONALLY ITS GRADIENT

CLONewithTHETA SELF THETA RETURNS A CLONE OF SELF WITH GIVEN HYPERPARAMETERS

THETA

DIAG SELF X RETURNS THE DIAGONAL OF THE KERNEL KX X

GETPARAMS SELF DEEP GET PARAMETERS OF THIS KERNEL

ISSTATIONARY SELF RETURNS WHETHER THE KERNEL IS STATIONARY

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS KERNEL

INIT SELFNOISELEVEL10 NOISELEVELBOUNDS1E05 1000000

CALL SELFXYNONE EVALGRADIENTFALSE

RETURN THE KERNEL KX Y AND OPTIONALLY ITS GRADIENT

PARAMETERS

XARRAY SHAPE NSAMPLESX NFEATURES LEFT ARGUMENT OF THE RETURNED KERNEL KX Y

YARRAY SHAPE NSAMPLESY NFEATURES OPTIONAL DEFAULTNONE RIGHT ARGUMENT OF THE RETURNED KERNEL KX Y IF NONE KX X IF EVALUATED INSTEAD

EVALGRADIENT BOOL OPTIONAL DEFAULTFALSE DETERMINES WHETHER THE GRADIENT WITH RESPECT TO THE KERNEL HYPERPARAMETER IS DETERMINED ONLY SUPPORTED WHEN Y IS NONE

RETURNS

KARRAY SHAPE NSAMPLESX NSAMPLESY KERNEL KX Y

KGRADIENT ARRAY OPT SHAPE NSAMPLESX NSAMPLESX NDIMS THE GRADIENT OF THE KERNEL KX X WITH RESPECT TO THE HYPERPARAMETER OF THE KERNEL ONLY RETURNED WHEN EVALGRADIENT IS TRUE

BOUNDS

RETURNS THE LOGTRANSFORMED BOUNDS ON THE THETA

RETURNS

BOUNDS ARRAY SHAPE NDIMS 2 THE LOGTRANSFORMED BOUNDS ON THE KERNEL’S HYPERPARAMETERS THETA

1836 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

CLONewithTHETA SELFTHETA

RETURNS A CLONE OF SELF WITH GIVEN HYPERPARAMETERS THETA

PARAMETERS

THETA ARRAY SHAPE NDIMS THE HYPERPARAMETERS

DIAGSELF

RETURNS THE DIAGONAL OF THE KERNEL KX X

THE RESULT OF THIS METHOD IS IDENTICAL TO NPDIAGSELF HOWEVER IT CAN BE EVALUATED MORE EFFICIENTLY

SINCE ONLY THE DIAGONAL IS EVALUATED

PARAMETERS

XARRAY SHAPE NSAMPLESX NFEATURES LEFT ARGUMENT OF THE RETURNED KERNEL KX Y

RETURNS

KDIAG ARRAY SHAPE NSAMPLESX DIAGONAL OF KERNEL KX X

GETPARAMS SELFDEEPTURE

GET PARAMETERS OF THIS KERNEL

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

HYPERPARAMETERS

RETURNS A LIST OF ALL HYPERPARAMETER SPECIFICATIONS

ISSTATIONARY SELF

RETURNS WHETHER THE KERNEL IS STATIONARY

NDIMS

RETURNS THE NUMBER OF NONFIXED HYPERPARAMETERS OF THE KERNEL

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS KERNEL

THE METHOD WORKS ON SIMPLE KERNELS AS WELL AS ON NESTED KERNELS THE LATTER HAVE PARAMETERS OF THE FORM

COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

THETA

RETURNS THE FLATTENED LOGTRANSFORMED NONFIXED HYPERPARAMETERS

NOTE THAT THETA ARE TYPICALLY THE LOGTRANSFORMED VALUES OF THE KERNEL'S HYPERPARAMETERS AS THIS REPRESENTA

TION OF THE SEARCH SPACE IS MORE AMENABLE FOR HYPERPARAMETER SEARCH AS HYPERPARAMETERS LIKE LENGTHSCALES

NATURALLY LIVE ON A LOGSCALE

RETURNS

THETA ARRAY SHAPE NDIMS THE NONFIXED LOGTRANSFORMED HYPERPARAMETERS OF THE KERNEL

617SKLEARNGAUSSIANPROCESS GAUSSIAN PROCESSES 1837

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES USING SKLEARN

- GAUSSIANPROCESSKERNELSWHITEKERNEL
- COMPARISON OF KERNEL RIDGE AND GAUSSIAN PROCESS REGRESSION
- GAUSSIAN PROCESS REGRESSION GPR WITH NOISELEVEL ESTIMATION
- GAUSSIAN PROCESS REGRESSION GPR ON MAUNA LOA CO2 DATA

618SKLEARNISOTONIC ISOTONIC REGRESSION

USER GUIDE SEE THE ISOTONIC REGRESSION SECTION FOR FURTHER DETAILS

ISOTONICISOTONICREGRESSION YMIN YMAX

ISOTONIC REGRESSION MODEL

6181SKLEARNISOTONIC ISOTONICREGRESSION

CLASSSSKLEARNISOTONIC ISOTONICREGRESSION YMINNONE YMAXNONE INCREASINGTRUE

OUTOFBOUNDS'NAN'

ISOTONIC REGRESSION MODEL

THE ISOTONIC REGRESSION OPTIMIZATION PROBLEM IS DEFINED BY

$$\min \sum w_i y_i - y_i^2$$

SUBJECT TO  $y_i \leq y_j$  WHENEVER  $x_i \leq x_j$

AND  $y_{\min} \leq y_i \leq y_{\max}$

WHERE

- $y_i$  ARE INPUTS REAL NUMBERS
- $y_i$  ARE FITTED
- $x_i$  SPECIFIES THE ORDER IF  $x_i$  IS NONDECREASING THEN  $y_i$  IS NONDECREASING
- $w_i$  ARE OPTIONAL STRICTLY POSITIVE WEIGHTS DEFAULT TO 10

READ MORE IN THE USER GUIDE

PARAMETERS

YMIN OPTIONAL DEFAULT NONE IF NOT NONE SET THE LOWEST VALUE OF THE FIT TO YMIN

YMAX OPTIONAL DEFAULT NONE IF NOT NONE SET THE HIGHEST VALUE OF THE FIT TO YMAX

INCREASING BOOLEAN OR STRING OPTIONAL DEFAULT TRUE IF BOOLEAN WHETHER OR NOT TO FIT THE ISOTONIC REGRESSION WITH Y INCREASING OR DECREASING

THE STRING VALUE "AUTO" DETERMINES WHETHER Y SHOULD INCREASE OR DECREASE BASED ON THE SPEARMAN CORRELATION ESTIMATE'S SIGN

OUTOFBOUNDS STRING OPTIONAL DEFAULT "NAN" THE OUTOFBOUNDS PARAMETER HANDLES HOW XVALUES OUTSIDE OF THE TRAINING DOMAIN ARE HANDLED WHEN SET TO "NAN" PREDICTED YVALUES WILL BE NAN WHEN SET TO "CLIP" PREDICTED YVALUES WILL BE SET TO THE VALUE CORRESPONDING TO THE NEAREST TRAIN INTERVAL ENDPOINT WHEN SET TO "RAISE" ALLOW INTERP1D TO THROW VALUEERROR

1838 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

ATTRIBUTES

XMIN FLOAT MINIMUM VALUE OF INPUT ARRAY XFOR LEFT BOUND

XMAX FLOAT MAXIMUM VALUE OF INPUT ARRAY XFOR RIGHT BOUND

FFUNCTION THE STEPWISE INTERPOLATING FUNCTION THAT COVERS THE INPUT DOMAIN X

NOTES

TIES ARE BROKEN USING THE SECONDARY METHOD FROM LEEUW 1977

REFERENCES

ISOTONIC MEDIAN REGRESSION A LINEAR PROGRAMMING APPROACH NILOTPAL CHAKRAVARTI MATHEMATICS OF OPERATIONS RESEARCH V OL 14 NO 2 MAY 1989 PP 303308

ISOTONE OPTIMIZATION IN R POOLADJACENTVIOLATORS ALGORITHM PA V A AND ACTIVE SET METHODS LEEUW HORNIK MAIR JOURNAL OF STATISTICAL SOFTWARE 2009

CORRECTNESS OF KRUSKAL'S ALGORITHMS FOR MONOTONE REGRESSION WITH TIES LEEUW PSYCHOMETRICA 1977

EXAMPLES

FROM SKLEARNDATASETS IMPORT MAKEREGRESSION

FROM SKLEARNISOTONIC IMPORT ISOTONICREGRESSION

X Y MAKEREGRESSIONNNSAMPLES10 NFEATURES1 RANDOMSTATE41

ISOREG ISOTONICREGRESSIONFITXFLATTEN Y

ISOREGPREDICT1 2

ARRAY18628 37256

METHODS

FITSELF X Y SAMPLEWEIGHT FIT THE MODEL USING X Y AS TRAINING DATA

FITTRANSFORM SELF X Y FIT TO DATA THEN TRANSFORM IT

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PREDICT SELF T PREDICT NEW DATA BY LINEAR INTERPOLATION

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PRE DITION

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF T TRANSFORM NEW DATA BY LINEAR INTERPOLATION

INIT SELFYMINNONE YMAXNONE INCREASINGTRUE OUTOFBOUNDS'NAN'

FITSELFXYSAMPLEWEIGHTNONE

FIT THE MODEL USING X Y AS TRAINING DATA

PARAMETERS

XARRAYLIKE SHAPENSAMPLES TRAINING DATA

YARRAYLIKE SHAPENSAMPLES TRAINING TARGET

618SKLEARNISOTONIC ISOTONIC REGRESSION 1839

SCIKITLEARN USER GUIDE RELEASE 0213

SAMPLEWEIGHT ARRAYLIKE SHAPENSAMPLES OPTIONAL DEFAULT NONE WEIGHTS IF SET TO NONE ALL WEIGHTS WILL BE SET TO 1 EQUAL WEIGHTS

RETURNS

SELF OBJECT RETURNS AN INSTANCE OF SELF

NOTES

X IS STORED FOR FUTURE USE AS TRANSFORM NEEDS X TO INTERPOLATE NEW INPUT DATA

FITTRANSFORM SELFXYNONE FITPARAMS

FIT TO DATA THEN TRANSFORM IT

FITS TRANSFORMER TO X AND Y WITH OPTIONAL PARAMETERS FITPARAMS AND RETURNS A TRANSFORMED VERSION OF X

PARAMETERS

XNUMPY ARRAY OF SHAPE NSAMPLES NFEATURES TRAINING SET

YNUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES

RETURNS

XNEW NUMPY ARRAY OF SHAPE NSAMPLES NFEATURESNEW TRANSFORMED ARRAY

GETPARAMS SELFDEEPTTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELF

PREDICT NEW DATA BY LINEAR INTERPOLATION

PARAMETERS

TARRAYLIKE SHAPENSAMPLES DATA TO TRANSFORM

RETURNS

TARRAY SHAPENSAMPLES TRANSFORMED DATA

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{U}{V}$  WHERE U IS THE RESIDUAL SUM OF SQUARES YTRUE YPRED

2SUM AND V IS THE TOTAL SUM OF SQUARES YTRUE YTRUEMEAN 2SUM THE BEST POSSIBLE SCORE

IS 10 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS

PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 00

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES FOR SOME ESTIMATORS THIS MAY

BE A PRECOMPUTED KERNEL MATRIX INSTEAD SHAPE NSAMPLES NSAMPLESFITTED WHERE

NSAMPLESFITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE VALUES FOR X

1840 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT R2 OF SELF PREDICTX WRT Y

NOTES

THE R2 SCORE USED WHEN CALLING SCORE ON A REGRESSOR WILL USE MULTIOUTPUTUNIFORMAVERAGE FROM VERSION 023 TO KEEP CONSISTENT WITH METRICSR2SCORE THIS WILL INFLUENCE THE SCORE METHOD OF ALL THE MULTIOUTPUT REGRESSORS EXCEPT FOR MULTIOUTPUTMULTIOUTPUTREGRESSOR TO SPECIFY THE DEFAULT VALUE MANUALLY AND AVOID THE WARNING PLEASE EITHER CALL METRICSR2SCORE DIRECTLY OR MAKE A CUSTOM SCORER WITH METRICSMAKESCORER THE BUILTIN SCORER R2 USES MULTIOUTPUTUNIFORMAVERAGE

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELF T

TRANSFORM NEW DATA BY LINEAR INTERPOLATION

PARAMETERS

TARRAYLIKE SHAPENSAMPLES DATA TO TRANSFORM

RETURNS

TARRAY SHAPENSAMPLES THE TRANSFORMED DATA

EXAMPLES USING SKLEARNISOTONICISOTONICREGRESSION

- ISOTONIC REGRESSION

ISOTONICCHECKINCREASING X Y DETERMINE WHETHER Y IS MONOTONICALLY CORRELATED WITH X

ISOTONICISOTONICREGRESSION Y SOLVE THE ISOTONIC REGRESSION MODEL

6182SKLEARNISOTONIC CHECKINCREASING

SKLEARNISOTONIC CHECKINCREASING XY

DETERMINE WHETHER Y IS MONOTONICALLY CORRELATED WITH X

Y IS FOUND INCREASING OR DECREASING WITH RESPECT TO X BASED ON A SPEARMAN CORRELATION TEST

PARAMETERS

XARRAYLIKE SHAPENSAMPLES TRAINING DATA

YARRAYLIKE SHAPENSAMPLES TRAINING TARGET

RETURNS

618SKLEARNISOTONIC ISOTONIC REGRESSION 1841

SCIKITLEARN USER GUIDE RELEASE 0213

INCREASINGBOOL BOOLEAN WHETHER THE RELATIONSHIP IS INCREASING OR DECREASING

NOTES

THE SPEARMAN CORRELATION COEFFICIENT IS ESTIMATED FROM THE DATA AND THE SIGN OF THE RESULTING ESTIMATE IS USED AS THE RESULT

IN THE EVENT THAT THE 95 CONFIDENCE INTERVAL BASED ON FISHER TRANSFORM SPANS ZERO A WARNING IS RAISED

REFERENCES

FISHER TRANSFORMATION WIKIPEDIA [HTTPS://ENWIKIPEDIA.ORG/WIKI/FISHERTRANSFORMATION](https://en.wikipedia.org/wiki/Fisher_transformation)

6183SKLEARNISOTONIC ISOTONICREGRESSION

SKLEARNISOTONIC ISOTONICREGRESSION YSAMPLEWEIGHTNONE YMINNONE YMAXNONE

INCREASINGTRUE

SOLVE THE ISOTONIC REGRESSION MODEL

MIN SUM WI YI YI 2

SUBJECT TO YMIN Y1 Y2 YN YMAX

WHERE

- YI ARE INPUTS REAL NUMBERS
- YI ARE FITTED
- WI ARE OPTIONAL STRICTLY POSITIVE WEIGHTS DEFAULT TO 10

READ MORE IN THE USER GUIDE

PARAMETERS

YITERABLE OF FLOATS THE DATA

SAMPLEWEIGHT ITERABLE OF FLOATS OPTIONAL DEFAULT NONE WEIGHTS ON EACH POINT OF THE REGRES

SION IF NONE WEIGHT IS SET TO 1 EQUAL WEIGHTS

YMIN OPTIONAL DEFAULT NONE IF NOT NONE SET THE LOWEST VALUE OF THE FIT TO YMIN

YMAX OPTIONAL DEFAULT NONE IF NOT NONE SET THE HIGHEST VALUE OF THE FIT TO YMAX

INCREASING BOOLEAN OPTIONAL DEFAULT TRUE WHETHER TO COMPUTE YIS INCREASING IF SET TO

TRUE OR DECREASING IF SET TO FALSE

RETURNS

YLIST OF FLOATS ISOTONIC FIT OF Y

REFERENCES

“ACTIVE SET ALGORITHMS FOR ISOTONIC REGRESSION A UNIFYING FRAMEWORK” BY MICHAEL J BEST AND NILOTPAL CHAKRAVARTI SECTION 3

1842 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213  
619SKLEARNIMPUTE IMPUTE  
TRANSFORMERS FOR MISSING VALUE IMPUTATION  
USER GUIDE SEE THE IMPUTATION OF MISSING VALUES SECTION FOR FURTHER DETAILS  
IMPUTESIMPLEIMPUTER MISSINGVALUES IMPUTATION TRANSFORMER FOR COMPLETING MISSING VALUES  
IMPUTEITERATIVEIMPUTER ESTIMATOR MULTIVARIATE IMPUTER THAT ESTIMATES EACH FEATURE FROM ALL  
THE OTHERS  
IMPUTEMISSINGINDICATOR MISSINGVALUES BINARY INDICATORS FOR MISSING VALUES  
6191SKLEARNIMPUTE SIMPLEIMPUTER  
CLASSSSKLEARNIMPUTE SIMPLEIMPUTER MISSINGVALUESNAN STRATEGY'MEAN' FILLVALUENONE VER  
BOSE0 COPYTRUE ADDINDICATORFALSE  
IMPUTATION TRANSFORMER FOR COMPLETING MISSING VALUES  
READ MORE IN THE USER GUIDE  
PARAMETERS  
MISSINGVALUES NUMBER STRING NPNAN DEFAULT OR NONE THE PLACEHOLDER FOR THE MISSING  
VALUES ALL OCCURRENCES OF MISSINGVALUES WILL BE IMPUTED  
STRATEGY STRING OPTIONAL DEFAULT"MEAN" THE IMPUTATION STRATEGY  
• IF "MEAN" THEN REPLACE MISSING VALUES USING THE MEAN ALONG EACH COLUMN CAN ONLY BE  
USED WITH NUMERIC DATA  
• IF "MEDIAN" THEN REPLACE MISSING VALUES USING THE MEDIAN ALONG EACH COLUMN CAN ONLY  
BE USED WITH NUMERIC DATA  
• IF "MOSTFREQUENT" THEN REPLACE MISSING USING THE MOST FREQUENT VALUE ALONG EACH COL  
UMN CAN BE USED WITH STRINGS OR NUMERIC DATA  
• IF "CONSTANT" THEN REPLACE MISSING VALUES WITH FILLVALUE CAN BE USED WITH STRINGS OR  
NUMERIC DATA  
NEW IN VERSION 020 STRATEGY"CONSTANT" FOR FIXED VALUE IMPUTATION  
FILLVALUE STRING OR NUMERICAL VALUE OPTIONAL DEFAULTNONE WHEN STRATEGY "CONSTANT"  
FILLVALUE IS USED TO REPLACE ALL OCCURRENCES OF MISSINGVALUES IF LEFT TO THE DEFAULT FILLVALUE  
WILL BE 0 WHEN IMPUTING NUMERICAL DATA AND "MISSINGVALUE" FOR STRINGS OR OBJECT DATA TYPES  
VERBOSE INTEGER OPTIONAL DEFAULT0 CONTROLS THE VERBOSITY OF THE IMPUTER  
COPY BOOLEAN OPTIONAL DEFAULTTRUE IF TRUE A COPY OF X WILL BE CREATED IF FALSE IMPUTATION  
WILL BE DONE INPLACE WHENEVER POSSIBLE NOTE THAT IN THE FOLLOWING CASES A NEW COPY WILL  
ALWAYS BE MADE EVEN IF COPYFALSE  
• IF X IS NOT AN ARRAY OF FLOATING VALUES  
• IF X IS ENCODED AS A CSR MATRIX  
• IF ADDINDICATORTRUE  
ADDINDICATOR BOOLEAN OPTIONAL DEFAULTFALSE IF TRUE A MISSINGINDICATOR TRANSFORM  
WILL STACK ONTO OUTPUT OF THE IMPUTER'S TRANSFORM THIS ALLOWS A PREDICTIVE ESTIMATOR TO  
ACCOUNT FOR MISSINGNESS DESPITE IMPUTATION IF A FEATURE HAS NO MISSING VALUES AT FITTRAIN  
619SKLEARNIMPUTE IMPUTE 1843

SCIKITLEARN USER GUIDE RELEASE 0213

TIME THE FEATURE WON'T APPEAR ON THE MISSING INDICATOR EVEN IF THERE ARE MISSING VALUES AT TRANSFORMTEST TIME

ATTRIBUTES

STATISTICS ARRAY OF SHAPE NFEATURES THE IMPUTATION FILL VALUE FOR EACH FEATURE

INDICATOR SKLEARNIMPUTEMISSINGINDICATOR INDICATOR USED TO ADD BINARY INDICATORS FOR MISSING VALUES NONE IF ADDINDICATOR IS FALSE

SEE ALSO

ITERATIVEIMPUTER MULTIVARIATE IMPUTATION OF MISSING VALUES

NOTES

COLUMNS WHICH ONLY CONTAINED MISSING VALUES AT FIT ARE DISCARDED UPON TRANSFORM IF STRATEGY IS NOT "CONSTANT"

EXAMPLES

```
import numpy as np
from sklearn.impute import SimpleImputer
imp_mean = SimpleImputer(missing_values=np.nan, strategy='mean')
imp_mean.fit(7 2 3 4 np.nan 6 10 5 9)
```

SIMPLEIMPUTERADDINDICATORFALSE COPYTRUE FILLVALUENONE

MISSINGVALUESNAN STRATEGYMEAN VERBOSE0

```
X = np.array([7, 2, 3, 4, np.nan, 6, 10, np.nan, 9])
print(imp_mean.transform(X))
```

```
7 2 3
4 35 6
10 35 9
```

METHODS

FITSELF X Y FIT THE IMPUTER ON X

FITTRANSFORM SELF X Y FIT TO DATA THEN TRANSFORM IT

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF X IMPUTE ALL MISSING VALUES IN X

INIT SELFMISSINGVALUESNAN STRATEGY'MEAN' FILLVALUENONE VERBOSE0 COPYTRUE

ADDINDICATORFALSE

FITSELFXYNONE

FIT THE IMPUTER ON X

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES INPUT DATA WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

1844 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

RETURNS

SELF SIMPLEIMPUTER

FITTRANSFORM SELFXYNONE FITPARAMS

FIT TO DATA THEN TRANSFORM IT

FITS TRANSFORMER TO X AND Y WITH OPTIONAL PARAMETERS FITPARAMS AND RETURNS A TRANSFORMED VERSION OF X

PARAMETERS

XNUMPY ARRAY OF SHAPE NSAMPLES NFEATURES TRAINING SET

YNUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES

RETURNS

XNEW NUMPY ARRAY OF SHAPE NSAMPLES NFEATURESNEW TRANSFORMED ARRAY

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE

PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT

OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELFXY

IMPUTE ALL MISSING VALUES IN X

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT DATA TO COMPLETE

EXAMPLES USING SKLEARNIMPUTESIMPLEIMPUTER

- COLUMN TRANSFORMER WITH MIXED TYPES
- IMPUTING MISSING VALUES WITH VARIANTS OF ITERATIVEIMPUTER
- IMPUTING MISSING VALUES BEFORE BUILDING AN ESTIMATOR

619SKLEARNIMPUTE IMPUTE 1845

SCIKITLEARN USER GUIDE RELEASE 0213  
61925KLEARNIMPUTE ITERATIVEIMPUTER  
CLASSSSKLEARNIMPUTE ITERATIVEIMPUTER ESTIMATORNONE MISSINGVALUESNAN SAM  
PLEPOSTERIORFALSE MAXITER10 TOL0001  
NNEARESTFEATURESNONE INITIALSTRATEGY'MEAN'  
IMPUTATIONORDER'ASCENDING' MINVALUENONE  
MAXVALUENONE VERBOSE0 RANDOMSTATENONE  
ADDINDICATORFALSE  
MULTIVARIATE IMPUTER THAT ESTIMATES EACH FEATURE FROM ALL THE OTHERS  
A STRATEGY FOR IMPUTING MISSING VALUES BY MODELING EACH FEATURE WITH MISSING VALUES AS A FUNCTION OF OTHER  
FEATURES IN A ROUNDROBIN FASHION  
READ MORE IN THE USER GUIDE  
NOTE THIS ESTIMATOR IS STILL EXPERIMENTAL FOR NOW THE PREDICTIONS AND THE API MIGHT CHANGE WITHOUT ANY  
DEPRECATION CYCLE TO USE IT YOU NEED TO EXPLICITLY IMPORT ENABLEITERATIVEIMPUTER  
EXPLICITLY REQUIRE THIS EXPERIMENTAL FEATURE  
FROM SKLEARNEXPERIMENTAL IMPORT ENABLEITERATIVEIMPUTER NOQA  
NOW YOU CAN IMPORT NORMALLY FROM SKLEARNIMPUTE  
FROM SKLEARNIMPUTE IMPORT ITERATIVEIMPUTER  
PARAMETERS  
ESTIMATOR ESTIMATOR OBJECT DEFAULTBAYESIANRIDGE THE ESTIMATOR TO USE AT EACH STEP OF  
THE ROUNDROBIN IMPUTATION IF SAMPLEPOSTERIOR IS TRUE THE ESTIMATOR MUST SUPPORT  
RETURNSTD IN ITSPREDICT METHOD  
MISSINGVALUES INT NPNAN OPTIONAL DEFAULTNPNAN THE PLACEHOLDER FOR THE MISSING VALUES  
ALL OCCURRENCES OF MISSINGVALUES WILL BE IMPUTED  
SAMPLEPOSTERIOR BOOLEAN DEFAULTFALSE WHETHER TO SAMPLE FROM THE GAUSSIAN PREDICTIVE  
POSTERIOR OF THE FITTED ESTIMATOR FOR EACH IMPUTATION ESTIMATOR MUST SUPPORT RETURNSTD  
IN ITSPREDICT METHOD IF SET TO TRUE SET TOTRUE IF USINGITERATIVEIMPUTER FOR  
MULTIPLE IMPUTATIONS  
MAXITER INT OPTIONAL DEFAULT10 MAXIMUM NUMBER OF IMPUTATION ROUNDS TO PERFORM  
BEFORE RETURNING THE IMPUTATIONS COMPUTED DURING THE FINAL ROUND A ROUND IS A SIN  
GLE IMPUTATION OF EACH FEATURE WITH MISSING VALUES THE STOPPING CRITERION IS MET  
ONCEABSMAXXT XT1ABSMAXXKNOWNVALS TOL WHERE  
XT ISXAT ITERATION T NOTE THAT EARLY STOPPING IS ONLY APPLIED IF  
SAMPLEPOSTERIORFALSE  
TOLFLOAT OPTIONAL DEFAULT1E3 TOLERANCE OF THE STOPPING CONDITION  
NNEARESTFEATURES INT OPTIONAL DEFAULTNONE NUMBER OF OTHER FEATURES TO USE TO ESTIMATE  
THE MISSING VALUES OF EACH FEATURE COLUMN NEARNESS BETWEEN FEATURES IS MEASURED USING  
THE ABSOLUTE CORRELATION COEFFICIENT BETWEEN EACH FEATURE PAIR AFTER INITIAL IMPUTATION TO  
ENSURE COVERAGE OF FEATURES THROUGHOUT THE IMPUTATION PROCESS THE NEIGHBOR FEATURES ARE  
NOT NECESSARILY NEAREST BUT ARE DRAWN WITH PROBABILITY PROPORTIONAL TO CORRELATION FOR EACH  
IMPUTED TARGET FEATURE CAN PROVIDE SIGNIFICANT SPEEDUP WHEN THE NUMBER OF FEATURES IS  
HUGE IFNONE ALL FEATURES WILL BE USED  
INITIALSTRATEGY STR OPTIONAL DEFAULT"MEAN" WHICH STRATEGY TO USE TO INITIALIZE THE MISSING  
VALUES SAME AS THE STRATEGY PARAMETER IN SKLEARNIMPUTESIMPLEIMPUTER  
VALID VALUES "MEAN" "MEDIAN" "MOSTFREQUENT" OR "CONSTANT"  
1846 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

IMPUTATIONORDER STR OPTIONAL DEFAULT"ASCENDING" THE ORDER IN WHICH THE FEATURES WILL BE IMPUTED POSSIBLE VALUES

"ASCENDING" FROM FEATURES WITH FEWEST MISSING VALUES TO MOST

"DESCENDING" FROM FEATURES WITH MOST MISSING VALUES TO FEWEST

"ROMAN" LEFT TO RIGHT

"ARABIC" RIGHT TO LEFT

"RANDOM" A RANDOM ORDER FOR EACH ROUND

MINVALUE FLOAT OPTIONAL DEFAULTNONE MINIMUM POSSIBLE IMPUTED VALUE DEFAULT OF NONE

WILL SET MINIMUM TO NEGATIVE INFINITY

MAXVALUE FLOAT OPTIONAL DEFAULTNONE MAXIMUM POSSIBLE IMPUTED VALUE DEFAULT OF

NONE WILL SET MAXIMUM TO POSITIVE INFINITY

VERBOSE INT OPTIONAL DEFAULT0 VERBOSITY FLAG CONTROLS THE DEBUG MESSAGES THAT ARE ISSUED

AS FUNCTIONS ARE EVALUATED THE HIGHER THE MORE VERBOSE CAN BE 0 1 OR 2

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE THE SEED OF THE

PSEUDO RANDOM NUMBER GENERATOR TO USE RANDOMIZES SELECTION OF ESTIMATOR FEATURES IF

NNEARESTFEATURES IS NOT NONE THE IMPUTATIONORDER IFRANDOM AND THE SAMPLING

FROM POSTERIOR IF SAMPLEPOSTERIOR IS TRUE USE AN INTEGER FOR DETERMINISM SEE THE

GLOSSARY

ADDINDICATOR BOOLEAN OPTIONAL DEFAULTFALSE IF TRUE A MISSINGINDICATOR TRANSFORM

WILL STACK ONTO OUTPUT OF THE IMPUTER'S TRANSFORM THIS ALLOWS A PREDICTIVE ESTIMATOR TO

ACCOUNT FOR MISSINGNESS DESPITE IMPUTATION IF A FEATURE HAS NO MISSING VALUES AT FITTRAIN

TIME THE FEATURE WON'T APPEAR ON THE MISSING INDICATOR EVEN IF THERE ARE MISSING VALUES AT

TRANSFORMTEST TIME

ATTRIBUTES

INITIALIMPUTER OBJECT OF TYPE SKLEARNIMPUTESIMPLEIMPUTER IMPUTER USED TO

INITIALIZE THE MISSING VALUES

IMPUTATIONSEQUENCE LIST OF TUPLES EACH TUPLE HAS FEATIDX

NEIGHBORFEATIDX ESTIMATOR WHEREFEATIDX IS THE CURRENT FEATURE

TO BE IMPUTED NEIGHBORFEATIDX IS THE ARRAY OF OTHER FEATURES USED TO IMPUTE THE

CURRENT FEATURE AND ESTIMATOR IS THE TRAINED ESTIMATOR USED FOR THE IMPUTATION LENGTH

ISSELFNFEATURESWITHMISSING SELFNITER

NITER INT NUMBER OF ITERATION ROUNDS THAT OCCURRED WILL BE LESS THAN SELFMAXITER IF

EARLY STOPPING CRITERION WAS REACHED

NFEATURESWITHMISSING INT NUMBER OF FEATURES WITH MISSING VALUES

INDICATOR SKLEARNIMPUTEMISSINGINDICATOR INDICATOR USED TO ADD BINARY INDI

CATORS FOR MISSING VALUES NONE IF ADDINDICATOR IS FALSE

SEE ALSO

SIMPLEIMPUTER UNIVARIATE IMPUTATION OF MISSING VALUES

NOTES

TO SUPPORT IMPUTATION IN INDUCTIVE MODE WE STORE EACH FEATURE'S ESTIMATOR DURING THE FIT PHASE AND PREDICT

WITHOUT REFITTING IN ORDER DURING THE TRANSFORM PHASE

619SKLEARNIMPUTE IMPUTE 1847

SCIKITLEARN USER GUIDE RELEASE 0213

FEATURES WHICH CONTAIN ALL MISSING VALUES AT FIT ARE DISCARDED UPON TRANSFORM

FEATURES WITH MISSING VALUES DURING TRANSFORM WHICH DID NOT HAVE ANY MISSING VALUES DURING FIT WILL BE IMPUTED WITH THE INITIAL IMPUTATION METHOD ONLY

REFERENCES

RCD31B817A31E1 RCD31B817A31E2

METHODS

FITSELF X Y FITS THE IMPUTER ON X AND RETURN SELF

FITTRANSFORM SELF X Y FITS THE IMPUTER ON X AND RETURN THE TRANSFORMED X

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF X IMPUTES ALL MISSING VALUES IN X

INIT SELFESTIMATORNONE MISSINGVALUESNAN SAMPLEPOSTERIORFALSE MAXITER10

TOL0001 NNEARESTFEATURESNONE INITIALSTRATEGY'MEAN' IMPUTA

TIONORDER'ASCENDING' MINVALUENONE MAXVALUENONE VERBOSE0 RAN

DOMSTATENONE ADDINDICATORFALSE

FITSELFXYNONE

FITS THE IMPUTER ON X AND RETURN SELF

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES INPUT DATA WHERE "NSAMPLES" IS THE NUMBER OF SAMPLES AND "NFEATURES" IS THE NUMBER OF FEATURES

YIGNORED

RETURNS

SELF OBJECT RETURNS SELF

FITTRANSFORM SELFXYNONE

FITS THE IMPUTER ON X AND RETURN THE TRANSFORMED X

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES INPUT DATA WHERE "NSAMPLES" IS THE NUMBER OF SAMPLES AND "NFEATURES" IS THE NUMBER OF FEATURES

YIGNORED

RETURNS

XTARRAYLIKE SHAPE NSAMPLES NFEATURES THE IMPUTED INPUT DATA

GETPARAMS SELFDEEPTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

1848 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELF

IMPUTES ALL MISSING VALUES IN X

NOTE THAT THIS IS STOCHASTIC AND THAT IF RANDOMSTATE IS NOT FIXED REPEATED CALLS OR PERMUTED INPUT WILL YIELD DIFFERENT RESULTS

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES THE INPUT DATA TO COMPLETE

RETURNS

XTARRAYLIKE SHAPE NSAMPLES NFEATURES THE IMPUTED INPUT DATA

EXAMPLES USING SKLEARNIMPUTEITERATIVEIMPUTER

- IMPUTING MISSING VALUES WITH VARIANTS OF ITERATIVEIMPUTER
- IMPUTING MISSING VALUES BEFORE BUILDING AN ESTIMATOR

6193SKLEARNIMPUTE MISSINGINDICATOR

CLASSSSKLEARNIMPUTE MISSINGINDICATOR MISSINGVALUESNAN FEATURES’MISSINGONLY’

SPARSE’AUTO’ ERRORONNEWTRUE

BINARY INDICATORS FOR MISSING VALUES

NOTE THAT THIS COMPONENT TYPICALLY SHOULD NOT BE USED IN A VANILLA PIPELINE CONSISTING OF TRANSFORMERS AND A CLASSIFIER BUT RATHER COULD BE ADDED USING A FEATUREUNION ORCOLUMNTRANSFORMER

READ MORE IN THE USER GUIDE

PARAMETERS

MISSINGVALUES NUMBER STRING NPNAN DEFAULT OR NONE THE PLACEHOLDER FOR THE MISSING

VALUES ALL OCCURRENCES OF MISSINGVALUES WILL BE INDICATED TRUE IN THE OUTPUT ARRAY

THE OTHER VALUES WILL BE MARKED AS FALSE

FEATURES STR OPTIONAL WHETHER THE IMPUTER MASK SHOULD REPRESENT ALL OR A SUBSET OF FEATURES

- IF “MISSINGONLY” DEFAULT THE IMPUTER MASK WILL ONLY REPRESENT FEATURES CONTAINING

MISSING VALUES DURING FIT TIME

- IF “ALL” THE IMPUTER MASK WILL REPRESENT ALL FEATURES

SPARSE BOOLEAN OR “AUTO” OPTIONAL WHETHER THE IMPUTER MASK FORMAT SHOULD BE SPARSE OR DENSE

- IF “AUTO” DEFAULT THE IMPUTER MASK WILL BE OF SAME TYPE AS INPUT
- IF TRUE THE IMPUTER MASK WILL BE A SPARSE MATRIX

6193SKLEARNIMPUTE IMPUTE 1849

SCIKITLEARN USER GUIDE RELEASE 0213

• IF FALSE THE IMPUTER MASK WILL BE A NUMPY ARRAY  
ERRORONNEW BOOLEAN OPTIONAL IF TRUE DEFAULT TRANSFORM WILL RAISE AN ERROR WHEN THERE  
ARE FEATURES WITH MISSING VALUES IN TRANSFORM THAT HAVE NO MISSING VALUES IN FIT THIS IS  
APPLICABLE ONLY WHEN FEATURESMISSINGONLY  
ATTRIBUTES  
FEATURES NDARRAY SHAPE NMISSINGFEATURES OR NFEATURES THE FEATURES INDICES WHICH  
WILL BE RETURNED WHEN CALLING TRANSFORM THEY ARE COMPUTED DURING FIT FOR  
FEATURESALL IT IS TORANGENFEATURES

EXAMPLES

IMPORT NUMPY AS NP  
FROM SKLEARNIMPUTE IMPORT MISSINGINDICATOR

X1 NPARRAYNPAN 1 3

4 0 NPNAN

8 1 0

X2 NPARRAY5 1 NPNAN

NPAN 2 3

2 4 0

INDICATOR MISSINGINDICATOR

INDICATORFITX1

MISSINGINDICATORERRORONNEWTRUE FEATURESMISSINGONLY

MISSINGVALUESNAN SPARSEAUTO

X2TR INDICATORTRANSFORMX2

X2TR

ARRAYFALSE TRUE

TRUE FALSE

FALSE FALSE

METHODS

FITSELF X Y FIT THE TRANSFORMER ON X

FITTRANSFORM SELF X Y GENERATE MISSING VALUES INDICATOR FOR X

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF X GENERATE MISSING VALUES INDICATOR FOR X

INIT SELFMISSINGVALUESNAN FEATURES'MISSINGONLY' SPARSE'AUTO' ERRORONNEWTRUE

FITSELFXYNONE

FIT THE TRANSFORMER ON X

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES INPUT DATA WHERE

NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

RETURNS

SELF OBJECT RETURNS SELF

FITTRANSFORM SELFXYNONE

GENERATE MISSING VALUES INDICATOR FOR X

1850 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT DATA TO COMPLETE

RETURNS

XTNDARRAY OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE MISSING INDICATOR FOR

INPUT DATA THE DATA TYPE OF XTWILL BE BOOLEAN

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE

PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT

OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELF

GENERATE MISSING VALUES INDICATOR FOR X

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT DATA TO COMPLETE

RETURNS

XTNDARRAY OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE MISSING INDICATOR FOR

INPUT DATA THE DATA TYPE OF XTWILL BE BOOLEAN

EXAMPLES USING SKLEARNIMPUTEMISSINGINDICATOR

- IMPUTING MISSING VALUES BEFORE BUILDING AN ESTIMATOR

620SKLEARNKERNELAPPROXIMATION KERNEL APPROXIMATION

THESKLEARNKERNELAPPROXIMATION MODULE IMPLEMENTS SEVERAL APPROXIMATE KERNEL FEATURE MAPS BASE ON

FOURIER TRANSFORMS

USER GUIDE SEE THE KERNEL APPROXIMATION SECTION FOR FURTHER DETAILS

KERNELAPPROXIMATION

ADDITIVECHI2SAMPLER APPROXIMATE FEATURE MAP FOR ADDITIVE CHI2 KERNEL

KERNELAPPROXIMATIONNNYSTROEM KERNEL APPROXIMATE A KERNEL MAP USING A SUBSET OF THE TRAINING

DATA

CONTINUED ON NEXT PAGE

620SKLEARNKERNELAPPROXIMATION KERNEL APPROXIMATION 1851

SCIKITLEARN USER GUIDE RELEASE 0213

TABLE 6141 – CONTINUED FROM PREVIOUS PAGE

KERNELAPPROXIMATIONRBFSSAMPLER GAMMA

APPROXIMATES FEATURE MAP OF AN RBF KERNEL BY MONTE CARLO APPROXIMATION OF ITS FOURIER TRANSFORM

KERNELAPPROXIMATION

SKEWEDCHI2SAMPLER APPROXIMATES FEATURE MAP OF THE “SKEWED CHISQUARED” KERNEL BY MONTE CARLO APPROXIMATION OF ITS FOURIER TRANSFORM

6201SKLEARNKERNELAPPROXIMATION ADDITIVECHI2SAMPLER

CLASSSSKLEARNKERNELAPPROXIMATION ADDITIVECHI2SAMPLER SAMPLESTEPS2 SAM

PLEINTERVALNONE

APPROXIMATE FEATURE MAP FOR ADDITIVE CHI2 KERNEL

USES SAMPLING THE FOURIER TRANSFORM OF THE KERNEL CHARACTERISTIC AT REGULAR INTERVALS

SINCE THE KERNEL THAT IS TO BE APPROXIMATED IS ADDITIVE THE COMPONENTS OF THE INPUT VECTORS CAN BE TREATED SEPARATELY EACH ENTRY IN THE ORIGINAL SPACE IS TRANSFORMED INTO 2SAMPLESTEPS1 FEATURES WHERE SAMPLESTEPS IS A PARAMETER OF THE METHOD TYPICAL VALUES OF SAMPLESTEPS INCLUDE 1 2 AND 3

OPTIMAL CHOICES FOR THE SAMPLING INTERVAL FOR CERTAIN DATA RANGES CAN BE COMPUTED SEE THE REFERENCE THE DEFAULT VALUES SHOULD BE REASONABLE

READ MORE IN THE USER GUIDE

PARAMETERS

SAMPLESTEPS INT OPTIONAL GIVES THE NUMBER OF COMPLEX SAMPLING POINTS

SAMPLEINTERVAL FLOAT OPTIONAL SAMPLING INTERVAL MUST BE SPECIFIED WHEN SAMPLESTEPS NOT IN 123

SEE ALSO

SKEWEDCHI2SAMPLER A FOURIERAPPROXIMATION TO A NONADDITIVE VARIANT OF THE CHI SQUARED KERNEL

SKLEARNMETRICSPAIRWISECHI2KERNEL THE EXACT CHI SQUARED KERNEL

SKLEARNMETRICSPAIRWISEADDITIVECHI2KERNEL THE EXACT ADDITIVE CHI SQUARED KERNEL

NOTES

THIS ESTIMATOR APPROXIMATES A SLIGHTLY DIFFERENT VERSION OF THE ADDITIVE CHI SQUARED KERNEL THEN METRIC ADDITIVECHI2 COMPUTES

REFERENCES

SEE “EFFICIENT ADDITIVE KERNELS VIA EXPLICIT FEATURE MAPS” A VEDALDI AND A ZISSERMAN PATTERN ANALYSIS AND MACHINE INTELLIGENCE 2011

EXAMPLES

FROM SKLEARNDATASETS IMPORT LOADDIGITS

FROM SKLEARNLINEARMODEL IMPORT SGDCLASSIFIER

FROM SKLEARNKERNELAPPROXIMATION IMPORT ADDITIVECHI2SAMPLER

X Y LOADDIGITSRETURNXY TRUE

1852 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213  
CHI2SAMPLER ADDITIVECHI2SAMPLERSAMPLESTEPS2  
XTRANSFORMED CHI2SAMPLERFITTRANSFORMX Y  
CLF SGDCLASSIFIERMAXITER5 RANDOMSTATE0 TOL1E3  
CLFFITXTRANSFORMED Y  
SGDCLASSIFIERALPHA00001 AVERAGEFALSE CLASSWEIGHTNONE  
EARLYSTOPPINGFALSE EPSILON01 ETA000 FITINTERCEPTTRUE  
L1RATIO015 LEARNINGRATEOPTIMAL LOSSHINGE MAXITER5  
NITERNOCHANGES5 NJOBSNONE PENALTYL2 POWERT05  
RANDOMSTATE0 SHUFFLETRUE TOL0001 VALIDATIONFRACTION01  
VERBOSE0 WARMSTARTFALSE  
CLFSCOREXTRANSFORMED Y  
09499  
METHODS  
FITSELF X Y SET THE PARAMETERS  
FITTRANSFORM SELF X Y FIT TO DATA THEN TRANSFORM IT  
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR  
SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR  
TRANSFORM SELF X APPLY APPROXIMATE FEATURE MAP TO X  
INIT SELFSAMPLESTEPS2 SAMPLEINTERVALNONE  
FITSELFXYNONE  
SET THE PARAMETERS  
PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA WHERE NSAMPLES IN THE NUMBER  
OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES  
RETURNS  
SELF OBJECT RETURNS THE TRANSFORMER  
FITTRANSFORM SELFXYNONE FITPARAMS  
FIT TO DATA THEN TRANSFORM IT  
FITS TRANSFORMER TO X AND Y WITH OPTIONAL PARAMETERS FITPARAMS AND RETURNS A TRANSFORMED VERSION OF X  
PARAMETERS  
XNUMPY ARRAY OF SHAPE NSAMPLES NFEATURES TRAINING SET  
YNUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES  
RETURNS  
XNEW NUMPY ARRAY OF SHAPE NSAMPLES NFEATURESNEW TRANSFORMED ARRAY  
GETPARAMS SELFDEEPTTRUE  
GET PARAMETERS FOR THIS ESTIMATOR  
PARAMETERS  
DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBJECTS THAT ARE ESTIMATORS  
RETURNS  
620SKLEARNKERNELAPPROXIMATION KERNEL APPROXIMATION 1853

SCIKITLEARN USER GUIDE RELEASE 0213

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELF

APPLY APPROXIMATE FEATURE MAP TO X

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES

RETURNS

XNEW ARRAY SPARSE MATRIX SHAPE NSAMPLES NFEATURES 2SAMPLESTEPS 1

WHETHER THE RETURN VALUE IS AN ARRAY OF SPARSE MATRIX DEPENDS ON THE TYPE OF THE INPUT X

6202SKLEARNKERNELAPPROXIMATION NYSTROEM

CLASSSSKLEARNKERNELAPPROXIMATION NYSTROEM KERNEL'RBF' GAMMANONE COEFONONE

DEGREENONE KERNELPARAMSNONE

NCOMPONENTS100 RANDOMSTATENONE

APPROXIMATE A KERNEL MAP USING A SUBSET OF THE TRAINING DATA

CONSTRUCTS AN APPROXIMATE FEATURE MAP FOR AN ARBITRARY KERNEL USING A SUBSET OF THE DATA AS BASIS

READ MORE IN THE USER GUIDE

PARAMETERS

KERNEL STRING OR CALLABLE DEFAULT"RBF" KERNEL MAP TO BE APPROXIMATED A CALLABLE SHOULD ACCEPT TWO ARGUMENTS AND THE KEYWORD ARGUMENTS PASSED TO THIS OBJECT AS KERNELPARAMS AND SHOULD RETURN A FLOATING POINT NUMBER

GAMMA FLOAT DEFAULTNONE GAMMA PARAMETER FOR THE RBF LAPLACIAN POLYNOMIAL EXPONENTIAL CHI2 AND SIGMOID KERNELS INTERPRETATION OF THE DEFAULT VALUE IS LEFT TO THE KERNEL SEE THE DOCUMENTATION FOR SKLEARNMETRICSPAIRWISE IGNORED BY OTHER KERNELS

COEF0 FLOAT DEFAULTNONE ZERO COEFFICIENT FOR POLYNOMIAL AND SIGMOID KERNELS IGNORED BY OTHER KERNELS

DEGREE FLOAT DEFAULTNONE DEGREE OF THE POLYNOMIAL KERNEL IGNORED BY OTHER KERNELS

KERNELPARAMS MAPPING OF STRING TO ANY OPTIONAL ADDITIONAL PARAMETERS KEYWORD ARGUMENTS FOR KERNEL FUNCTION PASSED AS CALLABLE OBJECT

NCOMPONENTS INT NUMBER OF FEATURES TO CONSTRUCT HOW MANY DATA POINTS WILL BE USED TO CONSTRUCT THE MAPPING

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RANDOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NRANDOM

ATTRIBUTES

1854 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

COMPONENTS ARRAY SHAPE NCOMPONENTS NFEATURES SUBSET OF TRAINING POINTS USED TO CONSTRUCT THE FEATURE MAP

COMPONENTINDICES ARRAY SHAPE NCOMPONENTS INDICES OF COMPONENTS IN THE TRAINING SET

NORMALIZATION ARRAY SHAPE NCOMPONENTS NCOMPONENTS NORMALIZATION MATRIX NEEDED FOR EMBEDDING SQUARE ROOT OF THE KERNEL MATRIX ON COMPONENTS

SEE ALSO

RBFSAMPLER AN APPROXIMATION TO THE RBF KERNEL USING RANDOM FOURIER FEATURES

SKLEARNMETRICSPAIRWISEKERNELMETRICS LIST OF BUILTIN KERNELS

REFERENCES

- WILLIAMS CKI AND SEEGER M “USING THE NYSTROEM METHOD TO SPEED UP KERNEL MACHINES” ADVANCES IN NEURAL INFORMATION PROCESSING SYSTEMS 2001
- T YANG Y LI M MAHDAVI R JIN AND Z ZHOU “NYSTROEM METHOD VS RANDOM FOURIER FEATURES A THEORETICAL AND EMPIRICAL COMPARISON” ADVANCES IN NEURAL INFORMATION PROCESSING SYSTEMS 2012

EXAMPLES

```
FROM SKLEARN IMPORT DATASETS SVM
FROM SKLEARNKERNELAPPROXIMATION IMPORT NYSTROEM
DIGITS DATASETSLOADDIGITSNCLASS9
DATA DIGITS DATA _16
CLF SVMLINEARSVC
FEATUREMAPNYSTROEM NYSTROEMGAMMA2
RANDOMSTATE1
NCOMPONENTS300
DATATransformed FEATUREMAPNYSTROEMFITTRANSFORMDATA
CLFFITDATATransformed DIGITSTARGET
```

LINEARSVCC10 CLASSWEIGHTNONE DUALTRUE FITINTERCEPTTRUE  
INTERCEPTSCALING1 LOSSSQUAREDHINGE MAXITER1000  
MULTICLASSOVR PENALTYL2 RANDOMSTATENONE TOL00001  
VERBOSE0  
CLFSCOREDATATransformed DIGITSTARGET  
09987

METHODS

```
FITSELF X Y FIT ESTIMATOR TO DATA
FITTRANSFORM SELF X Y FIT TO DATA THEN TRANSFORM IT
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR
SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR
TRANSFORM SELF X APPLY FEATURE MAP TO X
INIT SELFKERNEL'RBF' GAMMANONE COEFONONE DEGREEONE KERNELPARAMSNONE
NCOMPONENTS100 RANDOMSTATENONE
620SKLEARNKERNELAPPROXIMATION KERNEL APPROXIMATION 1855
```

SCIKITLEARN USER GUIDE RELEASE 0213

FITSELFXYNONE

FIT ESTIMATOR TO DATA

SAMPLES A SUBSET OF TRAINING POINTS COMPUTES KERNEL ON THESE AND COMPUTES NORMALIZATION MATRIX

PARAMETERS

XARRAYLIKE SHAPENSAMPLES NFEATURE TRAINING DATA

FITTRANSFORM SELFXYNONE FITPARAMS

FIT TO DATA THEN TRANSFORM IT

FITS TRANSFORMER TO X AND Y WITH OPTIONAL PARAMETERS FITPARAMS AND RETURNS A TRANSFORMED VERSION OF X

PARAMETERS

XNUMPY ARRAY OF SHAPE NSAMPLES NFEATURES TRAINING SET

YNUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES

RETURNS

XNEW NUMPY ARRAY OF SHAPE NSAMPLES NFEATURESNEW TRANSFORMED ARRAY

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE

PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT

OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELFXY

APPLY FEATURE MAP TO X

COMPUTES AN APPROXIMATE FEATURE MAP USING THE KERNEL BETWEEN SOME TRAINING POINTS AND X

PARAMETERS

XARRAYLIKE SHAPENSAMPLES NFEATURES DATA TO TRANSFORM

RETURNS

XTRANSFORMED ARRAY SHAPENSAMPLES NCOMPONENTS TRANSFORMED DATA

EXAMPLES USING SKLEARNKERNELAPPROXIMATIONNNYSTROEM

- EXPLICIT FEATURE MAP APPROXIMATION FOR RBF KERNELS

1856 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213  
6203SKLEARNKERNELAPPROXIMATION RBFSAMPLER  
CLASSSSKLEARNKERNELAPPROXIMATION RBFSAMPLER GAMMA10 NCOMPONENTS100 RAN  
DOMSTATENONE  
APPROXIMATES FEATURE MAP OF AN RBF KERNEL BY MONTE CARLO APPROXIMATION OF ITS FOURIER TRANSFORM  
IT IMPLEMENTS A VARIANT OF RANDOM KITCHEN SINKS1  
READ MORE IN THE USER GUIDE  
PARAMETERS  
GAMMA FLOAT PARAMETER OF RBF KERNEL EXPGAMMA X2  
NCOMPONENTS INT NUMBER OF MONTE CARLO SAMPLES PER ORIGINAL FEATURE EQUALS THE DIMEN  
SIONALITY OF THE COMPUTED FEATURE SPACE  
RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RAN  
DOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE  
RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS  
THE RANDOMSTATE INSTANCE USED BY NPRANDOM  
NOTES  
SEE “RANDOM FEATURES FOR LARGESCALE KERNEL MACHINES” BY A RAHIMI AND BENJAMIN RECHT  
1 “WEIGHTED SUMS OF RANDOM KITCHEN SINKS REPLACING MINIMIZATION WITH RANDOMIZATION IN LEARNING” BY A  
RAHIMI AND BENJAMIN RECHT [HTTPSPEOPLEEECSBERKELEYEDUBRECHTPAPERS08RAHRECNI.PDF](https://people.eecs.berkeley.edu/brecht/papers/08RAHRECNI.PDF)  
EXAMPLES  
FROM SKLEARNKERNELAPPROXIMATION IMPORT RBFSAMPLER  
FROM SKLEARNLINEARMODEL IMPORT SGDCLASSIFIER  
X 0 0 1 1 1 0 0 1  
Y 0 0 1 1  
RBFFEATURE RBFSAMPLERGAMMA1 RANDOMSTATE1  
XFEATURES RBFFEATUREFITTRANSFORMX  
CLF SGDCLASSIFIERMAXITER5 TOL1E3  
CLFFITXFEATURES Y  
  
SGDCLASSIFIERALPHA00001 AVERAGEFALSE CLASSWEIGHTNONE  
EARLYSTOPPINGFALSE EPSILON01 ETA000 FITINTERCEPTTRUE  
L1RATIO015 LEARNINGRATEOPTIMAL LOSSHINGE MAXITER5  
NITERNOCHANGE5 NJOBSNONE PENALTYL2 POWERT05  
RANDOMSTATENONE SHUFFLETRUE TOL0001 VALIDATIONFRACTION01  
VERBOSE0 WARMSTARTFALSE  
CLFSCOREXFEATURES Y  
10  
METHODS  
FITSELF X Y FIT THE MODEL WITH X  
FITTRANSFORM SELF X Y FIT TO DATA THEN TRANSFORM IT  
CONTINUED ON NEXT PAGE  
620SKLEARNKERNELAPPROXIMATION KERNEL APPROXIMATION 1857

SCIKITLEARN USER GUIDE RELEASE 0213

TABLE 6144 – CONTINUED FROM PREVIOUS PAGE

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF X APPLY THE APPROXIMATE FEATURE MAP TO X

INIT SELF GAMMA 10 NCOMPONENTS 100 RANDOM STATE NONE

FIT SELF X NONE

FIT THE MODEL WITH X

SAMPLES RANDOM PROJECTION ACCORDING TO NFEATURES

PARAMETERS

X ARRAY LIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES TRAINING DATA WHERE

NSAMPLES IN THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

RETURNS

SELF OBJECT RETURNS THE TRANSFORMER

FIT TRANSFORM SELF X NONE FIT PARAMS

FIT TO DATA THEN TRANSFORM IT

FITS TRANSFORMER TO X AND Y WITH OPTIONAL PARAMETERS FIT PARAMS AND RETURNS A TRANSFORMED VERSION OF X

PARAMETERS

X Numpy array of shape NSAMPLES NFEATURES TRAINING SET

Y Numpy array of shape NSAMPLES TARGET VALUES

RETURNS

X NEW Numpy array of shape NSAMPLES NFEATURES NEW TRANSFORMED ARRAY

GETPARAMS SELF DEEP TRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

SETPARAMS SELF PARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE

PARAMETERS OF THE FORM COMPONENT PARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT

OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELF X

APPLY THE APPROXIMATE FEATURE MAP TO X

PARAMETERS

X ARRAY LIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES NEW DATA WHERE NSAMPLES

IN THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

1858 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

RETURNS

XNEW ARRAYLIKE SHAPE NSAMPLES NCOMPONENTS

EXAMPLES USING SKLEARNKERNELAPPROXIMATIONRBFSSAMPLER

- EXPLICIT FEATURE MAP APPROXIMATION FOR RBF KERNELS

6204SKLEARNKERNELAPPROXIMATION SKEWEDCHI2SAMPLER

CLASSSSKLEARNKERNELAPPROXIMATION SKEWEDCHI2SAMPLER SKEWEDNESS10

NCOMPONENTS100 RAN

DOMSTATENONE

APPROXIMATES FEATURE MAP OF THE “SKEWED CHISQUARED” KERNEL BY MONTE CARLO APPROXIMATION OF ITS FOURIER TRANSFORM

READ MORE IN THE USER GUIDE

PARAMETERS

SKEWEDNESS FLOAT “SKEWEDNESS” PARAMETER OF THE KERNEL NEEDS TO BE CROSSVALIDATED

NCOMPONENTS INT NUMBER OF MONTE CARLO SAMPLES PER ORIGINAL FEATURE EQUALS THE DIMENSIONALITY OF THE COMPUTED FEATURE SPACE

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RAN

DOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE

RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS

THE RANDOMSTATE INSTANCE USED BY NPRANDOM

SEE ALSO

ADDITIVECHI2SAMPLER A DIFFERENT APPROACH FOR APPROXIMATING AN ADDITIVE VARIANT OF THE CHI SQUARED KERNEL

SKLEARNMETRICSPAIRWISECHI2KERNEL THE EXACT CHI SQUARED KERNEL

REFERENCES

SEE “RANDOM FOURIER APPROXIMATIONS FOR SKEWED MULTIPLICATIVE HISTOGRAM KERNELS” BY FUXIN LI CATALIN IONESCU AND CRISTIAN SMINCHISESCU

EXAMPLES

```
FROM SKLEARNKERNELAPPROXIMATION IMPORT SKEWEDCHI2SAMPLER
FROM SKLEARNLINEARMODEL IMPORT SGDCLASSIFIER
X 0 0 1 1 1 0 0 1
Y 0 0 1 1
CHI2FEATURE SKEWEDCHI2SAMPLERSKEWEDNESS01
NCOMPONENTS10
RANDOMSTATE0
XFEATURES CHI2FEATUREFITTRANSFORMX Y
CLF SGDCLASSIFIERMAXITER10 TOL1E3
CLFFITXFEATURES Y
SGDCLASSIFIERALPHA00001 AVERAGEFALSE CLASSWEIGHTNONE
620SKLEARNKERNELAPPROXIMATION KERNEL APPROXIMATION 1859
```

SCIKITLEARN USER GUIDE RELEASE 0213  
EARLYSTOPPINGFALSE EPSILON01 ETA000 FITINTERCEPTTRUE  
L1RATIO015 LEARNINGRATEOPTIMAL LOSSHINGE MAXITER10  
NITERNOCHANGES5 NJOBSNONE PENALTYL2 POWERT05  
RANDOMSTATENONE SHUFFLETRUE TOL0001 VALIDATIONFRACTION01  
VERBOSE0 WARMSTARTFALSE  
CLFSOREXFEATURES Y  
10  
METHODS  
FITSELF X Y FIT THE MODEL WITH X  
FITTRANSFORM SELF X Y FIT TO DATA THEN TRANSFORM IT  
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR  
SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR  
TRANSFORM SELF X APPLY THE APPROXIMATE FEATURE MAP TO X  
INIT SELF SKEWEDNESS10 NCOMPONENTS100 RANDOMSTATENONE  
FITSELFXYNONE  
FIT THE MODEL WITH X  
SAMPLES RANDOM PROJECTION ACCORDING TO NFEATURES  
PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA WHERE NSAMPLES IN THE NUMBER  
OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES  
RETURNS  
SELF OBJECT RETURNS THE TRANSFORMER  
FITTRANSFORM SELFXYNONE FITPARAMS  
FIT TO DATA THEN TRANSFORM IT  
FITS TRANSFORMER TO X AND Y WITH OPTIONAL PARAMETERS FITPARAMS AND RETURNS A TRANSFORMED VERSION OF X  
PARAMETERS  
XNUMPY ARRAY OF SHAPE NSAMPLES NFEATURES TRAINING SET  
YNUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES  
RETURNS  
XNEW NUMPY ARRAY OF SHAPE NSAMPLES NFEATURESNEW TRANSFORMED ARRAY  
GETPARAMS SELFDEEPTTRUE  
GET PARAMETERS FOR THIS ESTIMATOR  
PARAMETERS  
DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBOBJECTS THAT ARE ESTIMATORS  
RETURNS  
PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES  
1860 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELF

APPLY THE APPROXIMATE FEATURE MAP TO X

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES NEW DATA WHERE NSAMPLES IN THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES ALL VALUES OF X MUST BE STRICTLY GREATER THAN “SKEWEDNESS”

RETURNS

XNEW ARRAYLIKE SHAPE NSAMPLES NCOMPONENTS

621SKLEARNKERNELRIDGE KERNEL RIDGE REGRESSION

MODULESKLEARNKERNELRIDGE IMPLEMENTS KERNEL RIDGE REGRESSION

USER GUIDE SEE THE KERNEL RIDGE REGRESSION SECTION FOR FURTHER DETAILS

KERNELRIDGEKERNELRIDGE ALPHA KERNEL    KERNEL RIDGE REGRESSION

6211SKLEARNKERNELRIDGE KERNELRIDGE

CLASSSKLEARNKERNELRIDGE KERNELRIDGE ALPHA1 KERNEL’LINEAR’ GAMMANONE DEGREE3

COEF01 KERNELPARAMSNONE

KERNEL RIDGE REGRESSION

KERNEL RIDGE REGRESSION KRR COMBINES RIDGE REGRESSION LINEAR LEAST SQUARES WITH L2NORM REGULARIZATION WITH THE KERNEL TRICK IT THUS LEARNS A LINEAR FUNCTION IN THE SPACE INDUCED BY THE RESPECTIVE KERNEL AND THE DATA FOR NONLINEAR KERNELS THIS CORRESPONDS TO A NONLINEAR FUNCTION IN THE ORIGINAL SPACE

THE FORM OF THE MODEL LEARNED BY KRR IS IDENTICAL TO SUPPORT VECTOR REGRESSION SVR HOWEVER DIFFERENT LOSS FUNCTIONS ARE USED KRR USES SQUARED ERROR LOSS WHILE SUPPORT VECTOR REGRESSION USES EPSILONINSENSITIVE LOSS BOTH COMBINED WITH L2 REGULARIZATION IN CONTRAST TO SVR FITTING A KRR MODEL CAN BE DONE IN CLOSEDFORM AND IS TYPICALLY FASTER FOR MEDIUMSIZED DATASETS ON THE OTHER HAND THE LEARNED MODEL IS NONSPARSE AND THUS SLOWER THAN SVR WHICH LEARNS A SPARSE MODEL FOR EPSILON 0 AT PREDICTIONTIME

THIS ESTIMATOR HAS BUILTIN SUPPORT FOR MULTIVARIATE REGRESSION IE WHEN Y IS A 2DARRAY OF SHAPE NSAMPLES NTARGETS

READ MORE IN THE USER GUIDE

PARAMETERS

ALPHA FLOAT ARRAYLIKE SHAPE NTARGETS SMALL POSITIVE VALUES OF ALPHA IMPROVE THE CONDITIONING OF THE PROBLEM AND REDUCE THE VARIANCE OF THE ESTIMATES ALPHA CORRESPONDS TO 2C1 IN OTHER LINEAR MODELS SUCH AS LOGISTICREGRESSION OR LINEARSVC IF AN ARRAY IS

621SKLEARNKERNELRIDGE KERNEL RIDGE REGRESSION 1861

SCIKITLEARN USER GUIDE RELEASE 0213

PASSED PENALTIES ARE ASSUMED TO BE SPECIFIC TO THE TARGETS HENCE THEY MUST CORRESPOND IN NUMBER

KERNEL STRING OR CALLABLE DEFAULT"LINEAR" KERNEL MAPPING USED INTERNALLY A CALLABLE SHOULD ACCEPT TWO ARGUMENTS AND THE KEYWORD ARGUMENTS PASSED TO THIS OBJECT AS KERNELPARAMS AND SHOULD RETURN A FLOATING POINT NUMBER SET TO "PRECOMPUTED" IN ORDER TO PASS A PRECOMPUTED KERNEL MATRIX TO THE ESTIMATOR METHODS INSTEAD OF SAMPLES

GAMMA FLOAT DEFAULTNONE GAMMA PARAMETER FOR THE RBF LAPLACIAN POLYNOMIAL EXPONENTIAL CHI2 AND SIGMOID KERNELS INTERPRETATION OF THE DEFAULT VALUE IS LEFT TO THE KERNEL SEE THE DOCUMENTATION FOR SKLEARNMETRICSPAIRWISE IGNORED BY OTHER KERNELS

DEGREE FLOAT DEFAULT3 DEGREE OF THE POLYNOMIAL KERNEL IGNORED BY OTHER KERNELS

COEF0 FLOAT DEFAULT1 ZERO COEFFICIENT FOR POLYNOMIAL AND SIGMOID KERNELS IGNORED BY OTHER KERNELS

KERNELPARAMS MAPPING OF STRING TO ANY OPTIONAL ADDITIONAL PARAMETERS KEYWORD ARGUMENTS FOR KERNEL FUNCTION PASSED AS CALLABLE OBJECT

ATTRIBUTES

DUALCOEF ARRAY SHAPE NSAMPLES OR NSAMPLES NTARGETS REPRESENTATION OF WEIGHT VECTORS IN KERNEL SPACE

XFIT ARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES TRAINING DATA WHICH IS ALSO REQUIRED FOR PREDICTION IF KERNEL "PRECOMPUTED" THIS IS INSTEAD THE PRECOMPUTED TRAINING MATRIX SHAPE NSAMPLES NSAMPLES

SEE ALSO

SKLEARNLINEARMODELRIDGE LINEAR RIDGE REGRESSION

SKLEARNSVMSVR SUPPORT VECTOR REGRESSION IMPLEMENTED USING LIBSVM

REFERENCES

- KEVIN P MURPHY "MACHINE LEARNING A PROBABILISTIC PERSPECTIVE" THE MIT PRESS CHAPTER 1443 PP 492493

EXAMPLES

FROM SKLEARNKERNELRIDGE IMPORT KERNELRIDGE

IMPORT NUMPY AS NP

NSAMPLES NFEATURES 10 5

RNG NPRANDOMRANDOMSTATE0

Y RNGRANDNNSAMPLES

X RNGRANDNNSAMPLES NFEATURES

CLF KERNELRIDGEALPHA10

CLFFITX Y

KERNELRIDGEALPHA10 COEF01 DEGREE3 GAMMANONE KERNELLINEAR

KERNELPARAMSNONE

METHODS

1862 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

FITSELF X Y SAMPLEWEIGHT FIT KERNEL RIDGE REGRESSION MODEL

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PREDICT SELF X PREDICT USING THE KERNEL RIDGE MODEL

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELFALPHA1 KERNEL'LINEAR' GAMMANONE DEGREE3 COEF01 KERNELPARAMSNONE

FITSELFXYNONE SAMPLEWEIGHTNONE

FIT KERNEL RIDGE REGRESSION MODEL

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES TRAINING DATA IF KERNEL "PRECOMPUTED" THIS IS INSTEAD A PRECOMPUTED KERNEL MATRIX SHAPE NSAMPLES NSAMPLES

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NTARGETS TARGET VALUES

SAMPLEWEIGHT FLOAT OR ARRAYLIKE OF SHAPE NSAMPLES INDIVIDUAL WEIGHTS FOR EACH SAMPLE IGNORED IF NONE IS PASSED

RETURNS

SELF RETURNS AN INSTANCE OF SELF

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELF X

PREDICT USING THE KERNEL RIDGE MODEL

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES SAMPLES IF KERNEL "PRECOMPUTED" THIS IS INSTEAD A PRECOMPUTED KERNEL MATRIX SHAPE NSAMPLES NSAMPLESFITTED WHERE NSAMPLESFITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THIS ESTIMATOR

RETURNS

CARRAY SHAPE NSAMPLES OR NSAMPLES NTARGETS RETURNS PREDICTED VALUES

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{U}{V}$  WHERE U IS THE RESIDUAL SUM OF SQUARES  $Y_{TRUE} - Y_{PRED}$  2SUM AND V IS THE TOTAL SUM OF SQUARES  $Y_{TRUE} - Y_{TRUEMEAN}$  2SUM THE BEST POSSIBLE SCORE IS 10 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 00

PARAMETERS

621SKLEARNKERNELRIDGE KERNEL RIDGE REGRESSION 1863

SCIKITLEARN USER GUIDE RELEASE 0213

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES FOR SOME ESTIMATORS THIS MAY  
BE A PRECOMPUTED KERNEL MATRIX INSTEAD SHAPE NSAMPLES NSAMPLESFITTED WHERE  
NSAMPLESFITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR  
YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE VALUES FOR X  
SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT R2 OF SELF PREDICT X WRT Y

NOTES

THE R2 SCORE USED WHEN CALLING SCORE ON A REGRESSOR WILL USE MULTIOUTPUTUNIFORMAVERAGE  
FROM VERSION 023 TO KEEP CONSISTENT WITH METRICSR2SCORE THIS WILL INFLUENCE THE SCORE  
METHOD OF ALL THE MULTIOUTPUT REGRESSORS EXCEPT FOR MULTIOUTPUTMULTIOUTPUTREGRESSOR  
TO SPECIFY THE DEFAULT VALUE MANUALLY AND AVOID THE WARNING PLEASE EITHER CALL METRICSR2SCORE  
DIRECTLY OR MAKE A CUSTOM SCORER WITH METRICSMAKESCORER THE BUILTIN SCORER R2 USES  
MULTIOUTPUTUNIFORMAVERAGE

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE  
PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT  
OF A NESTED OBJECT

RETURNS

SELF

EXAMPLES USING SKLEARNKERNELRIDGEKERNELRIDGE

- COMPARISON OF KERNEL RIDGE REGRESSION AND SVR
- COMPARISON OF KERNEL RIDGE AND GAUSSIAN PROCESS REGRESSION

622SKLEARNLINEARMODEL GENERALIZED LINEAR MODELS

THESKLEARNLINEARMODEL MODULE IMPLEMENTS GENERALIZED LINEAR MODELS IT INCLUDES RIDGE REGRESSION  
BAYESIAN REGRESSION LASSO AND ELASTIC NET ESTIMATORS COMPUTED WITH LEAST ANGLE REGRESSION AND COORDINATE DE  
SCENT IT ALSO IMPLEMENTS STOCHASTIC GRADIENT DESCENT RELATED ALGORITHMS  
USER GUIDE SEE THE GENERALIZED LINEAR MODELS SECTION FOR FURTHER DETAILS

LINEARMODELARDREGRESSION NITER TOL BAYESIAN ARD REGRESSION

LINEARMODELBAYESIANRIDGE NITER TOL BAYESIAN RIDGE REGRESSION

LINEARMODELELASTICNET ALPHA L1RATIO LINEAR REGRESSION WITH COMBINED L1 AND L2 PRIORS AS REGU  
LARIZER

LINEARMODELELASTICNETCV L1RATIO EPS ELASTIC NET MODEL WITH ITERATIVE FITTING ALONG A REGULARIZA  
TION PATH

LINEARMODELHUBERREGRESSOR EPSILON LINEAR REGRESSION MODEL THAT IS ROBUST TO OUTLIERS

LINEARMODELLARS FITINTERCEPT VERBOSE LEAST ANGLE REGRESSION MODEL AKA

CONTINUED ON NEXT PAGE

1864 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

TABLE 6148 - CONTINUED FROM PREVIOUS PAGE

LINEARMODELLARSCV FITINTERCEPT CROSSVALIDATED LEAST ANGLE REGRESSION MODEL

LINEARMODELLASSO ALPHA FITINTERCEPT LINEAR MODEL TRAINED WITH L1 PRIOR AS REGULARIZER AKA THE LASSO

LINEARMODELLASSOCV EPS NALPHAS LASSO LINEAR MODEL WITH ITERATIVE FITTING ALONG A REGULARIZATION PATH

LINEARMODELLASSOLARS ALPHA LASSO MODEL FIT WITH LEAST ANGLE REGRESSION AKA

LINEARMODELLASSOLARSCV FITINTERCEPT CROSSVALIDATED LASSO USING THE LARS ALGORITHM

LINEARMODELLASSOLARSIC CRITERION LASSO MODEL FIT WITH LARS USING BIC OR AIC FOR MODEL SELECTION

LINEARMODELLINEARREGRESSION ORDINARY LEAST SQUARES LINEAR REGRESSION

LINEARMODELLOGISTICREGRESSION PENALTY

LOGISTIC REGRESSION AKA LOGIT MAXENT CLASSIFIER

LINEARMODELLOGISTICREGRESSIONCV CS

LOGISTIC REGRESSION CV AKA LOGIT MAXENT CLASSIFIER

LINEARMODELMULTITASKLASSO ALPHA MULTITASK LASSO MODEL TRAINED WITH L1L2 MIXEDNORM AS REGULARIZER

LINEARMODELMULTITASKELASTICNET ALPHA

MULTITASK ELASTICNET MODEL TRAINED WITH L1L2 MIXEDNORM AS REGULARIZER

LINEARMODELMULTITASKLASSOCV EPS MULTITASK LASSO MODEL TRAINED WITH L1L2 MIXEDNORM AS REGULARIZER

LINEARMODELMULTITASKELASTICNETCV MULTITASK L1L2 ELASTICNET WITH BUILTIN CROSSVALIDATION

LINEARMODELORTHOGONALMATCHINGPURSUIT ORTHOGONAL MATCHING PURSUIT MODEL OMP

LINEARMODELORTHOGONALMATCHINGPURSUITCV CROSSVALIDATED ORTHOGONAL MATCHING PURSUIT MODEL OMP

LINEARMODELPASSIVEAGGRESSIVECLASSIFIER PASSIVE AGGRESSIVE CLASSIFIER

LINEARMODELPASSIVEAGGRESSIVEREGRESSOR C

PASSIVE AGGRESSIVE REGRESSOR

LINEARMODELPERCEPTRON PENALTY ALPHA READ MORE IN THE USER GUIDE

LINEARMODELRANSACREGRESSOR RANSAC RANDOM SAMPLE CONSENSUS ALGORITHM

LINEARMODELRIDGE ALPHA FITINTERCEPT LINEAR LEAST SQUARES WITH L2 REGULARIZATION

LINEARMODELRIDGECLASSIFIER ALPHA CLASSIFIER USING RIDGE REGRESSION

LINEARMODELRIDGECLASSIFIERCV ALPHAS

RIDGE CLASSIFIER WITH BUILTIN CROSSVALIDATION

LINEARMODELRIDGECV ALPHAS RIDGE REGRESSION WITH BUILTIN CROSSVALIDATION

LINEARMODELSGDCLASSIFIER LOSS PENALTY

LINEAR CLASSIFIERS SVM LOGISTIC REGRESSION AO WITH SGD TRAINING

LINEARMODELSGDREGRESSOR LOSS PENALTY LINEAR MODEL FITTED BY MINIMIZING A REGULARIZED EMPIRICAL LOSS WITH SGD

LINEARMODELTHEILSENREGRESSOR THEILSEN ESTIMATOR ROBUST MULTIVARIATE REGRESSION MODEL

6221SKLEARNLINEARMODEL ARDREGRESSION

CLASSSKLEARNLINEARMODEL ARDREGRESSION NITER300 TOL0001 ALPHA11E06 ALPHA21E06 LAMBDA11E06

LAMBDA21E06 COMPUTESCOREFALSE THRESH

OLDLAMBDA100000 FITINTERCEPTTRUE NORMAL

IZEFALSE COPYXTRUE VERBOSEFALSE

BAYESIAN ARD REGRESSION

FIT THE WEIGHTS OF A REGRESSION MODEL USING AN ARD PRIOR THE WEIGHTS OF THE REGRESSION MODEL ARE ASSUMED TO BE IN GAUSSIAN DISTRIBUTIONS ALSO ESTIMATE THE PARAMETERS LAMBDA PRECISIONS OF THE DISTRIBUTIONS OF THE WEIGHTS AND ALPHA PRECISION OF THE DISTRIBUTION OF THE NOISE THE ESTIMATION IS DONE BY AN ITERATIVE PROCEDURES

622SKLEARNLINEARMODEL GENERALIZED LINEAR MODELS 1865

SCIKITLEARN USER GUIDE RELEASE 0213

EVIDENCE MAXIMIZATION

READ MORE IN THE USER GUIDE

PARAMETERS

NITER INT OPTIONAL MAXIMUM NUMBER OF ITERATIONS DEFAULT IS 300

TOLFLOAT OPTIONAL STOP THE ALGORITHM IF W HAS CONVERGED DEFAULT IS 1E3

ALPHA1 FLOAT OPTIONAL HYPERPARAMETER SHAPE PARAMETER FOR THE GAMMA DISTRIBUTION PRIOR OVER THE ALPHA PARAMETER DEFAULT IS 1E6

ALPHA2 FLOAT OPTIONAL HYPERPARAMETER INVERSE SCALE PARAMETER RATE PARAMETER FOR THE GAMMA DISTRIBUTION PRIOR OVER THE ALPHA PARAMETER DEFAULT IS 1E6

LAMBDA1 FLOAT OPTIONAL HYPERPARAMETER SHAPE PARAMETER FOR THE GAMMA DISTRIBUTION PRIOR OVER THE LAMBDA PARAMETER DEFAULT IS 1E6

LAMBDA2 FLOAT OPTIONAL HYPERPARAMETER INVERSE SCALE PARAMETER RATE PARAMETER FOR THE GAMMA DISTRIBUTION PRIOR OVER THE LAMBDA PARAMETER DEFAULT IS 1E6

COMPUTESCORE BOOLEAN OPTIONAL IF TRUE COMPUTE THE OBJECTIVE FUNCTION AT EACH STEP OF THE MODEL DEFAULT IS FALSE

THRESHOLDLAMBDA FLOAT OPTIONAL THRESHOLD FOR REMOVING PRUNING WEIGHTS WITH HIGH PRECISION FROM THE COMPUTATION DEFAULT IS 1E4

FITINTERCEPT BOOLEAN OPTIONAL WHETHER TO CALCULATE THE INTERCEPT FOR THIS MODEL IF SET TO FALSE NO INTERCEPT WILL BE USED IN CALCULATIONS EG DATA IS EXPECTED TO BE ALREADY CENTERED DEFAULT IS TRUE

NORMALIZE BOOLEAN OPTIONAL DEFAULT FALSE THIS PARAMETER IS IGNORED WHEN FITINTERCEPT IS SET TO FALSE IF TRUE THE REGRESSORS X WILL BE NORMALIZED BEFORE REGRESSION BY SUBTRACTING THE MEAN AND DIVIDING BY THE L2NORM IF YOU WISH TO STANDARDIZE PLEASE USE SKLEARNPREPROCESSINGSTANDARDSCALER BEFORE CALLINGFIT ON AN ESTIMATOR WITH NORMALIZEFALSE

COPYX BOOLEAN OPTIONAL DEFAULT TRUE IF TRUE X WILL BE COPIED ELSE IT MAY BE OVERWRITTEN

VERBOSE BOOLEAN OPTIONAL DEFAULT FALSE VERBOSE MODE WHEN FITTING THE MODEL

ATTRIBUTES

COEF ARRAY SHAPE NFEATURES COEFFICIENTS OF THE REGRESSION MODEL MEAN OF DISTRIBUTION

ALPHA FLOAT ESTIMATED PRECISION OF THE NOISE

LAMBDA ARRAY SHAPE NFEATURES ESTIMATED PRECISIONS OF THE WEIGHTS

SIGMA ARRAY SHAPE NFEATURES NFEATURES ESTIMATED VARIANCECOVARIANCE MATRIX OF THE WEIGHTS

SCORES FLOAT IF COMPUTED VALUE OF THE OBJECTIVE FUNCTION TO BE MAXIMIZED

NOTES

FOR AN EXAMPLE SEE EXAMPLESLINEARMODELPLOTARDPY

1866 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

REFERENCES

D J C MACKAY BAYESIAN NONLINEAR MODELING FOR THE PREDICTION COMPETITION ASHRAE TRANSACTIONS 1994

R SALAKHUTDINOV LECTURE NOTES ON STATISTICAL MACHINE LEARNING HTTPWWWUTSTATTORONTOEDURSALAKHU

STA4273NOTESLECTURE2PDFPAGE15 THEIR BETA IS OUR SELFALPHA THEIR ALPHA IS OUR SELF LAMBDA

ARD IS A LITTLE DIFFERENT THAN THE SLIDE ONLY DIMENSIONSFEATURES FOR WHICH SELF LAMBDA SELF

THRESHOLD LAMBDA ARE KEPT AND THE REST ARE DISCARDED

EXAMPLES

```
FROM SKLEARN IMPORT LINEAR MODEL
CLF LINEAR MODEL ARD REGRESSION
CLF FIT 00 1 1 2 2 0 1 2
```

ARD REGRESSION ALPHA 11E06 ALPHA 21E06 COMPUTE SCORE FALSE

COPY X TRUE FIT INTERCEPT TRUE LAMBDA 11E06 LAMBDA 21E06

NITER 300 NORMALIZE FALSE THRESHOLD LAMBDA 100000 TOL 0001

VERBOSE FALSE

```
CLF PREDICT 1 1
ARRAY 1
METHODS
FIT SELF X Y FIT THE ARD REGRESSION MODEL ACCORDING TO THE GIVEN
TRAINING DATA AND PARAMETERS
GET PARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR
PREDICT SELF X RETURN STD PREDICT USING THE LINEAR MODEL
SCORE SELF X Y SAMPLE WEIGHT RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PRE
DICTION
SET PARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR
INIT SELF NITER 300 TOL 0001 ALPHA 11E06 ALPHA 21E06 LAMBDA 11E06
LAMBDA 21E06 COMPUTE SCORE FALSE THRESHOLD LAMBDA 100000 FIT INTERCEPT TRUE
NORMALIZE FALSE COPY X TRUE VERBOSE FALSE
FIT SELF X Y
FIT THE ARD REGRESSION MODEL ACCORDING TO THE GIVEN TRAINING DATA AND PARAMETERS
ITERATIVE PROCEDURE TO MAXIMIZE THE EVIDENCE
PARAMETERS
X ARRAY LIKE SHAPE NSAMPLES NFEATURES TRAINING VECTOR WHERE NSAMPLES IN THE
NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES
Y ARRAY SHAPE NSAMPLES TARGET VALUES INTEGERS WILL BE CAST TO X'S DTYPE IF NECESSARY
RETURNS
SELF RETURNS AN INSTANCE OF SELF
GET PARAMS SELF DEEP TRUE
GET PARAMETERS FOR THIS ESTIMATOR
PARAMETERS
622 SKLEARN LINEAR MODEL GENERALIZED LINEAR MODELS 1867
```

SCIKITLEARN USER GUIDE RELEASE 0213

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELFRETURNSTDFALSE

PREDICT USING THE LINEAR MODEL

IN ADDITION TO THE MEAN OF THE PREDICTIVE DISTRIBUTION ALSO ITS STANDARD DEVIATION CAN BE RETURNED

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES SAMPLES

RETURNSTD BOOLEAN OPTIONAL WHETHER TO RETURN THE STANDARD DEVIATION OF POSTERIOR PREDICTION

RETURNS

YMEAN ARRAY SHAPE NSAMPLES MEAN OF PREDICTIVE DISTRIBUTION OF QUERY POINTS

YSTD ARRAY SHAPE NSAMPLES STANDARD DEVIATION OF PREDICTIVE DISTRIBUTION OF QUERY POINTS

SCORESELFYSAMPLEWEIGHTNONE

RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{U}{V}$  WHERE U IS THE RESIDUAL SUM OF SQUARES  $Y_{TRUE} - Y_{PRED}$  2SUM AND V IS THE TOTAL SUM OF SQUARES  $Y_{TRUE} - Y_{TRUEMEAN}$  2SUM THE BEST POSSIBLE SCORE IS 10 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 00

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES FOR SOME ESTIMATORS THIS MAY BE A PRECOMPUTED KERNEL MATRIX INSTEAD SHAPE NSAMPLES NSAMPLESFITTED WHERE NSAMPLESFITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE VALUES FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT R2 OF SELF PREDICTX WRT Y

NOTES

THE R2 SCORE USED WHEN CALLING SCORE ON A REGRESSOR WILL USE MULTIOUTPUTUNIFORMAVERAGE FROM VERSION 023 TO KEEP CONSISTENT WITH METRICSR2SCORE THIS WILL INFLUENCE THE SCORE METHOD OF ALL THE MULTIOUTPUT REGRESSORS EXCEPT FOR MULTIOUTPUTMULTIOUTPUTREGRESSOR TO SPECIFY THE DEFAULT VALUE MANUALLY AND AVOID THE WARNING PLEASE EITHER CALL METRICSR2SCORE DIRECTLY OR MAKE A CUSTOM SCORER WITH METRICSMAKESCORER THE BUILTIN SCORER R2 USES MULTIOUTPUTUNIFORMAVERAGE

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

1868 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

RETURNS  
SELF

EXAMPLES USING SKLEARNLINEARMODELARDREGRESSION

- AUTOMATIC RELEVANCE DETERMINATION REGRESSION ARD

6222SKLEARNLINEARMODEL BAYESIANRIDGE

CLASSSSKLEARNLINEARMODEL BAYESIANRIDGE NITER300 TOL0001 ALPHA11E06  
ALPHA21E06 LAMBDA11E06 LAMBDA21E

06 COMPUTESCOREFALSE FITINTERCEPTTRUE  
NORMALIZEFALSE COPYXTRUE VERBOSEFALSE

BAYESIAN RIDGE REGRESSION

FIT A BAYESIAN RIDGE MODEL SEE THE NOTES SECTION FOR DETAILS ON THIS IMPLEMENTATION AND THE OPTIMIZATION OF THE  
REGULARIZATION PARAMETERS LAMBDA PRECISION OF THE WEIGHTS AND ALPHA PRECISION OF THE NOISE

READ MORE IN THE USER GUIDE

PARAMETERS

NITER INT OPTIONAL MAXIMUM NUMBER OF ITERATIONS DEFAULT IS 300 SHOULD BE GREATER THAN OR  
EQUAL TO 1

TOLFLOAT OPTIONAL STOP THE ALGORITHM IF W HAS CONVERGED DEFAULT IS 1E3

ALPHA1 FLOAT OPTIONAL HYPERPARAMETER SHAPE PARAMETER FOR THE GAMMA DISTRIBUTION PRIOR  
OVER THE ALPHA PARAMETER DEFAULT IS 1E6

ALPHA2 FLOAT OPTIONAL HYPERPARAMETER INVERSE SCALE PARAMETER RATE PARAMETER FOR THE  
GAMMA DISTRIBUTION PRIOR OVER THE ALPHA PARAMETER DEFAULT IS 1E6

LAMBDA1 FLOAT OPTIONAL HYPERPARAMETER SHAPE PARAMETER FOR THE GAMMA DISTRIBUTION PRIOR  
OVER THE LAMBDA PARAMETER DEFAULT IS 1E6

LAMBDA2 FLOAT OPTIONAL HYPERPARAMETER INVERSE SCALE PARAMETER RATE PARAMETER FOR THE  
GAMMA DISTRIBUTION PRIOR OVER THE LAMBDA PARAMETER DEFAULT IS 1E6

COMPUTESCORE BOOLEAN OPTIONAL IF TRUE COMPUTE THE LOG MARGINAL LIKELIHOOD AT EACH ITERA  
TION OF THE OPTIMIZATION DEFAULT IS FALSE

FITINTERCEPT BOOLEAN OPTIONAL DEFAULT TRUE WHETHER TO CALCULATE THE INTERCEPT FOR THIS MODEL  
THE INTERCEPT IS NOT TREATED AS A PROBABILISTIC PARAMETER AND THUS HAS NO ASSOCIATED VARIANCE  
IF SET TO FALSE NO INTERCEPT WILL BE USED IN CALCULATIONS EG DATA IS EXPECTED TO BE ALREADY  
CENTERED

NORMALIZE BOOLEAN OPTIONAL DEFAULT FALSE THIS PARAMETER IS IGNORED WHEN  
FITINTERCEPT IS SET TO FALSE IF TRUE THE REGRESSORS X WILL BE NORMALIZED BE  
FORE REGRESSION BY SUBTRACTING THE MEAN AND DIVIDING BY THE L2NORM IF YOU WISH TO  
STANDARDIZE PLEASE USE SKLEARNPREPROCESSINGSTANDARDSCALER BEFORE  
CALLINGFIT ON AN ESTIMATOR WITH NORMALIZEFALSE

COPYX BOOLEAN OPTIONAL DEFAULT TRUE IF TRUE X WILL BE COPIED ELSE IT MAY BE OVERWRITTEN

VERBOSE BOOLEAN OPTIONAL DEFAULT FALSE VERBOSE MODE WHEN FITTING THE MODEL

ATTRIBUTES

COEF ARRAY SHAPE NFEATURES COEFFICIENTS OF THE REGRESSION MODEL MEAN OF DISTRIBUTION

622SKLEARNLINEARMODEL GENERALIZED LINEAR MODELS 1869

SCIKITLEARN USER GUIDE RELEASE 0213

INTERCEPT FLOAT INDEPENDENT TERM IN DECISION FUNCTION SET TO 00 IF FITINTERCEPT FALSE

ALPHA FLOAT ESTIMATED PRECISION OF THE NOISE

LAMBDA FLOAT ESTIMATED PRECISION OF THE WEIGHTS

SIGMA ARRAY SHAPE NFEATURES NFEATURES ESTIMATED VARIANCECOVARIANCE MATRIX OF THE WEIGHTS

SCORES ARRAY SHAPE NITER 1 IF COMPUTEDSCORE IS TRUE VALUE OF THE LOG MARGINAL LIKELIHOOD TO BE MAXIMIZED AT EACH ITERATION OF THE OPTIMIZATION THE ARRAY STARTS WITH THE VALUE OF THE LOG MARGINAL LIKELIHOOD OBTAINED FOR THE INITIAL VALUES OF ALPHA AND LAMBDA AND ENDS WITH THE VALUE OBTAINED FOR THE ESTIMATED ALPHA AND LAMBDA

NITER INT THE ACTUAL NUMBER OF ITERATIONS TO REACH THE STOPPING CRITERION

NOTES

THERE EXIST SEVERAL STRATEGIES TO PERFORM BAYESIAN RIDGE REGRESSION THIS IMPLEMENTATION IS BASED ON THE ALGORITHM DESCRIBED IN APPENDIX A OF TIPPING 2001 WHERE UPDATES OF THE REGULARIZATION PARAMETERS ARE DONE AS SUGGESTED IN MACKAY 1992 NOTE THAT ACCORDING TO A NEW VIEW OF AUTOMATIC RELEVANCE DETERMINATION WIPF AND NAGARAJAN 2008 THESE UPDATE RULES DO NOT GUARANTEE THAT THE MARGINAL LIKELIHOOD IS INCREASING BETWEEN TWO CONSECUTIVE ITERATIONS OF THE OPTIMIZATION

REFERENCES

D J C MACKAY BAYESIAN INTERPOLATION COMPUTATION AND NEURAL SYSTEMS V OL 4 NO 3 1992

M E TIPPING SPARSE BAYESIAN LEARNING AND THE RELEVANCE VECTOR MACHINE JOURNAL OF MACHINE LEARNING RESEARCH V OL 1 2001

EXAMPLES

```
FROM SKLEARN IMPORT LINEARMODEL
CLF LINEARMODELBAYESIANRIDGE
CLFFIT00 1 1 2 2 0 1 2
```

BAYESIANRIDGEALPHA11E06 ALPHA21E06 COMPUTESCOREFALSE

COPYXTRUE FITINTERCEPTTRUE LAMBDA11E06 LAMBDA21E06

NITER300 NORMALIZEFALSE TOL0001 VERBOSEFALSE

CLFPREDICT1 1

ARRAY1

METHODS

FITSELF X Y SAMPLEWEIGHT FIT THE MODEL

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PREDICT SELF X RETURNSTD PREDICT USING THE LINEAR MODEL

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

CONTINUED ON NEXT PAGE

1870 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

TABLE 6150 – CONTINUED FROM PREVIOUS PAGE

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELFNITER300 TOL0001 ALPHA11E06 ALPHA21E06 LAMBDA11E

06LAMBDA21E06 COMPUTESCOREFALSE FITINTERCEPTTRUE NORMALIZEFALSE

COPYXTRUE VERBOSEFALSE

FITSELFXYSAMPLEWEIGHTNONE

FIT THE MODEL

PARAMETERS

XNUMPY ARRAY OF SHAPE NSAMPLESNFEATURES TRAINING DATA

YNUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES WILL BE CAST TO X’S DTYPE IF NECESSARY

SAMPLEWEIGHT NUMPY ARRAY OF SHAPE NSAMPLES INDIVIDUAL WEIGHTS FOR EACH SAMPLE

NEW IN VERSION 020 PARAMETER SAMPLEWEIGHT SUPPORT TO BAYESIANRIDGE

RETURNS

SELF RETURNS AN INSTANCE OF SELF

GETPARAMS SELFDEEPTTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELFXTURNSTDFALSE

PREDICT USING THE LINEAR MODEL

IN ADDITION TO THE MEAN OF THE PREDICTIVE DISTRIBUTION ALSO ITS STANDARD DEVIATION CAN BE RETURNED

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES SAMPLES

RETURNSTD BOOLEAN OPTIONAL WHETHER TO RETURN THE STANDARD DEVIATION OF POSTERIOR PREDIC

TION

RETURNS

YMEAN ARRAY SHAPE NSAMPLES MEAN OF PREDICTIVE DISTRIBUTION OF QUERY POINTS

YSTD ARRAY SHAPE NSAMPLES STANDARD DEVIATION OF PREDICTIVE DISTRIBUTION OF QUERY

POINTS

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{U}{V}$  WHERE U IS THE RESIDUAL SUM OF SQUARES YTRUE YPRED

2SUM AND V IS THE TOTAL SUM OF SQUARES YTRUE YTRUEMEAN 2SUM THE BEST POSSIBLE SCORE

IS 10 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS

PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 00

PARAMETERS

622SKLEARNLINEARMODEL GENERALIZED LINEAR MODELS 1871

SCIKITLEARN USER GUIDE RELEASE 0213

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES FOR SOME ESTIMATORS THIS MAY BE A PRECOMPUTED KERNEL MATRIX INSTEAD SHAPE NSAMPLES NSAMPLESFITTED WHERE NSAMPLESFITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE VALUES FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT R2 OF SELF PREDICT X WRT Y

NOTES

THE R2 SCORE USED WHEN CALLING SCORE ON A REGRESSOR WILL USE MULTIOUTPUTUNIFORMAVERAGE FROM VERSION 023 TO KEEP CONSISTENT WITH METRICSR2SCORE THIS WILL INFLUENCE THE SCORE METHOD OF ALL THE MULTIOUTPUT REGRESSORS EXCEPT FOR MULTIOUTPUTMULTIOUTPUTREGRESSOR TO SPECIFY THE DEFAULT VALUE MANUALLY AND AVOID THE WARNING PLEASE EITHER CALL METRICSR2SCORE DIRECTLY OR MAKE A CUSTOM SCORER WITH METRICSMAKESCORER THE BUILTIN SCORER R2 USES MULTIOUTPUTUNIFORMAVERAGE

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

EXAMPLES USING SKLEARNLINEARMODELBAYESIANRIDGE

- FEATURE AGGLOMERATION VS UNIVARIATE SELECTION
- IMPUTING MISSING VALUES WITH VARIANTS OF ITERATIVEIMPUTER
- BAYESIAN RIDGE REGRESSION

6223SKLEARNLINEARMODEL ELASTICNET

CLASSSKLEARNLINEARMODEL ELASTICNET ALPHA10 L1RATIO05 FITINTERCEPTTRUE NOR

MALIZEFALSE PRECOMPUTEFALSE MAXITER1000

COPYXTRUE TOL00001 WARMSTARTFALSE POSITIVEFALSE RANDOMSTATENONE SELECTION'CYCLIC'

LINEAR REGRESSION WITH COMBINED L1 AND L2 PRIORS AS REGULARIZER

MINIMIZES THE OBJECTIVE FUNCTION

1 2NSAMPLES Y XW22

ALPHA L1RATIO W1

05ALPHA1 L1RATIO W22

IF YOU ARE INTERESTED IN CONTROLLING THE L1 AND L2 PENALTY SEPARATELY KEEP IN MIND THAT THIS IS EQUIVALENT TO

AL1 BL2

1872 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

WHERE

ALPHA A B AND L1RATIO A A B

THE PARAMETER L1RATIO CORRESPONDS TO ALPHA IN THE GLMNET R PACKAGE WHILE ALPHA CORRESPONDS TO THE LAMBDA PARAMETER IN GLMNET SPECIFICALLY L1RATIO 1 IS THE LASSO PENALTY CURRENTLY L1RATIO 001 IS NOT RELIABLE UNLESS YOU SUPPLY YOUR OWN SEQUENCE OF ALPHA

READ MORE IN THE USER GUIDE

PARAMETERS

ALPHA FLOAT OPTIONAL CONSTANT THAT MULTIPLIES THE PENALTY TERMS DEFAULTS TO 10 SEE THE NOTES FOR THE EXACT MATHEMATICAL MEANING OF THIS PARAMETER“ALPHA 0” IS EQUIVALENT TO AN ORDINARY LEAST SQUARE SOLVED BY THE LINEARREGRESSION OBJECT FOR NUMERICAL REASONS USINGALPHA 0 WITH THE LASSO OBJECT IS NOT ADVISED GIVEN THIS YOU SHOULD USE THE LINEARREGRESSION OBJECT

L1RATIO FLOAT THE ELASTICNET MIXING PARAMETER WITH 0 L1RATIO 1 FOR L1RATIO 0 THE PENALTY IS AN L2 PENALTY FOR L1RATIO 1 IT IS AN L1 PENALTY FOR 0 L1RATIO 1 THE PENALTY IS A COMBINATION OF L1 AND L2

FITINTERCEPT BOOL WHETHER THE INTERCEPT SHOULD BE ESTIMATED OR NOT IF FALSE THE DATA IS ASSUMED TO BE ALREADY CENTERED

NORMALIZE BOOLEAN OPTIONAL DEFAULT FALSE THIS PARAMETER IS IGNORED WHEN FITINTERCEPT IS SET TO FALSE IF TRUE THE REGRESSORS X WILL BE NORMALIZED BEFORE REGRESSION BY SUBTRACTING THE MEAN AND DIVIDING BY THE L2NORM IF YOU WISH TO STANDARDIZE PLEASE USE SKLEARNPREPROCESSINGSTANDARDSCALER BEFORE CALLING FIT ON AN ESTIMATOR WITH NORMALIZEFALSE

PRECOMPUTE TRUE FALSE ARRAYLIKE WHETHER TO USE A PRECOMPUTED GRAM MATRIX TO SPEED UP CALCULATIONS THE GRAM MATRIX CAN ALSO BE PASSED AS ARGUMENT FOR SPARSE INPUT THIS OPTION IS ALWAYS TRUE TO PRESERVE SPARSITY

MAXITER INT OPTIONAL THE MAXIMUM NUMBER OF ITERATIONS

COPYX BOOLEAN OPTIONAL DEFAULT TRUE IF TRUE X WILL BE COPIED ELSE IT MAY BE OVERWRITTEN

TOL FLOAT OPTIONAL THE TOLERANCE FOR THE OPTIMIZATION IF THE UPDATES ARE SMALLER THAN TOL THE OPTIMIZATION CODE CHECKS THE DUAL GAP FOR OPTIMALITY AND CONTINUES UNTIL IT IS SMALLER THAN TOL

WARMSTART BOOL OPTIONAL WHEN SET TO TRUE REUSE THE SOLUTION OF THE PREVIOUS CALL TO FIT AS INITIALIZATION OTHERWISE JUST ERASE THE PREVIOUS SOLUTION SEE THE GLOSSARY

POSITIVE BOOL OPTIONAL WHEN SET TO TRUE FORCES THE COEFFICIENTS TO BE POSITIVE

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULT NONE THE SEED OF THE PSEUDO RANDOM NUMBER GENERATOR THAT SELECTS A RANDOM FEATURE TO UPDATE IF INT RANDOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NRANDOM USED WHEN SELECTION ‘RANDOM’ SELECTION STR DEFAULT ‘CYCLIC’ IF SET TO ‘RANDOM’ A RANDOM COEFFICIENT IS UPDATED EVERY ITERATION RATHER THAN LOOPING OVER FEATURES SEQUENTIALLY BY DEFAULT THIS SETTING TO ‘RANDOM’ OFTEN LEADS TO SIGNIFICANTLY FASTER CONVERGENCE ESPECIALLY WHEN TOL IS HIGHER THAN 1E4

ATTRIBUTES

COEF ARRAY SHAPE NFEATURES NTARGETS NFEATURES PARAMETER VECTOR W IN THE COST FUNCTION FORMULA

622SKLEARNLINEARMODEL GENERALIZED LINEAR MODELS 1873

SCIKITLEARN USER GUIDE RELEASE 0213  
SPARSECOEF SCIPYSPARSE MATRIX SHAPE NFEATURES 1 NTARGETS NFEATURES SPARSE  
REPRESENTATION OF THE FITTED COEF  
INTERCEPT FLOAT ARRAY SHAPE NTARGETS INDEPENDENT TERM IN DECISION FUNCTION  
NITER ARRAYLIKE SHAPE NTARGETS NUMBER OF ITERATIONS RUN BY THE COORDINATE DESCENT SOLVER  
TO REACH THE SPECIFIED TOLERANCE  
SEE ALSO  
ELASTICNETCV ELASTIC NET MODEL WITH BEST MODEL SELECTION BY CROSSVALIDATION  
SGDREGRESSOR IMPLEMENTS ELASTIC NET REGRESSION WITH INCREMENTAL TRAINING  
SGDCLASSIFIER IMPLEMENTS LOGISTIC REGRESSION WITH ELASTIC NET PENALTY  
SGDCLASSIFIERLOSSLOG PENALTYELASTICNET  
NOTES  
TO AVOID UNNECESSARY MEMORY DUPLICATION THE X ARGUMENT OF THE FIT METHOD SHOULD BE DIRECTLY PASSED AS A  
FORTRANCONTIGUOUS NUMPY ARRAY  
EXAMPLES  
FROM SKLEARNLINEARMODEL IMPORT ELASTICNET  
FROM SKLEARNDATASETS IMPORT MAKEREGRESSION  
X Y MAKEREGRESSIONNFEATURES2 RANDOMSTATE0  
REGR ELASTICNETRANDOMSTATE0  
REGRFITX Y  
ELASTICNETALPHA10 COPYXTRUE FITINTERCEPTTRUE L1RATIO05  
MAXITER1000 NORMALIZEFALSE POSITIVEFALSE PRECOMPUTEFALSE  
RANDOMSTATE0 SELECTIONCYCLIC TOL00001 WARMSTARTFALSE  
PRINTREGRCOEF  
1883816048 6455968825  
PRINTREGRINTERCEPT  
1451  
PRINTREGRPREDICT0 0  
1451  
METHODS  
FITSELF X Y CHECKINPUT FIT MODEL WITH COORDINATE DESCENT  
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR  
PATH X Y L1RATIO EPS NALPHAS COMPUTE ELASTIC NET PATH WITH COORDINATE DESCENT  
PREDICT SELF X PREDICT USING THE LINEAR MODEL  
SCORE SELF X Y SAMPLEWEIGHT RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PRE  
DICTION  
SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR  
1874 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213  
INIT SELFALPHA10 L1RATIO05 FITINTERCEPTTRUE NORMALIZEFALSE PRECOMPUTEFALSE  
MAXITER1000 COPYXTRUE TOL00001 WARMSTARTFALSE POSITIVEFALSE RAN  
DOMSTATENONE SELECTION'CYCLIC'  
FITSELFXYCHECKINPUTTRUE  
FIT MODEL WITH COORDINATE DESCENT  
PARAMETERS  
XNDARRAY OR SCIPYSPARSE MATRIX NSAMPLES NFEATURES DATA  
YNDARRAY SHAPE NSAMPLES OR NSAMPLES NTARGETS TARGET WILL BE CAST TO X'S DTYPE  
IF NECESSARY  
CHECKINPUT BOOLEAN DEFAULTTRUE ALLOW TO BYPASS SEVERAL INPUT CHECKING DON'T USE  
THIS PARAMETER UNLESS YOU KNOW WHAT YOU DO  
NOTES  
COORDINATE DESCENT IS AN ALGORITHM THAT CONSIDERS EACH COLUMN OF DATA AT A TIME HENCE IT WILL AUTOMATICALLY  
CONVERT THE X INPUT AS A FORTRANCONTIGUOUS NUMPY ARRAY IF NECESSARY  
TO AVOID MEMORY REALLOCATION IT IS ADVISED TO ALLOCATE THE INITIAL DATA IN MEMORY DIRECTLY USING THAT FORMAT  
GETPARAMS SELFDEEPTTRUE  
GET PARAMETERS FOR THIS ESTIMATOR  
PARAMETERS  
DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBOBJECTS THAT ARE ESTIMATORS  
RETURNS  
PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES  
STATICPATHXYL1RATIO05 EPS0001 NALPHAS100 ALPHASNONE PRECOMPUTE'AUTO'  
XYNONE COPYXTRUE COEFINITNONE VERBOSEFALSE RETURNNNITERFALSE POSI  
TIVEFALSE CHECKINPUTTRUE PARAMS  
COMPUTE ELASTIC NET PATH WITH COORDINATE DESCENT  
THE ELASTIC NET OPTIMIZATION FUNCTION VARIES FOR MONO AND MULTIOUTPUTS  
FOR MONOOUTPUT TASKS IT IS  
1 2NSAMPLES Y XW22  
ALPHA L1RATIO W1  
05ALPHA1 L1RATIO W22  
FOR MULTIOUTPUT TASKS IT IS  
1 2 NSAMPLES Y XWFRO2  
ALPHA L1RATIO W21  
05ALPHA1 L1RATIO WFRO2  
WHERE  
W21 SUMI SQRTSUMJ WIJ2  
IE THE SUM OF NORM OF EACH ROW  
READ MORE IN THE USER GUIDE  
622SKLEARNLINEARMODEL GENERALIZED LINEAR MODELS 1875

SCIKITLEARN USER GUIDE RELEASE 0213

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA PASS DIRECTLY AS FORTRAN  
CONTIGUOUS DATA TO AVOID UNNECESSARY MEMORY DUPLICATION IF YIS MONOOUTPUT THEN X  
CAN BE SPARSE  
YNDARRAY SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TARGET VALUES  
L1RATIO FLOAT OPTIONAL FLOAT BETWEEN 0 AND 1 PASSED TO ELASTIC NET SCALING BETWEEN L1 AND  
L2 PENALTIES L1RATIO1 CORRESPONDS TO THE LASSO  
EPS FLOAT LENGTH OF THE PATH EPS1E3 MEANS THAT ALPHAMIN ALPHAMAX  
1E3  
NALPHAS INT OPTIONAL NUMBER OF ALPHAS ALONG THE REGULARIZATION PATH  
ALPHAS NDARRAY OPTIONAL LIST OF ALPHAS WHERE TO COMPUTE THE MODELS IF NONE ALPHAS ARE  
SET AUTOMATICALLY  
PRECOMPUTE TRUE FALSE 'AUTO' ARRAYLIKE WHETHER TO USE A PRECOMPUTED GRAM MATRIX  
TO SPEED UP CALCULATIONS IF SET TO AUTO LET US DECIDE THE GRAM MATRIX CAN ALSO BE  
PASSED AS ARGUMENT  
XYARRAYLIKE OPTIONAL XY NPDOTXT Y THAT CAN BE PRECOMPUTED IT IS USEFUL ONLY  
WHEN THE GRAM MATRIX IS PRECOMPUTED  
COPYX BOOLEAN OPTIONAL DEFAULT TRUE IF TRUE X WILL BE COPIED ELSE IT MAY BE OVER  
WRITTEN  
COEFINIT ARRAY SHAPE NFEATURES NONE THE INITIAL VALUES OF THE COEFFICIENTS  
VERBOSE BOOL OR INTEGER AMOUNT OF VERBOSITY  
RETURNNNITER BOOL WHETHER TO RETURN THE NUMBER OF ITERATIONS OR NOT  
POSITIVE BOOL DEFAULT FALSE IF SET TO TRUE FORCES COEFFICIENTS TO BE POSITIVE ONLY ALLOWED  
WHENYNDIM 1  
CHECKINPUT BOOL DEFAULT TRUE SKIP INPUT VALIDATION CHECKS INCLUDING THE GRAM MATRIX  
WHEN PROVIDED ASSUMING THERE ARE HANDLED BY THE CALLER WHEN CHECKINPUTFALSE  
PARAMS KWARGS KEYWORD ARGUMENTS PASSED TO THE COORDINATE DESCENT SOLVER  
RETURNS  
ALPHAS ARRAY SHAPE NALPHAS THE ALPHAS ALONG THE PATH WHERE MODELS ARE COMPUTED  
COEFS ARRAY SHAPE NFEATURES NALPHAS OR NOUTPUTS NFEATURES NALPHAS COEFFICIENTS  
ALONG THE PATH  
DUALGAPS ARRAY SHAPE NALPHAS THE DUAL GAPS AT THE END OF THE OPTIMIZATION FOR EACH  
ALPHA  
NITERS ARRAYLIKE SHAPE NALPHAS THE NUMBER OF ITERATIONS TAKEN BY THE COORDINATE  
DESCENT OPTIMIZER TO REACH THE SPECIFIED TOLERANCE FOR EACH ALPHA IS RETURNED WHEN  
RETURNNNITER IS SET TO TRUE  
SEE ALSO  
MULTITASKELASTICNET  
MULTITASKELASTICNETCV  
ELASTICNET  
ELASTICNETCV  
1876 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

NOTES

FOR AN EXAMPLE SEE EXAMPLESLINEARMODELPLOTLASSOCOORDINATEDESCENTPATHPY

PREDICTSELF

PREDICT USING THE LINEAR MODEL

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES SAMPLES

RETURNS

CARRAY SHAPE NSAMPLES RETURNS PREDICTED VALUES

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{U}{V}$  WHERE U IS THE RESIDUAL SUM OF SQUARES  $Y_{TRUE} - Y_{PRED}$

$\frac{2}{SUM}$  AND V IS THE TOTAL SUM OF SQUARES  $Y_{TRUE} - Y_{TRUEMEAN}$   $\frac{2}{SUM}$  THE BEST POSSIBLE SCORE

IS 10 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS

PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 00

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES FOR SOME ESTIMATORS THIS MAY

BE A PRECOMPUTED KERNEL MATRIX INSTEAD SHAPE NSAMPLES NSAMPLESFITTED WHERE

NSAMPLESFITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE VALUES FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT R2 OF SELF-PREDICTX WRT Y

NOTES

THE R2 SCORE USED WHEN CALLING SCORE ON A REGRESSOR WILL USE MULTIOUTPUTUNIFORMAVERAGE

FROM VERSION 023 TO KEEP CONSISTENT WITH METRICSR2SCORE THIS WILL INFLUENCE THE SCORE

METHOD OF ALL THE MULTIOUTPUT REGRESSORS EXCEPT FOR MULTIOUTPUTMULTIOUTPUTREGRESSOR

TO SPECIFY THE DEFAULT VALUE MANUALLY AND AVOID THE WARNING PLEASE EITHER CALL METRICSR2SCORE

DIRECTLY OR MAKE A CUSTOM SCORER WITH METRICSMAKESCORER THE BUILTIN SCORER R2 USES

MULTIOUTPUTUNIFORMAVERAGE

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE

PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT

OF A NESTED OBJECT

RETURNS

SELF

SPARSECOEF

SPARSE REPRESENTATION OF THE FITTED COEF

622SKLEARNLINEARMODEL GENERALIZED LINEAR MODELS 1877

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES USING SKLEARNLINEARMODELELASTICNET

- LASSO AND ELASTIC NET FOR SPARSE SIGNALS
- TRAIN ERROR VS TEST ERROR

6224SKLEARNLINEARMODEL HUBERREGRESSOR

CLASSSSKLEARNLINEARMODEL HUBERREGRESSOR EPSILON135 MAXITER100 ALPHA00001

WARMSTARTFALSE FITINTERCEPTTRUE TOL1E05

LINEAR REGRESSION MODEL THAT IS ROBUST TO OUTLIERS

THE HUBER REGRESSOR OPTIMIZES THE SQUARED LOSS FOR THE SAMPLES WHERE  $|Y - XW| < \text{SIGMA}$  EPSILON AND THE ABSOLUTE LOSS FOR THE SAMPLES WHERE  $|Y - XW| > \text{SIGMA}$  EPSILON WHERE W AND SIGMA ARE PARAMETERS TO BE OPTIMIZED THE PARAMETER SIGMA MAKES SURE THAT IF Y IS SCALED UP OR DOWN BY A CERTAIN FACTOR ONE DOES NOT NEED TO RESCALE EPSILON TO ACHIEVE THE SAME ROBUSTNESS NOTE THAT THIS DOES NOT TAKE INTO ACCOUNT THE FACT THAT THE DIFFERENT FEATURES OF X MAY BE OF DIFFERENT SCALES THIS MAKES SURE THAT THE LOSS FUNCTION IS NOT HEAVILY INFLUENCED BY THE OUTLIERS WHILE NOT COMPLETELY IGNORING THEIR EFFECT

READ MORE IN THE USER GUIDE

NEW IN VERSION 018

PARAMETERS

EPSILON FLOAT GREATER THAN 10 DEFAULT 135 THE PARAMETER EPSILON CONTROLS THE NUMBER OF SAMPLES THAT SHOULD BE CLASSIFIED AS OUTLIERS THE SMALLER THE EPSILON THE MORE ROBUST IT IS TO OUTLIERS

MAXITER INT DEFAULT 100 MAXIMUM NUMBER OF ITERATIONS THAT SCIPYOPTIMIZEFMINLBFGSB SHOULD RUN FOR

ALPHA FLOAT DEFAULT 00001 REGULARIZATION PARAMETER

WARMSTART BOOL DEFAULT FALSE THIS IS USEFUL IF THE STORED ATTRIBUTES OF A PREVIOUSLY USED MODEL HAS TO BE REUSED IF SET TO FALSE THEN THE COEFFICIENTS WILL BE REWRITTEN FOR EVERY CALL TO FIT SEE THE GLOSSARY

FITINTERCEPT BOOL DEFAULT TRUE WHETHER OR NOT TO FIT THE INTERCEPT THIS CAN BE SET TO FALSE IF THE DATA IS ALREADY CENTERED AROUND THE ORIGIN

TOLFLOAT DEFAULT 1E5 THE ITERATION WILL STOP WHEN  $\max |g_i| \leq \text{TOL}$

NTOL WHERE  $g_i$  IS THE ITH COMPONENT OF THE PROJECTED GRADIENT

ATTRIBUTES

COEF ARRAY SHAPE NFEATURES FEATURES GOT BY OPTIMIZING THE HUBER LOSS

INTERCEPT FLOAT BIAS

SCALE FLOAT THE VALUE BY WHICH  $|Y - XW|$  IS SCALED DOWN

NITER INT NUMBER OF ITERATIONS THAT FMINLBFGSB HAS RUN FOR

CHANGED IN VERSION 020 IN SCIPY 100 THE NUMBER OF LBFGS ITERATIONS MAY EXCEED MAXITER NITER WILL NOW REPORT AT MOST MAXITER

OUTLIERS ARRAY SHAPE NSAMPLES A BOOLEAN MASK WHICH IS SET TO TRUE WHERE THE SAMPLES ARE IDENTIFIED AS OUTLIERS

1878 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

REFERENCES

RE4616EF910FB1 RE4616EF910FB2

EXAMPLES

```
import numpy as np
from sklearn.linear_model import HuberRegressor, LinearRegression
from sklearn.datasets import make_regression
rng = np.random.RandomState(0)
X, y = make_regression(n_samples=200, n_features=2, noise=40, coef=True, random_state=0)
X4 = rng.uniform(10, 20, 4, 2)
y4 = rng.uniform(10, 20, 4)
huber = HuberRegressor(fit_x_y)
huber.score(X, y)
7284608623514573
huber.predict(X4)
array([8067200])
linear = LinearRegression(fit_x_y)
print('True coefficients:', coef)
true_coef = [204923, 341698]
print('Huber coefficients:', huber.coef_)
huber_coef = [177906, 310106]
print('Linear regression coefficients:', linear.coef_)
linear_coef = [19221, 70226]
```

METHODS

`fit(X, y, sample_weight=None)` Fit the model according to the given training data

`get_params()` Get parameters for this estimator

`predict(X)` Predict using the linear model

`score(X, y, sample_weight=None)` Returns the coefficient of determination  $R^2$  of the prediction

`set_params(**params)` Set the parameters of this estimator

`__init__(self, eps=1e-05, max_iter=100, alpha=0.0001, warm_start=False, fit_intercept=True, tol=1e-05, fit_x_y=False, sample_weight=None)`

`fit(X, y, sample_weight=None)` Fit the model according to the given training data

PARAMETERS

`X` : array-like, shape `(n_samples, n_features)` Training vector, where `n_samples` is the number of samples and `n_features` is the number of features

`y` : array-like, shape `(n_samples)` Target vector relative to `X`

`sample_weight` : array-like, shape `(n_samples)` Weight given to each sample

RETURNS

self

622SKLEARNLINEARMODEL GENERALIZED LINEAR MODELS 1879

SCIKITLEARN USER GUIDE RELEASE 0213

GETPARAMS SELFDEEPTREE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELF

PREDICT USING THE LINEAR MODEL

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES SAMPLES

RETURNS

CARRAY SHAPE NSAMPLES RETURNS PREDICTED VALUES

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{U}{V}$  WHERE U IS THE RESIDUAL SUM OF SQUARES YTRUE YPRED 2SUM AND V IS THE TOTAL SUM OF SQUARES YTRUE YTRUEMEAN 2SUM THE BEST POSSIBLE SCORE IS 10 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 00

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES FOR SOME ESTIMATORS THIS MAY BE A PRECOMPUTED KERNEL MATRIX INSTEAD SHAPE NSAMPLES NSAMPLESFITTED WHERE NSAMPLESFITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE VALUES FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT R2 OF SELF PREDICTX WRT Y

NOTES

THE R2 SCORE USED WHEN CALLING SCORE ON A REGRESSOR WILL USE MULTIOUTPUTUNIFORMAVERAGE FROM VERSION 023 TO KEEP CONSISTENT WITH METRICSR2SCORE THIS WILL INFLUENCE THE SCORE METHOD OF ALL THE MULTIOUTPUT REGRESSORS EXCEPT FOR MULTIOUTPUTMULTIOUTPUTREGRESSOR TO SPECIFY THE DEFAULT VALUE MANUALLY AND AVOID THE WARNING PLEASE EITHER CALL METRICSR2SCORE DIRECTLY OR MAKE A CUSTOM SCORER WITH METRICSMAKESCORER THE BUILTIN SCORER R2 USES MULTIOUTPUTUNIFORMAVERAGE

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

1880 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

SELF

EXAMPLES USING SKLEARNLINEARMODELHUBERREGRESSOR

•HUBERREGRESSOR VS RIDGE ON DATASET WITH STRONG OUTLIERS

•ROBUST LINEAR ESTIMATOR FITTING

6225SKLEARNLINEARMODEL LARS

CLASSSSKLEARNLINEARMODEL LARSFITINTERCEPTTRUE VERBOSEFALSE NORMALIZETRUE PRECOM

PUTE‘AUTO’ NNONZEROCOEF500 EPS2220446049250313E

16COPYXTRUE FITPATHTRUE POSITIVEFALSE

LEAST ANGLE REGRESSION MODEL AKA LAR

READ MORE IN THE USER GUIDE

PARAMETERS

FITINTERCEPT BOOLEAN WHETHER TO CALCULATE THE INTERCEPT FOR THIS MODEL IF SET TO FALSE NO

INTERCEPT WILL BE USED IN CALCULATIONS EG DATA IS EXPECTED TO BE ALREADY CENTERED

VERBOSE BOOLEAN OR INTEGER OPTIONAL SETS THE VERBOSITY AMOUNT

NORMALIZE BOOLEAN OPTIONAL DEFAULT TRUE THIS PARAMETER IS IGNORED WHEN FITINTERCEPT

IS SET TO FALSE IF TRUE THE REGRESSORS X WILL BE NORMALIZED BEFORE REGRESSION BY SUB

TRACTING THE MEAN AND DIVIDING BY THE L2NORM IF YOU WISH TO STANDARDIZE PLEASE USE

SKLEARNPREPROCESSINGSTANDARDSCALER BEFORE CALLING FIT ON AN ESTIMATOR

WITHNORMALIZEFALSE

PRECOMPUTE TRUE FALSE ‘AUTO’ ARRAYLIKE WHETHER TO USE A PRECOMPUTED GRAM MATRIX TO

SPEED UP CALCULATIONS IF SET TO AUTO LET US DECIDE THE GRAM MATRIX CAN ALSO BE PASSED

AS ARGUMENT

NNONZEROCOEF5 INT OPTIONAL TARGET NUMBER OF NONZERO COEFFICIENTS USE NPINF FOR NO

LIMIT

EPS FLOAT OPTIONAL THE MACHINEPRECISION REGULARIZATION IN THE COMPUTATION OF THE CHOLESKY

DIAGONAL FACTORS INCREASE THIS FOR VERY ILLCONDITIONED SYSTEMS UNLIKE THE TOL PARAMETER IN

SOME ITERATIVE OPTIMIZATIONBASED ALGORITHMS THIS PARAMETER DOES NOT CONTROL THE TOLERANCE

OF THE OPTIMIZATION

COPYX BOOLEAN OPTIONAL DEFAULT TRUE IF TRUE X WILL BE COPIED ELSE IT MAY BE OVERWRITTEN

FITPATH BOOLEAN IF TRUE THE FULL PATH IS STORED IN THE COEFPATH ATTRIBUTE IF YOU COMPUTE

THE SOLUTION FOR A LARGE PROBLEM OR MANY TARGETS SETTING FITPATH TOFALSE WILL LEAD TO

A SPEEDUP ESPECIALLY WITH A SMALL ALPHA

POSITIVE BOOLEAN DEFAULTFALSE RESTRICT COEFFICIENTS TO BE 0 BE AWARE THAT YOU MIGHT

WANT TO REMOVE FITINTERCEPT WHICH IS SET TRUE BY DEFAULT

DEPRECATED SINCE VERSION 020 THE OPTION IS BROKEN AND DEPRECATED IT WILL BE REMOVED IN

V022

ATTRIBUTES

ALPHAS ARRAY SHAPE NALPHAS 1 LIST OF NTARGETS SUCH ARRAYS MAXIMUM OF COVARI

ANCES IN ABSOLUTE VALUE AT EACH ITERATION NALPHAS IS EITHERNNONZEROCOEF5 OR

NFEATURES WHICHEVER IS SMALLER

6225SKLEARNLINEARMODEL GENERALIZED LINEAR MODELS 1881

SCIKITLEARN USER GUIDE RELEASE 0213

ACTIVE LIST LENGTH NALPHAS LIST OF NTARGETS SUCH LISTS INDICES OF ACTIVE VARIABLES AT THE END OF THE PATH

COEFPATH ARRAY SHAPE NFEATURES NALPHAS 1 LIST OF NTARGETS SUCH ARRAYS THE VARYING VALUES OF THE COEFFICIENTS ALONG THE PATH IT IS NOT PRESENT IF THE FITPATH PARAMETER IS FALSE

COEF ARRAY SHAPE NFEATURES OR NTARGETS NFEATURES PARAMETER VECTOR W IN THE FORMULATION FORMULA

INTERCEPT FLOAT ARRAY SHAPE NTARGETS INDEPENDENT TERM IN DECISION FUNCTION

NITER ARRAYLIKE OR INT THE NUMBER OF ITERATIONS TAKEN BY LARSPATH TO FIND THE GRID OF ALPHAS FOR EACH TARGET

SEE ALSO

LARSPATH LARSCV

SKLEARNDECOMPOSITIONSPARSEENCODE

EXAMPLES

FROM SKLEARN IMPORT LINEARMODEL

REG LINEARMODELLARSNNONZEROCOEFS1

REGFIT1 1 0 0 1 1 11111 0 11111

LARSCOPYXTRUE EPS FITINTERCEPTTRUE FITPATHTRUE

NNONZEROCOEFS1 NORMALIZETRUE POSITIVEFALSE PRECOMPUTEAUTO

VERBOSEFALSE

PRINTREGCOEF

0 111

METHODS

FITSELF X Y XY FIT THE MODEL USING X Y AS TRAINING DATA

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PREDICT SELF X PREDICT USING THE LINEAR MODEL

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELF FITINTERCEPTTRUE VERBOSEFALSE NORMALIZETRUE PRECOMPUTE'AUTO'

NNONZEROCOEFS500 EPS2220446049250313E16 COPYXTRUE FITPATHTRUE

POSITIVEFALSE

FITSELFXYXYNONE

FIT THE MODEL USING X Y AS TRAINING DATA

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NTARGETS TARGET VALUES

XYARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NTARGETS OPTIONAL XY NPDOTXT Y

1882 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

THAT CAN BE PRECOMPUTED IT IS USEFUL ONLY WHEN THE GRAM MATRIX IS PRECOMPUTED

RETURNS

SELF OBJECT RETURNS AN INSTANCE OF SELF

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELF

PREDICT USING THE LINEAR MODEL

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES SAMPLES

RETURNS

CARRAY SHAPE NSAMPLES RETURNS PREDICTED VALUES

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{U}{V}$  WHERE U IS THE RESIDUAL SUM OF SQUARES  $Y_{TRUE} - Y_{PRED}$

2SUM AND V IS THE TOTAL SUM OF SQUARES  $Y_{TRUE} - Y_{TRUEMEAN}$  2SUM THE BEST POSSIBLE SCORE

IS 10 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS

PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 00

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES FOR SOME ESTIMATORS THIS MAY

BE A PRECOMPUTED KERNEL MATRIX INSTEAD SHAPE NSAMPLES NSAMPLESFITTED WHERE

NSAMPLESFITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE VALUES FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT R2 OF SELF PREDICTX WRT Y

NOTES

THE R2 SCORE USED WHEN CALLING SCORE ON A REGRESSOR WILL USE MULTIOUTPUTUNIFORMAVERAGE

FROM VERSION 023 TO KEEP CONSISTENT WITH METRICSR2SCORE THIS WILL INFLUENCE THE SCORE

METHOD OF ALL THE MULTIOUTPUT REGRESSORS EXCEPT FOR MULTIOUTPUTMULTIOUTPUTREGRESSOR

TO SPECIFY THE DEFAULT VALUE MANUALLY AND AVOID THE WARNING PLEASE EITHER CALL METRICSR2SCORE

DIRECTLY OR MAKE A CUSTOM SCORER WITH METRICSMAKESCORER THE BUILTIN SCORER R2 USES

MULTIOUTPUTUNIFORMAVERAGE

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

622SKLEARNLINEARMODEL GENERALIZED LINEAR MODELS 1883

SCIKITLEARN USER GUIDE RELEASE 0213

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS  
SELF

62265KLEARNLINEARMODEL LASSO

CLASSSSKLEARNLINEARMODEL LASSOALPHA10 FITINTERCEPTTRUE NORMALIZEFALSE PRECOMPUTEFALSE COPYXTRUE MAXITER1000 TOL00001

WARMSTARTFALSE POSITIVEFALSE RANDOMSTATENONE SELECTION'CYCLIC'

LINEAR MODEL TRAINED WITH L1 PRIOR AS REGULARIZER AKA THE LASSO

THE OPTIMIZATION OBJECTIVE FOR LASSO IS

$\frac{1}{2} \sum_{i=1}^n (y_i - x_i^T w)^2 + \frac{\lambda}{2} \sum_{j=1}^p |w_j|$

TECHNICALLY THE LASSO MODEL IS OPTIMIZING THE SAME OBJECTIVE FUNCTION AS THE ELASTIC NET WITH L1RATIO10 NO L2 PENALTY

READ MORE IN THE USER GUIDE

PARAMETERS

ALPHA FLOAT OPTIONAL CONSTANT THAT MULTIPLIES THE L1 TERM DEFAULTS TO 10 ALPHA 0 IS EQUIVALENT TO AN ORDINARY LEAST SQUARE SOLVED BY THE LINEARREGRESSION OBJECT FOR NUMERICAL REASONS USING ALPHA 0 WITH THELASSO OBJECT IS NOT ADVISED GIVEN THIS YOU SHOULD USE THE LINEARREGRESSION OBJECT

FITINTERCEPT BOOLEAN OPTIONAL DEFAULT TRUE WHETHER TO CALCULATE THE INTERCEPT FOR THIS MODEL IF SET TO FALSE NO INTERCEPT WILL BE USED IN CALCULATIONS EG DATA IS EXPECTED TO BE ALREADY CENTERED

NORMALIZE BOOLEAN OPTIONAL DEFAULT FALSE THIS PARAMETER IS IGNORED WHEN FITINTERCEPT IS SET TO FALSE IF TRUE THE REGRESSORS X WILL BE NORMALIZED BEFORE REGRESSION BY SUBTRACTING THE MEAN AND DIVIDING BY THE L2NORM IF YOU WISH TO STANDARDIZE PLEASE USE SKLEARNPREPROCESSINGSTANDARDSCALER BEFORE CALLINGFIT ON AN ESTIMATOR WITH NORMALIZEFALSE

PRECOMPUTE TRUE FALSE ARRAYLIKE DEFAULTFALSE WHETHER TO USE A PRECOMPUTED GRAM MATRIX TO SPEED UP CALCULATIONS IF SET TO AUTO LET US DECIDE THE GRAM MATRIX CAN ALSO BE PASSED AS ARGUMENT FOR SPARSE INPUT THIS OPTION IS ALWAYS TRUE TO PRESERVE SPARSITY

COPYX BOOLEAN OPTIONAL DEFAULT TRUE IF TRUE X WILL BE COPIED ELSE IT MAY BE OVERWRITTEN

MAXITER INT OPTIONAL THE MAXIMUM NUMBER OF ITERATIONS

TOLFLOAT OPTIONAL THE TOLERANCE FOR THE OPTIMIZATION IF THE UPDATES ARE SMALLER THAN TOL THE OPTIMIZATION CODE CHECKS THE DUAL GAP FOR OPTIMALITY AND CONTINUES UNTIL IT IS SMALLER THAN TOL

WARMSTART BOOL OPTIONAL WHEN SET TO TRUE REUSE THE SOLUTION OF THE PREVIOUS CALL TO FIT AS INITIALIZATION OTHERWISE JUST ERASE THE PREVIOUS SOLUTION SEE THE GLOSSARY

POSITIVE BOOL OPTIONAL WHEN SET TO TRUE FORCES THE COEFFICIENTS TO BE POSITIVE

1884 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULT NONE THE SEED OF THE PSEUDO RANDOM NUMBER GENERATOR THAT SELECTS A RANDOM FEATURE TO UPDATE IF INT RANDOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NPRANDOM USED WHEN SELECTION 'RANDOM' SELECTION STR DEFAULT 'CYCLIC' IF SET TO 'RANDOM' A RANDOM COEFFICIENT IS UPDATED EVERY ITERATION RATHER THAN LOOPING OVER FEATURES SEQUENTIALLY BY DEFAULT THIS SETTING TO 'RANDOM' OFTEN LEADS TO SIGNIFICANTLY FASTER CONVERGENCE ESPECIALLY WHEN TOL IS HIGHER THAN 1E4

ATTRIBUTES

COEF ARRAY SHAPE NFEATURES NTARGETS NFEATURES PARAMETER VECTOR W IN THE COST FUNCTION FORMULA

SPARSECOEF SCIPYSPARSE MATRIX SHAPE NFEATURES 1 NTARGETS NFEATURES SPARSE REPRESENTATION OF THE FITTED COEF

INTERCEPT FLOAT ARRAY SHAPE NTARGETS INDEPENDENT TERM IN DECISION FUNCTION

NITER INT ARRAYLIKE SHAPE NTARGETS NUMBER OF ITERATIONS RUN BY THE COORDINATE DESCENT

SOLVER TO REACH THE SPECIFIED TOLERANCE

SEE ALSO

LARSPATH

LASSOPATH

LASSOLARS

LASSOCV

LASSOLARSCV

SKLEARNDECOMPOSITIONSPARSEENCODE

NOTES

THE ALGORITHM USED TO FIT THE MODEL IS COORDINATE DESCENT

TO AVOID UNNECESSARY MEMORY DUPLICATION THE X ARGUMENT OF THE FIT METHOD SHOULD BE DIRECTLY PASSED AS A

FORTRANCONTIGUOUS NUMPY ARRAY

EXAMPLES

FROM SKLEARN IMPORT LINEARMODEL

CLF LINEARMODELLASSOALPHA01

CLFFIT00 1 1 2 2 0 1 2

LASSOALPHA01 COPYXTRUE FITINTERCEPTTRUE MAXITER1000

NORMALIZEFALSE POSITIVEFALSE PRECOMPUTEFALSE RANDOMSTATENONE

SELECTIONCYCLIC TOL00001 WARMSTARTFALSE

PRINTCLFCOEF

085 0

PRINTCLFINTERCEPT

015

622SKLEARNLINEARMODEL GENERALIZED LINEAR MODELS 1885

SCIKITLEARN USER GUIDE RELEASE 0213

METHODS

FITSELF X Y CHECKINPUT FIT MODEL WITH COORDINATE DESCENT  
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR  
PATH X Y L1RATIO EPS NALPHAS COMPUTE ELASTIC NET PATH WITH COORDINATE DESCENT  
PREDICT SELF X PREDICT USING THE LINEAR MODEL  
SCORE SELF X Y SAMPLEWEIGHT RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PRE  
DICTION  
SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR  
INIT SELFALPHA10 FITINTERCEPTTRUE NORMALIZEFALSE PRECOMPUTEFALSE COPYXTRUE  
MAXITER1000 TOL00001 WARMSTARTFALSE POSITIVEFALSE RANDOMSTATENONE SE  
LECTION'CYCLIC'  
FITSELFXYCHECKINPUTTRUE  
FIT MODEL WITH COORDINATE DESCENT

PARAMETERS

XNDARRAY OR SCIPYSPARSE MATRIX NSAMPLES NFEATURES DATA  
YNDARRAY SHAPE NSAMPLES OR NSAMPLES NTARGETS TARGET WILL BE CAST TO X'S DTYPE  
IF NECESSARY  
CHECKINPUT BOOLEAN DEFAULTTRUE ALLOW TO BYPASS SEVERAL INPUT CHECKING DON'T USE  
THIS PARAMETER UNLESS YOU KNOW WHAT YOU DO

NOTES

COORDINATE DESCENT IS AN ALGORITHM THAT CONSIDERS EACH COLUMN OF DATA AT A TIME HENCE IT WILL AUTOMATICALLY  
CONVERT THE X INPUT AS A FORTRANCONTIGUOUS NUMPY ARRAY IF NECESSARY  
TO AVOID MEMORY REALLOCATION IT IS ADVISED TO ALLOCATE THE INITIAL DATA IN MEMORY DIRECTLY USING THAT FORMAT

GETPARAMS SELFDEEPTTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBJECTS THAT ARE ESTIMATORS  
RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

STATICPATHXYL1RATIO05 EPS0001 NALPHAS100 ALPHASNONE PRECOMPUTE'AUTO'

XYNONE COPYXTRUE COEFINITNONE VERBOSEFALSE RETURNNITERFALSE POSI

TIVEFALSE CHECKINPUTTRUE PARAMS

COMPUTE ELASTIC NET PATH WITH COORDINATE DESCENT

THE ELASTIC NET OPTIMIZATION FUNCTION VARIES FOR MONO AND MULTIOUTPUTS

FOR MONOOUTPUT TASKS IT IS

1 2NSAMPLES Y XW22

ALPHA L1RATIO W1

05ALPHA1 L1RATIO W22

1886 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

FOR MULTIOUTPUT TASKS IT IS

1 2 NSAMPLES Y XWFRO2

ALPHA L1RATIO W21

05ALPHA1 L1RATIO WFRO2

WHERE

W21 SUMI SQRTSUMJ WIJ2

IE THE SUM OF NORM OF EACH ROW

READ MORE IN THE USER GUIDE

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA PASS DIRECTLY AS FORTRAN

CONTIGUOUS DATA TO AVOID UNNECESSARY MEMORY DUPLICATION IF YIS MONOOUTPUT THEN X

CAN BE SPARSE

YNDARRAY SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TARGET VALUES

L1RATIO FLOAT OPTIONAL FLOAT BETWEEN 0 AND 1 PASSED TO ELASTIC NET SCALING BETWEEN L1 AND

L2 PENALTIES L1RATIO1 CORRESPONDS TO THE LASSO

EPS FLOAT LENGTH OF THE PATH EPS1E3 MEANS THAT ALPHAMIN ALPHAMAX

1E3

NALPHAS INT OPTIONAL NUMBER OF ALPHAS ALONG THE REGULARIZATION PATH

ALPHAS NDARRAY OPTIONAL LIST OF ALPHAS WHERE TO COMPUTE THE MODELS IF NONE ALPHAS ARE

SET AUTOMATICALLY

PRECOMPUTE TRUE FALSE 'AUTO' ARRAYLIKE WHETHER TO USE A PRECOMPUTED GRAM MATRIX

TO SPEED UP CALCULATIONS IF SET TO AUTO LET US DECIDE THE GRAM MATRIX CAN ALSO BE

PASSED AS ARGUMENT

XYARRAYLIKE OPTIONAL XY NPDOTXT Y THAT CAN BE PRECOMPUTED IT IS USEFUL ONLY

WHEN THE GRAM MATRIX IS PRECOMPUTED

COPYX BOOLEAN OPTIONAL DEFAULT TRUE IF TRUE X WILL BE COPIED ELSE IT MAY BE OVER

WRITTEN

COEFINIT ARRAY SHAPE NFEATURES NONE THE INITIAL VALUES OF THE COEFFICIENTS

VERBOSE BOOL OR INTEGER AMOUNT OF VERBOSITY

RETURNNITER BOOL WHETHER TO RETURN THE NUMBER OF ITERATIONS OR NOT

POSITIVE BOOL DEFAULT FALSE IF SET TO TRUE FORCES COEFFICIENTS TO BE POSITIVE ONLY ALLOWED

WHENYNDIM 1

CHECKINPUT BOOL DEFAULT TRUE SKIP INPUT VALIDATION CHECKS INCLUDING THE GRAM MATRIX

WHEN PROVIDED ASSUMING THERE ARE HANDLED BY THE CALLER WHEN CHECKINPUTFALSE

PARAMS KWARGS KEYWORD ARGUMENTS PASSED TO THE COORDINATE DESCENT SOLVER

RETURNS

ALPHAS ARRAY SHAPE NALPHAS THE ALPHAS ALONG THE PATH WHERE MODELS ARE COMPUTED

COEFS ARRAY SHAPE NFEATURES NALPHAS OR NOUTPUTS NFEATURES NALPHAS COEFFICIENTS

ALONG THE PATH

622SKLEARNLINEARMODEL GENERALIZED LINEAR MODELS 1887

SCIKITLEARN USER GUIDE RELEASE 0213

DUALGAPS ARRAY SHAPE NALPHAS THE DUAL GAPS AT THE END OF THE OPTIMIZATION FOR EACH ALPHA

NITERS ARRAYLIKE SHAPE NALPHAS THE NUMBER OF ITERATIONS TAKEN BY THE COORDINATE DESCENT OPTIMIZER TO REACH THE SPECIFIED TOLERANCE FOR EACH ALPHA IS RETURNED WHEN RETURNNITER IS SET TO TRUE

SEE ALSO

MULTITASKELASTICNET

MULTITASKELASTICNETCV

ELASTICNET

ELASTICNETCV

NOTES

FOR AN EXAMPLE SEE EXAMPLESLINEARMODELPLOTLASSOCOORDINATEDDESCENTPATHPY

PREDICTSELF

PREDICT USING THE LINEAR MODEL

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES SAMPLES

RETURNS

CARRAY SHAPE NSAMPLES RETURNS PREDICTED VALUES

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{U}{V}$  WHERE U IS THE RESIDUAL SUM OF SQUARES  $Y_{TRUE} - Y_{PRED}$  2SUM AND V IS THE TOTAL SUM OF SQUARES  $Y_{TRUE} - Y_{TRUEMEAN}$  2SUM THE BEST POSSIBLE SCORE IS 10 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 00

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES FOR SOME ESTIMATORS THIS MAY BE A PRECOMPUTED KERNEL MATRIX INSTEAD SHAPE NSAMPLES NSAMPLESFITTED WHERE NSAMPLESFITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE VALUES FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT R2 OF SELF-PREDICTX WRT Y

NOTES

THE R2 SCORE USED WHEN CALLING SCORE ON A REGRESSOR WILL USE MULTIOUTPUTUNIFORMAVERAGE FROM VERSION 023 TO KEEP CONSISTENT WITH METRICSR2SCORE THIS WILL INFLUENCE THE SCORE METHOD OF ALL THE MULTIOUTPUT REGRESSORS EXCEPT FOR MULTIOUTPUTMULTIOUTPUTREGRESSOR TO SPECIFY THE DEFAULT VALUE MANUALLY AND AVOID THE WARNING PLEASE EITHER CALL METRICSR2SCORE

1888 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213  
DIRECTLY OR MAKE A CUSTOM SCORER WITH METRICSMAKESCORER THE BUILTIN SCORER R2 USES  
MULTIOUTPUTUNIFORMAVERAGE  
SETPARAMS SELFPARAMS  
SET THE PARAMETERS OF THIS ESTIMATOR  
THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE  
PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT  
OF A NESTED OBJECT  
RETURNS  
SELF  
SPARSECOEF  
SPARSE REPRESENTATION OF THE FITTED COEF  
EXAMPLES USING SKLEARNLINEARMODELLASSO  
•COMPRESSIVE SENSING TOMOGRAPHY RECONSTRUCTION WITH L1 PRIOR LASSO  
•CROSSVALIDATION ON DIABETES DATASET EXERCISE  
•LASSO ON DENSE AND SPARSE DATA  
•JOINT FEATURE SELECTION WITH MULTITASK LASSO  
•LASSO AND ELASTIC NET FOR SPARSE SIGNALS  
6227SKLEARNLINEARMODEL LASSOLARS  
CLASSSKLEARNLINEARMODEL LASSOLARS ALPHA10 FITINTERCEPTTRUE VERBOSEFALSE NOR  
MALIZETTRUE PRECOMPUTE'AUTO' MAXITER500  
EPS2220446049250313E16 COPYXTRUE  
FITPATHTRUE POSITIVEFALSE  
LASSO MODEL FIT WITH LEAST ANGLE REGRESSION AKA LARS  
IT IS A LINEAR MODEL TRAINED WITH AN L1 PRIOR AS REGULARIZER  
THE OPTIMIZATION OBJECTIVE FOR LASSO IS  
$$\frac{1}{2} \sum_{i=1}^n (y_i - x_i^T w)^2 + \alpha \sum_{j=1}^p |w_j|$$
  
READ MORE IN THE USER GUIDE  
PARAMETERS  
ALPHA FLOAT CONSTANT THAT MULTIPLIES THE PENALTY TERM DEFAULTS TO 10 ALPHA 0 IS EQUIV  
ALENT TO AN ORDINARY LEAST SQUARE SOLVED BY LINEARREGRESSION FOR NUMERICAL REASONS  
USINGALPHA 0 WITH THE LASSOLARS OBJECT IS NOT ADVISED AND YOU SHOULD PREFER THE LIN  
EARREGRESSION OBJECT  
FITINTERCEPT BOOLEAN WHETHER TO CALCULATE THE INTERCEPT FOR THIS MODEL IF SET TO FALSE NO  
INTERCEPT WILL BE USED IN CALCULATIONS EG DATA IS EXPECTED TO BE ALREADY CENTERED  
VERBOSE BOOLEAN OR INTEGER OPTIONAL SETS THE VERBOSITY AMOUNT  
NORMALIZE BOOLEAN OPTIONAL DEFAULT TRUE THIS PARAMETER IS IGNORED WHEN FITINTERCEPT  
IS SET TO FALSE IF TRUE THE REGRESSORS X WILL BE NORMALIZED BEFORE REGRESSION BY SUB  
TRACTING THE MEAN AND DIVIDING BY THE L2NORM IF YOU WISH TO STANDARDIZE PLEASE USE  
622SKLEARNLINEARMODEL GENERALIZED LINEAR MODELS 1889

SCIKITLEARN USER GUIDE RELEASE 0213

SKLEARNPREPROCESSINGSTANDARDSCALER BEFORE CALLING FIT ON AN ESTIMATOR

WITHNORMALIZEFALSE

PRECOMPUTE TRUE FALSE ‘AUTO’ ARRAYLIKE WHETHER TO USE A PRECOMPUTED GRAM MATRIX TO SPEED UP CALCULATIONS IF SET TO AUTO LET US DECIDE THE GRAM MATRIX CAN ALSO BE PASSED AS ARGUMENT

MAXITER INTEGER OPTIONAL MAXIMUM NUMBER OF ITERATIONS TO PERFORM

EPS FLOAT OPTIONAL THE MACHINEPRECISION REGULARIZATION IN THE COMPUTATION OF THE CHOLESKY DIAGONAL FACTORS INCREASE THIS FOR VERY ILLCONDITIONED SYSTEMS UNLIKE THE TOL PARAMETER IN SOME ITERATIVE OPTIMIZATIONBASED ALGORITHMS THIS PARAMETER DOES NOT CONTROL THE TOLERANCE OF THE OPTIMIZATION

COPYX BOOLEAN OPTIONAL DEFAULT TRUE IF TRUE X WILL BE COPIED ELSE IT MAY BE OVERWRITTEN

FITPATH BOOLEAN IF TRUE THE FULL PATH IS STORED IN THE COEFPATH ATTRIBUTE IF YOU COMPUTE THE SOLUTION FOR A LARGE PROBLEM OR MANY TARGETS SETTING FITPATH TOFALSE WILL LEAD TO A SPEEDUP ESPECIALLY WITH A SMALL ALPHA

POSITIVE BOOLEAN DEFAULTFALSE RESTRICT COEFFICIENTS TO BE 0 BE AWARE THAT YOU MIGHT WANT TO REMOVE FITINTERCEPT WHICH IS SET TRUE BY DEFAULT UNDER THE POSITIVE RESTRICTION THE MODEL COEFFICIENTS WILL NOT CONVERGE TO THE ORDINARYLEASTSQUARES SOLUTION FOR SMALL VALUES OF ALPHA ONLY COEFFICIENTS UP TO THE SMALLEST ALPHA VALUE ALPHASALPHAS 0 MIN WHEN FITPATHTRUE REACHED BY THE STEPWISE LARSLASSO ALGORITHM ARE TYPICALLY IN CONGRUENCE WITH THE SOLUTION OF THE COORDINATE DESCENT LASSO ESTIMATOR

ATTRIBUTES

ALPHAS ARRAY SHAPE NALPHAS 1 LIST OF NTARGETS SUCH ARRAYS MAXIMUM OF COVARIANCES IN ABSOLUTE VALUE AT EACH ITERATION NALPHAS IS EITHERMAXITER NFEATURES OR THE NUMBER OF NODES IN THE PATH WITH CORRELATION GREATER THAN ALPHA WHICHEVER IS SMALLER

ACTIVE LIST LENGTH NALPHAS LIST OF NTARGETS SUCH LISTS INDICES OF ACTIVE VARIABLES AT THE END OF THE PATH

COEFPATH ARRAY SHAPE NFEATURES NALPHAS 1 OR LIST IF A LIST IS PASSED IT’S EXPECTED TO BE ONE OF NTARGETS SUCH ARRAYS THE VARYING VALUES OF THE COEFFICIENTS ALONG THE PATH IT IS NOT PRESENT IF THE FITPATH PARAMETER IS FALSE

COEF ARRAY SHAPE NFEATURES OR NTARGETS NFEATURES PARAMETER VECTOR W IN THE FORMULA TION FORMULA

INTERCEPT FLOAT ARRAY SHAPE NTARGETS INDEPENDENT TERM IN DECISION FUNCTION

NITER ARRAYLIKE OR INT THE NUMBER OF ITERATIONS TAKEN BY LARSPATH TO FIND THE GRID OF ALPHAS FOR EACH TARGET

SEE ALSO

LARSPATH

LASSOPATH

LASSO

LASSOCV

LASSOLARSCV

LASSOLARSIC

SKLEARNDECOMPOSITIONSPARSEENCODE

1890 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES

```
FROM SKLEARN IMPORT LINEARMODEL
REG LINEARMODELLASSOLARSALPHA001
REGFIT1 1 0 0 1 1 1 0 1
```

LASSOLARSALPHA001 COPYXTRUE EPS FITINTERCEPTTRUE  
FITPATHTRUE MAXITER500 NORMALIZETRUE POSITIVEFALSE  
PRECOMPUTEAUTO VERBOSEFALSE  
PRINTREGCOEF  
0 0963257

METHODS

FITSELF X Y XY FIT THE MODEL USING X Y AS TRAINING DATA  
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR  
PREDICT SELF X PREDICT USING THE LINEAR MODEL  
SCORE SELF X Y SAMPLEWEIGHT RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PRE  
DICTION  
SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR  
INIT SELFALPHA10 FITINTERCEPTTRUE VERBOSEFALSE NORMALIZETRUE PRECOMPUTE'AUTO'  
MAXITER500 EPS2220446049250313E16 COPYXTRUE FITPATHTRUE POSI  
TIVEFALSE  
FITSELFXYXYNONE  
FIT THE MODEL USING X Y AS TRAINING DATA  
PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA  
YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NTARGETS TARGET VALUES  
XYARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NTARGETS OPTIONAL XY NPDOTXT Y  
THAT CAN BE PRECOMPUTED IT IS USEFUL ONLY WHEN THE GRAM MATRIX IS PRECOMPUTED  
RETURNS  
SELF OBJECT RETURNS AN INSTANCE OF SELF  
GETPARAMS SELFDEEPTRUE  
GET PARAMETERS FOR THIS ESTIMATOR  
PARAMETERS  
DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBOBJECTS THAT ARE ESTIMATORS  
RETURNS  
PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES  
PREDICTSELF  
PREDICT USING THE LINEAR MODEL  
PARAMETERS  
XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES SAMPLES  
622SKLEARNLINEARMODEL GENERALIZED LINEAR MODELS 1891

SCIKITLEARN USER GUIDE RELEASE 0213

RETURNS

CARRAY SHAPE NSAMPLES RETURNS PREDICTED VALUES

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{U}{V}$  WHERE U IS THE RESIDUAL SUM OF SQUARES  $Y_{TRUE} - Y_{PRED}$  2SUM AND V IS THE TOTAL SUM OF SQUARES  $Y_{TRUE} - Y_{TRUEMEAN}$  2SUM THE BEST POSSIBLE SCORE IS 10 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 00

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES FOR SOME ESTIMATORS THIS MAY BE A PRECOMPUTED KERNEL MATRIX INSTEAD SHAPE NSAMPLES NSAMPLESFITTED WHERE NSAMPLESFITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE VALUES FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT R2 OF SELF PREDICTX WRT Y

NOTES

THE R2 SCORE USED WHEN CALLING SCORE ON A REGRESSOR WILL USE MULTIOUTPUTUNIFORMAVERAGE FROM VERSION 023 TO KEEP CONSISTENT WITH METRICSR2SCORE THIS WILL INFLUENCE THE SCORE METHOD OF ALL THE MULTIOUTPUT REGRESSORS EXCEPT FOR MULTIOUTPUTMULTIOUTPUTREGRESSOR TO SPECIFY THE DEFAULT VALUE MANUALLY AND AVOID THE WARNING PLEASE EITHER CALL METRICSR2SCORE DIRECTLY OR MAKE A CUSTOM SCORER WITH METRICSMAKESCORER THE BUILTIN SCORER R2 USES MULTIOUTPUTUNIFORMAVERAGE

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

6228SKLEARNLINEARMODEL LINEARREGRESSION

CLASSSKLEARNLINEARMODEL LINEARREGRESSION FITINTERCEPTTRUE NORMALIZEFALSE

COPYXTRUE NJOBSNONE

ORDINARY LEAST SQUARES LINEAR REGRESSION

PARAMETERS

FITINTERCEPT BOOLEAN OPTIONAL DEFAULT TRUE WHETHER TO CALCULATE THE INTERCEPT FOR THIS MODEL IF SET TO FALSE NO INTERCEPT WILL BE USED IN CALCULATIONS EG DATA IS EXPECTED TO BE ALREADY CENTERED

1892 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

NORMALIZE BOOLEAN OPTIONAL DEFAULT FALSE THIS PARAMETER IS IGNORED WHEN  
FITINTERCEPT IS SET TO FALSE IF TRUE THE REGRESSORS X WILL BE NORMALIZED BE  
FORE REGRESSION BY SUBTRACTING THE MEAN AND DIVIDING BY THE L2NORM IF YOU WISH TO  
STANDARDIZE PLEASE USE SKLEARNPREPROCESSINGSTANDARDSCALER BEFORE  
CALLINGFIT ON AN ESTIMATOR WITH NORMALIZEFALSE

COPYX BOOLEAN OPTIONAL DEFAULT TRUE IF TRUE X WILL BE COPIED ELSE IT MAY BE OVERWRITTEN  
NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF JOBS TO USE FOR THE COMPUTATION  
THIS WILL ONLY PROVIDE SPEEDUP FOR NTARGETS 1 AND SUFFICIENT LARGE PROBLEMS NONE MEANS  
1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL PROCESSORS SEE  
GLOSSARY FOR MORE DETAILS

ATTRIBUTES

COEF ARRAY SHAPE NFEATURES OR NTARGETS NFEATURES ESTIMATED COEFFICIENTS FOR THE LINEAR  
REGRESSION PROBLEM IF MULTIPLE TARGETS ARE PASSED DURING THE FIT Y 2D THIS IS A 2D ARRAY OF  
SHAPE NTARGETS NFEATURES WHILE IF ONLY ONE TARGET IS PASSED THIS IS A 1D ARRAY OF LENGTH  
NFEATURES

INTERCEPT ARRAY INDEPENDENT TERM IN THE LINEAR MODEL

NOTES

FROM THE IMPLEMENTATION POINT OF VIEW THIS IS JUST PLAIN ORDINARY LEAST SQUARES SCIPYLINALGLSTSQ WRAPPED AS  
A PREDICTOR OBJECT

EXAMPLES

```
import numpy as np
from sklearn.linear_model import LinearRegression
X = np.array([1, 1, 2, 2, 2, 2, 3])
Y = [1, 0, 2, 1, 3]
Y = np.dot(X, np.array([2, 3]))
reg = LinearRegression()
reg.fit(X, Y)
reg.score(X, Y)
10
reg.coef_
array([2, 3])
reg.intercept_
30000
reg.predict(np.array([3, 5]))
array([16, 16])
METHODS
fit(self, X, y, sample_weight=None) FIT LINEAR MODEL
get_params(self, deep=True) GET PARAMETERS FOR THIS ESTIMATOR
predict(self, X) PREDICT USING THE LINEAR MODEL
score(self, X, y, sample_weight=None) RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PRE  
DICTION
set_params(self, **kwargs) SET THE PARAMETERS OF THIS ESTIMATOR
```

622SKLEARNLINEARMODEL GENERALIZED LINEAR MODELS 1893

SCIKITLEARN USER GUIDE RELEASE 0213

INIT SELF FIT INTERCEPT TRUE NORMALIZE FALSE COPY X TRUE N JOBS NONE

FIT SELF X Y SAMPLE WEIGHT NONE

FIT LINEAR MODEL

PARAMETERS

X ARRAY LIKE OR SPARSE MATRIX SHAPE N SAMPLES N FEATURES TRAINING DATA

Y ARRAY LIKE SHAPE N SAMPLES N TARGETS TARGET VALUES WILL BE CAST TO X'S DTYPE IF NECESSARY

SAMPLE WEIGHT NUMPY ARRAY OF SHAPE N SAMPLES INDIVIDUAL WEIGHTS FOR EACH SAMPLE

NEW IN VERSION 0.17 PARAMETER SAMPLE WEIGHT SUPPORT TO LINEAR REGRESSION

RETURNS

SELF RETURNS AN INSTANCE OF SELF

GET PARAMS SELF DEEP TRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICT SELF X

PREDICT USING THE LINEAR MODEL

PARAMETERS

X ARRAY LIKE OR SPARSE MATRIX SHAPE N SAMPLES N FEATURES SAMPLES

RETURNS

C ARRAY SHAPE N SAMPLES RETURNS PREDICTED VALUES

SCORE SELF X Y SAMPLE WEIGHT NONE

RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{U}{V}$  WHERE U IS THE RESIDUAL SUM OF SQUARES  $Y_{TRUE} - Y_{PRED}$  2SUM AND V IS THE TOTAL SUM OF SQUARES  $Y_{TRUE} - Y_{TRUEMEAN}$  2SUM THE BEST POSSIBLE SCORE IS 1.0 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 0.0

PARAMETERS

X ARRAY LIKE SHAPE N SAMPLES N FEATURES TEST SAMPLES FOR SOME ESTIMATORS THIS MAY BE A PRECOMPUTED KERNEL MATRIX INSTEAD SHAPE N SAMPLES N SAMPLES FITTED WHERE N SAMPLES FITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR

Y ARRAY LIKE SHAPE N SAMPLES OR N SAMPLES N OUTPUTS TRUE VALUES FOR X

SAMPLE WEIGHT ARRAY LIKE SHAPE N SAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT R2 OF SELF PREDICT X WRT Y

1894 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

NOTES

THE R2 SCORE USED WHEN CALLING SCORE ON A REGRESSOR WILL USE MULTIOUTPUTUNIFORMAVERAGE FROM VERSION 023 TO KEEP CONSISTENT WITH METRICSR2SCORE THIS WILL INFLUENCE THE SCORE METHOD OF ALL THE MULTIOUTPUT REGRESSORS EXCEPT FOR MULTIOUTPUTMULTIOUTPUTREGRESSOR TO SPECIFY THE DEFAULT VALUE MANUALLY AND AVOID THE WARNING PLEASE EITHER CALL METRICSR2SCORE DIRECTLY OR MAKE A CUSTOM SCORER WITH METRICSMAKESCORER THE BUILTIN SCORER R2 USES MULTIOUTPUTUNIFORMAVERAGE  
SETPARAMS SELFPARAMS  
SET THE PARAMETERS OF THIS ESTIMATOR  
THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT  
RETURNS

SELF  
EXAMPLES USING SKLEARNLINEARMODELLINEARREGRESSION

- ISOTONIC REGRESSION
  - FACE COMPLETION WITH A MULTIOUTPUT ESTIMATORS
  - PLOT INDIVIDUAL AND VOTING REGRESSION PREDICTIONS
  - ORDINARY LEAST SQUARES AND RIDGE REGRESSION VARIANCE
  - LOGISTIC FUNCTION
  - LINEAR REGRESSION EXAMPLE
  - ROBUST LINEAR MODEL ESTIMATION USING RANSAC
  - SPARSITY EXAMPLE FITTING ONLY FEATURES 1 AND 2
  - THEILSEN REGRESSION
  - ROBUST LINEAR ESTIMATOR FITTING
  - AUTOMATIC RELEVANCE DETERMINATION REGRESSION ARD
  - BAYESIAN RIDGE REGRESSION
  - PLOTING CROSSVALIDATED PREDICTIONS
  - UNDERFITTING VS OVERFITTING
  - USING KBINSDISCRETIZER TO DISCRETIZE CONTINUOUS FEATURES
- 622SKLEARNLINEARMODEL GENERALIZED LINEAR MODELS 1895

SCIKITLEARN USER GUIDE RELEASE 0213  
62295KLEARNLINEARMODEL LOGISTICREGRESSION  
CLASSSSKLEARNLINEARMODEL LOGISTICREGRESSION PENALTY'L2' DUALFALSE TOL00001  
C10 FITINTERCEPTTRUE INTER  
CEPTSCALING1 CLASSWEIGHTNONE  
RANDOMSTATENONE SOLVER'WARN'  
MAXITER100 MULTICLASS'WARN' VER  
BOSE0 WARMSTARTFALSE NJOBSNONE  
L1RATIONONE  
LOGISTIC REGRESSION AKA LOGIT MAXENT CLASSIFIER  
IN THE MULTICLASS CASE THE TRAINING ALGORITHM USES THE ONEVSREST OVR SCHEME IF THE 'MULTICLASS' OPTION IS  
SET TO 'OVR' AND USES THE CROSSENTROPY LOSS IF THE 'MULTICLASS' OPTION IS SET TO 'MULTINOMIAL' CURRENTLY THE  
'MULTINOMIAL' OPTION IS SUPPORTED ONLY BY THE 'LBFGS' 'SAG' 'SAGA' AND 'NEWTONCG' SOLVERS  
THIS CLASS IMPLEMENTS REGULARIZED LOGISTIC REGRESSION USING THE 'LIBLINEAR' LIBRARY 'NEWTONCG' 'SAG' 'SAGA' AND  
'LBFGS' SOLVERS NOTE THAT REGULARIZATION IS APPLIED BY DEFAULT IT CAN HANDLE BOTH DENSE AND SPARSE INPUT USE  
CORDERED ARRAYS OR CSR MATRICES CONTAINING 64BIT FLOATS FOR OPTIMAL PERFORMANCE ANY OTHER INPUT FORMAT WILL  
BE CONVERTED AND COPIED  
THE 'NEWTONCG' 'SAG' AND 'LBFGS' SOLVERS SUPPORT ONLY L2 REGULARIZATION WITH PRIMAL FORMULATION OR NO REG  
ULARIZATION THE 'LIBLINEAR' SOLVER SUPPORTS BOTH L1 AND L2 REGULARIZATION WITH A DUAL FORMULATION ONLY FOR THE  
L2 PENALTY THE ELASTICNET REGULARIZATION IS ONLY SUPPORTED BY THE 'SAGA' SOLVER  
READ MORE IN THE USER GUIDE  
PARAMETERS  
PENALTY STR 'L1' 'L2' 'ELASTICNET' OR 'NONE' OPTIONAL DEFAULT'L2' USED TO SPECIFY THE NORM  
USED IN THE PENALIZATION THE 'NEWTONCG' 'SAG' AND 'LBFGS' SOLVERS SUPPORT ONLY L2 PENAL  
TIES 'ELASTICNET' IS ONLY SUPPORTED BY THE 'SAGA' SOLVER IF 'NONE' NOT SUPPORTED BY THE  
LIBLINEAR SOLVER NO REGULARIZATION IS APPLIED  
NEW IN VERSION 019 L1 PENALTY WITH SAGA SOLVER ALLOWING 'MULTINOMIAL' L1  
DUAL BOOL OPTIONAL DEFAULTFALSE DUAL OR PRIMAL FORMULATION DUAL FORMULATION IS ONLY  
IMPLEMENTED FOR L2 PENALTY WITH LIBLINEAR SOLVER PREFER DUALFALSE WHEN NSAMPLES  
NFEATURES  
TOLFLOAT OPTIONAL DEFAULT1E4 TOLERANCE FOR STOPPING CRITERIA  
CFLOAT OPTIONAL DEFAULT10 INVERSE OF REGULARIZATION STRENGTH MUST BE A POSITIVE FLOAT LIKE  
IN SUPPORT VECTOR MACHINES SMALLER VALUES SPECIFY STRONGER REGULARIZATION  
FITINTERCEPT BOOL OPTIONAL DEFAULTTRUE SPECIFIES IF A CONSTANT AKA BIAS OR INTERCEPT  
SHOULD BE ADDED TO THE DECISION FUNCTION  
INTERCEPTSCALING FLOAT OPTIONAL DEFAULT1 USEFUL ONLY WHEN THE SOLVER 'LIBLINEAR' IS USED  
AND SELFFITINTERCEPT IS SET TO TRUE IN THIS CASE X BECOMES X SELFINTERCEPTSCALING  
IE A "SYNTHETIC" FEATURE WITH CONSTANT VALUE EQUAL TO INTERCEPTSCALING IS AP  
PENDE TO THE INSTANCE VECTOR THE INTERCEPT BECOMES INTERCEPTSCALING  
SYNTHETICFEATUREWEIGHT  
NOTE THE SYNTHETIC FEATURE WEIGHT IS SUBJECT TO L1L2 REGULARIZATION AS ALL OTHER FEATURES TO  
LESSEN THE EFFECT OF REGULARIZATION ON SYNTHETIC FEATURE WEIGHT AND THEREFORE ON THE INTERCEPT  
INTERCEPTSCALING HAS TO BE INCREASED  
CLASSWEIGHT DICT OR 'BALANCED' OPTIONAL DEFAULTNONE WEIGHTS ASSOCIATED WITH CLASSES IN  
THE FORMCLASSLABEL WEIGHT IF NOT GIVEN ALL CLASSES ARE SUPPOSED TO HAVE  
WEIGHT ONE  
1896 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

THE “BALANCED” MODE USES THE VALUES OF Y TO AUTOMATICALLY ADJUST WEIGHTS INVERSELY PROPORTIONAL TO CLASS FREQUENCIES IN THE INPUT DATA AS NSAMPLES NCLASSES NP

BINCOUNTY

NOTE THAT THESE WEIGHTS WILL BE MULTIPLIED WITH SAMPLEWEIGHT PASSED THROUGH THE FIT METHOD IF SAMPLEWEIGHT IS SPECIFIED

NEW IN VERSION 017 CLASSWEIGHT‘BALANCED’

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE THE SEED OF THE PSEUDO RANDOM NUMBER GENERATOR TO USE WHEN SHUFFLING THE DATA IF INT RANDOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NPRANDOM USED WHEN SOLVER ‘SAG’ OR ‘LIBLINEAR’

SOLVER STR ‘NEWTONCG’ ‘LBFGS’ ‘LIBLINEAR’ ‘SAG’ ‘SAGA’ OPTIONAL DEFAULT‘LIBLINEAR’ ALGORITHM TO USE IN THE OPTIMIZATION PROBLEM

- FOR SMALL DATASETS ‘LIBLINEAR’ IS A GOOD CHOICE WHEREAS ‘SAG’ AND ‘SAGA’ ARE FASTER FOR LARGE ONES
- FOR MULTICLASS PROBLEMS ONLY ‘NEWTONCG’ ‘SAG’ ‘SAGA’ AND ‘LBFGS’ HANDLE MULTINOMIAL LOSS ‘LIBLINEAR’ IS LIMITED TO ONEVERSUSREST SCHEMES
- ‘NEWTONCG’ ‘LBFGS’ ‘SAG’ AND ‘SAGA’ HANDLE L2 OR NO PENALTY
- ‘LIBLINEAR’ AND ‘SAGA’ ALSO HANDLE L1 PENALTY
- ‘SAGA’ ALSO SUPPORTS ‘ELASTICNET’ PENALTY
- ‘LIBLINEAR’ DOES NOT HANDLE NO PENALTY

NOTE THAT ‘SAG’ AND ‘SAGA’ FAST CONVERGENCE IS ONLY GUARANTEED ON FEATURES WITH APPROXIMATELY THE SAME SCALE YOU CAN PREPROCESS THE DATA WITH A SCALER FROM SKLEARNPREPROCESSING

NEW IN VERSION 017 STOCHASTIC AVERAGE GRADIENT DESCENT SOLVER

NEW IN VERSION 019 SAGA SOLVER

CHANGED IN VERSION 020 DEFAULT WILL CHANGE FROM ‘LIBLINEAR’ TO ‘LBFGS’ IN 022

MAXITER INT OPTIONAL DEFAULT100 MAXIMUM NUMBER OF ITERATIONS TAKEN FOR THE SOLVERS TO CONVERGE

MULTICLASS STR ‘OVR’ ‘MULTINOMIAL’ ‘AUTO’ OPTIONAL DEFAULT‘OVR’ IF THE OPTION CHOSEN IS ‘OVR’ THEN A BINARY PROBLEM IS FIT FOR EACH LABEL FOR ‘MULTINOMIAL’ THE LOSS MINIMISED IS THE MULTINOMIAL LOSS FIT ACROSS THE ENTIRE PROBABILITY DISTRIBUTION EVEN WHEN THE DATA IS BINARY ‘MULTINOMIAL’ IS UNAVAILABLE WHEN SOLVER‘LIBLINEAR’ ‘AUTO’ SELECTS ‘OVR’ IF THE DATA IS BINARY OR IF SOLVER‘LIBLINEAR’ AND OTHERWISE SELECTS ‘MULTINOMIAL’

NEW IN VERSION 018 STOCHASTIC AVERAGE GRADIENT DESCENT SOLVER FOR ‘MULTINOMIAL’ CASE

CHANGED IN VERSION 020 DEFAULT WILL CHANGE FROM ‘OVR’ TO ‘AUTO’ IN 022

VERBOSE INT OPTIONAL DEFAULT0 FOR THE LIBLINEAR AND LBFGS SOLVERS SET VERBOSE TO ANY POSITIVE NUMBER FOR VERBOSITY

WARMSTART BOOL OPTIONAL DEFAULTFALSE WHEN SET TO TRUE REUSE THE SOLUTION OF THE PREVIOUS CALL TO FIT AS INITIALIZATION OTHERWISE JUST ERASE THE PREVIOUS SOLUTION USELESS FOR LIBLINEAR SOLVER SEE THE GLOSSARY

NEW IN VERSION 017 WARMSTART TO SUPPORT LBFGS NEWTONCG SAGSAGA SOLVERS

NJOBS INT OR NONE OPTIONAL DEFAULTNONE NUMBER OF CPU CORES USED WHEN PARALLELIZING OVER CLASSES IF MULTICLASS‘OVR’” THIS PARAMETER IS IGNORED WHEN THE SOLVER IS SET TO

622SKLEARNLINEARMODEL GENERALIZED LINEAR MODELS 1897

SCIKITLEARN USER GUIDE RELEASE 0213

‘LIBLINEAR’ REGARDLESS OF WHETHER ‘MULTICLASS’ IS SPECIFIED OR NOT NONE MEANS 1 UNLESS IN AJOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL PROCESSORS SEE GLOSSARY

FOR MORE DETAILS

L1RATIO FLOAT OR NONE OPTIONAL DEFAULTNONE THE ELASTICNET MIXING PARAMETER WITH 0 L1RATIO 1 ONLY USED IF PENALTYELASTICNET SETTING

L1RATIO0 IS EQUIVALENT TO USING PENALTYL2 WHILE SETTING L1RATIO1 IS EQUIVALENT TO USING PENALTYL1 FOR0 L1RATIO 1 THE PENALTY IS A COMBINATION OF L1 AND L2

ATTRIBUTES

CLASSES ARRAY SHAPE NCLASSES A LIST OF CLASS LABELS KNOWN TO THE CLASSIFIER

COEF ARRAY SHAPE 1 NFEATURES OR NCLASSES NFEATURES COEFFICIENT OF THE FEATURES IN THE DECISION FUNCTION

COEF IS OF SHAPE 1 NFEATURES WHEN THE GIVEN PROBLEM IS BINARY IN PARTICULAR WHEN MULTICLASSMULTINOMIAL COEF CORRESPONDS TO OUTCOME 1 TRUE AND COEF CORRESPONDS TO OUTCOME 0 FALSE

INTERCEPT ARRAY SHAPE 1 OR NCLASSES INTERCEPT AKA BIAS ADDED TO THE DECISION FUNCTION

IFFITINTERCEPT IS SET TO FALSE THE INTERCEPT IS SET TO ZERO INTERCEPT IS OF SHAPE 1 WHEN THE GIVEN PROBLEM IS BINARY IN PARTICULAR WHEN MULTICLASSMULTINOMIAL INTERCEPT CORRESPONDS TO OUTCOME 1 TRUE AND INTERCEPT CORRESPONDS TO OUTCOME 0 FALSE

NITER ARRAY SHAPE NCLASSES OR 1 ACTUAL NUMBER OF ITERATIONS FOR ALL CLASSES IF BINARY OR MULTINOMIAL IT RETURNS ONLY 1 ELEMENT FOR LIBLINEAR SOLVER ONLY THE MAXIMUM NUMBER OF ITERATION ACROSS ALL CLASSES IS GIVEN

CHANGED IN VERSION 020 IN SCIPY 100 THE NUMBER OF LBFGS ITERATIONS MAY EXCEED MAXITER NITER WILL NOW REPORT AT MOST MAXITER

SEE ALSO

SGDCCLASSIFIER INCREMENTALLY TRAINED LOGISTIC REGRESSION WHEN GIVEN THE PARAMETER LOSSLOG

LOGISTICREGRESSIONCV LOGISTIC REGRESSION WITH BUILTIN CROSS VALIDATION

NOTES

THE UNDERLYING C IMPLEMENTATION USES A RANDOM NUMBER GENERATOR TO SELECT FEATURES WHEN FITTING THE MODEL IT IS THUS NOT UNCOMMON TO HAVE SLIGHTLY DIFFERENT RESULTS FOR THE SAME INPUT DATA IF THAT HAPPENS TRY WITH A SMALLER TOL PARAMETER

PREDICT OUTPUT MAY NOT MATCH THAT OF STANDALONE LIBLINEAR IN CERTAIN CASES SEE DIFFERENCES FROM LIBLINEAR IN THE NARRATIVE DOCUMENTATION

REFERENCES

LIBLINEAR – A LIBRARY FOR LARGE LINEAR CLASSIFICATION [HTTPSWWWCSIENTUEDUTWCJLINLIBLINEAR](https://www.cs.tu.edu.tw/cjlin/liblinear)

SAG – MARK SCHMIDT NICOLAS LE ROUX AND FRANCIS BACH MINIMIZING FINITE SUMS WITH THE STOCHASTIC AVERAGE GRADIENT [HTTPSHALINRIAFRHAL00860051DOCUMENT](http://hal.inria.fr/hal-00860051/document)

1898 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

SAGA - DEFAZIO A BACH F LACOSTEJULIEN S 2014 SAGA A FAST INCREMENTAL GRADIENT METHOD WITH  
SUPPORT FOR NONSTRONGLY CONVEX COMPOSITE OBJECTIVES [HTTPSARXIVORGABS14070202](https://arxiv.org/abs/1407.0202)

HSIANGFU YU FANGLAN HUANG CHIHJEN LIN 2011 DUAL COORDINATE DESCENT METHODS FOR LOGISTIC RE  
GRESSION AND MAXIMUM ENTROPY MODELS MACHINE LEARNING 85124175 [HTTPSWWWCSIENTUEDUTW](https://www.cs.tu-dortmund.de/~tue/dutw)  
CJLINPAPERSMAXENTDUALPDF

EXAMPLES

```
FROM SKLEARNDATASETS IMPORT LOADIRIS
FROM SKLEARNLINEARMODEL IMPORT LOGISTICREGRESSION
X Y LOADIRISRETURNXY TRUE
CLF LOGISTICREGRESSIONRANDOMSTATE0 SOLVERLBFGS
MULTICLASSMULTINOMIALFITX Y
CLFPREDICTX2
ARRAY0 0
CLFPREDICTPROBAX2
ARRAY98E01 18E02 14E08
97E01 28E02 E08
CLFSCOREX Y
097
```

METHODS

DECISIONFUNCTION SELF X PREDICT CONFIDENCE SCORES FOR SAMPLES

DENSIFY SELF CONVERT COEFFICIENT MATRIX TO DENSE ARRAY FORMAT

FITSELF X Y SAMPLEWEIGHT FIT THE MODEL ACCORDING TO THE GIVEN TRAINING DATA

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PREDICT SELF X PREDICT CLASS LABELS FOR SAMPLES IN X

PREDICTLOGPROBA SELF X LOG OF PROBABILITY ESTIMATES

PREDICTPROBA SELF X PROBABILITY ESTIMATES

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND  
LABELS

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

SPARSIFY SELF CONVERT COEFFICIENT MATRIX TO SPARSE FORMAT

INIT SELFPENALTY'L2' DUALFALSE TOL00001 C10 FITINTERCEPTTRUE INTER  
CEPTSCALING1 CLASSWEIGHTNONE RANDOMSTATENONE SOLVER'WARN' MAXITER100

MULTICLASS'WARN' VERBOSE0 WARMSTARTFALSE NJOBSNONE L1RATIONONE

DECISIONFUNCTION SELF X

PREDICT CONFIDENCE SCORES FOR SAMPLES

THE CONFIDENCE SCORE FOR A SAMPLE IS THE SIGNED DISTANCE OF THAT SAMPLE TO THE HYPERPLANE

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES SAMPLES

RETURNS

ARRAY SHAPENSAMPLES IF NCLASSES > 2 ELSE NSAMPLES NCLASSES CONFIDENCE  
SCORES PER SAMPLE CLASS COMBINATION IN THE BINARY CASE CONFIDENCE SCORE FOR  
SELFCLASSES1 WHERE 0 MEANS THIS CLASS WOULD BE PREDICTED

622SKLEARNLINEARMODEL GENERALIZED LINEAR MODELS 1899

SCIKITLEARN USER GUIDE RELEASE 0213

DENSIFYSELF

CONVERT COEFFICIENT MATRIX TO DENSE ARRAY FORMAT

CONVERTS THE COEF MEMBER BACK TO A NUMPYNDARRAY THIS IS THE DEFAULT FORMAT OF COEF AND IS REQUIRED FOR FITTING SO CALLING THIS METHOD IS ONLY REQUIRED ON MODELS THAT HAVE PREVIOUSLY BEEN SPARSIFIED OTHERWISE IT IS A NOOP

RETURNS

SELF ESTIMATOR

FITSELFXYSAMPLEWEIGHTNONE

FIT THE MODEL ACCORDING TO THE GIVEN TRAINING DATA

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES TRAINING VECTOR WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

YARRAYLIKE SHAPE NSAMPLES TARGET VECTOR RELATIVE TO X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL ARRAY OF WEIGHTS THAT ARE ASSIGNED TO INDIVIDUAL SAMPLES IF NOT PROVIDED THEN EACH SAMPLE IS GIVEN UNIT WEIGHT

NEW IN VERSION 017 SAMPLEWEIGHT SUPPORT TO LOGISTICREGRESSION

RETURNS

SELF OBJECT

NOTES

THE SAGA SOLVER SUPPORTS BOTH FLOAT64 AND FLOAT32 BIT ARRAYS

GETPARAMS SELFDEEPTTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELF

PREDICT CLASS LABELS FOR SAMPLES IN X

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES SAMPLES

RETURNS

CARRAY SHAPE NSAMPLES PREDICTED CLASS LABEL PER SAMPLE

PREDICTLOGPROBA SELF

LOG OF PROBABILITY ESTIMATES

THE RETURNED ESTIMATES FOR ALL CLASSES ARE ORDERED BY THE LABEL OF CLASSES

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES

1900 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

RETURNS

TARRAYLIKE SHAPE NSAMPLES NCLASSES RETURNS THE LOGPROBABILITY OF THE SAMPLE FOR EACH CLASS IN THE MODEL WHERE CLASSES ARE ORDERED AS THEY ARE IN SELFCLASSES

PREDICTPROBA SELF

PROBABILITY ESTIMATES

THE RETURNED ESTIMATES FOR ALL CLASSES ARE ORDERED BY THE LABEL OF CLASSES

FOR A MULTICLASS PROBLEM IF MULTICLASS IS SET TO BE "MULTINOMIAL" THE SOFTMAX FUNCTION IS USED TO FIND THE PREDICTED PROBABILITY OF EACH CLASS ELSE USE A ONEVSREST APPROACH IE CALCULATE THE PROBABILITY OF EACH CLASS ASSUMING IT TO BE POSITIVE USING THE LOGISTIC FUNCTION AND NORMALIZE THESE VALUES ACROSS ALL THE CLASSES

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES

RETURNS

TARRAYLIKE SHAPE NSAMPLES NCLASSES RETURNS THE PROBABILITY OF THE SAMPLE FOR EACH CLASS IN THE MODEL WHERE CLASSES ARE ORDERED AS THEY ARE IN SELFCLASSES

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS

IN MULTILABEL CLASSIFICATION THIS IS THE SUBSET ACCURACY WHICH IS A HARSH METRIC SINCE YOU REQUIRE FOR EACH SAMPLE THAT EACH LABEL SET BE CORRECTLY PREDICTED

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE LABELS FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT MEAN ACCURACY OF SELF PREDICTX WRT Y

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

SPARSIFY SELF

CONVERT COEFFICIENT MATRIX TO SPARSE FORMAT

CONVERTS THE COEF MEMBER TO A SCIPYSPARSE MATRIX WHICH FOR L1REGULARIZED MODELS CAN BE MUCH MORE MEMORY AND STORAGEEFFICIENT THAN THE USUAL NUMPYNDARRAY REPRESENTATION

THEINTERCEPT MEMBER IS NOT CONVERTED

RETURNS

SELF ESTIMATOR

622SKLEARNLINEARMODEL GENERALIZED LINEAR MODELS 1901

SCIKITLEARN USER GUIDE RELEASE 0213

NOTES

FOR NONSPARSE MODELS IE WHEN THERE ARE NOT MANY ZEROS IN COEF THIS MAY ACTUALLY INCREASE MEMORY USAGE SO USE THIS METHOD WITH CARE A RULE OF THUMB IS THAT THE NUMBER OF ZERO ELEMENTS WHICH CAN BE COMPUTED WITH COEF\_0SUM MUST BE MORE THAN 50 FOR THIS TO PROVIDE SIGNIFICANT BENEFITS AFTER CALLING THIS METHOD FURTHER FITTING WITH THE PARTIALFIT METHOD IF ANY WILL NOT WORK UNTIL YOU CALL DENSIFY

EXAMPLES USING SKLEARNLINEARMODELLOGISTICREGRESSION

- COMPACT ESTIMATOR REPRESENTATIONS
- COMPARISON OF CALIBRATION OF CLASSIFIERS
- PROBABILITY CALIBRATION CURVES
- PLOT CLASSIFICATION PROBABILITY
- COLUMN TRANSFORMER WITH MIXED TYPES
- PLOT CLASS PROBABILITIES CALCULATED BY THE VOTINGCLASSIFIER
- FEATURE TRANSFORMATIONS WITH ENSEMBLES OF TREES
- DIGITS CLASSIFICATION EXERCISE
- REGULARIZATION PATH OF L1 LOGISTIC REGRESSION
- LOGISTIC FUNCTION
- LOGISTIC REGRESSION 3CLASS CLASSIFIER
- COMPARING VARIOUS ONLINE SOLVERS
- MNIST CLASSFICATION USING MULTINOMIAL LOGISTIC L1
- PLOT MULTINOMIAL AND ONEVSREST LOGISTIC REGRESSION
- L1 PENALTY AND SPARSITY IN LOGISTIC REGRESSION
- MULTICLASS SPARSE LOGISITIC REGRESSION ON NEWGROUPS20
- CLASSIFIER CHAIN
- RESTRICTED BOLTZMANN MACHINE FEATURES FOR DIGIT CLASSIFICATION
- FEATURE DISCRETIZATION

62210SKLEARNLINEARMODEL MULTITASKLASSO  
CLASSSSKLEARNLINEARMODEL MULTITASKLASSO ALPHA10 FITINTERCEPTTRUE NORMALIZEFALSE  
COPYXTRUE MAXITER1000 TOL00001  
WARMSTARTFALSE RANDOMSTATENONE SELEC  
TION'CYCLIC'  
MULTITASK LASSO MODEL TRAINED WITH L1L2 MIXEDNORM AS REGULARIZER  
THE OPTIMIZATION OBJECTIVE FOR LASSO IS  
1 2 NSAMPLES Y XW2FRO ALPHA W21  
WHERE  
1902 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

W21 SUMI SQRTSUMJ WJ2

IE THE SUM OF NORM OF EACH ROW

READ MORE IN THE USER GUIDE

PARAMETERS

ALPHA FLOAT OPTIONAL CONSTANT THAT MULTIPLIES THE L1L2 TERM DEFAULTS TO 10

FITINTERCEPT BOOLEAN WHETHER TO CALCULATE THE INTERCEPT FOR THIS MODEL IF SET TO FALSE NO

INTERCEPT WILL BE USED IN CALCULATIONS EG DATA IS EXPECTED TO BE ALREADY CENTERED

NORMALIZE BOOLEAN OPTIONAL DEFAULT FALSE THIS PARAMETER IS IGNORED WHEN

FITINTERCEPT IS SET TO FALSE IF TRUE THE REGRESSORS X WILL BE NORMALIZED BE

FORE REGRESSION BY SUBTRACTING THE MEAN AND DIVIDING BY THE L2NORM IF YOU WISH TO

STANDARDIZE PLEASE USE SKLEARNPREPROCESSINGSTANDARDSCALER BEFORE

CALLINGFIT ON AN ESTIMATOR WITH NORMALIZEFALSE

COPYX BOOLEAN OPTIONAL DEFAULT TRUE IF TRUE X WILL BE COPIED ELSE IT MAY BE OVERWRITTEN

MAXITER INT OPTIONAL THE MAXIMUM NUMBER OF ITERATIONS

TOLFLOAT OPTIONAL THE TOLERANCE FOR THE OPTIMIZATION IF THE UPDATES ARE SMALLER THAN TOL THE

OPTIMIZATION CODE CHECKS THE DUAL GAP FOR OPTIMALITY AND CONTINUES UNTIL IT IS SMALLER THAN

TOL

WARMSTART BOOL OPTIONAL WHEN SET TO TRUE REUSE THE SOLUTION OF THE PREVIOUS CALL TO FIT AS

INITIALIZATION OTHERWISE JUST ERASE THE PREVIOUS SOLUTION SEE THE GLOSSARY

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULT NONE THE SEED OF THE

PSEUDO RANDOM NUMBER GENERATOR THAT SELECTS A RANDOM FEATURE TO UPDATE IF INT RAN

DOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE

RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS

THE RANDOMSTATE INSTANCE USED BY NPRANDOM USED WHEN SELECTION 'RANDOM'

SELECTION STR DEFAULT 'CYCLIC' IF SET TO 'RANDOM' A RANDOM COEFFICIENT IS UPDATED EVERY ITERATION

RATHER THAN LOOPING OVER FEATURES SEQUENTIALLY BY DEFAULT THIS SETTING TO 'RANDOM' OFTEN

LEADS TO SIGNIFICANTLY FASTER CONVERGENCE ESPECIALLY WHEN TOL IS HIGHER THAN 1E4

ATTRIBUTES

COEF ARRAY SHAPE NTASKS NFEATURES PARAMETER VECTOR W IN THE COST FUNCTION FORMULA

NOTE THATCOEF STORES THE TRANSPOSE OF WWT

INTERCEPT ARRAY SHAPE NTASKS INDEPENDENT TERM IN DECISION FUNCTION

NITER INT NUMBER OF ITERATIONS RUN BY THE COORDINATE DESCENT SOLVER TO REACH THE SPECIFIED

TOLERANCE

SEE ALSO

MULTITASKLASSO MULTITASK L1L2 LASSO WITH BUILTIN CROSSVALIDATION

LASSO

MULTITASKELASTICNET

NOTES

THE ALGORITHM USED TO FIT THE MODEL IS COORDINATE DESCENT

622SKLEARNLINEARMODEL GENERALIZED LINEAR MODELS 1903

SCIKITLEARN USER GUIDE RELEASE 0213  
TO AVOID UNNECESSARY MEMORY DUPLICATION THE X ARGUMENT OF THE FIT METHOD SHOULD BE DIRECTLY PASSED AS A  
FORTRANCONTIGUOUS NUMPY ARRAY  
EXAMPLES  
FROM SKLEARN IMPORT LINEARMODEL  
CLF LINEARMODELMULTITASKLASSOALPHA01  
CLFFIT00 1 1 2 2 0 0 1 1 2 2

MULTITASKLASSOALPHA01 COPYXTRUE FITINTERCEPTTRUE MAXITER1000  
NORMALIZEFALSE RANDOMSTATENONE SELECTIONCYCLIC TOL00001  
WARMSTARTFALSE  
PRINTCLFCOEF  
089393398 0  
089393398 0  
PRINTCLFINTERCEPT  
010606602 010606602  
METHODS  
FITSELF X Y FIT MULTITASKELASTICNET MODEL WITH COORDINATE DESCENT  
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR  
PATH X Y L1RATIO EPS NALPHAS COMPUTE ELASTIC NET PATH WITH COORDINATE DESCENT  
PREDICT SELF X PREDICT USING THE LINEAR MODEL  
SCORE SELF X Y SAMPLEWEIGHT RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PRE  
DICTION  
SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR  
INIT SELFALPHA10 FITINTERCEPTTRUE NORMALIZEFALSE COPYXTRUE MAXITER1000  
TOL00001 WARMSTARTFALSE RANDOMSTATENONE SELECTION'CYCLIC'  
FITSELFXY  
FIT MULTITASKELASTICNET MODEL WITH COORDINATE DESCENT  
PARAMETERS  
XNDARRAY SHAPE NSAMPLES NFEATURES DATA  
YNDARRAY SHAPE NSAMPLES NTASKS TARGET WILL BE CAST TO X'S DTYPE IF NECESSARY  
NOTES  
COORDINATE DESCENT IS AN ALGORITHM THAT CONSIDERS EACH COLUMN OF DATA AT A TIME HENCE IT WILL AUTOMATICALLY  
CONVERT THE X INPUT AS A FORTRANCONTIGUOUS NUMPY ARRAY IF NECESSARY  
TO AVOID MEMORY REALLOCATION IT IS ADVISED TO ALLOCATE THE INITIAL DATA IN MEMORY DIRECTLY USING THAT FORMAT  
GETPARAMS SELFDEEPTTRUE  
GET PARAMETERS FOR THIS ESTIMATOR  
PARAMETERS  
DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBOBJECTS THAT ARE ESTIMATORS  
1904 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PATHXYL1RATIO05 EPS0001 NALPHAS100 ALPHASNONE PRECOMPUTE'AUTO' XYNONE

COPYXTRUE COEFINITNONE VERBOSEFALSE RETURNNITERFALSE POSITIVEFALSE

CHECKINPUTTRUE PARAMS

COMPUTE ELASTIC NET PATH WITH COORDINATE DESCENT

THE ELASTIC NET OPTIMIZATION FUNCTION VARIES FOR MONO AND MULTIOUTPUTS

FOR MONOOUTPUT TASKS IT IS

1 2NSAMPLES Y XW22

ALPHA L1RATIO W1

05ALPHA1 L1RATIO W22

FOR MULTIOUTPUT TASKS IT IS

1 2 NSAMPLES Y XWFRO2

ALPHA L1RATIO W21

05ALPHA1 L1RATIO WFRO2

WHERE

W21 SUMI SQRTSUMJ WIJ2

IE THE SUM OF NORM OF EACH ROW

READ MORE IN THE USER GUIDE

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA PASS DIRECTLY AS FORTRAN

CONTIGUOUS DATA TO AVOID UNNECESSARY MEMORY DUPLICATION IF YIS MONOOUTPUT THEN X

CAN BE SPARSE

YNDARRAY SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TARGET VALUES

L1RATIO FLOAT OPTIONAL FLOAT BETWEEN 0 AND 1 PASSED TO ELASTIC NET SCALING BETWEEN L1 AND

L2 PENALTIES L1RATIO1 CORRESPONDS TO THE LASSO

EPS FLOAT LENGTH OF THE PATH EPS1E3 MEANS THAT ALPHAMIN ALPHAMAX

1E3

NALPHAS INT OPTIONAL NUMBER OF ALPHAS ALONG THE REGULARIZATION PATH

ALPHAS NDARRAY OPTIONAL LIST OF ALPHAS WHERE TO COMPUTE THE MODELS IF NONE ALPHAS ARE

SET AUTOMATICALLY

PRECOMPUTE TRUE FALSE 'AUTO' ARRAYLIKE WHETHER TO USE A PRECOMPUTED GRAM MATRIX

TO SPEED UP CALCULATIONS IF SET TO AUTO LET US DECIDE THE GRAM MATRIX CAN ALSO BE

PASSED AS ARGUMENT

XYARRAYLIKE OPTIONAL XY NPDOTXT Y THAT CAN BE PRECOMPUTED IT IS USEFUL ONLY

WHEN THE GRAM MATRIX IS PRECOMPUTED

COPYX BOOLEAN OPTIONAL DEFAULT TRUE IF TRUE X WILL BE COPIED ELSE IT MAY BE OVER

WRITTEN

COEFINIT ARRAY SHAPE NFEATURES NONE THE INITIAL VALUES OF THE COEFFICIENTS

VERBOSE BOOL OR INTEGER AMOUNT OF VERBOSITY

622SKLEARNLINEARMODEL GENERALIZED LINEAR MODELS 1905

SCIKITLEARN USER GUIDE RELEASE 0213

RETURNNITER BOOL WHETHER TO RETURN THE NUMBER OF ITERATIONS OR NOT

POSITIVE BOOL DEFAULT FALSE IF SET TO TRUE FORCES COEFFICIENTS TO BE POSITIVE ONLY ALLOWED WHENYNDIM 1

CHECKINPUT BOOL DEFAULT TRUE SKIP INPUT VALIDATION CHECKS INCLUDING THE GRAM MATRIX WHEN PROVIDED ASSUMING THERE ARE HANDLED BY THE CALLER WHEN CHECKINPUTFALSE

PARAMS KWARGS KEYWORD ARGUMENTS PASSED TO THE COORDINATE DESCENT SOLVER

RETURNS

ALPHAS ARRAY SHAPE NALPHAS THE ALPHAS ALONG THE PATH WHERE MODELS ARE COMPUTED

COEFS ARRAY SHAPE NFEATURES NALPHAS OR NOUTPUTS NFEATURES NALPHAS COEFFICIENTS ALONG THE PATH

DUALGAPS ARRAY SHAPE NALPHAS THE DUAL GAPS AT THE END OF THE OPTIMIZATION FOR EACH ALPHA

NITERS ARRAYLIKE SHAPE NALPHAS THE NUMBER OF ITERATIONS TAKEN BY THE COORDINATE DESCENT OPTIMIZER TO REACH THE SPECIFIED TOLERANCE FOR EACH ALPHA IS RETURNED WHEN RETURNNITER IS SET TO TRUE

SEE ALSO

MULTITASKELASTICNET

MULTITASKELASTICNETCV

ELASTICNET

ELASTICNETCV

NOTES

FOR AN EXAMPLE SEE EXAMPLESLINEARMODELPLOTLASSOCOORDINATEDDESCENTPATHPY

PREDICTSELF

PREDICT USING THE LINEAR MODEL

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES SAMPLES

RETURNS

CARRAY SHAPE NSAMPLES RETURNS PREDICTED VALUES

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{UV}{2SUM}$  WHERE U IS THE RESIDUAL SUM OF SQUARES YTRUE YPRED 2SUM AND V IS THE TOTAL SUM OF SQUARES YTRUE YTRUEMEAN 2SUM THE BEST POSSIBLE SCORE IS 10 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 00

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES FOR SOME ESTIMATORS THIS MAY BE A PRECOMPUTED KERNEL MATRIX INSTEAD SHAPE NSAMPLES NSAMPLESFITTED WHERE NSAMPLESFITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE VALUES FOR X

1906 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213  
SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS  
RETURNS  
SCORE FLOAT R2 OF SELF PREDICTX WRT Y  
NOTES  
THE R2 SCORE USED WHEN CALLING SCORE ON A REGRESSOR WILL USE MULTIOUTPUTUNIFORMAVERAGE  
FROM VERSION 023 TO KEEP CONSISTENT WITH METRICSR2SCORE THIS WILL INFLUENCE THE SCORE  
METHOD OF ALL THE MULTIOUTPUT REGRESSORS EXCEPT FOR MULTIOUTPUTMULTIOUTPUTREGRESSOR  
TO SPECIFY THE DEFAULT VALUE MANUALLY AND AVOID THE WARNING PLEASE EITHER CALL METRICSR2SCORE  
DIRECTLY OR MAKE A CUSTOM SCORER WITH METRICSMAKESCORER THE BUILTIN SCORER R2 USES  
MULTIOUTPUTUNIFORMAVERAGE  
SETPARAMS SELFPARAMS  
SET THE PARAMETERS OF THIS ESTIMATOR  
THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE  
PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT  
OF A NESTED OBJECT  
RETURNS  
SELF  
SPARSECOEF  
SPARSE REPRESENTATION OF THE FITTED COEF  
EXAMPLES USING SKLEARNLINEARMODELMULTITASKLASSO  
•JOINT FEATURE SELECTION WITH MULTITASK LASSO  
62211SKLEARNLINEARMODEL MULTITASKELASTICNET  
CLASSSSKLEARNLINEARMODEL MULTITASKELASTICNET ALPHA10 L1RATIO05  
FITINTERCEPTTRUE NORMALIZEFALSE  
COPYXTRUE MAXITER1000 TOL00001  
WARMSTARTFALSE RANDOMSTATENONE  
SELECTION'CYCLIC'  
MULTITASK ELASTICNET MODEL TRAINED WITH L1L2 MIXEDNORM AS REGULARIZER  
THE OPTIMIZATION OBJECTIVE FOR MULTITASKELASTICNET IS  
1 2 NSAMPLES Y XWFR02  
ALPHA L1RATIO W21  
05ALPHA1 L1RATIO WFR02  
WHERE  
W21 SUMI SQRTSUMJ WJ 2  
IE THE SUM OF NORM OF EACH ROW  
READ MORE IN THE USER GUIDE  
PARAMETERS  
622SKLEARNLINEARMODEL GENERALIZED LINEAR MODELS 1907

SCIKITLEARN USER GUIDE RELEASE 0213

ALPHA FLOAT OPTIONAL CONSTANT THAT MULTIPLIES THE L1L2 TERM DEFAULTS TO 10  
L1RATIO FLOAT THE ELASTICNET MIXING PARAMETER WITH 0 L1RATIO 1 FOR L1RATIO 1 THE  
PENALTY IS AN L1L2 PENALTY FOR L1RATIO 0 IT IS AN L2 PENALTY FOR 0 L1RATIO  
1 THE PENALTY IS A COMBINATION OF L1L2 AND L2  
FITINTERCEPT BOOLEAN WHETHER TO CALCULATE THE INTERCEPT FOR THIS MODEL IF SET TO FALSE NO  
INTERCEPT WILL BE USED IN CALCULATIONS EG DATA IS EXPECTED TO BE ALREADY CENTERED  
NORMALIZE BOOLEAN OPTIONAL DEFAULT FALSE THIS PARAMETER IS IGNORED WHEN  
FITINTERCEPT IS SET TO FALSE IF TRUE THE REGRESSORS X WILL BE NORMALIZED BE  
FORE REGRESSION BY SUBTRACTING THE MEAN AND DIVIDING BY THE L2NORM IF YOU WISH TO  
STANDARDIZE PLEASE USE SKLEARNPREPROCESSINGSTANDARDSCALER BEFORE  
CALLINGFIT ON AN ESTIMATOR WITH NORMALIZEFALSE  
COPYX BOOLEAN OPTIONAL DEFAULT TRUE IF TRUE X WILL BE COPIED ELSE IT MAY BE OVERWRITTEN  
MAXITER INT OPTIONAL THE MAXIMUM NUMBER OF ITERATIONS  
TOLFLOAT OPTIONAL THE TOLERANCE FOR THE OPTIMIZATION IF THE UPDATES ARE SMALLER THAN TOL THE  
OPTIMIZATION CODE CHECKS THE DUAL GAP FOR OPTIMALITY AND CONTINUES UNTIL IT IS SMALLER THAN  
TOL  
WARMSTART BOOL OPTIONAL WHEN SET TO TRUE REUSE THE SOLUTION OF THE PREVIOUS CALL TO FIT AS  
INITIALIZATION OTHERWISE JUST ERASE THE PREVIOUS SOLUTION SEE THE GLOSSARY  
RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULT NONE THE SEED OF THE  
PSEUDO RANDOM NUMBER GENERATOR THAT SELECTS A RANDOM FEATURE TO UPDATE IF INT RAN  
DOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE  
RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS  
THE RANDOMSTATE INSTANCE USED BY NPRANDOM USED WHEN SELECTION 'RANDOM'  
SELECTION STR DEFAULT 'CYCLIC' IF SET TO 'RANDOM' A RANDOM COEFFICIENT IS UPDATED EVERY ITERATION  
RATHER THAN LOOPING OVER FEATURES SEQUENTIALLY BY DEFAULT THIS SETTING TO 'RANDOM' OFTEN  
LEADS TO SIGNIFICANTLY FASTER CONVERGENCE ESPECIALLY WHEN TOL IS HIGHER THAN 1E4  
ATTRIBUTES  
INTERCEPT ARRAY SHAPE NTASKS INDEPENDENT TERM IN DECISION FUNCTION  
COEF ARRAY SHAPE NTASKS NFEATURES PARAMETER VECTOR W IN THE COST FUNCTION FORMULA IF  
A 1D Y IS PASSED IN AT FIT NON MULTITASK USAGE COEF IS THEN A 1D ARRAY NOTE THAT COEF  
STORES THE TRANSPOSE OF WWT  
NITER INT NUMBER OF ITERATIONS RUN BY THE COORDINATE DESCENT SOLVER TO REACH THE SPECIFIED  
TOLERANCE  
SEE ALSO  
MULTITASKELASTICNET MULTITASK L1L2 ELASTICNET WITH BUILTIN CROSSVALIDATION  
ELASTICNET  
MULTITASKLASSO  
NOTES  
THE ALGORITHM USED TO FIT THE MODEL IS COORDINATE DESCENT  
TO AVOID UNNECESSARY MEMORY DUPLICATION THE X ARGUMENT OF THE FIT METHOD SHOULD BE DIRECTLY PASSED AS A  
FORTRANCONTIGUOUS NUMPY ARRAY  
1908 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES

```
FROM SKLEARN IMPORT LINEARMODEL
CLF LINEARMODELMULTITASKELASTICNETALPHA01
CLFFIT00 1 1 2 2 0 0 1 1 2 2
```

```
MULTITASKELASTICNETALPHA01 COPYXTRUE FITINTERCEPTTRUE
L1RATIO05 MAXITER1000 NORMALIZEFALSE RANDOMSTATENONE
SELECTIONCYCLIC TOL00001 WARMSTARTFALSE
PRINTCLFCOEF
045663524 045612256
045663524 045612256
PRINTCLFINTERCEPT
00872422 00872422
METHODS
FITSELF X Y FIT MULTITASKELASTICNET MODEL WITH COORDINATE DESCENT
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR
PATH X Y L1RATIO EPS NALPHAS COMPUTE ELASTIC NET PATH WITH COORDINATE DESCENT
PREDICT SELF X PREDICT USING THE LINEAR MODEL
SCORE SELF X Y SAMPLEWEIGHT RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PRE
DICTION
SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR
INIT SELFALPHA10 L1RATIO05 FITINTERCEPTTRUE NORMALIZEFALSE COPYXTRUE
MAXITER1000 TOL00001 WARMSTARTFALSE RANDOMSTATENONE SELECTION'CYCLIC'
FITSELFXY
FIT MULTITASKELASTICNET MODEL WITH COORDINATE DESCENT
PARAMETERS
XNDARRAY SHAPE NSAMPLES NFEATURES DATA
YNDARRAY SHAPE NSAMPLES NTASKS TARGET WILL BE CAST TO X'S DTYPE IF NECESSARY
NOTES
COORDINATE DESCENT IS AN ALGORITHM THAT CONSIDERS EACH COLUMN OF DATA AT A TIME HENCE IT WILL AUTOMATICALLY
CONVERT THE X INPUT AS A FORTRANCONTIGUOUS NUMPY ARRAY IF NECESSARY
TO AVOID MEMORY REALLOCATION IT IS ADVISED TO ALLOCATE THE INITIAL DATA IN MEMORY DIRECTLY USING THAT FORMAT
GETPARAMS SELFDEEPTURE
GET PARAMETERS FOR THIS ESTIMATOR
PARAMETERS
DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED
SUBOBJECTS THAT ARE ESTIMATORS
RETURNS
PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES
622SKLEARNLINEARMODEL GENERALIZED LINEAR MODELS 1909
```

SCIKITLEARN USER GUIDE RELEASE 0213

PATHXYL1RATIO05 EPS0001 NALPHAS100 ALPHASNONE PRECOMPUTE'AUTO' XYNONE  
COPYXTRUE COEFINITNONE VERBOSEFALSE RETURNNITERFALSE POSITIVEFALSE  
CHECKINPUTTRUE PARAMS  
COMPUTE ELASTIC NET PATH WITH COORDINATE DESCENT  
THE ELASTIC NET OPTIMIZATION FUNCTION VARIES FOR MONO AND MULTIOUTPUTS  
FOR MONOOUTPUT TASKS IT IS  
1 2NSAMPLES Y XW22  
ALPHA L1RATIO W1  
05ALPHA1 L1RATIO W22  
FOR MULTIOUTPUT TASKS IT IS  
1 2 NSAMPLES Y XWFRO2  
ALPHA L1RATIO W21  
05ALPHA1 L1RATIO W21  
WHERE  
W21 SUMI SQRTSUMJ WIJ2  
IE THE SUM OF NORM OF EACH ROW  
READ MORE IN THE USER GUIDE  
PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA PASS DIRECTLY AS FORTRAN  
CONTIGUOUS DATA TO AVOID UNNECESSARY MEMORY DUPLICATION IF YIS MONOOUTPUT THEN X  
CAN BE SPARSE  
YNDARRAY SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TARGET VALUES  
L1RATIO FLOAT OPTIONAL FLOAT BETWEEN 0 AND 1 PASSED TO ELASTIC NET SCALING BETWEEN L1 AND  
L2 PENALTIES L1RATIO1 CORRESPONDS TO THE LASSO  
EPS FLOAT LENGTH OF THE PATH EPS1E3 MEANS THAT ALPHAMIN ALPHAMAX  
1E3  
NALPHAS INT OPTIONAL NUMBER OF ALPHAS ALONG THE REGULARIZATION PATH  
ALPHAS NDARRAY OPTIONAL LIST OF ALPHAS WHERE TO COMPUTE THE MODELS IF NONE ALPHAS ARE  
SET AUTOMATICALLY  
PRECOMPUTE TRUE FALSE 'AUTO' ARRAYLIKE WHETHER TO USE A PRECOMPUTED GRAM MATRIX  
TO SPEED UP CALCULATIONS IF SET TO AUTO LET US DECIDE THE GRAM MATRIX CAN ALSO BE  
PASSED AS ARGUMENT  
XYARRAYLIKE OPTIONAL XY NPDOTXT Y THAT CAN BE PRECOMPUTED IT IS USEFUL ONLY  
WHEN THE GRAM MATRIX IS PRECOMPUTED  
COPYX BOOLEAN OPTIONAL DEFAULT TRUE IF TRUE X WILL BE COPIED ELSE IT MAY BE OVER  
WRITTEN  
COEFINIT ARRAY SHAPE NFEATURES NONE THE INITIAL VALUES OF THE COEFFICIENTS  
VERBOSE BOOL OR INTEGER AMOUNT OF VERBOSITY  
RETURNNITER BOOL WHETHER TO RETURN THE NUMBER OF ITERATIONS OR NOT  
POSITIVE BOOL DEFAULT FALSE IF SET TO TRUE FORCES COEFFICIENTS TO BE POSITIVE ONLY ALLOWED  
WHENYNDIM 1  
1910 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

CHECKINPUT BOOL DEFAULT TRUE SKIP INPUT VALIDATION CHECKS INCLUDING THE GRAM MATRIX WHEN PROVIDED ASSUMING THERE ARE HANDLED BY THE CALLER WHEN CHECKINPUTFALSE

PARAMS KWARGS KEYWORD ARGUMENTS PASSED TO THE COORDINATE DESCENT SOLVER

RETURNS

ALPHAS ARRAY SHAPE NALPHAS THE ALPHAS ALONG THE PATH WHERE MODELS ARE COMPUTED

COEFS ARRAY SHAPE NFEATURES NALPHAS OR NOUTPUTS NFEATURES NALPHAS COEFFICIENTS ALONG THE PATH

DUALGAPS ARRAY SHAPE NALPHAS THE DUAL GAPS AT THE END OF THE OPTIMIZATION FOR EACH ALPHA

NITERS ARRAYLIKE SHAPE NALPHAS THE NUMBER OF ITERATIONS TAKEN BY THE COORDINATE DESCENT OPTIMIZER TO REACH THE SPECIFIED TOLERANCE FOR EACH ALPHA IS RETURNED WHEN RETURNNITER IS SET TO TRUE

SEE ALSO

MULTITASKELASTICNET

MULTITASKELASTICNETCV

ELASTICNET

ELASTICNETCV

NOTES

FOR AN EXAMPLE SEE EXAMPLESLINEARMODELPLOTLASSOCOORDINATEDESCENTPATHPY

PREDICTSELF

PREDICT USING THE LINEAR MODEL

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES SAMPLES

RETURNS

CARRAY SHAPE NSAMPLES RETURNS PREDICTED VALUES

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{U}{V}$  WHERE U IS THE RESIDUAL SUM OF SQUARES  $Y_{TRUE} - Y_{PRED}$  2SUM AND V IS THE TOTAL SUM OF SQUARES  $Y_{TRUE} - Y_{TRUEMEAN}$  2SUM THE BEST POSSIBLE SCORE IS 10 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 00

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES FOR SOME ESTIMATORS THIS MAY BE A PRECOMPUTED KERNEL MATRIX INSTEAD SHAPE NSAMPLES NSAMPLESFITTED WHERE NSAMPLESFITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE VALUES FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

6225KLEARNLINEARMODEL GENERALIZED LINEAR MODELS 1911

SCIKITLEARN USER GUIDE RELEASE 0213

SCORE FLOAT R2 OF SELF PREDICT X WRT Y

NOTES

THE R2 SCORE USED WHEN CALLING SCORE ON A REGRESSOR WILL USE MULTIOUTPUTUNIFORMAVERAGE FROM VERSION 023 TO KEEP CONSISTENT WITH METRICSR2SCORE THIS WILL INFLUENCE THE SCORE METHOD OF ALL THE MULTIOUTPUT REGRESSORS EXCEPT FOR MULTIOUTPUTMULTIOUTPUTREGRESSOR TO SPECIFY THE DEFAULT VALUE MANUALLY AND AVOID THE WARNING PLEASE EITHER CALL METRICSR2SCORE DIRECTLY OR MAKE A CUSTOM SCORER WITH METRICSMAKESCORER THE BUILTIN SCORER R2 USES MULTIOUTPUTUNIFORMAVERAGE

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

SPARSECOEF

SPARSE REPRESENTATION OF THE FITTED COEF

62212SKLEARNLINEARMODEL ORTHOGONALMATCHINGPURSUIT

CLASSSKLEARNLINEARMODEL ORTHOGONALMATCHINGPURSUIT NNONZEROCOEFSSNONE

TOLNONE FITINTERCEPTTRUE NOR

MALIZETRUE PRECOMPUTE'AUTO'

ORTHOGONAL MATCHING PURSUIT MODEL OMP

READ MORE IN THE USER GUIDE

PARAMETERS

NNONZEROCOEFSS INT OPTIONAL DESIRED NUMBER OF NONZERO ENTRIES IN THE SOLUTION IF NONE BY DEFAULT THIS VALUE IS SET TO 10 OF NFEATURES

TOLFLOAT OPTIONAL MAXIMUM NORM OF THE RESIDUAL IF NOT NONE OVERRIDES NNONZEROCOEFSS

FITINTERCEPT BOOLEAN OPTIONAL WHETHER TO CALCULATE THE INTERCEPT FOR THIS MODEL IF SET TO FALSE NO INTERCEPT WILL BE USED IN CALCULATIONS EG DATA IS EXPECTED TO BE ALREADY CENTERED

NORMALIZE BOOLEAN OPTIONAL DEFAULT TRUE THIS PARAMETER IS IGNORED WHEN FITINTERCEPT IS SET TO FALSE IF TRUE THE REGRESSORS X WILL BE NORMALIZED BEFORE REGRESSION BY SUBTRACTING THE MEAN AND DIVIDING BY THE L2NORM IF YOU WISH TO STANDARDIZE PLEASE USE SKLEARNPREPROCESSINGSTANDARDSCALER BEFORE CALLING FIT ON AN ESTIMATOR

WITHNORMALIZEFALSE

PRECOMPUTE TRUE FALSE 'AUTO' DEFAULT 'AUTO' WHETHER TO USE A PRECOMPUTED GRAM AND XY MATRIX TO SPEED UP CALCULATIONS IMPROVES PERFORMANCE WHEN NTARGETS OR NSAMPLES IS VERY LARGE NOTE THAT IF YOU ALREADY HAVE SUCH MATRICES YOU CAN PASS THEM DIRECTLY TO THE FIT METHOD

ATTRIBUTES

COEF ARRAY SHAPE NFEATURES OR NTARGETS NFEATURES PARAMETER VECTOR W IN THE FORMULA

INTERCEPT FLOAT OR ARRAY SHAPE NTARGETS INDEPENDENT TERM IN DECISION FUNCTION

1912 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213  
NITER INT OR ARRAYLIKE NUMBER OF ACTIVE FEATURES ACROSS EVERY TARGET  
SEE ALSO  
ORTHOGONALMP  
ORTHOGONALMPGRAM  
LARSPATH  
LARS  
LASSOLARS  
DECOMPOSITIONSPARSEENCODE  
ORTHOGONALMATCHINGPURSUITCV  
NOTES  
ORTHOGONAL MATCHING PURSUIT WAS INTRODUCED IN G MALLAT Z ZHANG MATCHING PURSUITS WITH TIMEFREQUENCY  
DICTIONARIES IEEE TRANSACTIONS ON SIGNAL PROCESSING V OL 41 NO 12 DECEMBER 1993 PP 33973415  
HTTPBLANCHEPOLYTECHNIQUEFRMALLATPAPIERSMALLATPURSUIT93PDF  
THIS IMPLEMENTATION IS BASED ON RUBINSTEIN R ZIBULEVSKY M AND ELAD M EFFICIENT IMPLEMENTATION OF  
THE KSVD ALGORITHM USING BATCH ORTHOGONAL MATCHING PURSUIT TECHNICAL REPORT CS TECHNION APRIL 2008  
HTTPSWWWCSTECHNIONACILRONRUBINPUBLICATIONSKSVDOMPV2PDF  
EXAMPLES  
FROM SKLEARNLINEARMODEL IMPORT ORTHOGONALMATCHINGPURSUIT  
FROM SKLEARNDATASETS IMPORT MAKEREGRESSION  
X Y MAKEREGRESSIONNOISE4 RANDOMSTATE0  
REG ORTHOGONALMATCHINGPURSUITFITX Y  
REGSCOREX Y  
09991  
REGPREDICTX1  
ARRAY783854  
METHODS  
FITSELF X Y FIT THE MODEL USING X Y AS TRAINING DATA  
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR  
PREDICT SELF X PREDICT USING THE LINEAR MODEL  
SCORE SELF X Y SAMPLEWEIGHT RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PRE  
DICTION  
SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR  
INIT SELFNNONZEROCOEFSSNONE TOLNONE FITINTERCEPTTRUE NORMALIZETRUE PRECOM  
PUTE'AUTO'  
FITSELFXY  
FIT THE MODEL USING X Y AS TRAINING DATA  
PARAMETERS  
622SKLEARNLINEARMODEL GENERALIZED LINEAR MODELS 1913

SCIKITLEARN USER GUIDE RELEASE 0213

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NTARGETS TARGET VALUES WILL BE CAST TO X'S DTYPE IF NECESSARY

RETURNS

SELF OBJECT RETURNS AN INSTANCE OF SELF

GETPARAMS SELFDEEPTREE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELF

PREDICT USING THE LINEAR MODEL

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES SAMPLES

RETURNS

CARRAY SHAPE NSAMPLES RETURNS PREDICTED VALUES

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{U}{V}$  WHERE U IS THE RESIDUAL SUM OF SQUARES  $Y_{TRUE} - Y_{PRED}$   $2 \sum$  AND V IS THE TOTAL SUM OF SQUARES  $Y_{TRUE} - Y_{TRUEMEAN}$   $2 \sum$  THE BEST POSSIBLE SCORE IS 1.0 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 0.0

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES FOR SOME ESTIMATORS THIS MAY BE A PRECOMPUTED KERNEL MATRIX INSTEAD SHAPE NSAMPLES NSAMPLESFITTED WHERE NSAMPLESFITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE VALUES FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT R2 OF SELF PREDICT X WRT Y

NOTES

THE R2 SCORE USED WHEN CALLING SCORE ON A REGRESSOR WILL USE MULTIOUTPUTUNIFORMAVERAGE FROM VERSION 0.23 TO KEEP CONSISTENT WITH METRICS R2 SCORE THIS WILL INFLUENCE THE SCORE METHOD OF ALL THE MULTIOUTPUT REGRESSORS EXCEPT FOR MULTIOUTPUTMULTIOUTPUTREGRESSOR TO SPECIFY THE DEFAULT VALUE MANUALLY AND AVOID THE WARNING PLEASE EITHER CALL METRICS R2 SCORE DIRECTLY OR MAKE A CUSTOM SCORER WITH METRICSMAKESCORER THE BUILT IN SCORER R2 USES MULTIOUTPUTUNIFORMAVERAGE

1914 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

EXAMPLES USING SKLEARNLINEARMODELORTHOAGONALMATCHINGPURSUIT

- ORTHOAGONAL MATCHING PURSUIT

62213SKLEARNLINEARMODEL PASSIVEAGGRESSIVECLASSIFIER

CLASSSSKLEARNLINEARMODEL PASSIVEAGGRESSIVECLASSIFIER C10 FITINTERCEPTTRUE

MAXITER1000 TOL0001

EARLYSTOPPINGFALSE

VALIDATIONFRACTION01

NITERNOCHANGE5

SHUFFLETRUE VERBOSE0

LOSS'HINGE' NJOBSNONE

RANDOMSTATENONE

WARMSTARTFALSE

CLASSWEIGHTNONE AVER

AGEFALSE

PASSIVE AGGRESSIVE CLASSIFIER

READ MORE IN THE USER GUIDE

PARAMETERS

CFLOAT MAXIMUM STEP SIZE REGULARIZATION DEFAULTS TO 10

FITINTERCEPT BOOL DEFAULTFALSE WHETHER THE INTERCEPT SHOULD BE ESTIMATED OR NOT IF FALSE THE DATA IS ASSUMED TO BE ALREADY CENTERED

MAXITER INT OPTIONAL DEFAULT1000 THE MAXIMUM NUMBER OF PASSES OVER THE TRAINING DATA AKA EPOCHS IT ONLY IMPACTS THE BEHAVIOR IN THE FIT METHOD AND NOT THE PARTIALFIT

NEW IN VERSION 019

TOLFLOAT OR NONE OPTIONAL DEFAULT1E3 THE STOPPING CRITERION IF IT IS NOT NONE THE ITERATIONS WILL STOP WHEN LOSS PREVIOUSLOSS TOL

NEW IN VERSION 019

EARLYSTOPPING BOOL DEFAULTFALSE WHETHER TO USE EARLY STOPPING TO TERMINATE TRAINING WHEN VALIDATION SCORE IS NOT IMPROVING IF SET TO TRUE IT WILL AUTOMATICALLY SET ASIDE A STRATIFIED FRACTION OF TRAINING DATA AS VALIDATION AND TERMINATE TRAINING WHEN VALIDATION SCORE IS NOT IMPROVING BY AT LEAST TOL FOR NITERNOCHANGE CONSECUTIVE EPOCHS

NEW IN VERSION 020

VALIDATIONFRACTION FLOAT DEFAULT01 THE PROPORTION OF TRAINING DATA TO SET ASIDE AS VALIDATION SET FOR EARLY STOPPING MUST BE BETWEEN 0 AND 1 ONLY USED IF EARLYSTOPPING IS TRUE

NEW IN VERSION 020

622SKLEARNLINEARMODEL GENERALIZED LINEAR MODELS 1915

SCIKITLEARN USER GUIDE RELEASE 0213

NITERNOCHANGE INT DEFAULT5 NUMBER OF ITERATIONS WITH NO IMPROVEMENT TO WAIT BEFORE  
EARLY STOPPING

NEW IN VERSION 020

SHUFFLE BOOL DEFAULTTRUE WHETHER OR NOT THE TRAINING DATA SHOULD BE SHUFFLED AFTER EACH  
EPOCH

VERBOSE INTEGER OPTIONAL THE VERBOSITY LEVEL

LOSS STRING OPTIONAL THE LOSS FUNCTION TO BE USED HINGE EQUIVALENT TO PAI IN THE REFERENCE  
PAPER SQUAREDHINGE EQUIVALENT TO PAII IN THE REFERENCE PAPER

NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF CPUS TO USE TO DO THE OV A ONE  
VERSUS ALL FOR MULTICLASS PROBLEMS COMPUTATION NONE MEANS 1 UNLESS IN A JOBLIB  
PARALLELBACKEND CONTEXT1MEANS USING ALL PROCESSORS SEE GLOSSARY FOR MORE  
DETAILS

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE THE SEED OF THE  
PSEUDO RANDOM NUMBER GENERATOR TO USE WHEN SHUFFLING THE DATA IF INT RANDOMSTATE IS  
THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS  
THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE  
INSTANCE USED BY NPRANDOM

WARMSTART BOOL OPTIONAL WHEN SET TO TRUE REUSE THE SOLUTION OF THE PREVIOUS CALL TO FIT AS  
INITIALIZATION OTHERWISE JUST ERASE THE PREVIOUS SOLUTION SEE THE GLOSSARY  
REPEATEDLY CALLING FIT OR PARTIALFIT WHEN WARMSTART IS TRUE CAN RESULT IN A DIFFERENT SOLUTION  
THAN WHEN CALLING FIT A SINGLE TIME BECAUSE OF THE WAY THE DATA IS SHUFFLED

CLASSWEIGHT DICT CLASSLABEL WEIGHT OR “BALANCED” OR NONE OPTIONAL PRESET FOR THE  
CLASSWEIGHT FIT PARAMETER

WEIGHTS ASSOCIATED WITH CLASSES IF NOT GIVEN ALL CLASSES ARE SUPPOSED TO HAVE WEIGHT ONE  
THE “BALANCED” MODE USES THE VALUES OF Y TO AUTOMATICALLY ADJUST WEIGHTS INVERSELY PRO  
PORTIONAL TO CLASS FREQUENCIES IN THE INPUT DATA AS NSAMPLES NCLASSES NP

BINCOUNTY

NEW IN VERSION 017 PARAMETER CLASSWEIGHT TO AUTOMATICALLY WEIGHT SAMPLES

AVERAGE BOOL OR INT OPTIONAL WHEN SET TO TRUE COMPUTES THE AVERAGED SGD WEIGHTS AND  
STORES THE RESULT IN THE COEF ATTRIBUTE IF SET TO AN INT GREATER THAN 1 AVERAGING WILL BEGIN  
ONCE THE TOTAL NUMBER OF SAMPLES SEEN REACHES AVERAGE SO AVERAGE10 WILL BEGIN AVERAGING  
AFTER SEEING 10 SAMPLES

NEW IN VERSION 019 PARAMETER AVERAGE TO USE WEIGHTS AVERAGING IN SGD

ATTRIBUTES

COEF ARRAY SHAPE 1 NFEATURES IF NCLASSES 2 ELSE NCLASSES NFEATURES WEIGHTS  
ASSIGNED TO THE FEATURES

INTERCEPT ARRAY SHAPE 1 IF NCLASSES 2 ELSE NCLASSES CONSTANTS IN DECISION FUNCTION

NITER INT THE ACTUAL NUMBER OF ITERATIONS TO REACH THE STOPPING CRITERION FOR MULTICLASS FITS  
IT IS THE MAXIMUM OVER EVERY BINARY FIT

SEE ALSO

SGDClassifier

Perceptron

1916 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

REFERENCES

ONLINE PASSIVEAGGRESSIVE ALGORITHMS HTTPJMLRCSAILMITEDUPAPERSVOLUME7CRAMMER06ACRAMMER06A  
PDF K CRAMMER O DEKEL J KESHAT S SHALEVSHWARTZ Y SINGER JMLR 2006

EXAMPLES

```
FROM SKLEARNLINEARMODEL IMPORT PASSIVEAGGRESSIVECLASSIFIER
FROM SKLEARNDATASETS IMPORT MAKECLASSIFICATION
X Y MAKECLASSIFICATIONNFEATURES4 RANDOMSTATE0
CLF PASSIVEAGGRESSIVECLASSIFIERMAXITER1000 RANDOMSTATE0
TOL1E3
CLFFITX Y
PASSIVEAGGRESSIVECLASSIFIERC10 AVERAGEFALSE CLASSWEIGHTNONE
EARLYSTOPPINGFALSE FITINTERCEPTTRUE LOSSHINGE
MAXITER1000 NITERNOCHANGE5 NJOBSNONE
RANDOMSTATE0 SHUFFLETRUE TOL0001
VALIDATIONFRACTION01 VERBOSE0 WARMSTARTFALSE
PRINTCLFCOEF
026642044 045070924 067251877 064185414
PRINTCLFINTERCEPT
184127814
PRINTCLFPREDICT0 0 0 0
1
```

METHODS

DECISIONFUNCTION SELF X PREDICT CONFIDENCE SCORES FOR SAMPLES  
DENSIFY SELF CONVERT COEFFICIENT MATRIX TO DENSE ARRAY FORMAT  
FITSELF X Y COEFINIT INTERCEPTINIT FIT LINEAR MODEL WITH PASSIVE AGGRESSIVE ALGORITHM  
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR  
PARTIALFIT SELF X Y CLASSES FIT LINEAR MODEL WITH PASSIVE AGGRESSIVE ALGORITHM  
PREDICT SELF X PREDICT CLASS LABELS FOR SAMPLES IN X  
SCORE SELF X Y SAMPLEWEIGHT RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS  
SETPARAMS SELF ARGS KWARGS  
SPARSIFY SELF CONVERT COEFFICIENT MATRIX TO SPARSE FORMAT  
INIT SELF C10 FITINTERCEPTTRUE MAXITER1000 TOL0001 EARLYSTOPPINGFALSE VALIDATIONFRACTION01 NITERNOCHANGE5 SHUFFLETRUE VERBOSE0 LOSS'HINGE' NJOBSNONE RANDOMSTATENONE WARMSTARTFALSE CLASSWEIGHTNONE AVERAGEFALSE  
DECISIONFUNCTION SELF X  
PREDICT CONFIDENCE SCORES FOR SAMPLES  
THE CONFIDENCE SCORE FOR A SAMPLE IS THE SIGNED DISTANCE OF THAT SAMPLE TO THE HYPERPLANE  
PARAMETERS  
XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES SAMPLES  
622SKLEARNLINEARMODEL GENERALIZED LINEAR MODELS 1917

SCIKITLEARN USER GUIDE RELEASE 0213

RETURNS

ARRAY SHAPENSAMPLES IF NCLASSES > 2 ELSE NSAMPLES NCLASSES CONFIDENCE SCORES PER SAMPLE CLASS COMBINATION IN THE BINARY CASE CONFIDENCE SCORE FOR SELFCLASSES1 WHERE 0 MEANS THIS CLASS WOULD BE PREDICTED

DENSIFYSELF

CONVERT COEFFICIENT MATRIX TO DENSE ARRAY FORMAT

CONVERTS THE COEF MEMBER BACK TO A NUMPYNDARRAY THIS IS THE DEFAULT FORMAT OF COEF AND IS REQUIRED FOR FITTING SO CALLING THIS METHOD IS ONLY REQUIRED ON MODELS THAT HAVE PREVIOUSLY BEEN SPARSIFIED OTHERWISE IT IS A NOOP

RETURNS

SELF ESTIMATOR

FITSELFXYCOEFINITNONE INTERCEPTINITNONE

FIT LINEAR MODEL WITH PASSIVE AGGRESSIVE ALGORITHM

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES TRAINING DATA

YNUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES

COEFINIT ARRAY SHAPE NCLASSESNFEATURES THE INITIAL COEFFICIENTS TO WARMSTART THE OPTIMIZATION

INTERCEPTINIT ARRAY SHAPE NCLASSES THE INITIAL INTERCEPT TO WARMSTART THE OPTIMIZATION

RETURNS

SELF RETURNS AN INSTANCE OF SELF

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PARTIALFIT SELFXYCLASSESNONE

FIT LINEAR MODEL WITH PASSIVE AGGRESSIVE ALGORITHM

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES SUBSET OF THE TRAINING DATA

YNUMPY ARRAY OF SHAPE NSAMPLES SUBSET OF THE TARGET VALUES

CLASSES ARRAY SHAPE NCLASSES CLASSES ACROSS ALL CALLS TO PARTIALFIT CAN BE OBTAINED BY VIANPUNIQUEYALL WHERE YALL IS THE TARGET VECTOR OF THE ENTIRE DATASET THIS ARGUMENT IS REQUIRED FOR THE FIRST CALL TO PARTIALFIT AND CAN BE OMITTED IN THE SUBSEQUENT CALLS NOTE THAT Y DOESN'T NEED TO CONTAIN ALL LABELS IN CLASSES

RETURNS

SELF RETURNS AN INSTANCE OF SELF

1918 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

PREDICTSELF

PREDICT CLASS LABELS FOR SAMPLES IN X

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES SAMPLES

RETURNS

CARRAY SHAPE NSAMPLES PREDICTED CLASS LABEL PER SAMPLE

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS

IN MULTILABEL CLASSIFICATION THIS IS THE SUBSET ACCURACY WHICH IS A HARSH METRIC SINCE YOU REQUIRE FOR EACH SAMPLE THAT EACH LABEL SET BE CORRECTLY PREDICTED

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE LABELS FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT MEAN ACCURACY OF SELF PREDICTX WRT Y

SPARSIFY SELF

CONVERT COEFFICIENT MATRIX TO SPARSE FORMAT

CONVERTS THE COEF MEMBER TO A SCIPYSPARSE MATRIX WHICH FOR L1REGULARIZED MODELS CAN BE MUCH MORE MEMORY AND STORAGEEFFICIENT THAN THE USUAL NUMPYNDARRAY REPRESENTATION

THEINTERCEPT MEMBER IS NOT CONVERTED

RETURNS

SELF ESTIMATOR

NOTES

FOR NONSPARSE MODELS IE WHEN THERE ARE NOT MANY ZEROS IN COEF THIS MAY ACTUALLY INCREASE MEMORY USAGE SO USE THIS METHOD WITH CARE A RULE OF THUMB IS THAT THE NUMBER OF ZERO ELEMENTS WHICH CAN BE COMPUTED WITH COEF\_0SUM\_ MUST BE MORE THAN 50 FOR THIS TO PROVIDE SIGNIFICANT BENEFITS

AFTER CALLING THIS METHOD FURTHER FITTING WITH THE PARTIALFIT METHOD IF ANY WILL NOT WORK UNTIL YOU CALL DENSIFY

EXAMPLES USING SKLEARNLINEARMODELPASSIVEAGGRESSIVECLASSIFIER

- OUTOFCORE CLASSIFICATION OF TEXT DOCUMENTS
- COMPARING VARIOUS ONLINE SOLVERS
- CLASSIFICATION OF TEXT DOCUMENTS USING SPARSE FEATURES

622SKLEARNLINEARMODEL GENERALIZED LINEAR MODELS 1919

SCIKITLEARN USER GUIDE RELEASE 0213  
62214SKLEARNLINEARMODEL PASSIVEAGGRESSIVEREGRESSOR  
CLASSSSKLEARNLINEARMODEL PASSIVEAGGRESSIVEREGRESSOR C10 FITINTERCEPTTRUE  
MAXITER1000 TOL0001  
EARLYSTOPPINGFALSE  
VALIDATIONFRACTION01  
NITERNOCHANGE5 SHUF  
FLETRUE VERBOSE0  
LOSS'EPSILONINSENSITIVE' EP  
SILON01 RANDOMSTATENONE  
WARMSTARTFALSE AVER  
AGEFALSE  
PASSIVE AGGRESSIVE REGRESSOR  
READ MORE IN THE USER GUIDE  
PARAMETERS  
CFLOAT MAXIMUM STEP SIZE REGULARIZATION DEFAULTS TO 10  
FITINTERCEPT BOOL WHETHER THE INTERCEPT SHOULD BE ESTIMATED OR NOT IF FALSE THE DATA IS  
ASSUMED TO BE ALREADY CENTERED DEFAULTS TO TRUE  
MAXITER INT OPTIONAL DEFAULT1000 THE MAXIMUM NUMBER OF PASSES OVER THE TRAINING DATA  
AKA EPOCHS IT ONLY IMPACTS THE BEHAVIOR IN THE FIT METHOD AND NOT THE PARTIALFIT  
NEW IN VERSION 019  
TOLFLOAT OR NONE OPTIONAL DEFAULT1E3 THE STOPPING CRITERION IF IT IS NOT NONE THE ITERA  
TIONS WILL STOP WHEN LOSS PREVIOUSLOSS TOL  
NEW IN VERSION 019  
EARLYSTOPPING BOOL DEFAULTFALSE WHETHER TO USE EARLY STOPPING TO TERMINATE TRAINING WHEN  
VALIDATION SCORE IS NOT IMPROVING IF SET TO TRUE IT WILL AUTOMATICALLY SET ASIDE A FRACTION OF  
TRAINING DATA AS VALIDATION AND TERMINATE TRAINING WHEN VALIDATION SCORE IS NOT IMPROVING BY  
AT LEAST TOL FOR NITERNOCHANGE CONSECUTIVE EPOCHS  
NEW IN VERSION 020  
VALIDATIONFRACTION FLOAT DEFAULT01 THE PROPORTION OF TRAINING DATA TO SET ASIDE AS VALIDATION  
SET FOR EARLY STOPPING MUST BE BETWEEN 0 AND 1 ONLY USED IF EARLYSTOPPING IS TRUE  
NEW IN VERSION 020  
NITERNOCHANGE INT DEFAULT5 NUMBER OF ITERATIONS WITH NO IMPROVEMENT TO WAIT BEFORE  
EARLY STOPPING  
NEW IN VERSION 020  
SHUFFLE BOOL DEFAULTTRUE WHETHER OR NOT THE TRAINING DATA SHOULD BE SHUFFLED AFTER EACH  
EPOCH  
VERBOSE INTEGER OPTIONAL THE VERBOSITY LEVEL  
LOSS STRING OPTIONAL THE LOSS FUNCTION TO BE USED EPSILONINSENSITIVE EQUIVALENT TO PAI IN  
THE REFERENCE PAPER SQUAREDEPSILONINSENSITIVE EQUIVALENT TO PAII IN THE REFERENCE PAPER  
EPSILON FLOAT IF THE DIFFERENCE BETWEEN THE CURRENT PREDICTION AND THE CORRECT LABEL IS BELOW  
THIS THRESHOLD THE MODEL IS NOT UPDATED  
RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE THE SEED OF THE  
PSEUDO RANDOM NUMBER GENERATOR TO USE WHEN SHUFFLING THE DATA IF INT RANDOMSTATE IS  
THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS  
1920 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE  
INSTANCE USED BY NPRANDOM

WARMSTART BOOL OPTIONAL WHEN SET TO TRUE REUSE THE SOLUTION OF THE PREVIOUS CALL TO FIT AS  
INITIALIZATION OTHERWISE JUST ERASE THE PREVIOUS SOLUTION SEE THE GLOSSARY

REPEATEDLY CALLING FIT OR PARTIALFIT WHEN WARMSTART IS TRUE CAN RESULT IN A DIFFERENT SOLUTION  
THAN WHEN CALLING FIT A SINGLE TIME BECAUSE OF THE WAY THE DATA IS SHUFFLED

AVERAGE BOOL OR INT OPTIONAL WHEN SET TO TRUE COMPUTES THE AVERAGED SGD WEIGHTS AND  
STORES THE RESULT IN THE COEF ATTRIBUTE IF SET TO AN INT GREATER THAN 1 AVERAGING WILL BEGIN  
ONCE THE TOTAL NUMBER OF SAMPLES SEEN REACHES AVERAGE SO AVERAGE10 WILL BEGIN AVERAGING  
AFTER SEEING 10 SAMPLES

NEW IN VERSION 019 PARAMETER AVERAGE TO USE WEIGHTS AVERAGING IN SGD

ATTRIBUTES

COEF ARRAY SHAPE 1 NFEATURES IF NCLASSES 2 ELSE NCLASSES NFEATURES WEIGHTS  
ASSIGNED TO THE FEATURES

INTERCEPT ARRAY SHAPE 1 IF NCLASSES 2 ELSE NCLASSES CONSTANTS IN DECISION FUNCTION

NITER INT THE ACTUAL NUMBER OF ITERATIONS TO REACH THE STOPPING CRITERION

SEE ALSO

SGDREGRESSOR

REFERENCES

ONLINE PASSIVEAGGRESSIVE ALGORITHMS [HTTPJMLRCSAILMITEDUPAPERSVOLUME7CRAMMER06ACRAMMER06A](http://jmlr.csail.mit.edu/papers/volume7/crammer06a/crammer06a.pdf)  
PDF K CRAMMER O DEKEL J KESHAT S SHALEVSHWARTZ Y SINGER JMLR 2006

EXAMPLES

```
FROM SKLEARNLINEARMODEL IMPORT PASSIVEAGGRESSIVEREGRESSOR
FROM SKLEARNDATASETS IMPORT MAKEREGRESSION
X Y MAKEREGRESSIONNFEATURES4 RANDOMSTATE0
REGR PASSIVEAGGRESSIVEREGRESSORMAXITER100 RANDOMSTATE0
TOL1E3
REGRFITX Y
PASSIVEAGGRESSIVEREGRESSORC10 AVERAGEFALSE EARLYSTOPPINGFALSE
EPSILON01 FITINTERCEPTTRUE LOSSEPSILONINSENSITIVE
MAXITER100 NITERNOCHANGE5 RANDOMSTATE0
SHUFFLETRUE TOL0001 VALIDATIONFRACTION01
VERBOSE0 WARMSTARTFALSE
PRINTREGRCOEF
2048736655 3418818427 6759122734 8794731329
PRINTREGRINTERCEPT
002306214
PRINTREGRPREDICT0 0 0 0
002306214
622SKLEARNLINEARMODEL GENERALIZED LINEAR MODELS 1921
```

SCIKITLEARN USER GUIDE RELEASE 0213

METHODS

DENSIFY SELF CONVERT COEFFICIENT MATRIX TO DENSE ARRAY FORMAT

FITSELF X Y COEFINIT INTERCEPTINIT FIT LINEAR MODEL WITH PASSIVE AGGRESSIVE ALGORITHM

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PARTIALFIT SELF X Y FIT LINEAR MODEL WITH PASSIVE AGGRESSIVE ALGORITHM

PREDICT SELF X PREDICT USING THE LINEAR MODEL

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

SETPARAMS SELF ARGS KWARGS

SPARSIFY SELF CONVERT COEFFICIENT MATRIX TO SPARSE FORMAT

INIT SELF C10 FITINTERCEPTTRUE MAXITER1000 TOL0001 EARLYSTOPPINGFALSE

VALIDATIONFRACTION01 NITERNOCHANGE5 SHUFFLETRUE VERBOSE0

LOSS'EPSILONINSENSITIVE' EPSILON01 RANDOMSTATENONE WARMSTARTFALSE AVERAGEFALSE

DENSIFYSELF

CONVERT COEFFICIENT MATRIX TO DENSE ARRAY FORMAT

CONVERTS THE COEF MEMBER BACK TO A NUMPYNDARRAY THIS IS THE DEFAULT FORMAT OF COEF AND IS REQUIRED FOR FITTING SO CALLING THIS METHOD IS ONLY REQUIRED ON MODELS THAT HAVE PREVIOUSLY BEEN SPARSIFIED OTHERWISE IT IS A NOOP

RETURNS

SELF ESTIMATOR

FITSELFXYCOEFINITNONE INTERCEPTINITNONE

FIT LINEAR MODEL WITH PASSIVE AGGRESSIVE ALGORITHM

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES TRAINING DATA

Y NUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES

COEFINIT ARRAY SHAPE NFEATURES THE INITIAL COEFFICIENTS TO WARMSTART THE OPTIMIZATION

INTERCEPTINIT ARRAY SHAPE 1 THE INITIAL INTERCEPT TO WARMSTART THE OPTIMIZATION

RETURNS

SELF RETURNS AN INSTANCE OF SELF

GETPARAMS SELFDEEPTTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PARTIALFIT SELFXY

FIT LINEAR MODEL WITH PASSIVE AGGRESSIVE ALGORITHM

PARAMETERS

1922 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES SUBSET OF TRAINING DATA

YNUMPY ARRAY OF SHAPE NSAMPLES SUBSET OF TARGET VALUES

RETURNS

SELF RETURNS AN INSTANCE OF SELF

PREDICTSELF

PREDICT USING THE LINEAR MODEL

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES

RETURNS

ARRAY SHAPE NSAMPLES PREDICTED TARGET VALUES PER ELEMENT IN X

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{U}{V}$  WHERE U IS THE RESIDUAL SUM OF SQUARES  $Y_{TRUE} - Y_{PRED}$

$2 \sum$  AND V IS THE TOTAL SUM OF SQUARES  $Y_{TRUE} - Y_{TRUEMEAN}$   $2 \sum$  THE BEST POSSIBLE SCORE

IS 10 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS

PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 00

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES FOR SOME ESTIMATORS THIS MAY

BE A PRECOMPUTED KERNEL MATRIX INSTEAD SHAPE NSAMPLES NSAMPLESFITTED WHERE

NSAMPLESFITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE VALUES FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT R2 OF SELF PREDICTX WRT Y

NOTES

THE R2 SCORE USED WHEN CALLING SCORE ON A REGRESSOR WILL USE MULTIOUTPUTUNIFORMAVERAGE

FROM VERSION 023 TO KEEP CONSISTENT WITH METRICSR2SCORE THIS WILL INFLUENCE THE SCORE

METHOD OF ALL THE MULTIOUTPUT REGRESSORS EXCEPT FOR MULTIOUTPUTMULTIOUTPUTREGRESSOR

TO SPECIFY THE DEFAULT VALUE MANUALLY AND AVOID THE WARNING PLEASE EITHER CALL METRICSR2SCORE

DIRECTLY OR MAKE A CUSTOM SCORER WITH METRICSMAKESCORER THE BUILTIN SCORER R2 USES

MULTIOUTPUTUNIFORMAVERAGE

SPARSIFY SELF

CONVERT COEFFICIENT MATRIX TO SPARSE FORMAT

CONVERTS THE COEF MEMBER TO A SCIPYSPARSE MATRIX WHICH FOR L1REGULARIZED MODELS CAN BE MUCH MORE

MEMORY AND STORAGEEFFICIENT THAN THE USUAL NUMPYNDARRAY REPRESENTATION

THEINTERCEPT MEMBER IS NOT CONVERTED

RETURNS

SELF ESTIMATOR

622SKLEARNLINEARMODEL GENERALIZED LINEAR MODELS 1923

SCIKITLEARN USER GUIDE RELEASE 0213

NOTES

FOR NONSPARSE MODELS IE WHEN THERE ARE NOT MANY ZEROS IN COEF THIS MAY ACTUALLY INCREASE MEMORY USAGE SO USE THIS METHOD WITH CARE A RULE OF THUMB IS THAT THE NUMBER OF ZERO ELEMENTS WHICH CAN BE COMPUTED WITH COEF\_0SUM\_ MUST BE MORE THAN 50 FOR THIS TO PROVIDE SIGNIFICANT BENEFITS AFTER CALLING THIS METHOD FURTHER FITTING WITH THE PARTIALFIT METHOD IF ANY WILL NOT WORK UNTIL YOU CALL DENSIFY

62215SKLEARNLINEARMODEL PERCEPTRON  
CLASSSSKLEARNLINEARMODEL PERCEPTRON PENALTYNONE ALPHA00001 FITINTERCEPTTRUE  
MAXITER1000 TOL0001 SHUFFLETRUE VER  
BOSE0 ETA010 NJOBSNONE RANDOMSTATE0  
EARLYSTOPPINGFALSE VALIDATIONFRACTION01  
NITERNOCHANGE5 CLASSWEIGHTNONE  
WARMSTARTFALSE

READ MORE IN THE USER GUIDE

PARAMETERS

PENALTY NONE 'L2' OR 'L1' OR 'ELASTICNET' THE PENALTY AKA REGULARIZATION TERM TO BE USED  
DEFAULTS TO NONE

ALPHA FLOAT CONSTANT THAT MULTIPLIES THE REGULARIZATION TERM IF REGULARIZATION IS USED DEFAULTS  
TO 00001

FITINTERCEPT BOOL WHETHER THE INTERCEPT SHOULD BE ESTIMATED OR NOT IF FALSE THE DATA IS  
ASSUMED TO BE ALREADY CENTERED DEFAULTS TO TRUE

MAXITER INT OPTIONAL DEFAULT1000 THE MAXIMUM NUMBER OF PASSES OVER THE TRAINING DATA  
AKA EPOCHS IT ONLY IMPACTS THE BEHAVIOR IN THE FIT METHOD AND NOT THE PARTIALFIT  
NEW IN VERSION 019

TOLFLOAT OR NONE OPTIONAL DEFAULT1E3 THE STOPPING CRITERION IF IT IS NOT NONE THE ITERA  
TIONS WILL STOP WHEN LOSS\_ PREVIOUSLOSS\_ TOL  
NEW IN VERSION 019

SHUFFLE BOOL OPTIONAL DEFAULT TRUE WHETHER OR NOT THE TRAINING DATA SHOULD BE SHUFFLED AFTER  
EACH EPOCH

VERBOSE INTEGER OPTIONAL THE VERBOSITY LEVEL

ETA0 DOUBLE CONSTANT BY WHICH THE UPDATES ARE MULTIPLIED DEFAULTS TO 1

NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF CPUS TO USE TO DO THE OV A ONE

VERSUS ALL FOR MULTICLASS PROBLEMS COMPUTATION NONE MEANS 1 UNLESS IN A JOBLIB

PARALLELBACKEND CONTEXT1MEANS USING ALL PROCESSORS SEE GLOSSARY FOR MORE

DETAILS

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULT NONE THE SEED OF THE  
PSEUDO RANDOM NUMBER GENERATOR TO USE WHEN SHUFFLING THE DATA IF INT RANDOMSTATE IS  
THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS  
THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE  
INSTANCE USED BY NPRANDOM

EARLYSTOPPING BOOL DEFAULTFALSE WHETHER TO USE EARLY STOPPING TO TERMINATE TRAINING WHEN  
VALIDATION SCORE IS NOT IMPROVING IF SET TO TRUE IT WILL AUTOMATICALLY SET ASIDE A STRATIFIED

1924 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

FRACTION OF TRAINING DATA AS VALIDATION AND TERMINATE TRAINING WHEN VALIDATION SCORE IS NOT IMPROVING BY AT LEAST TOL FOR NITERNOCHANGE CONSECUTIVE EPOCHS

NEW IN VERSION 020

VALIDATIONFRACTION FLOAT DEFAULT01 THE PROPORTION OF TRAINING DATA TO SET ASIDE AS VALIDATION SET FOR EARLY STOPPING MUST BE BETWEEN 0 AND 1 ONLY USED IF EARLYSTOPPING IS TRUE

NEW IN VERSION 020

NITERNOCHANGE INT DEFAULT5 NUMBER OF ITERATIONS WITH NO IMPROVEMENT TO WAIT BEFORE EARLY STOPPING

NEW IN VERSION 020

CLASSWEIGHT DICT CLASSLABEL WEIGHT OR “BALANCED” OR NONE OPTIONAL PRESET FOR THE CLASSWEIGHT FIT PARAMETER

WEIGHTS ASSOCIATED WITH CLASSES IF NOT GIVEN ALL CLASSES ARE SUPPOSED TO HAVE WEIGHT ONE THE “BALANCED” MODE USES THE VALUES OF Y TO AUTOMATICALLY ADJUST WEIGHTS INVERSELY PROPORTIONAL TO CLASS FREQUENCIES IN THE INPUT DATA AS NSAMPLES NCLASSES NP

BINCOUNTY

WARMSTART BOOL OPTIONAL WHEN SET TO TRUE REUSE THE SOLUTION OF THE PREVIOUS CALL TO FIT AS INITIALIZATION OTHERWISE JUST ERASE THE PREVIOUS SOLUTION SEE THE GLOSSARY

ATTRIBUTES

COEF ARRAY SHAPE 1 NFEATURES IF NCLASSES 2 ELSE NCLASSES NFEATURES WEIGHTS ASSIGNED TO THE FEATURES

INTERCEPT ARRAY SHAPE 1 IF NCLASSES 2 ELSE NCLASSES CONSTANTS IN DECISION FUNCTION

NITER INT THE ACTUAL NUMBER OF ITERATIONS TO REACH THE STOPPING CRITERION FOR MULTICLASS FITS IT IS THE MAXIMUM OVER EVERY BINARY FIT

SEE ALSO

SGDCCLASSIFIER

NOTES

PERCEPTRON IS A CLASSIFICATION ALGORITHM WHICH SHARES THE SAME UNDERLYING IMPLEMENTATION WITH SGDCCLASSIFIER IN FACTPERCEPTRON IS EQUIVALENT TO SGDCCLASSIFIERLOSSPERCEPTRON

ETA01 LEARNINGRATECONSTANT PENALTYNONE

REFERENCES

HTTPSENWIKIPEDIAORGWIKIPERCEPTRON AND REFERENCES THEREIN

EXAMPLES

FROM SKLEARNDATASETS IMPORT LOADDIGITS

FROM SKLEARNLINEARMODEL IMPORT PERCEPTRON

X Y LOADDIGITSRETURNXY TRUE

CLF PERCEPTRONTOL1E3 RANDOMSTATE0

CLFFITX Y

622SKLEARNLINEARMODEL GENERALIZED LINEAR MODELS 1925

SCIKITLEARN USER GUIDE RELEASE 0213  
PERCEPTRONALPHA00001 CLASSWEIGHTNONE EARLYSTOPPINGFALSE ETA010  
FITINTERCEPTTRUE MAXITER1000 NITERNOCHANGE5 NJOBSNONE  
PENALTYNONE RANDOMSTATE0 SHUFFLETRUE TOL0001  
VALIDATIONFRACTION01 VERBOSE0 WARMSTARTFALSE  
CLFSCOREX Y  
0939  
METHODS  
DECISIONFUNCTION SELF X PREDICT CONFIDENCE SCORES FOR SAMPLES  
DENSIFY SELF CONVERT COEFFICIENT MATRIX TO DENSE ARRAY FORMAT  
FITSELF X Y COEFINIT INTERCEPTINIT FIT LINEAR MODEL WITH STOCHASTIC GRADIENT DESCENT  
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR  
PARTIALFIT SELF X Y CLASSES SAMPLEWEIGHT PERFORM ONE EPOCH OF STOCHASTIC GRADIENT DESCENT ON  
GIVEN SAMPLES  
PREDICT SELF X PREDICT CLASS LABELS FOR SAMPLES IN X  
SCORE SELF X Y SAMPLEWEIGHT RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND  
LABELS  
SETPARAMS SELF ARGS KWARGS  
SPARSIFY SELF CONVERT COEFFICIENT MATRIX TO SPARSE FORMAT  
INIT SELFPENALTYNONE ALPHA00001 FITINTERCEPTTRUE MAXITER1000 TOL0001 SHUF  
FLETRUE VERBOSE0 ETA010 NJOBSNONE RANDOMSTATE0 EARLYSTOPPINGFALSE VAL  
IDATIONFRACTION01 NITERNOCHANGE5 CLASSWEIGHTNONE WARMSTARTFALSE  
DECISIONFUNCTION SELF X  
PREDICT CONFIDENCE SCORES FOR SAMPLES  
THE CONFIDENCE SCORE FOR A SAMPLE IS THE SIGNED DISTANCE OF THAT SAMPLE TO THE HYPERPLANE  
PARAMETERS  
XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES SAMPLES  
RETURNS  
ARRAY SHAPENSAMPLES IF NCLASSES > 2 ELSE NSAMPLES NCLASSES CONFIDENCE  
SCORES PER SAMPLE CLASS COMBINATION IN THE BINARY CASE CONFIDENCE SCORE FOR  
SELFCLASSES1 WHERE 0 MEANS THIS CLASS WOULD BE PREDICTED  
DENSIFYSELF  
CONVERT COEFFICIENT MATRIX TO DENSE ARRAY FORMAT  
CONVERTS THE COEF MEMBER BACK TO A NUMPYNDARRAY THIS IS THE DEFAULT FORMAT OF COEF AND IS  
REQUIRED FOR FITTING SO CALLING THIS METHOD IS ONLY REQUIRED ON MODELS THAT HAVE PREVIOUSLY BEEN SPARSIFIED  
OTHERWISE IT IS A NOOP  
RETURNS  
SELF ESTIMATOR  
FITSELFXYCOEFINITNONE INTERCEPTINITNONE SAMPLEWEIGHTNONE  
FIT LINEAR MODEL WITH STOCHASTIC GRADIENT DESCENT  
PARAMETERS  
XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES TRAINING DATA  
1926 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

YNUMPY ARRAY SHAPE NSAMPLES TARGET VALUES

COEFINIT ARRAY SHAPE NCLASSES NFEATURES THE INITIAL COEFFICIENTS TO WARMSTART THE OPTIMIZATION

INTERCEPTINIT ARRAY SHAPE NCLASSES THE INITIAL INTERCEPT TO WARMSTART THE OPTIMIZATION

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL WEIGHTS APPLIED TO INDIVIDUAL SAMPLES IF NOT PROVIDED UNIFORM WEIGHTS ARE ASSUMED THESE WEIGHTS WILL BE MULTIPLIED WITH CLASSWEIGHT PASSED THROUGH THE CONSTRUCTOR IF CLASSWEIGHT IS SPECIFIED

RETURNS

SELF RETURNS AN INSTANCE OF SELF

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PARTIALFIT SELFXYCLASSESNONE SAMPLEWEIGHTNONE

PERFORM ONE EPOCH OF STOCHASTIC GRADIENT DESCENT ON GIVEN SAMPLES

INTERNALLY THIS METHOD USES MAXITER 1 THEREFORE IT IS NOT GUARANTEED THAT A MINIMUM OF THE COST FUNCTION IS REACHED AFTER CALLING IT ONCE MATTERS SUCH AS OBJECTIVE CONVERGENCE AND EARLY STOPPING SHOULD BE HANDLED BY THE USER

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES SUBSET OF THE TRAINING DATA

YNUMPY ARRAY SHAPE NSAMPLES SUBSET OF THE TARGET VALUES

CLASSES ARRAY SHAPE NCLASSES CLASSES ACROSS ALL CALLS TO PARTIALFIT CAN BE OBTAINED BY VIANPUNIQUEYALL WHERE YALL IS THE TARGET VECTOR OF THE ENTIRE DATASET THIS ARGUMENT IS REQUIRED FOR THE FIRST CALL TO PARTIALFIT AND CAN BE OMITTED IN THE SUBSEQUENT CALLS NOTE THAT Y DOESN'T NEED TO CONTAIN ALL LABELS IN CLASSES

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL WEIGHTS APPLIED TO INDIVIDUAL SAMPLES IF NOT PROVIDED UNIFORM WEIGHTS ARE ASSUMED

RETURNS

SELF RETURNS AN INSTANCE OF SELF

PREDICTSELF

PREDICT CLASS LABELS FOR SAMPLES IN X

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES SAMPLES

RETURNS

CARRAY SHAPE NSAMPLES PREDICTED CLASS LABEL PER SAMPLE

622SKLEARNLINEARMODEL GENERALIZED LINEAR MODELS 1927

SCIKITLEARN USER GUIDE RELEASE 0213

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS

IN MULTILABEL CLASSIFICATION THIS IS THE SUBSET ACCURACY WHICH IS A HARSH METRIC SINCE YOU REQUIRE FOR EACH SAMPLE THAT EACH LABEL SET BE CORRECTLY PREDICTED

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE LABELS FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT MEAN ACCURACY OF SELF PREDICTX WRT Y

SPARSIFY SELF

CONVERT COEFFICIENT MATRIX TO SPARSE FORMAT

CONVERTS THE COEF MEMBER TO A SCIPYSPARSE MATRIX WHICH FOR L1REGULARIZED MODELS CAN BE MUCH MORE MEMORY AND STORAGEEFFICIENT THAN THE USUAL NUMPYNDARRAY REPRESENTATION

THEINTERCEPT MEMBER IS NOT CONVERTED

RETURNS

SELF ESTIMATOR

NOTES

FOR NONSPARSE MODELS IE WHEN THERE ARE NOT MANY ZEROS IN COEF THIS MAY ACTUALLY INCREASE MEMORY USAGE SO USE THIS METHOD WITH CARE A RULE OF THUMB IS THAT THE NUMBER OF ZERO ELEMENTS WHICH CAN BE COMPUTED WITH COEF 0SUM MUST BE MORE THAN 50 FOR THIS TO PROVIDE SIGNIFICANT BENEFITS

AFTER CALLING THIS METHOD FURTHER FITTING WITH THE PARTIALFIT METHOD IF ANY WILL NOT WORK UNTIL YOU CALL DENSIFY

EXAMPLES USING SKLEARNLINEARMODELPERCEPTRON

- OUTOFCORE CLASSIFICATION OF TEXT DOCUMENTS
- COMPARING VARIOUS ONLINE SOLVERS
- CLASSIFICATION OF TEXT DOCUMENTS USING SPARSE FEATURES

62216SKLEARNLINEARMODEL RANSACREGRESSOR

CLASSSSKLEARNLINEARMODEL RANSACREGRESSOR BASEESTIMATORNONE MINSAMPLESNONE

RESIDUALTHRESHOLDNONE ISDATAVALIDNONE

ISMODELVALIDNONE MAXTRIALS100

MAXSKIPSINF STOPNINLIERSINF

STOPSCOREINF STOPPROBABILITY099

LOSS'ABSOLUTELOSS' RANDOMSTATENONE

RANSAC RANDOM SAMPLE CONSENSUS ALGORITHM

RANSAC IS AN ITERATIVE ALGORITHM FOR THE ROBUST ESTIMATION OF PARAMETERS FROM A SUBSET OF INLIERS FROM THE COMPLETE DATA SET MORE INFORMATION CAN BE FOUND IN THE GENERAL DOCUMENTATION OF LINEAR MODELS

1928 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

A DETAILED DESCRIPTION OF THE ALGORITHM CAN BE FOUND IN THE DOCUMENTATION OF THE LINEARMODEL SUBPACKAGE  
READ MORE IN THE USER GUIDE

PARAMETERS

BASEESTIMATOR OBJECT OPTIONAL BASE ESTIMATOR OBJECT WHICH IMPLEMENTS THE FOLLOWING METHODS

- FITX Y FIT MODEL TO GIVEN TRAINING DATA AND TARGET VALUES
- SCOREX Y RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA WHICH IS USED FOR THE STOP CRITERION DEFINED BY STOPSCORE ADDITIONALLY THE SCORE IS USED TO DECIDE WHICH OF TWO EQUALLY LARGE CONSENSUS SETS IS CHOSEN AS THE BETTER ONE
- PREDICTX RETURNS PREDICTED VALUES USING THE LINEAR MODEL WHICH IS USED TO COMPUTE RESIDUAL ERROR USING LOSS FUNCTION

IFBASEESTIMATOR IS NONE THEN BASEESTIMATORSKLEARNLINEARMODEL LINEARREGRESSION IS USED FOR TARGET VALUES OF DTYPE FLOAT

NOTE THAT THE CURRENT IMPLEMENTATION ONLY SUPPORTS REGRESSION ESTIMATORS

MINSAMPLES INT 1 OR FLOAT 0 1 OPTIONAL MINIMUM NUMBER OF SAMPLES CHOSEN RANDOMLY FROM ORIGINAL DATA TREATED AS AN ABSOLUTE NUMBER OF SAMPLES FOR MINSAMPLES 1 TREATED AS A RELATIVE NUMBER CEILMINSAMPLES X SHAPE0 FORMINSAMPLES 1 THIS IS TYPICALLY CHOSEN AS THE MINIMAL NUMBER OF SAMPLES NECESSARY TO ESTIMATE THE GIVEN BASEESTIMATOR BY DEFAULT A SKLEARN LINEARMODELLINEARREGRESSION ESTIMATOR IS ASSUMED AND MINSAMPLES IS CHOSEN ASXSHAPE1 1

RESIDUALTHRESHOLD FLOAT OPTIONAL MAXIMUM RESIDUAL FOR A DATA SAMPLE TO BE CLASSIFIED AS AN INLIER BY DEFAULT THE THRESHOLD IS CHOSEN AS THE MAD MEDIAN ABSOLUTE DEVIATION OF THE TARGET VALUES Y

ISDATAVALID CALLABLE OPTIONAL THIS FUNCTION IS CALLED WITH THE RANDOMLY SELECTED DATA BEFORE THE MODEL IS FITTED TO IT ISDATAVALIDX Y IF ITS RETURN VALUE IS FALSE THE CURRENT RANDOMLY CHOSEN SUBSAMPLE IS SKIPPED

ISMODELVALID CALLABLE OPTIONAL THIS FUNCTION IS CALLED WITH THE ESTIMATED MODEL AND THE RANDOMLY SELECTED DATA ISMODELVALIDMODEL X Y IF ITS RETURN VALUE IS FALSE THE CURRENT RANDOMLY CHOSEN SUBSAMPLE IS SKIPPED REJECTING SAMPLES WITH THIS FUNCTION IS COMPUTATIONALLY COSTLIER THAN WITH ISDATAVALID ISMODELVALID SHOULD THEREFORE ONLY BE USED IF THE ESTIMATED MODEL IS NEEDED FOR MAKING THE REJECTION DECISION

MAXTRIALS INT OPTIONAL MAXIMUM NUMBER OF ITERATIONS FOR RANDOM SAMPLE SELECTION

MAXSKIPS INT OPTIONAL MAXIMUM NUMBER OF ITERATIONS THAT CAN BE SKIPPED DUE TO FINDING ZERO INLIERS OR INVALID DATA DEFINED BY ISDATAVALID OR INVALID MODELS DEFINED BY ISMODELVALID

NEW IN VERSION 019

STOPNINLIERS INT OPTIONAL STOP ITERATION IF AT LEAST THIS NUMBER OF INLIERS ARE FOUND

STOPSCORE FLOAT OPTIONAL STOP ITERATION IF SCORE IS GREATER EQUAL THAN THIS THRESHOLD

STOPPROBABILITY FLOAT IN RANGE 0 1 OPTIONAL RANSAC ITERATION STOPS IF AT LEAST ONE OUTLIER FREE SET OF THE TRAINING DATA IS SAMPLED IN RANSAC THIS REQUIRES TO GENERATE AT LEAST N SAMPLES ITERATIONS

N LOG1 PROBABILITY LOG1 E M

622SKLEARNLINEARMODEL GENERALIZED LINEAR MODELS 1929

SCIKITLEARN USER GUIDE RELEASE 0213

WHERE THE PROBABILITY CONFIDENCE IS TYPICALLY SET TO HIGH VALUE SUCH AS 0.99 THE DEFAULT  
AND E IS THE CURRENT FRACTION OF INLIERS WRT THE TOTAL NUMBER OF SAMPLES  
LOSS STRING CALLABLE OPTIONAL DEFAULT "ABSOLUTELOSS" STRING INPUTS "ABSOLUTELOSS" AND  
"SQUAREDLOSS" ARE SUPPORTED WHICH FIND THE ABSOLUTE LOSS AND SQUARED LOSS PER SAMPLE RE  
SPECTIVELY  
IFLOSS IS A CALLABLE THEN IT SHOULD BE A FUNCTION THAT TAKES TWO ARRAYS AS INPUTS THE TRUE  
AND PREDICTED VALUE AND RETURNS A 1D ARRAY WITH THE ITH VALUE OF THE ARRAY CORRESPONDING TO  
THE LOSS ON XI  
IF THE LOSS ON A SAMPLE IS GREATER THAN THE RESIDUALTHRESHOLD THEN THIS SAMPLE IS  
CLASSIFIED AS AN OUTLIER  
RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULT NONE THE GENERATOR USED  
TO INITIALIZE THE CENTERS IF INT RANDOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENER  
ATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE  
RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NPRANDOM  
ATTRIBUTES  
ESTIMATOR OBJECT BEST FITTED MODEL COPY OF THE BASEESTIMATOR OBJECT  
NTRIALS INT NUMBER OF RANDOM SELECTION TRIALS UNTIL ONE OF THE STOP CRITERIA IS MET IT IS  
ALWAYS MAXTRIALS  
INLIERMASK BOOL ARRAY OF SHAPE NSAMPLES BOOLEAN MASK OF INLIERS CLASSIFIED AS TRUE  
NSKIPSNINLIERS INT NUMBER OF ITERATIONS SKIPPED DUE TO FINDING ZERO INLIERS  
NEW IN VERSION 0.19  
NSKIPSINVALIDDATA INT NUMBER OF ITERATIONS SKIPPED DUE TO INVALID DATA DEFINED BY  
ISDATAVALID  
NEW IN VERSION 0.19  
NSKIPSINVALIDMODEL INT NUMBER OF ITERATIONS SKIPPED DUE TO AN INVALID MODEL DEFINED BY  
ISMODELVALID  
NEW IN VERSION 0.19  
REFERENCES  
R80CE5B25CF9D1 R80CE5B25CF9D2 R80CE5B25CF9D3  
EXAMPLES  
FROM SKLEARNLINEARMODEL IMPORT RANSACREGRESSOR  
FROM SKLEARNDATASETS IMPORT MAKEREGRESSION  
X Y MAKEREGRESSION  
NSAMPLES200 NFEATURES2 NOISE40 RANDOMSTATE0  
REG RANSACREGRESSORRANDOMSTATE0FITX Y  
REGSCOREX Y  
0.9885  
REGPREDICTX1  
ARRAY319417  
1930 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

METHODS

FITSELF X Y SAMPLEWEIGHT FIT ESTIMATOR USING RANSAC ALGORITHM

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PREDICT SELF X PREDICT USING THE ESTIMATED MODEL

SCORE SELF X Y RETURNS THE SCORE OF THE PREDICTION

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELF BASEESTIMATORNONE MINSAMPLESNONE RESIDUALTHRESHOLDNONE

ISDATAVALIDNONE ISMODELVALIDNONE MAXTRIALS100 MAXSKIPSINF

STOPNINLIERSINF STOPSCOREINF STOPPROBABILITY099 LOSS'ABSOLUTELOSS' RAN

DOMSTATENONE

FITSELFXYSAMPLEWEIGHTNONE

FIT ESTIMATOR USING RANSAC ALGORITHM

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES TRAINING DATA

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NTARGETS TARGET VALUES

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES INDIVIDUAL WEIGHTS FOR EACH SAMPLE RAISES

ERROR IF SAMPLEWEIGHT IS PASSED AND BASEESTIMATOR FIT METHOD DOES NOT SUPPORT IT

RAISES

VALUEERROR IF NO VALID CONSENSUS SET COULD BE FOUND THIS OCCURS IF ISDATAVALID AND

ISMODELVALID RETURN FALSE FOR ALL MAXTRIALS RANDOMLY CHOSEN SUBSAMPLES

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELF

PREDICT USING THE ESTIMATED MODEL

THIS IS A WRAPPER FOR ESTIMATORPREDICTX

PARAMETERS

XNUMPY ARRAY OF SHAPE NSAMPLES NFEATURES

RETURNS

YARRAY SHAPE NSAMPLES OR NSAMPLES NTARGETS RETURNS PREDICTED VALUES

SCORESELFXY

RETURNS THE SCORE OF THE PREDICTION

THIS IS A WRAPPER FOR ESTIMATORSCOREX Y

PARAMETERS

XNUMPY ARRAY OR SPARSE MATRIX OF SHAPE NSAMPLES NFEATURES TRAINING DATA

622SKLEARNLINEARMODEL GENERALIZED LINEAR MODELS 1931

SCIKITLEARN USER GUIDE RELEASE 0213

YARRAY SHAPE NSAMPLES OR NSAMPLES NTARGETS TARGET VALUES

RETURNS

ZFLOAT SCORE OF THE PREDICTION

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

EXAMPLES USING SKLEARNLINEARMODELRANSACREGRESSOR

- ROBUST LINEAR MODEL ESTIMATION USING RANSAC
- THEILSEN REGRESSION
- ROBUST LINEAR ESTIMATOR FITTING

62217SKLEARNLINEARMODEL RIDGE

CLASSSKLEARNLINEARMODEL RIDGEALPHA10 FITINTERCEPTTRUE NORMALIZEFALSE

COPYXTRUE MAXITERNONE TOL0001 SOLVER’AUTO’

RANDOMSTATENONE

LINEAR LEAST SQUARES WITH L2 REGULARIZATION

MINIMIZES THE OBJECTIVE FUNCTION

Y XW22 ALPHA W22

THIS MODEL SOLVES A REGRESSION MODEL WHERE THE LOSS FUNCTION IS THE LINEAR LEAST SQUARES FUNCTION AND REGULARIZATION IS GIVEN BY THE L2NORM ALSO KNOWN AS RIDGE REGRESSION OR TIKHONOV REGULARIZATION THIS ESTIMATOR HAS BUILTIN SUPPORT FOR MULTIVARIATE REGRESSION IE WHEN Y IS A 2DARRAY OF SHAPE NSAMPLES NTARGETS

READ MORE IN THE USER GUIDE

PARAMETERS

ALPHA FLOAT ARRAYLIKE SHAPE NTARGETS REGULARIZATION STRENGTH MUST BE A POSITIVE FLOAT

REGULARIZATION IMPROVES THE CONDITIONING OF THE PROBLEM AND REDUCES THE VARIANCE OF THE ESTIMATES LARGER VALUES SPECIFY STRONGER REGULARIZATION ALPHA CORRESPONDS TO C1 IN OTHER LINEAR MODELS SUCH AS LOGISTICREGRESSION OR LINEARSVC IF AN ARRAY IS PASSED PENALTIES ARE ASSUMED TO BE SPECIFIC TO THE TARGETS HENCE THEY MUST CORRESPOND IN NUMBER

FITINTERCEPT BOOLEAN WHETHER TO CALCULATE THE INTERCEPT FOR THIS MODEL IF SET TO FALSE NO INTERCEPT WILL BE USED IN CALCULATIONS EG DATA IS EXPECTED TO BE ALREADY CENTERED

NORMALIZE BOOLEAN OPTIONAL DEFAULT FALSE THIS PARAMETER IS IGNORED WHEN FITINTERCEPT IS SET TO FALSE IF TRUE THE REGRESSORS X WILL BE NORMALIZED BEFORE REGRESSION BY SUBTRACTING THE MEAN AND DIVIDING BY THE L2NORM IF YOU WISH TO STANDARDIZE PLEASE USE SKLEARNPREPROCESSINGSTANDARDSCALER BEFORE

CALLINGFIT ON AN ESTIMATOR WITH NORMALIZEFALSE

COPYX BOOLEAN OPTIONAL DEFAULT TRUE IF TRUE X WILL BE COPIED ELSE IT MAY BE OVERWRITTEN

1932 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

MAXITER INT OPTIONAL MAXIMUM NUMBER OF ITERATIONS FOR CONJUGATE GRADIENT SOLVER FOR 'SPARSECG' AND 'LSQR' SOLVERS THE DEFAULT VALUE IS DETERMINED BY SCIPYSPARSELINALG FOR 'SAG' SOLVER THE DEFAULT VALUE IS 1000

TOLFLOAT PRECISION OF THE SOLUTION

SOLVER 'AUTO' 'SVD' 'CHOLESKY' 'LSQR' 'SPARSECG' 'SAG' 'SAGA' SOLVER TO USE IN THE COMPUTATIONAL ROUTINES

- 'AUTO' CHOOSES THE SOLVER AUTOMATICALLY BASED ON THE TYPE OF DATA
- 'SVD' USES A SINGULAR VALUE DECOMPOSITION OF X TO COMPUTE THE RIDGE COEFFICIENTS MORE STABLE FOR SINGULAR MATRICES THAN 'CHOLESKY'
- 'CHOLESKY' USES THE STANDARD SCIPYLINALGSOLVE FUNCTION TO OBTAIN A CLOSEDFORM SOLUTION
- 'SPARSECG' USES THE CONJUGATE GRADIENT SOLVER AS FOUND IN SCIPYSPARSELINALGCG AS AN ITERATIVE ALGORITHM THIS SOLVER IS MORE APPROPRIATE THAN 'CHOLESKY' FOR LARGESCALE DATA POSSIBILITY TO SET TOL ANDMAXITER
- 'LSQR' USES THE DEDICATED REGULARIZED LEASTSQUARES ROUTINE SCIPYSPARSELINALGLSQR IT IS THE FASTEST AND USES AN ITERATIVE PROCEDURE
- 'SAG' USES A STOCHASTIC AVERAGE GRADIENT DESCENT AND 'SAGA' USES ITS IMPROVED UNBIASED VERSION NAMED SAGA BOTH METHODS ALSO USE AN ITERATIVE PROCEDURE AND ARE OFTEN FASTER THAN OTHER SOLVERS WHEN BOTH NSAMPLES AND NFEATURES ARE LARGE NOTE THAT 'SAG' AND 'SAGA' FAST CONVERGENCE IS ONLY GUARANTEED ON FEATURES WITH APPROXIMATELY THE SAME SCALE YOU CAN PREPROCESS THE DATA WITH A SCALER FROM SKLEARNPREPROCESSING

ALL LAST FIVE SOLVERS SUPPORT BOTH DENSE AND SPARSE DATA HOWEVER ONLY 'SAG' AND 'SPARSECG' SUPPORTS SPARSE INPUT WHEN FITINTERCEPT IS TRUE

NEW IN VERSION 017 STOCHASTIC AVERAGE GRADIENT DESCENT SOLVER

NEW IN VERSION 019 SAGA SOLVER

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULT NONE THE SEED OF THE PSEUDO RANDOM NUMBER GENERATOR TO USE WHEN SHUFFLING THE DATA IF INT RANDOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NPRANDOM USED WHEN SOLVER 'SAG'

NEW IN VERSION 017 RANDOMSTATE TO SUPPORT STOCHASTIC AVERAGE GRADIENT

ATTRIBUTES

COEF ARRAY SHAPE NFEATURES OR NTARGETS NFEATURES WEIGHT VECTORS

INTERCEPT FLOAT ARRAY SHAPE NTARGETS INDEPENDENT TERM IN DECISION FUNCTION SET TO 00

IFFITINTERCEPT FALSE

NITER ARRAY OR NONE SHAPE NTARGETS ACTUAL NUMBER OF ITERATIONS FOR EACH TARGET AVAILABLE ONLY FOR SAG AND LSQR SOLVERS OTHER SOLVERS WILL RETURN NONE

NEW IN VERSION 017

SEE ALSO

RIDGECLASSIFIER RIDGE CLASSIFIER

RIDGECV RIDGE REGRESSION WITH BUILTIN CROSS VALIDATION

SKLEARNKERNELRIDGEKERNELRIDGE KERNEL RIDGE REGRESSION COMBINES RIDGE REGRESSION WITH THE KERNEL TRICK

622SKLEARNLINEARMODEL GENERALIZED LINEAR MODELS 1933

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES

FROM SKLEARNLINEARMODEL IMPORT RIDGE

IMPORT NUMPY AS NP

NSAMPLES NFEATURES 10 5

RNG NPRANDOMRANDOMSTATE0

Y RNGRANDNNSAMPLES

X RNGRANDNNSAMPLES NFEATURES

CLF RIDGEALPHA10

CLFFITX Y

RIDGEALPHA10 COPYXTRUE FITINTERCEPTTRUE MAXITERNONE

NORMALIZEFALSE RANDOMSTATENONE SOLVERAUTO TOL0001

METHODS

FITSELF X Y SAMPLEWEIGHT FIT RIDGE REGRESSION MODEL

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PREDICT SELF X PREDICT USING THE LINEAR MODEL

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PRE

DICTION

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELFALPHA10 FITINTERCEPTTRUE NORMALIZEFALSE COPYXTRUE MAXITERNONE

TOL0001 SOLVER'AUTO' RANDOMSTATENONE

FITSELFXYSAMPLEWEIGHTNONE

FIT RIDGE REGRESSION MODEL

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES TRAINING DATA

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NTARGETS TARGET VALUES

SAMPLEWEIGHT FLOAT OR NUMPY ARRAY OF SHAPE NSAMPLES INDIVIDUAL WEIGHTS FOR EACH

SAMPLE

RETURNS

SELF RETURNS AN INSTANCE OF SELF

GETPARAMS SELFDEEPTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELF

PREDICT USING THE LINEAR MODEL

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES SAMPLES

1934 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

RETURNS

CARRAY SHAPE NSAMPLES RETURNS PREDICTED VALUES

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{U}{V}$  WHERE U IS THE RESIDUAL SUM OF SQUARES  $Y_{TRUE} - Y_{PRED}$

2SUM AND V IS THE TOTAL SUM OF SQUARES  $Y_{TRUE} - Y_{TRUEMEAN}$  2SUM THE BEST POSSIBLE SCORE

IS 10 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS

PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 00

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES FOR SOME ESTIMATORS THIS MAY

BE A PRECOMPUTED KERNEL MATRIX INSTEAD SHAPE NSAMPLES NSAMPLESFITTED WHERE

NSAMPLESFITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE VALUES FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT R2 OF SELF-PREDICTX WRT Y

NOTES

THE R2 SCORE USED WHEN CALLING SCORE ON A REGRESSOR WILL USE MULTIOUTPUTUNIFORMAVERAGE

FROM VERSION 023 TO KEEP CONSISTENT WITH METRICSR2SCORE THIS WILL INFLUENCE THE SCORE

METHOD OF ALL THE MULTIOUTPUT REGRESSORS EXCEPT FOR MULTIOUTPUTMULTIOUTPUTREGRESSOR

TO SPECIFY THE DEFAULT VALUE MANUALLY AND AVOID THE WARNING PLEASE EITHER CALL METRICSR2SCORE

DIRECTLY OR MAKE A CUSTOM SCORER WITH METRICSMAKESCORER THE BUILTIN SCORER R2 USES

MULTIOUTPUTUNIFORMAVERAGE

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE

PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT

OF A NESTED OBJECT

RETURNS

SELF

EXAMPLES USING SKLEARNLINEARMODELRIDGE

- COMPRESSIVE SENSING TOMOGRAPHY RECONSTRUCTION WITH L1 PRIOR LASSO
  - PREDICTION LATENCY
  - PLOT RIDGE COEFFICIENTS AS A FUNCTION OF THE REGULARIZATION
  - ORDINARY LEAST SQUARES AND RIDGE REGRESSION VARIANCE
  - PLOT RIDGE COEFFICIENTS AS A FUNCTION OF THE L2 REGULARIZATION
  - POLYNOMIAL INTERPOLATION
  - HUBERREGRESSOR VS RIDGE ON DATASET WITH STRONG OUTLIERS
- 622SKLEARNLINEARMODEL GENERALIZED LINEAR MODELS 1935

SCIKITLEARN USER GUIDE RELEASE 0213  
62218SKLEARNLINEARMODEL RIDGECLASSIFIER  
CLASSSSKLEARNLINEARMODEL RIDGECLASSIFIER ALPHA10 FITINTERCEPTTRUE NORMALIZEFALSE  
COPYXTRUE MAXITERNONE TOL0001  
CLASSWEIGHTNONE SOLVER'AUTO' RAN  
DOMSTATENONE  
CLASSIFIER USING RIDGE REGRESSION  
READ MORE IN THE USER GUIDE  
PARAMETERS  
ALPHA FLOAT REGULARIZATION STRENGTH MUST BE A POSITIVE FLOAT REGULARIZATION IMPROVES THE CON  
DITITIONING OF THE PROBLEM AND REDUCES THE VARIANCE OF THE ESTIMATES LARGER VALUES SPECIFY  
STRONGER REGULARIZATION ALPHA CORRESPONDS TO C1 IN OTHER LINEAR MODELS SUCH AS LOGISTI  
CREGRESSION OR LINEARSVC  
FITINTERCEPT BOOLEAN WHETHER TO CALCULATE THE INTERCEPT FOR THIS MODEL IF SET TO FALSE NO  
INTERCEPT WILL BE USED IN CALCULATIONS EG DATA IS EXPECTED TO BE ALREADY CENTERED  
NORMALIZE BOOLEAN OPTIONAL DEFAULT FALSE THIS PARAMETER IS IGNORED WHEN  
FITINTERCEPT IS SET TO FALSE IF TRUE THE REGRESSORS X WILL BE NORMALIZED BE  
FORE REGRESSION BY SUBTRACTING THE MEAN AND DIVIDING BY THE L2NORM IF YOU WISH TO  
STANDARDIZE PLEASE USE SKLEARNPREPROCESSINGSTANDARDSCALER BEFORE  
CALLINGFIT ON AN ESTIMATOR WITH NORMALIZEFALSE  
COPYX BOOLEAN OPTIONAL DEFAULT TRUE IF TRUE X WILL BE COPIED ELSE IT MAY BE OVERWRITTEN  
MAXITER INT OPTIONAL MAXIMUM NUMBER OF ITERATIONS FOR CONJUGATE GRADIENT SOLVER THE  
DEFAULT VALUE IS DETERMINED BY SCIPYSPARSELINALG  
TOLFLOAT PRECISION OF THE SOLUTION  
CLASSWEIGHT DICT OR 'BALANCED' OPTIONAL WEIGHTS ASSOCIATED WITH CLASSES IN THE FORM  
CLASSLABEL WEIGHT IF NOT GIVEN ALL CLASSES ARE SUPPOSED TO HAVE WEIGHT ONE  
THE "BALANCED" MODE USES THE VALUES OF Y TO AUTOMATICALLY ADJUST WEIGHTS INVERSELY PRO  
PORTIONAL TO CLASS FREQUENCIES IN THE INPUT DATA AS NSAMPLES NCLASSES NP  
BINCOUNTY  
SOLVER 'AUTO' 'SVD' 'CHOLESKY' 'LSQR' 'SPARSECG' 'SAG' 'SAGA' SOLVER TO USE IN THE COM  
PUTATIONAL ROUTINES  
• 'AUTO' CHOOSES THE SOLVER AUTOMATICALLY BASED ON THE TYPE OF DATA  
• 'SVD' USES A SINGULAR VALUE DECOMPOSITION OF X TO COMPUTE THE RIDGE COEFFICIENTS MORE  
STABLE FOR SINGULAR MATRICES THAN 'CHOLESKY'  
• 'CHOLESKY' USES THE STANDARD SCIPYLINALGSOLVE FUNCTION TO OBTAIN A CLOSEDFORM SOLUTION  
• 'SPARSECG' USES THE CONJUGATE GRADIENT SOLVER AS FOUND IN SCIPYSPARSELINALGCG AS AN  
ITERATIVE ALGORITHM THIS SOLVER IS MORE APPROPRIATE THAN 'CHOLESKY' FOR LARGESCALE DATA  
POSSIBILITY TO SET TOL ANDMAXITER  
• 'LSQR' USES THE DEDICATED REGULARIZED LEASTSQUARES ROUTINE SCIPYSPARSELINALGLSQR IT IS THE  
FASTEST AND USES AN ITERATIVE PROCEDURE  
• 'SAG' USES A STOCHASTIC AVERAGE GRADIENT DESCENT AND 'SAGA' USES ITS UNBIASED AND MORE  
FLEXIBLE VERSION NAMED SAGA BOTH METHODS USE AN ITERATIVE PROCEDURE AND ARE OFTEN  
FASTER THAN OTHER SOLVERS WHEN BOTH NSAMPLES AND NFEATURES ARE LARGE NOTE THAT 'SAG'  
AND 'SAGA' FAST CONVERGENCE IS ONLY GUARANTEED ON FEATURES WITH APPROXIMATELY THE SAME  
SCALE YOU CAN PREPROCESS THE DATA WITH A SCALER FROM SKLEARNPREPROCESSING  
1936 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213  
 NEW IN VERSION 017 STOCHASTIC AVERAGE GRADIENT DESCENT SOLVER  
 NEW IN VERSION 019 SAGA SOLVER  
 RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULT NONE THE SEED OF THE  
 PSEUDO RANDOM NUMBER GENERATOR TO USE WHEN SHUFFLING THE DATA IF INT RANDOMSTATE IS  
 THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS  
 THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE  
 INSTANCE USED BY NPRANDOM USED WHEN SOLVER 'SAG'  
 ATTRIBUTES  
 COEF ARRAY SHAPE 1 NFEATURES OR NCLASSES NFEATURES COEFFICIENT OF THE FEATURES IN THE  
 DECISION FUNCTION  
 COEF IS OF SHAPE 1 NFEATURES WHEN THE GIVEN PROBLEM IS BINARY  
 INTERCEPT FLOAT ARRAY SHAPE NTARGETS INDEPENDENT TERM IN DECISION FUNCTION SET TO 00  
 IFFITINTERCEPT FALSE  
 NITER ARRAY OR NONE SHAPE NTARGETS ACTUAL NUMBER OF ITERATIONS FOR EACH TARGET AVAIL  
 ABLE ONLY FOR SAG AND LSQR SOLVERS OTHER SOLVERS WILL RETURN NONE  
 SEE ALSO  
 RIDGE RIDGE REGRESSION  
 RIDGECLASSIFIERCV RIDGE CLASSIFIER WITH BUILTIN CROSS VALIDATION  
 NOTES  
 FOR MULTICLASS CLASSIFICATION NCLASS CLASSIFIERS ARE TRAINED IN A ONEVERSUSALL APPROACH CONCRETELY THIS IS  
 IMPLEMENTED BY TAKING ADVANTAGE OF THE MULTIVARIATE RESPONSE SUPPORT IN RIDGE  
 EXAMPLES  
 FROM SKLEARNDATASETS IMPORT LOADBREASTCANCER  
 FROM SKLEARNLINEARMODEL IMPORT RIDGECLASSIFIER  
 X Y LOADBREASTCANCERRETURNXY TRUE  
 CLF RIDGECLASSIFIERFITX Y  
 CLFScoreX Y  
 09595  
 METHODS  
 DECISIONFUNCTION SELF X PREDICT CONFIDENCE SCORES FOR SAMPLES  
 FITSELF X Y SAMPLEWEIGHT FIT RIDGE REGRESSION MODEL  
 GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR  
 PREDICT SELF X PREDICT CLASS LABELS FOR SAMPLES IN X  
 SCORE SELF X Y SAMPLEWEIGHT RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND  
 LABELS  
 SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR  
 622SKLEARNLINEARMODEL GENERALIZED LINEAR MODELS 1937

SCIKITLEARN USER GUIDE RELEASE 0213  
INIT SELFALPHA10 FITINTERCEPTTRUE NORMALIZEFALSE COPYXTRUE MAXITERNONE  
TOL0001 CLASSWEIGHTNONE SOLVER'AUTO' RANDOMSTATENONE  
DECISIONFUNCTION SELF  
PREDICT CONFIDENCE SCORES FOR SAMPLES  
THE CONFIDENCE SCORE FOR A SAMPLE IS THE SIGNED DISTANCE OF THAT SAMPLE TO THE HYPERPLANE  
PARAMETERS  
XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES SAMPLES  
RETURNS  
ARRAY SHAPENSAMPLES IF NCLASSES 2 ELSE NSAMPLES NCLASSES CONFIDENCE  
SCORES PER SAMPLE CLASS COMBINATION IN THE BINARY CASE CONFIDENCE SCORE FOR  
SELFCLASSES1 WHERE 0 MEANS THIS CLASS WOULD BE PREDICTED  
FITSELFXYSAMPLEWEIGHTNONE  
FIT RIDGE REGRESSION MODEL  
PARAMETERS  
XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLESNFEATURES TRAINING DATA  
YARRAYLIKE SHAPE NSAMPLES TARGET VALUES  
SAMPLEWEIGHT FLOAT OR NUMPY ARRAY OF SHAPE NSAMPLES SAMPLE WEIGHT  
NEW IN VERSION 017 SAMPLEWEIGHT SUPPORT TO CLASSIFIER  
RETURNS  
SELF RETURNS AN INSTANCE OF SELF  
GETPARAMS SELFDEEPTTRUE  
GET PARAMETERS FOR THIS ESTIMATOR  
PARAMETERS  
DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBOBJECTS THAT ARE ESTIMATORS  
RETURNS  
PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES  
PREDICTSELF  
PREDICT CLASS LABELS FOR SAMPLES IN X  
PARAMETERS  
XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES SAMPLES  
RETURNS  
CARRAY SHAPE NSAMPLES PREDICTED CLASS LABEL PER SAMPLE  
SCORESELFXYSAMPLEWEIGHTNONE  
RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS  
IN MULTILABEL CLASSIFICATION THIS IS THE SUBSET ACCURACY WHICH IS A HARSH METRIC SINCE YOU REQUIRE FOR EACH  
SAMPLE THAT EACH LABEL SET BE CORRECTLY PREDICTED  
PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES  
YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE LABELS FOR X  
1938 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT MEAN ACCURACY OF SELF PREDICT X WRT Y

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

EXAMPLES USING SKLEARN LINEAR MODEL RIDGE CLASSIFIER

- CLASSIFICATION OF TEXT DOCUMENTS USING SPARSE FEATURES

62219 SKLEARN LINEAR MODEL SGD CLASSIFIER

CLASS SKLEARN LINEAR MODEL SGD CLASSIFIER LOSS 'HINGE' PENALTY 'L2' ALPHA 00001

L1 RATIO 0015 FIT INTERCEPT TRUE MAX ITER 1000

TOL 0001 SHUFFLE TRUE VERBOSE 0 EPSILON 01

NJOBS NONE RANDOM STATE NONE LEARN

IN RATE 'OPTIMAL' ETA 000 POWER T 05

EARLY STOPPING FALSE VALIDATION FRACTION 01

NITER NO CHANGE 5 CLASS WEIGHT NONE

WARM START FALSE AVERAGE FALSE

LINEAR CLASSIFIERS SVM LOGISTIC REGRESSION AO WITH SGD TRAINING

THIS ESTIMATOR IMPLEMENTS REGULARIZED LINEAR MODELS WITH STOCHASTIC GRADIENT DESCENT SGD LEARNING THE GRADIENT OF THE LOSS IS ESTIMATED EACH SAMPLE AT A TIME AND THE MODEL IS UPDATED ALONG THE WAY WITH A DECREASING STRENGTH SCHEDULE AKA LEARNING RATE SGD ALLOWS MINIBATCH ONLINE OUT OF CORE LEARNING SEE THE PARTIAL FIT METHOD FOR BEST RESULTS USING THE DEFAULT LEARNING RATE SCHEDULE THE DATA SHOULD HAVE ZERO MEAN AND UNIT VARIANCE

THIS IMPLEMENTATION WORKS WITH DATA REPRESENTED AS DENSE OR SPARSE ARRAYS OF FLOATING POINT VALUES FOR THE FEATURES THE MODEL IT FITS CAN BE CONTROLLED WITH THE LOSS PARAMETER BY DEFAULT IT FITS A LINEAR SUPPORT VECTOR MACHINE SVM

THE REGULARIZER IS A PENALTY ADDED TO THE LOSS FUNCTION THAT SHRINKS MODEL PARAMETERS TOWARDS THE ZERO VECTOR USING EITHER THE SQUARED EUCLIDEAN NORM L2 OR THE ABSOLUTE NORM L1 OR A COMBINATION OF BOTH ELASTIC NET IF THE PARAMETER UPDATE CROSSES THE 00 VALUE BECAUSE OF THE REGULARIZER THE UPDATE IS TRUNCATED TO 00 TO ALLOW FOR LEARNING SPARSE MODELS AND ACHIEVE ONLINE FEATURE SELECTION

READ MORE IN THE USER GUIDE

PARAMETERS

LOSS STR DEFAULT 'HINGE' THE LOSS FUNCTION TO BE USED DEFAULTS TO 'HINGE' WHICH GIVES A LINEAR SVM

THE POSSIBLE OPTIONS ARE 'HINGE' 'LOG' 'MODIFIED HUBER' 'SQUARED HINGE' 'PERCEPTRON' OR A REGRESSION LOSS 'SQUARED LOSS' 'HUBER' 'EPSILON INSENSITIVE' OR 'SQUARED EPSILON INSENSITIVE'

622 SKLEARN LINEAR MODEL GENERALIZED LINEAR MODELS 1939

SCIKITLEARN USER GUIDE RELEASE 0213

THE 'LOG' LOSS GIVES LOGISTIC REGRESSION A PROBABILISTIC CLASSIFIER 'MODIFIEDHUBER' IS ANOTHER SMOOTH LOSS THAT BRINGS TOLERANCE TO OUTLIERS AS WELL AS PROBABILITY ESTIMATES 'SQUAREDHINGE' IS LIKE HINGE BUT IS QUADRATICALLY PENALIZED 'PERCEPTRON' IS THE LINEAR LOSS USED BY THE PERCEPTRON ALGORITHM THE OTHER LOSSES ARE DESIGNED FOR REGRESSION BUT CAN BE USEFUL IN CLASSIFICATION AS WELL SEE SGDREGRESSOR FOR A DESCRIPTION

PENALTY STR 'NONE' 'L2' 'L1' OR 'ELASTICNET' THE PENALTY AKA REGULARIZATION TERM TO BE USED DEFAULTS TO 'L2' WHICH IS THE STANDARD REGULARIZER FOR LINEAR SVM MODELS 'L1' AND 'ELASTICNET' MIGHT BRING SPARSITY TO THE MODEL FEATURE SELECTION NOT ACHIEVABLE WITH 'L2'

ALPHA FLOAT CONSTANT THAT MULTIPLIES THE REGULARIZATION TERM DEFAULTS TO 00001 ALSO USED TO COMPUTE LEARNINGRATE WHEN SET TO 'OPTIMAL'

L1RATIO FLOAT THE ELASTIC NET MIXING PARAMETER WITH 0 L1RATIO 1 L1RATIO0 CORRESPONDS TO L2 PENALTY L1RATIO1 TO L1 DEFAULTS TO 015

FITINTERCEPT BOOL WHETHER THE INTERCEPT SHOULD BE ESTIMATED OR NOT IF FALSE THE DATA IS ASSUMED TO BE ALREADY CENTERED DEFAULTS TO TRUE

MAXITER INT OPTIONAL DEFAULT1000 THE MAXIMUM NUMBER OF PASSES OVER THE TRAINING DATA AKA EPOCHS IT ONLY IMPACTS THE BEHAVIOR IN THE FIT METHOD AND NOT THE PARTIALFIT

NEW IN VERSION 019

TOLFLOAT OR NONE OPTIONAL DEFAULT1E3 THE STOPPING CRITERION IF IT IS NOT NONE THE ITERATIONS WILL STOP WHEN LOSS - BESTLOSS - TOL FOR NITERNOCHANGE CONSECUTIVE EPOCHS

NEW IN VERSION 019

SHUFFLE BOOL OPTIONAL WHETHER OR NOT THE TRAINING DATA SHOULD BE SHUFFLED AFTER EACH EPOCH DEFAULTS TO TRUE

VERBOSE INTEGER OPTIONAL THE VERBOSITY LEVEL

EPSILON FLOAT EPSILON IN THE EPSILONINSENSITIVE LOSS FUNCTIONS ONLY IF LOSS IS 'HUBER' 'EPSILONINSENSITIVE' OR 'SQUAREDEPSILONINSENSITIVE' FOR 'HUBER' DETERMINES THE THRESHOLD AT WHICH IT BECOMES LESS IMPORTANT TO GET THE PREDICTION EXACTLY RIGHT FOR EPSILONINSENSITIVE ANY DIFFERENCES BETWEEN THE CURRENT PREDICTION AND THE CORRECT LABEL ARE IGNORED IF THEY ARE LESS THAN THIS THRESHOLD

NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF CPUS TO USE TO DO THE OVA ONE VERSUS ALL FOR MULTICLASS PROBLEMS COMPUTATION NONE MEANS 1 UNLESS IN A JOBLIB PARALLELBACKEND CONTEXT1MEANS USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE THE SEED OF THE PSEUDO RANDOM NUMBER GENERATOR TO USE WHEN SHUFFLING THE DATA IF INT RANDOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NPRANDOM

LEARNINGRATE STRING OPTIONAL THE LEARNING RATE SCHEDULE

'CONSTANT' ETA - ETA0

'OPTIMAL' DEFAULT ETA - 10 \* ALPHA \* T / T0 WHERE T0 IS CHOSEN BY A HEURISTIC PROPOSED BY LEON BOTTOU

'INVSCALING' ETA - ETA0 \* POWT / POWT

'ADAPTIVE' ETA - ETA0 AS LONG AS THE TRAINING KEEPS DECREASING EACH TIME

NITERNOCHANGE CONSECUTIVE EPOCHS FAIL TO DECREASE THE TRAINING LOSS BY TOL OR FAIL TO

1940 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

INCREASE VALIDATION SCORE BY TOL IF EARLYSTOPPING IS TRUE THE CURRENT LEARNING RATE IS DIVIDED BY 5

ETA0 DOUBLE THE INITIAL LEARNING RATE FOR THE ‘CONSTANT’ ‘INVSCALING’ OR ‘ADAPTIVE’ SCHEDULES THE DEFAULT VALUE IS 00 AS ETA0 IS NOT USED BY THE DEFAULT SCHEDULE ‘OPTIMAL’

POWERT DOUBLE THE EXPONENT FOR INVERSE SCALING LEARNING RATE DEFAULT 05

EARLYSTOPPING BOOL DEFAULTFALSE WHETHER TO USE EARLY STOPPING TO TERMINATE TRAINING WHEN VALIDATION SCORE IS NOT IMPROVING IF SET TO TRUE IT WILL AUTOMATICALLY SET ASIDE A STRATIFIED FRACTION OF TRAINING DATA AS VALIDATION AND TERMINATE TRAINING WHEN VALIDATION SCORE IS NOT IMPROVING BY AT LEAST TOL FOR NITERNOCHANGE CONSECUTIVE EPOCHS

NEW IN VERSION 020

VALIDATIONFRACTION FLOAT DEFAULT01 THE PROPORTION OF TRAINING DATA TO SET ASIDE AS VALIDATION SET FOR EARLY STOPPING MUST BE BETWEEN 0 AND 1 ONLY USED IF EARLYSTOPPING IS TRUE

NEW IN VERSION 020

NITERNOCHANGE INT DEFAULT5 NUMBER OF ITERATIONS WITH NO IMPROVEMENT TO WAIT BEFORE EARLY STOPPING

NEW IN VERSION 020

CLASSWEIGHT DICT CLASSLABEL WEIGHT OR “BALANCED” OR NONE OPTIONAL PRESET FOR THE CLASSWEIGHT FIT PARAMETER

WEIGHTS ASSOCIATED WITH CLASSES IF NOT GIVEN ALL CLASSES ARE SUPPOSED TO HAVE WEIGHT ONE THE “BALANCED” MODE USES THE VALUES OF Y TO AUTOMATICALLY ADJUST WEIGHTS INVERSELY PROPORTIONAL TO CLASS FREQUENCIES IN THE INPUT DATA AS NSAMPLES NCLASSES NP

BINCOUNTY

WARMSTART BOOL OPTIONAL WHEN SET TO TRUE REUSE THE SOLUTION OF THE PREVIOUS CALL TO FIT AS INITIALIZATION OTHERWISE JUST ERASE THE PREVIOUS SOLUTION SEE THE GLOSSARY

REPEATEDLY CALLING FIT OR PARTIALFIT WHEN WARMSTART IS TRUE CAN RESULT IN A DIFFERENT SOLUTION THAN WHEN CALLING FIT A SINGLE TIME BECAUSE OF THE WAY THE DATA IS SHUFFLED IF A DYNAMIC LEARNING RATE IS USED THE LEARNING RATE IS ADAPTED DEPENDING ON THE NUMBER OF SAMPLES ALREADY SEEN CALLING FIT RESETS THIS COUNTER WHILE PARTIALFIT WILL RESULT IN INCREASING THE EXISTING COUNTER

AVERAGE BOOL OR INT OPTIONAL WHEN SET TO TRUE COMPUTES THE AVERAGED SGD WEIGHTS AND STORES THE RESULT IN THE COEF ATTRIBUTE IF SET TO AN INT GREATER THAN 1 AVERAGING WILL BEGIN ONCE THE TOTAL NUMBER OF SAMPLES SEEN REACHES AVERAGE SO AVERAGE10 WILL BEGIN AVERAGING AFTER SEEING 10 SAMPLES

ATTRIBUTES

COEF ARRAY SHAPE 1 NFEATURES IF NCLASSES 2 ELSE NCLASSES NFEATURES WEIGHTS AS SIGNED TO THE FEATURES

INTERCEPT ARRAY SHAPE 1 IF NCLASSES 2 ELSE NCLASSES CONSTANTS IN DECISION FUNCTION

NITER INT THE ACTUAL NUMBER OF ITERATIONS TO REACH THE STOPPING CRITERION FOR MULTICLASS FITS IT IS THE MAXIMUM OVER EVERY BINARY FIT

LOSSFUNCTION CONCRETELOSSFUNCTION

SEE ALSO

SKLEARNVMLINEARSVC LOGISTICREGRESSION PERCEPTRON

622SKLEARNLINEARMODEL GENERALIZED LINEAR MODELS 1941

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES

```
import numpy as np
from sklearn import linear_model
X = np.array([1, 2, 1, 1, 1, 2, 1])
Y = np.array([1, 2, 2])
clf = linear_model.SGDClassifier(max_iter=1000, tol=1e-3)
clf.fit(X, Y)
```

```
sgd_classifier(alpha=0.0001, average=False, class_weight=None,
early_stopping=False, epsilon=0.1, eta=0.0001, fit_intercept=True,
l1_ratio=0.15, learning_rate='optimal', loss='hinge', max_iter=1000,
n_iter_no_change=5, n_jobs=None, penalty='l2', power_t=0.5,
random_state=None, shuffle=True, tol=0.0001, validation_fraction=0.1,
verbose=0, warm_start=False)
print(clf.predict(0.8))
```

```
1
METHODS
DECISIONFUNCTION SELF X PREDICT CONFIDENCE SCORES FOR SAMPLES
DENSIFY SELF CONVERT COEFFICIENT MATRIX TO DENSE ARRAY FORMAT
FITSELF X Y COEFINIT INTERCEPTINIT FIT LINEAR MODEL WITH STOCHASTIC GRADIENT DESCENT
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR
PARTIALFIT SELF X Y CLASSES SAMPLEWEIGHT PERFORM ONE EPOCH OF STOCHASTIC GRADIENT DESCENT ON
GIVEN SAMPLES
PREDICT SELF X PREDICT CLASS LABELS FOR SAMPLES IN X
SCORE SELF X Y SAMPLEWEIGHT RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND
LABELS
SETPARAMS SELF ARGS KWARGS
SPARSIFY SELF CONVERT COEFFICIENT MATRIX TO SPARSE FORMAT
INIT SELFLOSS'HINGE' PENALTY'L2' ALPHA00001 L1RATIO015 FITINTERCEPTTRUE
MAXITER1000 TOL0001 SHUFFLETRUE VERBOSE0 EPSILON01 NJOBSNONE RAN
DOMSTATENONE LEARNINGRATE'OPTIMAL' ETA000 POWER05 EARLYSTOPPINGFALSE
VALIDATIONFRACTION01 NITERNOCHANGE5 CLASSWEIGHTNONE WARMSTARTFALSE AV
ERAGEFALSE
DECISIONFUNCTION SELF X
PREDICT CONFIDENCE SCORES FOR SAMPLES
THE CONFIDENCE SCORE FOR A SAMPLE IS THE SIGNED DISTANCE OF THAT SAMPLE TO THE HYPERPLANE
PARAMETERS
XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES SAMPLES
RETURNS
ARRAY SHAPENSAMPLES IF NCLASSES > 2 ELSE NSAMPLES NCLASSES CONFIDENCE
SCORES PER SAMPLE CLASS COMBINATION IN THE BINARY CASE CONFIDENCE SCORE FOR
SELFCLASSES1 WHERE 0 MEANS THIS CLASS WOULD BE PREDICTED
1942 CHAPTER 6 API REFERENCE
```

SCIKITLEARN USER GUIDE RELEASE 0213

DENSIFYSELF

CONVERT COEFFICIENT MATRIX TO DENSE ARRAY FORMAT

CONVERTS THE COEF MEMBER BACK TO A NUMPYNDARRAY THIS IS THE DEFAULT FORMAT OF COEF AND IS REQUIRED FOR FITTING SO CALLING THIS METHOD IS ONLY REQUIRED ON MODELS THAT HAVE PREVIOUSLY BEEN SPARSIFIED OTHERWISE IT IS A NOOP

RETURNS

SELF ESTIMATOR

FITSELFXYCOEFINITNONE INTERCEPTINITNONE SAMPLEWEIGHTNONE

FIT LINEAR MODEL WITH STOCHASTIC GRADIENT DESCENT

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES TRAINING DATA

YNUMPY ARRAY SHAPE NSAMPLES TARGET VALUES

COEFINIT ARRAY SHAPE NCLASSES NFEATURES THE INITIAL COEFFICIENTS TO WARMSTART THE OPTIMIZATION

INTERCEPTINIT ARRAY SHAPE NCLASSES THE INITIAL INTERCEPT TO WARMSTART THE OPTIMIZATION

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL WEIGHTS APPLIED TO INDIVIDUAL SAMPLES IF NOT PROVIDED UNIFORM WEIGHTS ARE ASSUMED THESE WEIGHTS WILL BE MULTIPLIED WITH CLASSWEIGHT PASSED THROUGH THE CONSTRUCTOR IF CLASSWEIGHT IS SPECIFIED

RETURNS

SELF RETURNS AN INSTANCE OF SELF

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PARTIALFIT SELFXYCLASSESNONE SAMPLEWEIGHTNONE

PERFORM ONE EPOCH OF STOCHASTIC GRADIENT DESCENT ON GIVEN SAMPLES

INTERNALLY THIS METHOD USES MAXITER 1 THEREFORE IT IS NOT GUARANTEED THAT A MINIMUM OF THE COST FUNCTION IS REACHED AFTER CALLING IT ONCE MATTERS SUCH AS OBJECTIVE CONVERGENCE AND EARLY STOPPING SHOULD BE HANDLED BY THE USER

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES SUBSET OF THE TRAINING DATA

YNUMPY ARRAY SHAPE NSAMPLES SUBSET OF THE TARGET VALUES

CLASSES ARRAY SHAPE NCLASSES CLASSES ACROSS ALL CALLS TO PARTIALFIT CAN BE OBTAINED BY VIANPUNIQUEYALL WHERE YALL IS THE TARGET VECTOR OF THE ENTIRE DATASET THIS ARGUMENT IS REQUIRED FOR THE FIRST CALL TO PARTIALFIT AND CAN BE OMITTED IN THE SUBSEQUENT CALLS NOTE THAT Y DOESN'T NEED TO CONTAIN ALL LABELS IN CLASSES

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL WEIGHTS APPLIED TO INDIVIDUAL SAMPLES IF NOT PROVIDED UNIFORM WEIGHTS ARE ASSUMED

622SKLEARNLINEARMODEL GENERALIZED LINEAR MODELS 1943

SCIKITLEARN USER GUIDE RELEASE 0213

**RETURNS**  
SELF RETURNS AN INSTANCE OF SELF

**PREDICTSELF**  
PREDICT CLASS LABELS FOR SAMPLES IN X

**PARAMETERS**  
XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES SAMPLES

**RETURNS**  
CARRAY SHAPE NSAMPLES PREDICTED CLASS LABEL PER SAMPLE

**PREDICTLOGPROBA**  
LOG OF PROBABILITY ESTIMATES  
THIS METHOD IS ONLY AVAILABLE FOR LOG LOSS AND MODIFIED HUBER LOSS  
WHEN LOSS"MODIFIEDHUBER" PROBABILITY ESTIMATES MAY BE HARD ZEROS AND ONES SO TAKING THE LOGARITHM IS NOT POSSIBLE  
SEEPREDICTPROBA FOR DETAILS

**PARAMETERS**  
XARRAYLIKE SHAPE NSAMPLES NFEATURES

**RETURNS**  
TARRAYLIKE SHAPE NSAMPLES NCLASSES RETURNS THE LOGPROBABILITY OF THE SAMPLE FOR EACH CLASS IN THE MODEL WHERE CLASSES ARE ORDERED AS THEY ARE IN SELFCLASSES

**PREDICTPROBA**  
PROBABILITY ESTIMATES  
THIS METHOD IS ONLY AVAILABLE FOR LOG LOSS AND MODIFIED HUBER LOSS  
MULTICLASS PROBABILITY ESTIMATES ARE DERIVED FROM BINARY ONEVSREST ESTIMATES BY SIMPLE NORMALIZATION AS RECOMMENDED BY ZADROZNY AND ELKAN  
BINARY PROBABILITY ESTIMATES FOR LOSS"MODIFIEDHUBER" ARE GIVEN BY CLIPDECISIONFUNCTIONX 1 1  
1 2 FOR OTHER LOSS FUNCTIONS IT IS NECESSARY TO PERFORM PROPER PROBABILITY CALIBRATION BY WRAPPING THE CLASSIFIER WITH SKLEARNCALIBRATIONCALIBRATEDCLASSIFIERCV INSTEAD

**PARAMETERS**  
XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES

**RETURNS**  
ARRAY SHAPE NSAMPLES NCLASSES RETURNS THE PROBABILITY OF THE SAMPLE FOR EACH CLASS IN THE MODEL WHERE CLASSES ARE ORDERED AS THEY ARE IN SELFCLASSES

**REFERENCES**  
ZADROZNY AND ELKAN "TRANSFORMING CLASSIFIER SCORES INTO MULTICLASS PROBABILITY ESTIMATES" SIGKDD'02  
HTTPWWWRESEARCHIBMCOMPEOPLEZZADROZNYKDD2002TRANSFPDF  
THE JUSTIFICATION FOR THE FORMULA IN THE LOSS"MODIFIEDHUBER" CASE IS IN THE APPENDIX B IN HTTPJMLR  
CSAILMITEDUPAPERSVOLUME2ZHANG02CZHANG02CPDF  
1944 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS

IN MULTILABEL CLASSIFICATION THIS IS THE SUBSET ACCURACY WHICH IS A HARSH METRIC SINCE YOU REQUIRE FOR EACH SAMPLE THAT EACH LABEL SET BE CORRECTLY PREDICTED

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE LABELS FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT MEAN ACCURACY OF SELF PREDICT X WRT Y

SPARSIFY SELF

CONVERT COEFFICIENT MATRIX TO SPARSE FORMAT

CONVERTS THE COEF MEMBER TO A SCIPYSPARSE MATRIX WHICH FOR L1REGULARIZED MODELS CAN BE MUCH MORE MEMORY AND STORAGEEFFICIENT THAN THE USUAL NUMPYNDARRAY REPRESENTATION

THEINTERCEPT MEMBER IS NOT CONVERTED

RETURNS

SELF ESTIMATOR

NOTES

FOR NONSPARSE MODELS IE WHEN THERE ARE NOT MANY ZEROS IN COEF THIS MAY ACTUALLY INCREASE MEMORY USAGE SO USE THIS METHOD WITH CARE A RULE OF THUMB IS THAT THE NUMBER OF ZERO ELEMENTS WHICH CAN BE COMPUTED WITH COEF 0SUM MUST BE MORE THAN 50 FOR THIS TO PROVIDE SIGNIFICANT BENEFITS

AFTER CALLING THIS METHOD FURTHER FITTING WITH THE PARTIALFIT METHOD IF ANY WILL NOT WORK UNTIL YOU CALL DENSIFY

EXAMPLES USING SKLEARNLINEARMODELSGDCLASSIFIER

- MODEL COMPLEXITY INFLUENCE
- OUTOFCORE CLASSIFICATION OF TEXT DOCUMENTS
- PIPELINING CHAINING A PCA AND A LOGISTIC REGRESSION
- SGD MAXIMUM MARGIN SEPARATING HYPERPLANE
- SGD WEIGHTED SAMPLES
- COMPARING VARIOUS ONLINE SOLVERS
- PLOT MULTICLASS SGD ON THE IRIS DATASET
- EARLY STOPPING OF STOCHASTIC GRADIENT DESCENT
- SAMPLE PIPELINE FOR TEXT FEATURE EXTRACTION AND EVALUATION
- CLASSIFICATION OF TEXT DOCUMENTS USING SPARSE FEATURES

622SKLEARNLINEARMODEL GENERALIZED LINEAR MODELS 1945

SCIKITLEARN USER GUIDE RELEASE 0213

62220SKLEARNLINEARMODEL SGDREGRESSOR  
CLASSSSKLEARNLINEARMODEL SGDREGRESSOR LOSS‘SQUAREDLOSS’ PENALTY‘L2’ ALPHA00001  
L1RATIO015 FITINTERCEPTTRUE MAXITER1000  
TOL0001 SHUFFLETRUE VERBOSE0 EPSILON01  
RANDOMSTATENONE LEARNINGRATE‘INVSCALING’  
ETA0001 POWERT025 EARLYSTOPPINGFALSE  
VALIDATIONFRACTION01 NITERNOCHANGE5  
WARMSTARTFALSE AVERAGEFALSE

LINEAR MODEL FITTED BY MINIMIZING A REGULARIZED EMPIRICAL LOSS WITH SGD  
SGD STANDS FOR STOCHASTIC GRADIENT DESCENT THE GRADIENT OF THE LOSS IS ESTIMATED EACH SAMPLE AT A TIME AND THE  
MODEL IS UPDATED ALONG THE WAY WITH A DECREASING STRENGTH SCHEDULE AKA LEARNING RATE  
THE REGULARIZER IS A PENALTY ADDED TO THE LOSS FUNCTION THAT SHRINKS MODEL PARAMETERS TOWARDS THE ZERO VECTOR  
USING EITHER THE SQUARED EUCLIDEAN NORM L2 OR THE ABSOLUTE NORM L1 OR A COMBINATION OF BOTH ELASTIC NET IF  
THE PARAMETER UPDATE CROSSES THE 00 VALUE BECAUSE OF THE REGULARIZER THE UPDATE IS TRUNCATED TO 00 TO ALLOW FOR  
LEARNING SPARSE MODELS AND ACHIEVE ONLINE FEATURE SELECTION  
THIS IMPLEMENTATION WORKS WITH DATA REPRESENTED AS DENSE NUMPY ARRAYS OF FLOATING POINT VALUES FOR THE FEATURES  
READ MORE IN THE USER GUIDE

PARAMETERS

LOSS STR DEFAULT ‘SQUAREDLOSS’ THE LOSS FUNCTION TO BE USED THE POSSIBLE VALUES ARE  
‘SQUAREDLOSS’ ‘HUBER’ ‘EPSILONINSENSITIVE’ OR ‘SQUAREDEPSILONINSENSITIVE’  
THE ‘SQUAREDLOSS’ REFERS TO THE ORDINARY LEAST SQUARES FIT ‘HUBER’ MODIFIES ‘SQUAREDLOSS’ TO  
FOCUS LESS ON GETTING OUTLIERS CORRECT BY SWITCHING FROM SQUARED TO LINEAR LOSS PAST A DISTANCE  
OF EPSILON ‘EPSILONINSENSITIVE’ IGNORES ERRORS LESS THAN EPSILON AND IS LINEAR PAST THAT THIS  
IS THE LOSS FUNCTION USED IN SVR ‘SQUAREDEPSILONINSENSITIVE’ IS THE SAME BUT BECOMES  
SQUARED LOSS PAST A TOLERANCE OF EPSILON  
PENALTY STR ‘NONE’ ‘L2’ ‘L1’ OR ‘ELASTICNET’ THE PENALTY AKA REGULARIZATION TERM TO BE USED  
DEFAULTS TO ‘L2’ WHICH IS THE STANDARD REGULARIZER FOR LINEAR SVM MODELS ‘L1’ AND ‘ELASTICNET’  
MIGHT BRING SPARSITY TO THE MODEL FEATURE SELECTION NOT ACHIEVABLE WITH ‘L2’  
ALPHA FLOAT CONSTANT THAT MULTIPLIES THE REGULARIZATION TERM DEFAULTS TO 00001 ALSO USED TO  
COMPUTE LEARNINGRATE WHEN SET TO ‘OPTIMAL’  
L1RATIO FLOAT THE ELASTIC NET MIXING PARAMETER WITH 0 L1RATIO 1 L1RATIO0 CORRE  
SPONDS TO L2 PENALTY L1RATIO1 TO L1 DEFAULTS TO 015  
FITINTERCEPT BOOL WHETHER THE INTERCEPT SHOULD BE ESTIMATED OR NOT IF FALSE THE DATA IS  
ASSUMED TO BE ALREADY CENTERED DEFAULTS TO TRUE  
MAXITER INT OPTIONAL DEFAULT1000 THE MAXIMUM NUMBER OF PASSES OVER THE TRAINING DATA  
AKA EPOCHS IT ONLY IMPACTS THE BEHAVIOR IN THE FIT METHOD AND NOT THE PARTIALFIT  
NEW IN VERSION 019  
TOLFLOAT OR NONE OPTIONAL DEFAULT1E3 THE STOPPING CRITERION IF IT IS NOT NONE THE ITERA  
TIONS WILL STOP WHEN LOSS BESTLOSS TOL FOR NITERNOCHANGE CONSECUTIVE EPOCHS  
NEW IN VERSION 019  
SHUFFLE BOOL OPTIONAL WHETHER OR NOT THE TRAINING DATA SHOULD BE SHUFFLED AFTER EACH EPOCH  
DEFAULTS TO TRUE  
VERBOSE INTEGER OPTIONAL THE VERBOSITY LEVEL

1946 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

EPSILON FLOAT EPSILON IN THE EPSILONINSENSITIVE LOSS FUNCTIONS ONLY IF LOSS IS 'HUBER' 'EP  
SILONINSENSITIVE' OR 'SQUAREDEPSILONINSENSITIVE' FOR 'HUBER' DETERMINES THE THRESH  
OLD AT WHICH IT BECOMES LESS IMPORTANT TO GET THE PREDICTION EXACTLY RIGHT FOR EPSILON  
INSENSITIVE ANY DIFFERENCES BETWEEN THE CURRENT PREDICTION AND THE CORRECT LABEL ARE IGNORED  
IF THEY ARE LESS THAN THIS THRESHOLD

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE THE SEED OF THE  
PSEUDO RANDOM NUMBER GENERATOR TO USE WHEN SHUFFLING THE DATA IF INT RANDOMSTATE IS  
THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS  
THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE  
INSTANCE USED BY NPRANDOM

LEARNINGRATE STRING OPTIONAL THE LEARNING RATE SCHEDULE

'CONSTANT' ETA ETA0

'OPTIMAL' ETA 10 ALPHA T T0 WHERE T0 IS CHOSEN BY A HEURISTIC PROPOSED BY LEON  
BOTTOU

'INVSCALING' DEFAULT ETA ETA0 POWT POWER

'ADAPTIVE' ETA ETA0 AS LONG AS THE TRAINING KEEPS DECREASING EACH TIME

NITERNOCHANGE CONSECUTIVE EPOCHS FAIL TO DECREASE THE TRAINING LOSS BY TOL OR FAIL TO  
INCREASE VALIDATION SCORE BY TOL IF EARLYSTOPPING IS TRUE THE CURRENT LEARNING RATE IS DI  
VIDED BY 5

ETA0 DOUBLE THE INITIAL LEARNING RATE FOR THE 'CONSTANT' 'INVSCALING' OR 'ADAPTIVE' SCHEDULES  
THE DEFAULT VALUE IS 001

POWERT DOUBLE THE EXPONENT FOR INVERSE SCALING LEARNING RATE DEFAULT 05

EARLYSTOPPING BOOL DEFAULTFALSE WHETHER TO USE EARLY STOPPING TO TERMINATE TRAINING WHEN  
VALIDATION SCORE IS NOT IMPROVING IF SET TO TRUE IT WILL AUTOMATICALLY SET ASIDE A FRACTION OF  
TRAINING DATA AS VALIDATION AND TERMINATE TRAINING WHEN VALIDATION SCORE IS NOT IMPROVING BY  
AT LEAST TOL FOR NITERNOCHANGE CONSECUTIVE EPOCHS

NEW IN VERSION 020

VALIDATIONFRACTION FLOAT DEFAULT01 THE PROPORTION OF TRAINING DATA TO SET ASIDE AS VALIDATION  
SET FOR EARLY STOPPING MUST BE BETWEEN 0 AND 1 ONLY USED IF EARLYSTOPPING IS TRUE

NEW IN VERSION 020

NITERNOCHANGE INT DEFAULT5 NUMBER OF ITERATIONS WITH NO IMPROVEMENT TO WAIT BEFORE  
EARLY STOPPING

NEW IN VERSION 020

WARMSTART BOOL OPTIONAL WHEN SET TO TRUE REUSE THE SOLUTION OF THE PREVIOUS CALL TO FIT AS  
INITIALIZATION OTHERWISE JUST ERASE THE PREVIOUS SOLUTION SEE THE GLOSSARY

REPEATEDLY CALLING FIT OR PARTIALFIT WHEN WARMSTART IS TRUE CAN RESULT IN A DIFFERENT SOLUTION  
THAN WHEN CALLING FIT A SINGLE TIME BECAUSE OF THE WAY THE DATA IS SHUFFLED IF A DYNAMIC  
LEARNING RATE IS USED THE LEARNING RATE IS ADAPTED DEPENDING ON THE NUMBER OF SAMPLES AL  
READY SEEN CALLING FIT RESETS THIS COUNTER WHILE PARTIALFIT WILL RESULT IN INCREASING  
THE EXISTING COUNTER

AVERAGE BOOL OR INT OPTIONAL WHEN SET TO TRUE COMPUTES THE AVERAGED SGD WEIGHTS AND  
STORES THE RESULT IN THE COEF ATTRIBUTE IF SET TO AN INT GREATER THAN 1 AVERAGING WILL BE  
GIN ONCE THE TOTAL NUMBER OF SAMPLES SEEN REACHES AVERAGE SO AVERAGE10 WILL BEGIN  
AVERAGING AFTER SEEING 10 SAMPLES

ATTRIBUTES

622SKLEARNLINEARMODEL GENERALIZED LINEAR MODELS 1947

SCIKITLEARN USER GUIDE RELEASE 0213

COEF ARRAY SHAPE NFEATURES WEIGHTS ASSIGNED TO THE FEATURES

INTERCEPT ARRAY SHAPE 1 THE INTERCEPT TERM

AVERAGECOEF ARRAY SHAPE NFEATURES AVERAGED WEIGHTS ASSIGNED TO THE FEATURES

AVERAGEINTERCEPT ARRAY SHAPE 1 THE AVERAGED INTERCEPT TERM

NITER INT THE ACTUAL NUMBER OF ITERATIONS TO REACH THE STOPPING CRITERION

SEE ALSO

RIDGE ELASTICNET LASSO SKLEARN SVMSVR

EXAMPLES

```
import numpy as np
from sklearn import linear_model
n_samples, n_features = 10, 5
rng = np.random.RandomState(0)
y = rng.randn(n_samples)
X = rng.randn(n_samples, n_features)
clf = linear_model.SGDRegressor(max_iter=1000, tol=1e-3,
                                clf_fit_x=y)
```

SGDRegressor(alpha=0.0001, average=False, early\_stopping=False, epsilon=0.1, eta=0.001, fit\_intercept=True, l1\_ratio=0.15, learning\_rate='invscaling', loss='squared\_loss', max\_iter=1000, n\_iter\_no\_change=5, penalty='l2', power\_t=0.25, random\_state=None, shuffle=True, tol=0.001, validation\_fraction=0.1, verbose=0, warm\_start=False)

METHODS

densify\_self: CONVERT COEFFICIENT MATRIX TO DENSE ARRAY FORMAT

fit\_self(x, y, coef\_init=None, intercept\_init=None): FIT LINEAR MODEL WITH STOCHASTIC GRADIENT DESCENT

get\_params(self, deep=True): GET PARAMETERS FOR THIS ESTIMATOR

partial\_fit(self, x, y, sample\_weight=None): PERFORM ONE EPOCH OF STOCHASTIC GRADIENT DESCENT ON GIVEN SAMPLES

predict(self, x): PREDICT USING THE LINEAR MODEL

score(self, x, y, sample\_weight=None): RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

set\_params(self, \*\*kwargs)

sparsify\_self: CONVERT COEFFICIENT MATRIX TO SPARSE FORMAT

init(self, loss='squared\_loss', penalty='l2', alpha=0.0001, l1\_ratio=0.15, fit\_intercept=True, max\_iter=1000, tol=0.001, shuffle=True, verbose=0, epsilon=0.1, random\_state=None, learning\_rate='invscaling', eta=0.001, power\_t=0.25, early\_stopping=False, validation\_fraction=0.1, n\_iter\_no\_change=5, warm\_start=False, average=False)

densify\_self

CONVERT COEFFICIENT MATRIX TO DENSE ARRAY FORMAT

converts the coef member back to a numpy ndarray. This is the default format of coef and is required for fitting so calling this method is only required on models that have previously been sparsified

1948 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213  
OTHERWISE IT IS A NOOP  
RETURNS  
SELF ESTIMATOR  
FITSELFXYCOEFINITNONE INTERCEPTINITNONE SAMPLEWEIGHTNONE  
FIT LINEAR MODEL WITH STOCHASTIC GRADIENT DESCENT  
PARAMETERS  
XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES TRAINING DATA  
YNUMPY ARRAY SHAPE NSAMPLES TARGET VALUES  
COEFINIT ARRAY SHAPE NFEATURES THE INITIAL COEFFICIENTS TO WARMSTART THE OPTIMIZATION  
INTERCEPTINIT ARRAY SHAPE 1 THE INITIAL INTERCEPT TO WARMSTART THE OPTIMIZATION  
SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL WEIGHTS APPLIED TO INDIVIDUAL  
SAMPLES 1 FOR UNWEIGHTED  
RETURNS  
SELF RETURNS AN INSTANCE OF SELF  
GETPARAMS SELFDEEPTURE  
GET PARAMETERS FOR THIS ESTIMATOR  
PARAMETERS  
DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBOBJECTS THAT ARE ESTIMATORS  
RETURNS  
PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES  
PARTIALFIT SELFXYSAMPLEWEIGHTNONE  
PERFORM ONE EPOCH OF STOCHASTIC GRADIENT DESCENT ON GIVEN SAMPLES  
INTERNALLY THIS METHOD USES MAXITER 1 THEREFORE IT IS NOT GUARANTEED THAT A MINIMUM OF THE COST  
FUNCTION IS REACHED AFTER CALLING IT ONCE MATTERS SUCH AS OBJECTIVE CONVERGENCE AND EARLY STOPPING SHOULD  
BE HANDLED BY THE USER  
PARAMETERS  
XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES SUBSET OF TRAINING DATA  
YNUMPY ARRAY OF SHAPE NSAMPLES SUBSET OF TARGET VALUES  
SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL WEIGHTS APPLIED TO INDIVIDUAL  
SAMPLES IF NOT PROVIDED UNIFORM WEIGHTS ARE ASSUMED  
RETURNS  
SELF RETURNS AN INSTANCE OF SELF  
PREDICTSELF  
PREDICT USING THE LINEAR MODEL  
PARAMETERS  
XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES  
RETURNS  
ARRAY SHAPE NSAMPLES PREDICTED TARGET VALUES PER ELEMENT IN X  
622SKLEARNLINEARMODEL GENERALIZED LINEAR MODELS 1949

SCIKITLEARN USER GUIDE RELEASE 0213

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{U}{V}$  WHERE U IS THE RESIDUAL SUM OF SQUARES  $Y_{TRUE} - Y_{PRED}$   $2SUM$  AND V IS THE TOTAL SUM OF SQUARES  $Y_{TRUE} - Y_{TRUEMEAN}$   $2SUM$  THE BEST POSSIBLE SCORE IS 10 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 00

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES FOR SOME ESTIMATORS THIS MAY BE A PRECOMPUTED KERNEL MATRIX INSTEAD SHAPE NSAMPLES NSAMPLESFITTED WHERE NSAMPLESFITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE VALUES FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT R2 OF SELF PREDICT X WRT Y

NOTES

THE R2 SCORE USED WHEN CALLING SCORE ON A REGRESSOR WILL USE MULTIOUTPUTUNIFORMAVERAGE FROM VERSION 023 TO KEEP CONSISTENT WITH METRICSR2SCORE THIS WILL INFLUENCE THE SCORE METHOD OF ALL THE MULTIOUTPUT REGRESSORS EXCEPT FOR MULTIOUTPUTMULTIOUTPUTREGRESSOR TO SPECIFY THE DEFAULT VALUE MANUALLY AND AVOID THE WARNING PLEASE EITHER CALL METRICSR2SCORE DIRECTLY OR MAKE A CUSTOM SCORER WITH METRICSMAKESCORER THE BUILTIN SCORER R2 USES MULTIOUTPUTUNIFORMAVERAGE

SPARSIFY SELF

CONVERT COEFFICIENT MATRIX TO SPARSE FORMAT

CONVERTS THE COEF MEMBER TO A SCIPYSPARSE MATRIX WHICH FOR L1REGULARIZED MODELS CAN BE MUCH MORE MEMORY AND STORAGEEFFICIENT THAN THE USUAL NUMPYNDARRAY REPRESENTATION

THEINTERCEPT MEMBER IS NOT CONVERTED

RETURNS

SELF ESTIMATOR

NOTES

FOR NONSPARSE MODELS IE WHEN THERE ARE NOT MANY ZEROS IN COEF THIS MAY ACTUALLY INCREASE MEMORY USAGE SO USE THIS METHOD WITH CARE A RULE OF THUMB IS THAT THE NUMBER OF ZERO ELEMENTS WHICH CAN BE COMPUTED WITH COEF 0SUM MUST BE MORE THAN 50 FOR THIS TO PROVIDE SIGNIFICANT BENEFITS

AFTER CALLING THIS METHOD FURTHER FITTING WITH THE PARTIALFIT METHOD IF ANY WILL NOT WORK UNTIL YOU CALL DENSIFY

EXAMPLES USING SKLEARNLINEARMODELSGREGRESSOR

- PREDICTION LATENCY

1950 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213  
62221SKLEARNLINEARMODEL THEILSENREGRESSOR  
CLASSSSKLEARNLINEARMODEL THEILSENREGRESSOR FITINTERCEPTTRUE COPYXTRUE  
MAXSUBPOPULATION100000  
NSUBSAMPLESNONE MAXITER300  
TOL0001 RANDOMSTATENONE  
NJOBSNONE VERBOSEFALSE  
THEILSEN ESTIMATOR ROBUST MULTIVARIATE REGRESSION MODEL  
THE ALGORITHM CALCULATES LEAST SQUARE SOLUTIONS ON SUBSETS WITH SIZE NSUBSAMPLES OF THE SAMPLES IN X ANY VALUE  
OF NSUBSAMPLES BETWEEN THE NUMBER OF FEATURES AND SAMPLES LEADS TO AN ESTIMATOR WITH A COMPROMISE BETWEEN  
ROBUSTNESS AND EFFICIENCY SINCE THE NUMBER OF LEAST SQUARE SOLUTIONS IS “NSAMPLES CHOOSE NSUBSAMPLES” IT  
CAN BE EXTREMELY LARGE AND CAN THEREFORE BE LIMITED WITH MAXSUBPOPULATION IF THIS LIMIT IS REACHED THE SUBSETS  
ARE CHOSEN RANDOMLY IN A FINAL STEP THE SPATIAL MEDIAN OR L1 MEDIAN IS CALCULATED OF ALL LEAST SQUARE SOLUTIONS  
READ MORE IN THE USER GUIDE  
PARAMETERS  
FITINTERCEPT BOOLEAN OPTIONAL DEFAULT TRUE WHETHER TO CALCULATE THE INTERCEPT FOR THIS MODEL  
IF SET TO FALSE NO INTERCEPT WILL BE USED IN CALCULATIONS  
COPYX BOOLEAN OPTIONAL DEFAULT TRUE IF TRUE X WILL BE COPIED ELSE IT MAY BE OVERWRITTEN  
MAXSUBPOPULATION INT OPTIONAL DEFAULT 1E4 INSTEAD OF COMPUTING WITH A SET OF CARDINALITY  
‘N CHOOSE K’ WHERE N IS THE NUMBER OF SAMPLES AND K IS THE NUMBER OF SUBSAMPLES AT  
LEAST NUMBER OF FEATURES CONSIDER ONLY A STOCHASTIC SUBPOPULATION OF A GIVEN MAXIMAL SIZE  
IF ‘N CHOOSE K’ IS LARGER THAN MAXSUBPOPULATION FOR OTHER THAN SMALL PROBLEM SIZES THIS  
PARAMETER WILL DETERMINE MEMORY USAGE AND RUNTIME IF NSUBSAMPLES IS NOT CHANGED  
NSUBSAMPLES INT OPTIONAL DEFAULT NONE NUMBER OF SAMPLES TO CALCULATE THE PARAMETERS  
THIS IS AT LEAST THE NUMBER OF FEATURES PLUS 1 IF FITINTERCEPTTRUE AND THE NUMBER OF SAM  
PLES AS A MAXIMUM A LOWER NUMBER LEADS TO A HIGHER BREAKDOWN POINT AND A LOW EFFICIENCY  
WHILE A HIGH NUMBER LEADS TO A LOW BREAKDOWN POINT AND A HIGH EFFICIENCY IF NONE TAKE THE  
MINIMUM NUMBER OF SUBSAMPLES LEADING TO MAXIMAL ROBUSTNESS IF NSUBSAMPLES IS SET TO  
NSAMPLES THEILSEN IS IDENTICAL TO LEAST SQUARES  
MAXITER INT OPTIONAL DEFAULT 300 MAXIMUM NUMBER OF ITERATIONS FOR THE CALCULATION OF SPA  
TIAL MEDIAN  
TOLFLOAT OPTIONAL DEFAULT 1E3 TOLERANCE WHEN CALCULATING SPATIAL MEDIAN  
RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULT NONE A RANDOM NUMBER  
GENERATOR INSTANCE TO DEFINE THE STATE OF THE RANDOM PERMUTATIONS GENERATOR IF INT RAN  
DOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RAN  
DOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE  
RANDOMSTATE INSTANCE USED BY NPRANDOM  
NJOBS INT OR NONE OPTIONAL DEFAULTNONE NUMBER OF CPUS TO USE DURING THE CROSS VALIDA  
TIONNONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING  
ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS  
VERBOSE BOOLEAN OPTIONAL DEFAULT FALSE VERBOSE MODE WHEN FITTING THE MODEL  
ATTRIBUTES  
COEF ARRAY SHAPE NFEATURES COEFFICIENTS OF THE REGRESSION MODEL MEDIAN OF DISTRIBUTION  
INTERCEPT FLOAT ESTIMATED INTERCEPT OF REGRESSION MODEL  
BREAKDOWN FLOAT APPROXIMATED BREAKDOWN POINT  
NITER INT NUMBER OF ITERATIONS NEEDED FOR THE SPATIAL MEDIAN  
622SKLEARNLINEARMODEL GENERALIZED LINEAR MODELS 1951

SCIKITLEARN USER GUIDE RELEASE 0213

NSUBPOPULATION INT NUMBER OF COMBINATIONS TAKEN INTO ACCOUNT FROM ‘N CHOOSE K’ WHERE  
N IS THE NUMBER OF SAMPLES AND K IS THE NUMBER OF SUBSAMPLES

REFERENCES

- THEILSEN ESTIMATORS IN A MULTIPLE LINEAR REGRESSION MODEL 2009 XIN DANG HANXIANG PENG XUEQIN  
WANG AND HEPING ZHANG HTTPHOMEOLEMISSEDUXDANGPAPERSMTSEP

EXAMPLES

```
FROM SKLEARNLINEARMODEL IMPORT THEILSENREGRESSOR
FROM SKLEARNDATASETS IMPORT MAKEREGRESSION
X Y MAKEREGRESSION
NSAMPLES200 NFEATURES2 NOISE40 RANDOMSTATE0
REG THEILSENREGRESSORRANDOMSTATE0FITX Y
REGSCOREX Y
09884
REGPREDICTX1
ARRAY315871
METHODS
FITSELF X Y FIT LINEAR MODEL
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR
PREDICT SELF X PREDICT USING THE LINEAR MODEL
SCORE SELF X Y SAMPLEWEIGHT RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PRE
DICTION
SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR
INIT SELF FITINTERCEPTTRUE COPYXTRUE MAXSUBPOPULATION100000 NSUBSAMPLESNONE
MAXITER300 TOL0001 RANDOMSTATENONE NJOBSNONE VERBOSEFALSE
FITSELFXY
FIT LINEAR MODEL
PARAMETERS
XNUMPY ARRAY OF SHAPE NSAMPLES NFEATURES TRAINING DATA
YNUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES
RETURNS
SELF RETURNS AN INSTANCE OF SELF
GETPARAMS SELFDEEPTTRUE
GET PARAMETERS FOR THIS ESTIMATOR
PARAMETERS
DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED
SUBOBJECTS THAT ARE ESTIMATORS
RETURNS
```

1952 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELF

PREDICT USING THE LINEAR MODEL

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES SAMPLES

RETURNS

CARRAY SHAPE NSAMPLES RETURNS PREDICTED VALUES

SCORESELF

SCORESELF SAMPLEWEIGHT NONE

RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{U}{V}$  WHERE U IS THE RESIDUAL SUM OF SQUARES  $\frac{1}{n} \sum (y - \hat{y})^2$  YTRUE YPRED 2SUM AND V IS THE TOTAL SUM OF SQUARES  $\frac{1}{n} \sum (y - \bar{y})^2$  YTRUE YTRUEMEAN 2SUM THE BEST POSSIBLE SCORE IS 1.0 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 0.0

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES FOR SOME ESTIMATORS THIS MAY BE A PRECOMPUTED KERNEL MATRIX INSTEAD SHAPE NSAMPLES NSAMPLESFITTED WHERE NSAMPLESFITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE VALUES FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT R2 OF SELF PREDICT X WRT Y

NOTES

THE R2 SCORE USED WHEN CALLING SCORE ON A REGRESSOR WILL USE MULTIOUTPUT UNIFORM AVERAGE FROM VERSION 0.23 TO KEEP CONSISTENT WITH METRICS R2 SCORE THIS WILL INFLUENCE THE SCORE METHOD OF ALL THE MULTIOUTPUT REGRESSORS EXCEPT FOR MULTIOUTPUT MULTICLASS REGRESSOR TO SPECIFY THE DEFAULT VALUE MANUALLY AND AVOID THE WARNING PLEASE EITHER CALL METRICS R2 SCORE DIRECTLY OR MAKE A CUSTOM SCORER WITH METRICS MAKE SCORER THE BUILT IN SCORER R2 USES MULTIOUTPUT UNIFORM AVERAGE

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENT PARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

EXAMPLES USING SKLEARN LINEAR MODEL THEILSEN REGRESSOR

- THEILSEN REGRESSION
- ROBUST LINEAR ESTIMATOR FITTING

622 SKLEARN LINEAR MODEL GENERALIZED LINEAR MODELS 1953

SCIKITLEARN USER GUIDE RELEASE 0213

LINEARMODELENETPATH X Y L1RATIO    COMPUTE ELASTIC NET PATH WITH COORDINATE DESCENT

LINEARMODELLARSPATH X Y XY GRAM    COMPUTE LEAST ANGLE REGRESSION OR LASSO PATH USING

LARS ALGORITHM 1

LINEARMODELLARSPATHGRAM XY GRAM

NSAMPLESLARSPATH IN THE SUFFICIENT STATS MODE 1

LINEARMODELLASSOPATH X Y EPS    COMPUTE LASSO PATH WITH COORDINATE DESCENT

LINEARMODELORTHOGONALMP X Y    ORTHOGONAL MATCHING PURSUIT OMP

LINEARMODELORTHOGONALMPGRAM GRAM XY

GRAM ORTHOGONAL MATCHING PURSUIT OMP

LINEARMODELRIDGEREGRESSION X Y ALPHA

SOLVE THE RIDGE EQUATION BY THE METHOD OF NORMAL EQUATIONS

62222SKLEARNLINEARMODEL ENETPATH

SKLEARNLINEARMODEL ENETPATH XYL1RATIO05 EPS0001 NALPHAS100 ALPHASNONE

PRECOMPUTE'AUTO' XYNONE COPYXTRUE COEFINITNONE

VERBOSEFALSE RETURNNITERFALSE POSITIVEFALSE

CHECKINPUTTRUE PARAMS

COMPUTE ELASTIC NET PATH WITH COORDINATE DESCENT

THE ELASTIC NET OPTIMIZATION FUNCTION VARIES FOR MONO AND MULTIOUTPUTS

FOR MONOOUTPUT TASKS IT IS

1 2NSAMPLES Y XW22

ALPHA L1RATIO W1

05ALPHA1 L1RATIO W22

FOR MULTIOUTPUT TASKS IT IS

1 2 NSAMPLES Y XWFRO2

ALPHA L1RATIO W21

05ALPHA1 L1RATIO WFO2

WHERE

W21 SUMI SQRTSUMJ WIJ2

IE THE SUM OF NORM OF EACH ROW

READ MORE IN THE USER GUIDE

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA PASS DIRECTLY AS FORTRAN

CONTIGUOUS DATA TO AVOID UNNECESSARY MEMORY DUPLICATION IF YIS MONOOUTPUT THEN XCAN

BE SPARSE

YNDARRAY SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TARGET VALUES

L1RATIO FLOAT OPTIONAL FLOAT BETWEEN 0 AND 1 PASSED TO ELASTIC NET SCALING BETWEEN L1 AND L2

PENALTIESL1RATIO1 CORRESPONDS TO THE LASSO

EPS FLOAT LENGTH OF THE PATH EPS1E3 MEANS THAT ALPHAMIN ALPHAMAX

1E3

NALPHAS INT OPTIONAL NUMBER OF ALPHAS ALONG THE REGULARIZATION PATH

1954 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

ALPHAS NDARRAY OPTIONAL LIST OF ALPHAS WHERE TO COMPUTE THE MODELS IF NONE ALPHAS ARE SET AUTOMATICALLY

PRECOMPUTE TRUE FALSE 'AUTO' ARRAYLIKE WHETHER TO USE A PRECOMPUTED GRAM MATRIX TO SPEED UP CALCULATIONS IF SET TO AUTO LET US DECIDE THE GRAM MATRIX CAN ALSO BE PASSED AS ARGUMENT

XYARRAYLIKE OPTIONAL XY NPDOTXT Y THAT CAN BE PRECOMPUTED IT IS USEFUL ONLY WHEN THE GRAM MATRIX IS PRECOMPUTED

COPYX BOOLEAN OPTIONAL DEFAULT TRUE IF TRUE X WILL BE COPIED ELSE IT MAY BE OVERWRITTEN

COEFINIT ARRAY SHAPE NFEATURES NONE THE INITIAL VALUES OF THE COEFFICIENTS

VERBOSE BOOL OR INTEGER AMOUNT OF VERBOSITY

RETURNNITER BOOL WHETHER TO RETURN THE NUMBER OF ITERATIONS OR NOT

POSITIVE BOOL DEFAULT FALSE IF SET TO TRUE FORCES COEFFICIENTS TO BE POSITIVE ONLY ALLOWED WHENYNDIM 1

CHECKINPUT BOOL DEFAULT TRUE SKIP INPUT VALIDATION CHECKS INCLUDING THE GRAM MATRIX WHEN PROVIDED ASSUMING THERE ARE HANDLED BY THE CALLER WHEN CHECKINPUTFALSE

PARAMS KWARGS KEYWORD ARGUMENTS PASSED TO THE COORDINATE DESCENT SOLVER

RETURNS

ALPHAS ARRAY SHAPE NALPHAS THE ALPHAS ALONG THE PATH WHERE MODELS ARE COMPUTED

COEFS ARRAY SHAPE NFEATURES NALPHAS OR NOUTPUTS NFEATURES NALPHAS COEFFICIENTS ALONG THE PATH

DUALGAPS ARRAY SHAPE NALPHAS THE DUAL GAPS AT THE END OF THE OPTIMIZATION FOR EACH ALPHA

NITERS ARRAYLIKE SHAPE NALPHAS THE NUMBER OF ITERATIONS TAKEN BY THE COORDINATE DESCENT OPTIMIZER TO REACH THE SPECIFIED TOLERANCE FOR EACH ALPHA IS RETURNED WHEN RETURNNITER IS SET TO TRUE

SEE ALSO

MULTITASKELASTICNET

MULTITASKELASTICNETCV

ELASTICNET

ELASTICNETCV

NOTES

FOR AN EXAMPLE SEE EXAMPLESLINEARMODELPLOTLASSOCOORDINATEDESCENTPATHPY

EXAMPLES USING SKLEARNLINEARMODELENETPATH

•LASSO AND ELASTIC NET

622SKLEARNLINEARMODEL GENERALIZED LINEAR MODELS 1955

SCIKITLEARN USER GUIDE RELEASE 0213  
62223SKLEARNLINEARMODEL LARSPATH  
SKLEARNLINEARMODEL LARSPATH XYXYNONE GRAMNONE MAXITER500 ALPHAMINO  
METHOD'LAR' COPYXTRUE EPS2220446049250313E  
16COPYGRAMTRUE VERBOSE0 RETURNPATHTRUE RE  
TURNNNITERFALSE POSITIVEFALSE  
COMPUTE LEAST ANGLE REGRESSION OR LASSO PATH USING LARS ALGORITHM 1  
THE OPTIMIZATION OBJECTIVE FOR THE CASE METHOD'LASSO' IS  
1 2 NSAMPLES Y XW22 ALPHA W1  
IN THE CASE OF METHOD'LARS' THE OBJECTIVE FUNCTION IS ONLY KNOWN IN THE FORM OF AN IMPLICIT EQUATION SEE  
DISCUSSION IN 1  
READ MORE IN THE USER GUIDE  
PARAMETERS  
XNONE OR ARRAY SHAPE NSAMPLES NFEATURES INPUT DATA NOTE THAT IF X IS NONE THEN THE  
GRAM MATRIX MUST BE SPECIFIED IE CANNOT BE NONE OR FALSE  
DEPRECATED SINCE VERSION 021 THE USE OF XISNONE IN COMBINATION WITH GRAM IS NOT  
NONE WILL BE REMOVED IN V023 USE LARSPATHGRAM INSTEAD  
YNONE OR ARRAY SHAPE NSAMPLES INPUT TARGETS  
XYARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NTARGETS OPTIONAL XY NPDOTXT Y THAT  
CAN BE PRECOMPUTED IT IS USEFUL ONLY WHEN THE GRAM MATRIX IS PRECOMPUTED  
GRAM NONE 'AUTO' ARRAY SHAPE NFEATURES NFEATURES OPTIONAL PRECOMPUTED GRAM MATRIX  
X' X IF AUTO THE GRAM MATRIX IS PRECOMPUTED FROM THE GIVEN X IF THERE ARE MORE  
SAMPLES THAN FEATURES  
DEPRECATED SINCE VERSION 021 THE USE OF XISNONE IN COMBINATION WITH GRAM IS NOT NONE  
WILL BE REMOVED IN V023 USE LARSPATHGRAM INSTEAD  
MAXITER INTEGER OPTIONAL DEFAULT500 MAXIMUM NUMBER OF ITERATIONS TO PERFORM SET TO  
INFINITY FOR NO LIMIT  
ALPHAMIN FLOAT OPTIONAL DEFAULT0 MINIMUM CORRELATION ALONG THE PATH IT CORRESPONDS TO  
THE REGULARIZATION PARAMETER ALPHA PARAMETER IN THE LASSO  
METHOD 'LAR' 'LASSO' OPTIONAL DEFAULT'LAR' SPECIFIES THE RETURNED MODEL SELECT LAR  
FOR LEAST ANGLE REGRESSION LASSO FOR THE LASSO  
COPYX BOOL OPTIONAL DEFAULTTRUE IF FALSE XIS OVERWRITTEN  
EPS FLOAT OPTIONAL DEFAULT''NPFINFONPFLOATEPS'' THE MACHINEPRECISION REGULARIZATION IN  
THE COMPUTATION OF THE CHOLESKY DIAGONAL FACTORS INCREASE THIS FOR VERY ILLCONDITIONED  
SYSTEMS  
COPYGRAM BOOL OPTIONAL DEFAULTTRUE IF FALSE GRAM IS OVERWRITTEN  
VERBOSE INT DEFAULT0 CONTROLS OUTPUT VERBOSITY  
RETURNPATH BOOL OPTIONAL DEFAULTTRUE IF RETURNPATHTRUE RETURNS THE ENTIRE PATH  
ELSE RETURNS ONLY THE LAST POINT OF THE PATH  
RETURNNNITER BOOL OPTIONAL DEFAULTFALSE WHETHER TO RETURN THE NUMBER OF ITERATIONS  
POSITIVE BOOLEAN DEFAULTFALSE RESTRICT COEFFICIENTS TO BE 0 THIS OPTION IS ONLY ALLOWED  
WITH METHOD 'LASSO' NOTE THAT THE MODEL COEFFICIENTS WILL NOT CONVERGE TO THE ORDINARY  
LEASTSQUARES SOLUTION FOR SMALL VALUES OF ALPHA ONLY COEFFICIENTS UP TO THE SMALLEST ALPHA  
1956 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

VALUE ALPHASALPHAS 0MIN WHEN FITPATHTRUE REACHED BY THE STEP

WISE LARSLASSO ALGORITHM ARE TYPICALLY IN CONGRUENCE WITH THE SOLUTION OF THE COORDINATE

DESCENT LASSOPATH FUNCTION

RETURNS

ALPHAS ARRAY SHAPE NALPHAS 1 MAXIMUM OF COVARIANCES IN ABSOLUTE VALUE AT EACH ITER

ATIONNALPHAS IS EITHERMAXITER NFEATURES OR THE NUMBER OF NODES IN THE PATH

WITHALPHA ALPHAMIN WHICHEVER IS SMALLER

ACTIVE ARRAY SHAPE NALPHAS INDICES OF ACTIVE VARIABLES AT THE END OF THE PATH

COEFS ARRAY SHAPE NFEATURES NALPHAS 1 COEFFICIENTS ALONG THE PATH

NITER INT NUMBER OF ITERATIONS RUN RETURNED ONLY IF RETURNNNITER IS SET TO TRUE

SEE ALSO

LARSPATHGRAM

LASSOPATH

LASSOPATHGRAM

LASSOLARS

LARS

LASSOLARSCV

LARSCV

SKLEARNDECOMPOSITIONSPARSEENCODE

REFERENCES

123

EXAMPLES USING SKLEARNLINEARMODELLARSPATH

- LASSO PATH USING LARS

62224SKLEARNLINEARMODEL LARSPATHGRAM

SKLEARNLINEARMODEL LARSPATHGRAM XY GRAM NSAMPLES MAXITER500 AL

PHAMINO METHOD'LAR' COPYXTRUE

EPS2220446049250313E16 COPYGRAMTRUE

VERBOSE0 RETURNPATHTRUE RETURNNNITERFALSE

POSITIVEFALSE

LARSPATH IN THE SUFFICIENT STATS MODE 1

THE OPTIMIZATION OBJECTIVE FOR THE CASE METHOD'LASSO' IS

1 2 NSAMPLES Y XW22 ALPHA W1

IN THE CASE OF METHOD'LARS' THE OBJECTIVE FUNCTION IS ONLY KNOWN IN THE FORM OF AN IMPLICIT EQUATION SEE

DISCUSSION IN 1

READ MORE IN THE USER GUIDE

622SKLEARNLINEARMODEL GENERALIZED LINEAR MODELS 1957

SCIKITLEARN USER GUIDE RELEASE 0213

PARAMETERS

XYARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NTARGETS XY NPDOTXT Y  
GRAM ARRAY SHAPE NFEATURES NFEATURES GRAM NPDOTXT X  
NSAMPLES INTEGER OR FLOAT EQUIVALENT SIZE OF SAMPLE  
MAXITER INTEGER OPTIONAL DEFAULT500 MAXIMUM NUMBER OF ITERATIONS TO PERFORM SET TO  
INFINITY FOR NO LIMIT  
ALPHAMIN FLOAT OPTIONAL DEFAULT0 MINIMUM CORRELATION ALONG THE PATH IT CORRESPONDS TO  
THE REGULARIZATION PARAMETER ALPHA PARAMETER IN THE LASSO  
METHOD 'LAR' 'LASSO' OPTIONAL DEFAULT'LAR' SPECIFIES THE RETURNED MODEL SELECT LAR  
FOR LEAST ANGLE REGRESSION LASSO FOR THE LASSO  
COPYX BOOL OPTIONAL DEFAULTTRUE IF FALSE XIS OVERWRITTEN  
EPS FLOAT OPTIONAL DEFAULT''NPFINFONPFLOATEPS'' THE MACHINEPRECISION REGULARIZATION IN  
THE COMPUTATION OF THE CHOLSKY DIAGONAL FACTORS INCREASE THIS FOR VERY ILLCONDITIONED  
SYSTEMS  
COPYGRAM BOOL OPTIONAL DEFAULTTRUE IF FALSE GRAM IS OVERWRITTEN  
VERBOSE INT DEFAULT0 CONTROLS OUTPUT VERBOSITY  
RETURNPATH BOOL OPTIONAL DEFAULTTRUE IF RETURNPATHTRUE RETURNS THE ENTIRE PATH  
ELSE RETURNS ONLY THE LAST POINT OF THE PATH  
RETURNNITER BOOL OPTIONAL DEFAULTFALSE WHETHER TO RETURN THE NUMBER OF ITERATIONS  
POSITIVE BOOLEAN DEFAULTFALSE RESTRICT COEFFICIENTS TO BE 0 THIS OPTION IS ONLY ALLOWED  
WITH METHOD 'LASSO' NOTE THAT THE MODEL COEFFICIENTS WILL NOT CONVERGE TO THE ORDINARY  
LEASTSQUARES SOLUTION FOR SMALL VALUES OF ALPHA ONLY COEFFICIENTS UP TO THE SMALLEST ALPHA  
VALUE ALPHASALPHAS 0MIN WHEN FITPATHTRUE REACHED BY THE STEP  
WISE LARSLASSO ALGORITHM ARE TYPICALLY IN CONGRUENCE WITH THE SOLUTION OF THE COORDINATE  
DESCENT LASSOPATH FUNCTION  
RETURNS  
ALPHAS ARRAY SHAPE NALPHAS 1 MAXIMUM OF COVARIANCES IN ABSOLUTE VALUE AT EACH ITER  
ATIONNNALPHAS IS EITHERMAXITER NFEATURES OR THE NUMBER OF NODES IN THE PATH  
WITHALPHA ALPHAMIN WHICHEVER IS SMALLER  
ACTIVE ARRAY SHAPE NALPHAS INDICES OF ACTIVE VARIABLES AT THE END OF THE PATH  
COEFS ARRAY SHAPE NFEATURES NALPHAS 1 COEFFICIENTS ALONG THE PATH  
NITER INT NUMBER OF ITERATIONS RUN RETURNED ONLY IF RETURNNITER IS SET TO TRUE  
SEE ALSO  
LARSPATH  
LASSOPATH  
LASSOPATHGRAM  
LASSOLARS  
LARS  
LASSOLARSCV  
LARSCV  
1958 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213  
SKLEARNDECOMPOSITIONSPARSEENCODE  
REFERENCES  
123  
62225SKLEARNLINEARMODEL LASSOPATH  
SKLEARNLINEARMODEL LASSOPATH XYEPS0001 NALPHAS100 ALPHASNONE PRECOM  
PUTE'AUTO' XYNONE COPYXTRUE COEFINITNONE VER  
BOSEFALSE RETURNNITERFALSE POSITIVEFALSE PARAMS  
COMPUTE LASSO PATH WITH COORDINATE DESCENT  
THE LASSO OPTIMIZATION FUNCTION VARIES FOR MONO AND MULTIOUTPUTS  
FOR MONOOUTPUT TASKS IT IS  
1 2 NSAMPLES Y XW22 ALPHA W1  
FOR MULTIOUTPUT TASKS IT IS  
1 2 NSAMPLES Y XW2FRO ALPHA W21  
WHERE  
W21 SUMI SQRTSUMJ WIJ2  
IE THE SUM OF NORM OF EACH ROW  
READ MORE IN THE USER GUIDE  
PARAMETERS  
XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES TRAINING DATA PASS DIRECTLY AS  
FORTRANCONTIGUOUS DATA TO AVOID UNNECESSARY MEMORY DUPLICATION IF YIS MONOOUTPUT THEN  
XCAN BE SPARSE  
YNDARRAY SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TARGET VALUES  
EPS FLOAT OPTIONAL LENGTH OF THE PATH EPS1E3 MEANS THATALPHAMIN ALPHAMAX  
1E3  
NALPHAS INT OPTIONAL NUMBER OF ALPHAS ALONG THE REGULARIZATION PATH  
ALPHAS NDARRAY OPTIONAL LIST OF ALPHAS WHERE TO COMPUTE THE MODELS IF NONE ALPHAS ARE SET  
AUTOMATICALLY  
PRECOMPUTE TRUE FALSE 'AUTO' ARRAYLIKE WHETHER TO USE A PRECOMPUTED GRAM MATRIX TO  
SPEED UP CALCULATIONS IF SET TO AUTO LET US DECIDE THE GRAM MATRIX CAN ALSO BE PASSED  
AS ARGUMENT  
XYARRAYLIKE OPTIONAL XY NPDOTXT Y THAT CAN BE PRECOMPUTED IT IS USEFUL ONLY WHEN  
THE GRAM MATRIX IS PRECOMPUTED  
COPYX BOOLEAN OPTIONAL DEFAULT TRUE IF TRUE X WILL BE COPIED ELSE IT MAY BE OVERWRITTEN  
COEFINIT ARRAY SHAPE NFEATURES NONE THE INITIAL VALUES OF THE COEFFICIENTS  
VERBOSE BOOL OR INTEGER AMOUNT OF VERBOSITY  
RETURNNITER BOOL WHETHER TO RETURN THE NUMBER OF ITERATIONS OR NOT  
622SKLEARNLINEARMODEL GENERALIZED LINEAR MODELS 1959

SCIKITLEARN USER GUIDE RELEASE 0213

POSITIVE BOOL DEFAULT FALSE IF SET TO TRUE FORCES COEFFICIENTS TO BE POSITIVE ONLY ALLOWED WHEN YNDIM = 1

PARAMS KWARGS KEYWORD ARGUMENTS PASSED TO THE COORDINATE DESCENT SOLVER

RETURNS

ALPHAS ARRAY SHAPE NALPHAS THE ALPHAS ALONG THE PATH WHERE MODELS ARE COMPUTED

COEFS ARRAY SHAPE NFEATURES NALPHAS OR NOUTPUTS NFEATURES NALPHAS COEFFICIENTS ALONG THE PATH

DUALGAPS ARRAY SHAPE NALPHAS THE DUAL GAPS AT THE END OF THE OPTIMIZATION FOR EACH ALPHA

NITERS ARRAYLIKE SHAPE NALPHAS THE NUMBER OF ITERATIONS TAKEN BY THE COORDINATE DESCENT OPTIMIZER TO REACH THE SPECIFIED TOLERANCE FOR EACH ALPHA

SEE ALSO

LARSPATH

LASSO

LASSOLARS

LASSOCV

LASSOLARSCV

SKLEARNDECOMPOSITIONSPARSEENCODE

NOTES

FOR AN EXAMPLE SEE EXAMPLESLINEARMODELPLOTLASSOCOORDINATEDDESCENTPATHPY

TO AVOID UNNECESSARY MEMORY DUPLICATION THE X ARGUMENT OF THE FIT METHOD SHOULD BE DIRECTLY PASSED AS A FORTRANCONTIGUOUS NUMPY ARRAY

NOTE THAT IN CERTAIN CASES THE LARS SOLVER MAY BE SIGNIFICANTLY FASTER TO IMPLEMENT THIS FUNCTIONALITY IN PARTICULAR LINEAR INTERPOLATION CAN BE USED TO RETRIEVE MODEL COEFFICIENTS BETWEEN THE VALUES OUTPUT BY LARSPATH

EXAMPLES

COMPARING LASSOPATH AND LARSPATH WITH INTERPOLATION

```
X = np.array([2, 31, 23, 54, 43])
Y = np.array([2, 31])

# USE LASSOPATH TO COMPUTE A COEFFICIENT PATH
coef_path = lassopath(X, Y, alphas=5)
print(coef_path)
# 0.0 0.46874778
# 0.2159048 0.4425765 0.23689075

# NOW USE LARSPATH AND 1D LINEAR INTERPOLATION TO COMPUTE THE SAME PATH
from sklearn.linear_model import LARSPath
alphas_active, coef_path_lars = larspath(X, Y, method='lasso')
from scipy import interpolate
coef_path_continuous = interpolate.interp1d(alphas_active, coef_path_lars)
```

1960 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213  
COEFFPATHLARS 1  
PRINTCOEFFPATHCONTINUOUS 1 5  
0 0 046915237  
02159048 04425765 023668876  
EXAMPLES USING SKLEARNLINEARMODELLASSOPATH  
•LASSO AND ELASTIC NET  
62226SKLEARNLINEARMODEL ORTHOGONALMP  
SKLEARNLINEARMODEL ORTHOGONALMP XYNONZEROCOEFSSNONE TOLNONE PRECOM  
PUTEFALSE COPYXTRUE RETURNPATHFALSE RE  
TURNNITERFALSE  
ORTHOGONAL MATCHING PURSUIT OMP  
SOLVES NTARGETS ORTHOGONAL MATCHING PURSUIT PROBLEMS AN INSTANCE OF THE PROBLEM HAS THE FORM  
WHEN PARAMETRIZED BY THE NUMBER OF NONZERO COEFFICIENTS USING NNONZEROCOEFSS ARGMIN Y  
XGAMMA2 SUBJECT TO GAMMA0 NNONZERO COEFS  
WHEN PARAMETRIZED BY ERROR USING THE PARAMETER TOL ARGMIN GAMMA0 SUBJECT TO Y XGAMMA2 TOL  
READ MORE IN THE USER GUIDE  
PARAMETERS  
XARRAY SHAPE NSAMPLES NFEATURES INPUT DATA COLUMNS ARE ASSUMED TO HAVE UNIT NORM  
YARRAY SHAPE NSAMPLES OR NSAMPLES NTARGETS INPUT TARGETS  
NNONZEROCOEFSS INT DESIRED NUMBER OF NONZERO ENTRIES IN THE SOLUTION IF NONE BY DEFAULT  
THIS VALUE IS SET TO 10 OF NFEATURES  
TOLFLOAT MAXIMUM NORM OF THE RESIDUAL IF NOT NONE OVERRIDES NNONZEROCOEFSS  
PRECOMPUTE TRUE FALSE 'AUTO' WHETHER TO PERFORM PRECOMPUTATIONS IMPROVES PERFOR  
MANCE WHEN NTARGETS OR NSAMPLES IS VERY LARGE  
COPYX BOOL OPTIONAL WHETHER THE DESIGN MATRIX X MUST BE COPIED BY THE ALGORITHM A FALSE  
VALUE IS ONLY HELPFUL IF X IS ALREADY FORTRANORDERED OTHERWISE A COPY IS MADE ANYWAY  
RETURNPATH BOOL OPTIONAL DEFAULT FALSE WHETHER TO RETURN EVERY VALUE OF THE NONZERO  
COEFFICIENTS ALONG THE FORWARD PATH USEFUL FOR CROSSVALIDATION  
RETURNNITER BOOL OPTIONAL DEFAULT FALSE WHETHER OR NOT TO RETURN THE NUMBER OF ITERATIONS  
RETURNS  
COEF ARRAY SHAPE NFEATURES OR NFEATURES NTARGETS COEFFICIENTS OF THE OMP SOLUTION IF  
RETURNPATHTRUE THIS CONTAINS THE WHOLE COEFFICIENT PATH IN THIS CASE ITS SHAPE IS  
NFEATURES NFEATURES OR NFEATURES NTARGETS NFEATURES AND ITERATING OVER THE LAST AXIS  
YIELDS COEFFICIENTS IN INCREASING ORDER OF ACTIVE FEATURES  
NITERS ARRAYLIKE OR INT NUMBER OF ACTIVE FEATURES ACROSS EVERY TARGET RETURNED ONLY IF  
RETURNNITER IS SET TO TRUE  
SEE ALSO  
ORTHOGONALMATCHINGPURSUIT  
622SKLEARNLINEARMODEL GENERALIZED LINEAR MODELS 1961

SCIKITLEARN USER GUIDE RELEASE 0213  
ORTHOGONALMPGRAM  
LARSPATH  
DECOMPOSITIONSPARSEENCODE  
NOTES  
ORTHOGONAL MATCHING PURSUIT WAS INTRODUCED IN S MALLAT Z ZHANG MATCHING PURSUITS WITH TIMEFREQUENCY  
DICTIONARIES IEEE TRANSACTIONS ON SIGNAL PROCESSING V OL 41 NO 12 DECEMBER 1993 PP 33973415  
HTTPBLANCHEPOLYTECHNIQUEFRMALLATPAPIERSMALLATPURSUIT93PDF  
THIS IMPLEMENTATION IS BASED ON RUBINSTEIN R ZIBULEVSKY M AND ELAD M EFFICIENT IMPLEMENTATION OF  
THE KSVD ALGORITHM USING BATCH ORTHOGONAL MATCHING PURSUIT TECHNICAL REPORT CS TECHNION APRIL 2008  
HTTPSWWWCSTECHNIONACILRONRUBINPUBLICATIONSKSVDOMPV2PDF  
62227SKLEARNLINEARMODEL ORTHOGONALMPGRAM  
SKLEARNLINEARMODEL ORTHOGONALMPGRAM GRAM XYNONZEROCOEFSSNONE TOLNONE  
NORMSSQUAREDNONE COPYGRAMTRUE  
COPYXYTRUE RETURNPATHFALSE RE  
TURNNNITERFALSE  
GRAM ORTHOGONAL MATCHING PURSUIT OMP  
SOLVES NTARGETS ORTHOGONAL MATCHING PURSUIT PROBLEMS USING ONLY THE GRAM MATRIX XT X AND THE PRODUCT  
XT Y  
READ MORE IN THE USER GUIDE  
PARAMETERS  
GRAM ARRAY SHAPE NFEATURES NFEATURES GRAM MATRIX OF THE INPUT DATA XT X  
XYARRAY SHAPE NFEATURES OR NFEATURES NTARGETS INPUT TARGETS MULTIPLIED BY X XT Y  
NNONZEROCOEFSS INT DESIRED NUMBER OF NONZERO ENTRIES IN THE SOLUTION IF NONE BY DEFAULT  
THIS VALUE IS SET TO 10 OF NFEATURES  
TOLFLOAT MAXIMUM NORM OF THE RESIDUAL IF NOT NONE OVERRIDES NNONZEROCOEFSS  
NORMSSQUARED ARRAYLIKE SHAPE NTARGETS SQUARED L2 NORMS OF THE LINES OF Y REQUIRED  
IF TOL IS NOT NONE  
COPYGRAM BOOL OPTIONAL WHETHER THE GRAM MATRIX MUST BE COPIED BY THE ALGORITHM A FALSE  
VALUE IS ONLY HELPFUL IF IT IS ALREADY FORTRANORDERED OTHERWISE A COPY IS MADE ANYWAY  
COPYXY BOOL OPTIONAL WHETHER THE COVARIANCE VECTOR XY MUST BE COPIED BY THE ALGORITHM  
IF FALSE IT MAY BE OVERWRITTEN  
RETURNPATH BOOL OPTIONAL DEFAULT FALSE WHETHER TO RETURN EVERY VALUE OF THE NONZERO  
COEFFICIENTS ALONG THE FORWARD PATH USEFUL FOR CROSSVALIDATION  
RETURNNNITER BOOL OPTIONAL DEFAULT FALSE WHETHER OR NOT TO RETURN THE NUMBER OF ITERATIONS  
RETURNS  
COEF ARRAY SHAPE NFEATURES OR NFEATURES NTARGETS COEFFICIENTS OF THE OMP SOLUTION IF  
RETURNPATHTRUE THIS CONTAINS THE WHOLE COEFFICIENT PATH IN THIS CASE ITS SHAPE IS  
NFEATURES NFEATURES OR NFEATURES NTARGETS NFEATURES AND ITERATING OVER THE LAST AXIS  
YIELDS COEFFICIENTS IN INCREASING ORDER OF ACTIVE FEATURES  
1962 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213  
NITERS ARRAYLIKE OR INT NUMBER OF ACTIVE FEATURES ACROSS EVERY TARGET RETURNED ONLY IF  
RETURNNITER IS SET TO TRUE  
SEE ALSO  
ORTHOGONALMATCHINGPURSUIT  
ORTHOGONALMP  
LARSPATH  
DECOMPOSITIONSPARSEENCODE  
NOTES  
ORTHOGONAL MATCHING PURSUIT WAS INTRODUCED IN G MALLAT Z ZHANG MATCHING PURSUITS WITH TIMEFREQUENCY  
DICTIONARIES IEEE TRANSACTIONS ON SIGNAL PROCESSING V OL 41 NO 12 DECEMBER 1993 PP 33973415  
HTTPBLANCHEPOLYTECHNIQUEFRMALLATPAPIERSMALLATPURSUIT93PDF  
THIS IMPLEMENTATION IS BASED ON RUBINSTEIN R ZIBULEVSKY M AND ELAD M EFFICIENT IMPLEMENTATION OF  
THE KSVD ALGORITHM USING BATCH ORTHOGONAL MATCHING PURSUIT TECHNICAL REPORT CS TECHNION APRIL 2008  
HTTPWWWCSTECHNIONACILRONRUBINPUBLICATIONSKSVDOMPV2PDF  
62228SKLEARNLINEARMODEL RIDGEREGRESSION  
SKLEARNLINEARMODEL RIDGEREGRESSION XYALPHA SAMPLEWEIGHTNONE SOLVER'AUTO'  
MAXITERNONE TOL0001 VERBOSE0 RAN  
DOMSTATENONE RETURNNITERFALSE RE  
TURNINTERCEPTFALSE CHECKINPUTTRUE  
SOLVE THE RIDGE EQUATION BY THE METHOD OF NORMAL EQUATIONS  
READ MORE IN THE USER GUIDE  
PARAMETERS  
XARRAYLIKE SPARSE MATRIX LINEAROPERATOR SHAPE NSAMPLES NFEATURES TRAINING DATA  
YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NTARGETS TARGET VALUES  
ALPHA FLOAT ARRAYLIKE SHAPE NTARGETS IF ARRAYLIKE REGULARIZATION STRENGTH MUST BE  
A POSITIVE FLOAT REGULARIZATION IMPROVES THE CONDITIONING OF THE PROBLEM AND REDUCES THE  
VARIANCE OF THE ESTIMATES LARGER VALUES SPECIFY STRONGER REGULARIZATION ALPHA CORRESPONDS  
TOC1 IN OTHER LINEAR MODELS SUCH AS LOGISTICREGRESSION OR LINEARSVC IF AN ARRAY IS  
PASSED PENALTIES ARE ASSUMED TO BE SPECIFIC TO THE TARGETS HENCE THEY MUST CORRESPOND IN  
NUMBER  
SAMPLEWEIGHT FLOAT OR NUMPY ARRAY OF SHAPE NSAMPLES INDIVIDUAL WEIGHTS FOR EACH SAM  
PLE IF SAMPLEWEIGHT IS NOT NONE AND SOLVER'AUTO' THE SOLVER WILL BE SET TO 'CHOLESKY'  
NEW IN VERSION 017  
SOLVER 'AUTO' 'SVD' 'CHOLESKY' 'LSQR' 'SPARSECG' 'SAG' 'SAGA' SOLVER TO USE IN THE COM  
PUTATIONAL ROUTINES  
• 'AUTO' CHOOSES THE SOLVER AUTOMATICALLY BASED ON THE TYPE OF DATA  
• 'SVD' USES A SINGULAR VALUE DECOMPOSITION OF X TO COMPUTE THE RIDGE COEFFICIENTS MORE  
STABLE FOR SINGULAR MATRICES THAN 'CHOLESKY'  
622SKLEARNLINEARMODEL GENERALIZED LINEAR MODELS 1963

SCIKITLEARN USER GUIDE RELEASE 0213

- ‘CHOLESKY’ USES THE STANDARD SCIPYLINALGSOLVE FUNCTION TO OBTAIN A CLOSEDFORM SOLUTION VIA A CHOLESKY DECOMPOSITION OF DOTXT X
  - ‘SPARSECG’ USES THE CONJUGATE GRADIENT SOLVER AS FOUND IN SCIPYSPARSELINALGCG AS AN ITERATIVE ALGORITHM THIS SOLVER IS MORE APPROPRIATE THAN ‘CHOLESKY’ FOR LARGESCALE DATA POSSIBILITY TO SET TOL ANDMAXITER
  - ‘LSQR’ USES THE DEDICATED REGULARIZED LEASTSQUARES ROUTINE SCIPYSPARSELINALGLSQR IT IS THE FASTEST AND USES AN ITERATIVE PROCEDURE
  - ‘SAG’ USES A STOCHASTIC AVERAGE GRADIENT DESCENT AND ‘SAGA’ USES ITS IMPROVED UNBIASED VERSION NAMED SAGA BOTH METHODS ALSO USE AN ITERATIVE PROCEDURE AND ARE OFTEN FASTER THAN OTHER SOLVERS WHEN BOTH NSAMPLES AND NFEATURES ARE LARGE NOTE THAT ‘SAG’ AND ‘SAGA’ FAST CONVERGENCE IS ONLY GUARANTEED ON FEATURES WITH APPROXIMATELY THE SAME SCALE YOU CAN PREPROCESS THE DATA WITH A SCALER FROM SKLEARNPREPROCESSING
- ALL LAST FIVE SOLVERS SUPPORT BOTH DENSE AND SPARSE DATA HOWEVER ONLY ‘SAG’ AND ‘SPARSECG’ SUPPORTS SPARSE INPUT WHEN‘FITINTERCEPT’ IS TRUE
- NEW IN VERSION 017 STOCHASTIC AVERAGE GRADIENT DESCENT SOLVER
- NEW IN VERSION 019 SAGA SOLVER
- MAXITER INT OPTIONAL MAXIMUM NUMBER OF ITERATIONS FOR CONJUGATE GRADIENT SOLVER FOR THE ‘SPARSECG’ AND ‘LSQR’ SOLVERS THE DEFAULT VALUE IS DETERMINED BY SCIPYSPARSELINALG FOR ‘SAG’ AND SAGA SOLVER THE DEFAULT VALUE IS 1000
- TOLFLOAT PRECISION OF THE SOLUTION
- VERBOSE INT VERBOSITY LEVEL SETTING VERBOSE 0 WILL DISPLAY ADDITIONAL INFORMATION DEPENDING ON THE SOLVER USED
- RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULT NONE THE SEED OF THE PSEUDO RANDOM NUMBER GENERATOR TO USE WHEN SHUFFLING THE DATA IF INT RANDOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NPRANDOM USED WHEN SOLVER ‘SAG’
- RETURNNITER BOOLEAN DEFAULT FALSE IF TRUE THE METHOD ALSO RETURNS NITER THE ACTUAL NUMBER OF ITERATION PERFORMED BY THE SOLVER
- NEW IN VERSION 017
- RETURNINTERCEPT BOOLEAN DEFAULT FALSE IF TRUE AND IF X IS SPARSE THE METHOD ALSO RETURNS THE INTERCEPT AND THE SOLVER IS AUTOMATICALLY CHANGED TO ‘SAG’ THIS IS ONLY A TEMPORARY FIX FOR FITTING THE INTERCEPT WITH SPARSE DATA FOR DENSE DATA USE SKLEARNLINEARMODELPREPROCESSDATA BEFORE YOUR REGRESSION
- NEW IN VERSION 017
- CHECKINPUT BOOLEAN DEFAULT TRUE IF FALSE THE INPUT ARRAYS X AND Y WILL NOT BE CHECKED
- NEW IN VERSION 021
- RETURNS
- COEF ARRAY SHAPE NFEATURES OR NTARGETS NFEATURES WEIGHT VECTORS
- NITER INT OPTIONAL THE ACTUAL NUMBER OF ITERATION PERFORMED BY THE SOLVER ONLY RETURNED IF RETURNNITER IS TRUE
- INTERCEPT FLOAT OR ARRAY SHAPE NTARGETS THE INTERCEPT OF THE MODEL ONLY RETURNED IF RETURNINTERCEPT IS TRUE AND IF X IS A SCIPY SPARSE ARRAY

SCIKITLEARN USER GUIDE RELEASE 0213

NOTES

THIS FUNCTION WON'T COMPUTE THE INTERCEPT

623SKLEARNMANIFOLD MANIFOLD LEARNING

THESKLEARNMANIFOLD MODULE IMPLEMENTS DATA EMBEDDING TECHNIQUES

USER GUIDE SEE THE MANIFOLD LEARNING SECTION FOR FURTHER DETAILS

MANIFOLDISOMAP NNEIGHBORS NCOMPONENTS ISOMAP EMBEDDING

MANIFOLDLOCALLYLINEAREMBEDDING LOCALLY LINEAR EMBEDDING

MANIFOLDMDS NCOMPONENTS METRIC NINIT MULTIDIMENSIONAL SCALING

MANIFOLDSPECTRALEMBEDDING NCOMPONENTS

SPECTRAL EMBEDDING FOR NONLINEAR DIMENSIONALITY REDUC

TION

MANIFOLDTSNE NCOMPONENTS PERPLEXITY TDISTRIBUTED STOCHASTIC NEIGHBOR EMBEDDING

6231SKLEARNMANIFOLD ISOMAP

CLASSSKLEARNMANIFOLD ISOMAPNNEIGHBORS5 NCOMPONENTS2 EIGENSOLVER'AUTO' TOL0

MAXITERNONE PATHMETHOD'AUTO' NEIGHBORSALGORITHM'AUTO'

NJOBSNONE

ISOMAP EMBEDDING

NONLINEAR DIMENSIONALITY REDUCTION THROUGH ISOMETRIC MAPPING

READ MORE IN THE USER GUIDE

PARAMETERS

NNEIGHBORS INTEGER NUMBER OF NEIGHBORS TO CONSIDER FOR EACH POINT

NCOMPONENTS INTEGER NUMBER OF COORDINATES FOR THE MANIFOLD

EIGENSOLVER 'AUTO''ARPACK''DENSE' 'AUTO' ATTEMPT TO CHOOSE THE MOST EFFICIENT SOLVER FOR

THE GIVEN PROBLEM

'ARPACK' USE ARNOLDI DECOMPOSITION TO FIND THE EIGENVALUES AND EIGENVECTORS

'DENSE' USE A DIRECT SOLVER IE LAPACK FOR THE EIGENVALUE DECOMPOSITION

TOLFLOAT CONVERGENCE TOLERANCE PASSED TO ARPACK OR LOBPCG NOT USED IF EIGENSOLVER

'DENSE'

MAXITER INTEGER MAXIMUM NUMBER OF ITERATIONS FOR THE ARPACK SOLVER NOT USED IF

EIGENSOLVER 'DENSE'

PATHMETHOD STRING 'AUTO''FW''D' METHOD TO USE IN FINDING SHORTEST PATH

'AUTO' ATTEMPT TO CHOOSE THE BEST ALGORITHM AUTOMATICALLY

'FW' FLOYDWARSHALL ALGORITHM

'D' DIJKSTRA'S ALGORITHM

NEIGHBORSALGORITHM STRING 'AUTO''BRUTE''KDTree''BALLTree' ALGORITHM TO USE FOR NEAREST

NEIGHBORS SEARCH PASSED TO NEIGHBORSNEARESTNEIGHBORS INSTANCE

623SKLEARNMANIFOLD MANIFOLD LEARNING 1965

SCIKITLEARN USER GUIDE RELEASE 0213

NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF PARALLEL JOBS TO RUN NONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

ATTRIBUTES

EMBEDDING ARRAYLIKE SHAPE NSAMPLES NCOMPONENTS STORES THE EMBEDDING VECTORS

KERNELPCA OBJECTKERNELPCA OBJECT USED TO IMPLEMENT THE EMBEDDING

TRAININGDATA ARRAYLIKE SHAPE NSAMPLES NFEATURES STORES THE TRAINING DATA

NBRS SKLEARNNEIGHBORSNEARESTNEIGHBORS INSTANCE STORES NEAREST NEIGHBORS INSTANCE IN CLUDING BALLTREE OR KDTREE IF APPLICABLE

DISTMATRIX ARRAYLIKE SHAPE NSAMPLES NSAMPLES STORES THE GEODESIC DISTANCE MATRIX OF TRAINING DATA

REFERENCES

R7F4D308F50541

EXAMPLES

FROM SKLEARNDATASETS IMPORT LOADDIGITS

FROM SKLEARNMANIFOLD IMPORT ISOMAP

X LOADDIGITSRETURNXY TRUE

XSHAPE

1797 64

EMBEDDING ISOMAPNCOMPONENTS2

XTRANSFORMED EMBEDDINGFITTRANSFORMX100

XTRANSFORMEDSHAPE

100 2

METHODS

FITSELF X Y COMPUTE THE EMBEDDING VECTORS FOR DATA X

FITTRANSFORM SELF X Y FIT THE MODEL FROM DATA IN X AND TRANSFORM X

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

RECONSTRUCTIONERROR SELF COMPUTE THE RECONSTRUCTION ERROR FOR THE EMBEDDING

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF X TRANSFORM X

INIT SELFNNEIGHBORS5 NCOMPONENTS2 EIGENSOLVER'AUTO' TOL0 MAXITERNONE

PATHMETHOD'AUTO' NEIGHBORSALGORITHM'AUTO' NJOBSNONE

FITSELFXYNONE

COMPUTE THE EMBEDDING VECTORS FOR DATA X

PARAMETERS

XARRAYLIKE SPARSE MATRIX BALLTREE KDTREE NEARESTNEIGHBORS SAMPLE DATA SHAPE NSAMPLES NFEATURES IN THE FORM OF A NUMPY ARRAY PRECOMPUTED TREE OR NEAREST NEIGHBORS OBJECT

1966 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

YIGNORED

RETURNS

SELF RETURNS AN INSTANCE OF SELF

FITTRANSFORM SELFXYNONE

FIT THE MODEL FROM DATA IN X AND TRANSFORM X

PARAMETERS

XARRAYLIKE SPARSE MATRIX BALLTREE KDTRREE TRAINING VECTOR WHERE NSAMPLES IN THE

NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

YIGNORED

RETURNS

XNEW ARRAYLIKE SHAPE NSAMPLES NCOMPONENTS

GETPARAMS SELFDEEPTREE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

RECONSTRUCTIONERROR SELF

COMPUTE THE RECONSTRUCTION ERROR FOR THE EMBEDDING

RETURNS

RECONSTRUCTIONERROR FLOAT

NOTES

THE COST FUNCTION OF AN ISOMAP EMBEDDING IS

$E = \frac{1}{n} \sum_{i=1}^n \text{FROBENIUSNORM}(D_i - KDFIT_i)^2$  WHERE D IS THE MATRIX OF DISTANCES FOR THE INPUT DATA X DFIT IS THE MATRIX OF DISTANCES FOR THE OUTPUT

EMBEDDING XFIT AND K IS THE ISOMAP KERNEL

$KD = \frac{1}{n} \sum_{i=1}^n \text{FROBENIUSNORM}(D_i - KDFIT_i)^2$

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE

PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT

OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELFXY

TRANSFORM X

623SKLEARNMANIFOLD MANIFOLD LEARNING 1967

SCIKITLEARN USER GUIDE RELEASE 0213

THIS IS IMPLEMENTED BY LINKING THE POINTS X INTO THE GRAPH OF GEODESIC DISTANCES OF THE TRAINING DATA FIRST THENNEIGHBORS NEAREST NEIGHBORS OF X ARE FOUND IN THE TRAINING DATA AND FROM THESE THE SHORTEST GEODESIC DISTANCES FROM EACH POINT IN X TO EACH POINT IN THE TRAINING DATA ARE COMPUTED IN ORDER TO CONSTRUCT THE KERNEL THE EMBEDDING OF X IS THE PROJECTION OF THIS KERNEL ONTO THE EMBEDDING VECTORS OF THE TRAINING SET

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES

RETURNS

XNEW ARRAYLIKE SHAPE NSAMPLES NCOMPONENTS

EXAMPLES USING SKLEARNMANIFOLDISOMAP

- COMPARISON OF MANIFOLD LEARNING METHODS
- MANIFOLD LEARNING METHODS ON A SEVERED SPHERE
- MANIFOLD LEARNING ON HANDWRITTEN DIGITS LOCALLY LINEAR EMBEDDING ISOMAP

6232SKLEARNMANIFOLD LOCALLYLINEAREMBEDDING

CLASSSSKLEARNMANIFOLD LOCALLYLINEAREMBEDDING NNEIGHBORS5 NCOMPONENTS2

REG0001 EIGENSOLVER'AUTO' TOL1E

06MAXITER100 METHOD'STANDARD'

HESSIAN TOL00001 MODIFIEDTOL1E

12 NEIGHBORSALGORITHM'AUTO' RAN

DOMSTATENONE NJOBSNONE

LOCALLY LINEAR EMBEDDING

READ MORE IN THE USER GUIDE

PARAMETERS

NNEIGHBORS INTEGER NUMBER OF NEIGHBORS TO CONSIDER FOR EACH POINT

NCOMPONENTS INTEGER NUMBER OF COORDINATES FOR THE MANIFOLD

REGFLOAT REGULARIZATION CONSTANT MULTIPLIES THE TRACE OF THE LOCAL COVARIANCE MATRIX OF THE DISTANCES

EIGENSOLVER STRING 'AUTO' 'ARPACK' 'DENSE' AUTO ALGORITHM WILL ATTEMPT TO CHOOSE THE BEST METHOD FOR INPUT DATA

ARPACK USE ARNOLDI ITERATION IN SHIFTINVERT MODE FOR THIS METHOD M MAY BE A DENSE MATRIX SPARSE MATRIX OR GENERAL LINEAR OPERATOR WARNING ARPACK CAN BE UNSTABLE FOR SOME PROBLEMS IT IS BEST TO TRY SEVERAL RANDOM SEEDS IN ORDER TO CHECK RESULTS

DENSE USE STANDARD DENSE MATRIX OPERATIONS FOR THE EIGENVALUE DECOMPOSITION FOR THIS METHOD M MUST BE AN ARRAY OR MATRIX TYPE THIS METHOD SHOULD BE AVOIDED FOR LARGE PROBLEMS

TOLFLOAT OPTIONAL TOLERANCE FOR 'ARPACK' METHOD NOT USED IF EIGENSOLVER'DENSE'

MAXITER INTEGER MAXIMUM NUMBER OF ITERATIONS FOR THE ARPACK SOLVER NOT USED IF EIGENSOLVER'DENSE'

METHOD STRING 'STANDARD' 'HESSIAN' 'MODIFIED' OR 'LTSA'

STANDARD USE THE STANDARD LOCALLY LINEAR EMBEDDING ALGORITHM SEE REFERENCE 1

1968 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

HESSIAN USE THE HESSIAN EIGENMAP METHOD THIS METHOD REQUIRES NNEIGHBORS  
NCOMPONENTS 1 NCOMPONENTS 1 2 SEE REFERENCE 2

MODIFIED USE THE MODIFIED LOCALLY LINEAR EMBEDDING ALGORITHM SEE REFERENCE 3

LTSA USE LOCAL TANGENT SPACE ALIGNMENT ALGORITHM SEE REFERENCE 4

HESSIAN\_TOL FLOAT OPTIONAL TOLERANCE FOR HESSIAN EIGENMAPPING METHOD ONLY USED IF  
METHOD HESSIAN

MODIFIED\_TOL FLOAT OPTIONAL TOLERANCE FOR MODIFIED LLE METHOD ONLY USED IF METHOD  
MODIFIED

NEIGHBORS\_ALGORITHM STRING 'AUTO''BRUTE''KDTREE''BALLTREE' ALGORITHM TO USE FOR NEAREST  
NEIGHBORS SEARCH PASSED TO NEIGHBORS\_NEAREST\_NEIGHBORS INSTANCE

RANDOM\_STATE INT RANDOM\_STATE INSTANCE OR NONE OPTIONAL DEFAULT NONE IF INT RAN  
DOM\_STATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOM\_STATE INSTANCE RAN  
DOM\_STATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE  
RANDOM\_STATE INSTANCE USED BY NPRANDOM USED WHEN EIGENSOLVER 'ARPACK'

NJOBS INT OR NONE OPTIONAL DEFAULT NONE THE NUMBER OF PARALLEL JOBS TO RUN NONE MEANS  
1 UNLESS IN A JOBLIBPARALLEL\_BACKEND CONTEXT 1 MEANS USING ALL PROCESSORS SEE  
GLOSSARY FOR MORE DETAILS

ATTRIBUTES

EMBEDDING\_ARRAY LIKE SHAPE NSAMPLES NCOMPONENTS STORES THE EMBEDDING VECTORS

RECONSTRUCTION\_ERROR FLOAT RECONSTRUCTION ERROR ASSOCIATED WITH EMBEDDING

NBRS\_NEAREST\_NEIGHBORS OBJECT STORES NEAREST NEIGHBORS INSTANCE INCLUDING BALLTREE OR  
KDTREE IF APPLICABLE

REFERENCES

R62E36DD1B0561 R62E36DD1B0562 R62E36DD1B0563 R62E36DD1B0564

EXAMPLES

FROM SKLEARN\_DATASETS IMPORT LOADDIGITS

FROM SKLEARN\_MANIFOLD IMPORT LOCALLYLINEAREMBEDDING

X LOADDIGITS\_RETURN\_XY TRUE

X\_SHAPE

1797 64

EMBEDDING LOCALLYLINEAREMBEDDING\_NCOMPONENTS 2

X\_TRANSFORMED EMBEDDING\_FIT\_TRANSFORM\_X 100

X\_TRANSFORMED\_SHAPE

100 2

METHODS

FIT\_SELF X Y COMPUTE THE EMBEDDING VECTORS FOR DATA X

FIT\_TRANSFORM\_SELF X Y COMPUTE THE EMBEDDING VECTORS FOR DATA X AND TRANS  
FORM X

CONTINUED ON NEXT PAGE

623 SKLEARNMANIFOLD MANIFOLD LEARNING 1969

SCIKITLEARN USER GUIDE RELEASE 0213

TABLE 6173 – CONTINUED FROM PREVIOUS PAGE

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF X TRANSFORM NEW POINTS INTO EMBEDDING SPACE

INIT SELFNNNEIGHBORS5 NCOMPONENTS2 REG0001 EIGENSOLVER'AUTO' TOL1E06

MAXITER100 METHOD'STANDARD' HESSIANTOL00001 MODIFIEDTOL1E12 NEIGH

BORSALGORITHM'AUTO' RANDOMSTATENONE NJOBSNONE

FITSELFXYNONE

COMPUTE THE EMBEDDING VECTORS FOR DATA X

PARAMETERS

XARRAYLIKE OF SHAPE NSAMPLES NFEATURES TRAINING SET

YIGNORED

RETURNS

SELF RETURNS AN INSTANCE OF SELF

FITTRANSFORM SELFXYNONE

COMPUTE THE EMBEDDING VECTORS FOR DATA X AND TRANSFORM X

PARAMETERS

XARRAYLIKE OF SHAPE NSAMPLES NFEATURES TRAINING SET

YIGNORED

RETURNS

XNEW ARRAYLIKE SHAPE NSAMPLES NCOMPONENTS

GETPARAMS SELFDEEPTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE

PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT

OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELF X

TRANSFORM NEW POINTS INTO EMBEDDING SPACE

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES

RETURNS

1970 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213  
XNEW ARRAY SHAPE NSAMPLES NCOMPONENTS  
NOTES

BECAUSE OF SCALING PERFORMED BY THIS METHOD IT IS DISCOURAGED TO USE IT TOGETHER WITH METHODS THAT ARE NOT  
SCALEINVARIANT LIKE SVMs

EXAMPLES USING SKLEARNMANIFOLDLOCALLYLINEAREMBEDDING

- VISUALIZING THE STOCK MARKET STRUCTURE
- COMPARISON OF MANIFOLD LEARNING METHODS
- MANIFOLD LEARNING METHODS ON A SEVERED SPHERE
- MANIFOLD LEARNING ON HANDWRITTEN DIGITS LOCALLY LINEAR EMBEDDING ISOMAP

6233SKLEARNMANIFOLD MDS

CLASSSSKLEARNMANIFOLD MDSNCOMPONENTS2 METRICTRUE NINIT4 MAXITER300 VER

BOSE0 EPS0001 NJOBSNONE RANDOMSTATENONE DISSIMILAR

ITY'EUCLIDEAN'

MULTIDIMENSIONAL SCALING

READ MORE IN THE USER GUIDE

PARAMETERS

NCOMPONENTS INT OPTIONAL DEFAULT 2 NUMBER OF DIMENSIONS IN WHICH TO IMMERSE THE DIS  
SIMILARITIES

METRIC BOOLEAN OPTIONAL DEFAULT TRUE IF TRUE PERFORM METRIC MDS OTHERWISE PERFORM  
NONMETRIC MDS

NINIT INT OPTIONAL DEFAULT 4 NUMBER OF TIMES THE SMACOF ALGORITHM WILL BE RUN WITH  
DIFFERENT INITIALIZATIONS THE FINAL RESULTS WILL BE THE BEST OUTPUT OF THE RUNS DETERMINED BY  
THE RUN WITH THE SMALLEST FINAL STRESS

MAXITER INT OPTIONAL DEFAULT 300 MAXIMUM NUMBER OF ITERATIONS OF THE SMACOF ALGO  
RITHM FOR A SINGLE RUN

VERBOSE INT OPTIONAL DEFAULT 0 LEVEL OF VERBOSITY

EPS FLOAT OPTIONAL DEFAULT 1E3 RELATIVE TOLERANCE WITH RESPECT TO STRESS AT WHICH TO DECLARE  
CONVERGENCE

NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF JOBS TO USE FOR THE COMPUTATION  
IF MULTIPLE INITIALIZATIONS ARE USED NINIT EACH RUN OF THE ALGORITHM IS COMPUTED IN  
PARALLEL

NONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL  
PROCESSORS SEE GLOSSARY FOR MORE DETAILS

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULT NONE THE GENERATOR  
USED TO INITIALIZE THE CENTERS IF INT RANDOMSTATE IS THE SEED USED BY THE RANDOM NUMBER  
GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE  
THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NPRANDOM  
6233SKLEARNMANIFOLD MANIFOLD LEARNING 1971

SCIKITLEARN USER GUIDE RELEASE 0213

DISSIMILARITY ‘EUCLIDEAN’ ‘PRECOMPUTED’ OPTIONAL DEFAULT ‘EUCLIDEAN’ DISSIMILARITY MEASURE TO USE

- ‘EUCLIDEAN’ PAIRWISE EUCLIDEAN DISTANCES BETWEEN POINTS IN THE DATASET
- ‘PRECOMPUTED’ PRECOMPUTED DISSIMILARITIES ARE PASSED DIRECTLY TO FIT AND FITTRANSFORM

ATTRIBUTES

EMBEDDING ARRAYLIKE SHAPE NSAMPLES NCOMPONENTS STORES THE POSITION OF THE DATASET IN THE EMBEDDING SPACE

STRESS FLOAT THE FINAL VALUE OF THE STRESS SUM OF SQUARED DISTANCE OF THE DISPARITIES AND THE DISTANCES FOR ALL CONSTRAINED POINTS

REFERENCES

“MODERN MULTIDIMENSIONAL SCALING THEORY AND APPLICATIONS” BORG I GROENEN P SPRINGER SERIES IN STATISTICS 1997

“NONMETRIC MULTIDIMENSIONAL SCALING A NUMERICAL METHOD” KRUSKAL J PSYCHOMETRIKA 29 1964

“MULTIDIMENSIONAL SCALING BY OPTIMIZING GOODNESS OF FIT TO A NONMETRIC HYPOTHESIS” KRUSKAL J PSYCHOMETRIKA 29 1964

EXAMPLES

```
FROM SKLEARNDATASETS IMPORT LOADDIGITS
FROM SKLEARNMANIFOLD IMPORT MDS
X LOADDIGITSRETURNXY TRUE
XSHAPE
1797 64
EMBEDDING MDSNCOMPONENTS2
XTRANSFORMED EMBEDDINGFITTRANSFORMX100
XTRANSFORMEDSHAPE
100 2
```

METHODS

FITSELF X Y INIT COMPUTES THE POSITION OF THE POINTS IN THE EMBEDDING SPACE

FITTRANSFORM SELF X Y INIT FIT THE DATA FROM X AND RETURNS THE EMBEDDED COORDINATES

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELFNCOMPONENTS2 METRICTRUE NINIT4 MAXITER300 VERBOSE0 EPS0001

NJOBSNONE RANDOMSTATENONE DISSIMILARITY‘EUCLIDEAN’

FITSELFXYNONE INITNONE

COMPUTES THE POSITION OF THE POINTS IN THE EMBEDDING SPACE

PARAMETERS

1972 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

XARRAY SHAPE NSAMPLES NFEATURES OR NSAMPLES NSAMPLES INPUT DATA IF  
DISSIMILARITYPRECOMPUTED THE INPUT SHOULD BE THE DISSIMILARITY MATRIX  
YIGNORED

INIT NDARRAY SHAPE NSAMPLES OPTIONAL DEFAULT NONE STARTING CONFIGURATION OF THE EM  
BEDDING TO INITIALIZE THE SMACOF ALGORITHM BY DEFAULT THE ALGORITHM IS INITIALIZED WITH  
A RANDOMLY CHOSEN ARRAY

FITTRANSFORM SELFXYNONE INITNONE

FIT THE DATA FROM X AND RETURNS THE EMBEDDED COORDINATES

PARAMETERS

XARRAY SHAPE NSAMPLES NFEATURES OR NSAMPLES NSAMPLES INPUT DATA IF  
DISSIMILARITYPRECOMPUTED THE INPUT SHOULD BE THE DISSIMILARITY MATRIX  
YIGNORED

INIT NDARRAY SHAPE NSAMPLES OPTIONAL DEFAULT NONE STARTING CONFIGURATION OF THE EM  
BEDDING TO INITIALIZE THE SMACOF ALGORITHM BY DEFAULT THE ALGORITHM IS INITIALIZED WITH  
A RANDOMLY CHOSEN ARRAY

GETPARAMS SELFDEEPTTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE  
PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT  
OF A NESTED OBJECT

RETURNS

SELF

EXAMPLES USING SKLEARNMANIFOLDMDS

- MULTIDIMENSIONAL SCALING
- COMPARISON OF MANIFOLD LEARNING METHODS
- MANIFOLD LEARNING METHODS ON A SEVERED SPHERE
- MANIFOLD LEARNING ON HANDWRITTEN DIGITS LOCALLY LINEAR EMBEDDING ISOMAP

623SKLEARNMANIFOLD MANIFOLD LEARNING 1973

SCIKITLEARN USER GUIDE RELEASE 0213  
62345KLEARNMANIFOLD SPECTRALEMMBEDDING  
CLASSSSKLEARNMANIFOLD SPECTRALEMMBEDDING NCOMPONENTS2 AFFINITY'NEARESTNEIGHBORS'  
GAMMANONE RANDOMSTATENONE  
EIGENSOLVERNONE NNEIGHBORSNONE  
NJOBSNONE  
SPECTRAL EMBEDDING FOR NONLINEAR DIMENSIONALITY REDUCTION  
FORMS AN AFFINITY MATRIX GIVEN BY THE SPECIFIED FUNCTION AND APPLIES SPECTRAL DECOMPOSITION TO THE CORRESPONDING  
GRAPH LAPLACIAN THE RESULTING TRANSFORMATION IS GIVEN BY THE VALUE OF THE EIGENVECTORS FOR EACH DATA POINT  
NOTE LAPLACIAN EIGENMAPS IS THE ACTUAL ALGORITHM IMPLEMENTED HERE  
READ MORE IN THE USER GUIDE  
PARAMETERS  
NCOMPONENTS INTEGER DEFAULT 2 THE DIMENSION OF THE PROJECTED SUBSPACE  
AFFINITY STRING OR CALLABLE DEFAULT  
HOW TO CONSTRUCT THE AFFINITY MATRIX  
• 'NEARESTNEIGHBORS' CONSTRUCT AFFINITY MATRIX BY KNN GRAPH  
• 'RBF' CONSTRUCT AFFINITY MATRIX BY RBF KERNEL  
• 'PRECOMPUTED' INTERPRET X AS PRECOMPUTED AFFINITY MATRIX  
• CALLABLE USE PASSED IN FUNCTION AS AFFINITY THE FUNCTION TAKES IN DATA MATRIX NSAMPLES  
NFEATURES AND RETURN AFFINITY MATRIX NSAMPLES NSAMPLES  
GAMMA FLOAT OPTIONAL DEFAULT KERNEL COEFFICIENT FOR RBF KERNEL  
RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULT NONE A PSEUDO RANDOM  
NUMBER GENERATOR USED FOR THE INITIALIZATION OF THE LOBPCG EIGENVECTORS IF INT RANDOMSTATE  
IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS  
THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE  
INSTANCE USED BY NRANDOM USED WHEN SOLVER 'AMG'  
EIGENSOLVER NONE 'ARPACK' 'LOBPCG' OR 'AMG' THE EIGENVALUE DECOMPOSITION STRATEGY TO  
USE AMG REQUIRES PYAMG TO BE INSTALLED IT CAN BE FASTER ON VERY LARGE SPARSE PROBLEMS  
BUT MAY ALSO LEAD TO INSTABILITIES  
NNEIGHBORS INT DEFAULT NUMBER OF NEAREST NEIGHBORS FOR NEARESTNEIGHBORS GRAPH BUILDING  
NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF PARALLEL JOBS TO RUN NONE MEANS  
1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL PROCESSORS SEE  
GLOSSARY FOR MORE DETAILS  
ATTRIBUTES  
EMBEDDING ARRAY SHAPE NSAMPLES NCOMPONENTS SPECTRAL EMBEDDING OF THE TRAINING  
MATRIX  
AFFINITYMATRIX ARRAY SHAPE NSAMPLES NSAMPLES AFFINITYMATRIX CONSTRUCTED FROM  
SAMPLES OR PRECOMPUTED  
REFERENCES  
• A TUTORIAL ON SPECTRAL CLUSTERING 2007 ULRIKE VON LUXBURG HTTPCITESEERXISTPSUEDUVIEWDOC  
SUMMARYDOI10111659323  
1974 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

- ON SPECTRAL CLUSTERING ANALYSIS AND AN ALGORITHM 2001 ANDREW Y NG MICHAEL I JORDAN YAIR WEISS  
HTTPCITESEERXISTPSUEDUVIEWDOC SUMMARYDOI1011198100
- NORMALIZED CUTS AND IMAGE SEGMENTATION 2000 JIANBO SHI JITENDRA MALIK HTTPCITESEERISTPSUEDU  
VIEWDOC SUMMARYDOI10111602324

EXAMPLES

```
FROM SKLEARN DATASETS IMPORT LOADDIGITS
FROM SKLEARN MANIFOLD IMPORT SPECTRALEMBEDDING
X = LOADDIGITS RETURN XY TRUE
```

X SHAPE

1797 64

EMBEDDING = SPECTRALEMBEDDING(NCOMPONENTS=2

X TRANSFORMED = EMBEDDING(FIT TRANSFORM X)100

X TRANSFORMED SHAPE

100 2

METHODS

FIT SELF X Y FIT THE MODEL FROM DATA IN X

FIT TRANSFORM SELF X Y FIT THE MODEL FROM DATA IN X AND TRANSFORM X

GET PARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

SET PARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELF(NCOMPONENTS=2 AFFINITY='NEAREST NEIGHBORS' GAMMA=None RANDOM STATE=None

EIGEN SOLVER=None N NEIGHBORS=None N JOBS=None

FIT SELF(X=None)

FIT THE MODEL FROM DATA IN X

PARAMETERS

X: ARRAYLIKE SHAPE (NSAMPLES, NFEATURES) TRAINING VECTOR WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

IF AFFINITY IS "PRECOMPUTED" X: ARRAYLIKE SHAPE (NSAMPLES, NSAMPLES) INTERPRET X AS

PRECOMPUTED ADJACENCY GRAPH COMPUTED FROM SAMPLES

RETURNS

SELF OBJECT RETURNS THE INSTANCE ITSELF

FIT TRANSFORM SELF(X=None)

FIT THE MODEL FROM DATA IN X AND TRANSFORM X

PARAMETERS

X: ARRAYLIKE SHAPE (NSAMPLES, NFEATURES) TRAINING VECTOR WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

IF AFFINITY IS "PRECOMPUTED" X: ARRAYLIKE SHAPE (NSAMPLES, NSAMPLES) INTERPRET X AS

PRECOMPUTED ADJACENCY GRAPH COMPUTED FROM SAMPLES

RETURNS

X NEW: ARRAYLIKE SHAPE (NSAMPLES, NCOMPONENTS)

623 SKLEARN MANIFOLD MANIFOLD LEARNING 1975

SCIKITLEARN USER GUIDE RELEASE 0213

GETPARAMS SELFDEEPTREE  
GET PARAMETERS FOR THIS ESTIMATOR  
PARAMETERS  
DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBOBJECTS THAT ARE ESTIMATORS  
RETURNS  
PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES  
SETPARAMS SELFPARAMS  
SET THE PARAMETERS OF THIS ESTIMATOR  
THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE  
PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT  
OF A NESTED OBJECT  
RETURNS  
SELF

EXAMPLES USING SKLEARNMANIFOLDSPECTRALEMBEDDING

- VARIOUS AGGLOMERATIVE CLUSTERING ON A 2D EMBEDDING OF DIGITS
- COMPARISON OF MANIFOLD LEARNING METHODS
- MANIFOLD LEARNING METHODS ON A SEVERED SPHERE
- MANIFOLD LEARNING ON HANDWRITTEN DIGITS LOCALLY LINEAR EMBEDDING ISOMAP

6235SKLEARNMANIFOLD TSNE  
CLASSSSKLEARNMANIFOLD TSNENCOMPONENTS2 PERPLEXITY300 EARLYEXAGGERATION120  
LEARNINGRATE2000 NITER1000 NITERWITHOUTPROGRESS300  
MINGRADNORM1E07 METRIC'EUCLIDEAN' INIT'RANDOM' VERBOSE0  
RANDOMSTATENONE METHOD'BARNESHUT' ANGLE05  
TDISTRIBUTED STOCHASTIC NEIGHBOR EMBEDDING

TSNE 1 IS A TOOL TO VISUALIZE HIGHDIMENSIONAL DATA IT CONVERTS SIMILARITIES BETWEEN DATA POINTS TO JOINT  
PROBABILITIES AND TRIES TO MINIMIZE THE KULLBACKLEIBLER DIVERGENCE BETWEEN THE JOINT PROBABILITIES OF THE LOW  
DIMENSIONAL EMBEDDING AND THE HIGHDIMENSIONAL DATA TSNE HAS A COST FUNCTION THAT IS NOT CONVEX IE WITH  
DIFFERENT INITIALIZATIONS WE CAN GET DIFFERENT RESULTS  
IT IS HIGHLY RECOMMENDED TO USE ANOTHER DIMENSIONALITY REDUCTION METHOD EG PCA FOR DENSE DATA OR TRUNCAT  
EDSVD FOR SPARSE DATA TO REDUCE THE NUMBER OF DIMENSIONS TO A REASONABLE AMOUNT EG 50 IF THE NUMBER OF  
FEATURES IS VERY HIGH THIS WILL SUPPRESS SOME NOISE AND SPEED UP THE COMPUTATION OF PAIRWISE DISTANCES BETWEEN  
SAMPLES FOR MORE TIPS SEE LAURENS VAN DER MAATEN'S FAQ 2  
READ MORE IN THE USER GUIDE

PARAMETERS  
NCOMPONENTS INT OPTIONAL DEFAULT 2 DIMENSION OF THE EMBEDDED SPACE  
PERPLEXITY FLOAT OPTIONAL DEFAULT 30 THE PERPLEXITY IS RELATED TO THE NUMBER OF NEAREST  
NEIGHBORS THAT IS USED IN OTHER MANIFOLD LEARNING ALGORITHMS LARGER DATASETS USUALLY REQUIRE  
A LARGER PERPLEXITY CONSIDER SELECTING A VALUE BETWEEN 5 AND 50 DIFFERENT VALUES CAN RESULT  
IN SIGNIFICANLTY DIFFERENT RESULTS  
1976 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

EARLYEXAGGERATION FLOAT OPTIONAL DEFAULT 120 CONTROLS HOW TIGHT NATURAL CLUSTERS IN THE ORIGINAL SPACE ARE IN THE EMBEDDED SPACE AND HOW MUCH SPACE WILL BE BETWEEN THEM FOR LARGER VALUES THE SPACE BETWEEN NATURAL CLUSTERS WILL BE LARGER IN THE EMBEDDED SPACE AGAIN THE CHOICE OF THIS PARAMETER IS NOT VERY CRITICAL IF THE COST FUNCTION INCREASES DURING INITIAL OPTIMIZATION THE EARLY EXAGGERATION FACTOR OR THE LEARNING RATE MIGHT BE TOO HIGH LEARNINGRATE FLOAT OPTIONAL DEFAULT 2000 THE LEARNING RATE FOR TSNE IS USUALLY IN THE RANGE 100 10000 IF THE LEARNING RATE IS TOO HIGH THE DATA MAY LOOK LIKE A 'BALL' WITH ANY POINT APPROXIMATELY EQUIDISTANT FROM ITS NEAREST NEIGHBOURS IF THE LEARNING RATE IS TOO LOW MOST POINTS MAY LOOK COMPRESSED IN A DENSE CLOUD WITH FEW OUTLIERS IF THE COST FUNCTION GETS STUCK IN A BAD LOCAL MINIMUM INCREASING THE LEARNING RATE MAY HELP NITER INT OPTIONAL DEFAULT 1000 MAXIMUM NUMBER OF ITERATIONS FOR THE OPTIMIZATION SHOULD BE AT LEAST 250

NITERWITHOUTPROGRESS INT OPTIONAL DEFAULT 300 MAXIMUM NUMBER OF ITERATIONS WITHOUT PROGRESS BEFORE WE ABORT THE OPTIMIZATION USED AFTER 250 INITIAL ITERATIONS WITH EARLY EXAGGERATION NOTE THAT PROGRESS IS ONLY CHECKED EVERY 50 ITERATIONS SO THIS VALUE IS ROUNDED TO THE NEXT MULTIPLE OF 50

NEW IN VERSION 017 PARAMETER NITERWITHOUTPROGRESS TO CONTROL STOPPING CRITERIA MINGRADNORM FLOAT OPTIONAL DEFAULT 1E7 IF THE GRADIENT NORM IS BELOW THIS THRESHOLD THE OPTIMIZATION WILL BE STOPPED

METRIC STRING OR CALLABLE OPTIONAL THE METRIC TO USE WHEN CALCULATING DISTANCE BETWEEN INSTANCES IN A FEATURE ARRAY IF METRIC IS A STRING IT MUST BE ONE OF THE OPTIONS ALLOWED BY SCIPYSPATIALDISTANCEPDIST FOR ITS METRIC PARAMETER OR A METRIC LISTED IN PAIRWISEPAIRWISEDISTANCEFUNCTIONS IF METRIC IS "PRECOMPUTED" X IS ASSUMED TO BE A DISTANCE MATRIX ALTERNATIVELY IF METRIC IS A CALLABLE FUNCTION IT IS CALLED ON EACH PAIR OF INSTANCES ROWS AND THE RESULTING VALUE RECORDED THE CALLABLE SHOULD TAKE TWO ARRAYS FROM X AS INPUT AND RETURN A VALUE INDICATING THE DISTANCE BETWEEN THEM THE DEFAULT IS "EUCLIDEAN" WHICH IS INTERPRETED AS SQUARED EUCLIDEAN DISTANCE INIT STRING OR NUMPY ARRAY OPTIONAL DEFAULT "RANDOM" INITIALIZATION OF EMBEDDING POSSIBLE OPTIONS ARE 'RANDOM' 'PCA' AND A NUMPY ARRAY OF SHAPE NSAMPLES NCOMPONENTS PCA INITIALIZATION CANNOT BE USED WITH PRECOMPUTED DISTANCES AND IS USUALLY MORE GLOBALLY STABLE THAN RANDOM INITIALIZATION

VERBOSE INT OPTIONAL DEFAULT 0 VERBOSITY LEVEL RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULT NONE IF INT RANDOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NPRANDOM NOTE THAT DIFFERENT INITIALIZATIONS MIGHT RESULT IN DIFFERENT LOCAL MINIMA OF THE COST FUNCTION METHOD STRING DEFAULT 'BARNESHUT' BY DEFAULT THE GRADIENT CALCULATION ALGORITHM USES BARNESHUT APPROXIMATION RUNNING IN ONLOGN TIME METHOD'EXACT' WILL RUN ON THE SLOWER BUT EXACT ALGORITHM IN ON2 TIME THE EXACT ALGORITHM SHOULD BE USED WHEN NEARESTNEIGHBOR ERRORS NEED TO BE BETTER THAN 3 HOWEVER THE EXACT METHOD CANNOT SCALE TO MILLIONS OF EXAMPLES

NEW IN VERSION 017 APPROXIMATE OPTIMIZATION METHOD VIA THE BARNESHUT ANGLE FLOAT DEFAULT 05 ONLY USED IF METHOD'BARNESHUT' THIS IS THE TRADEOFF BETWEEN SPEED AND ACCURACY FOR BARNESHUT TSNE 'ANGLE' IS THE ANGULAR SIZE REFERRED TO AS THETA IN 3 OF A DISTANT NODE AS MEASURED FROM A POINT IF THIS SIZE IS BELOW 'ANGLE' THEN IT IS USED AS A SUMMARY NODE OF ALL POINTS CONTAINED WITHIN IT THIS METHOD IS NOT VERY SENSITIVE TO 623SKLEARNMANIFOLD MANIFOLD LEARNING 1977

SCIKITLEARN USER GUIDE RELEASE 0213

CHANGES IN THIS PARAMETER IN THE RANGE OF 02 08 ANGLE LESS THAN 02 HAS QUICKLY INCREASING  
COMPUTATION TIME AND ANGLE GREATER 08 HAS QUICKLY INCREASING ERROR

ATTRIBUTES

EMBEDDING ARRAYLIKE SHAPE NSAMPLES NCOMPONENTS STORES THE EMBEDDING VECTORS

KLDIVERGENCE FLOAT KULLBACKLEIBLER DIVERGENCE AFTER OPTIMIZATION

NITER INT NUMBER OF ITERATIONS RUN

REFERENCES

1 VAN DER MAATEN LJP HINTON GE VISUALIZING HIGHDIMENSIONAL DATA USING TSNE JOURNAL OF MA  
CHINE LEARNING RESEARCH 925792605 2008

2 VAN DER MAATEN LJP TDISTRIBUTED STOCHASTIC NEIGHBOR EMBEDDING HTTPSLVDMATENGITHUBIOTSNE

3 LJP VAN DER MAATEN ACCELERATING TSNE USING TREEBASED ALGORITHMS JOURNAL OF MACHINE LEARN  
ING RESEARCH 15OCT32213245 2014 HTTPSLVDMATENGITHUBIOPUBLICATIONSPAPERSJMLR2014PDF

EXAMPLES

```
import numpy as np
from sklearn.manifold import TSNE
X = np.array([0, 0, 0, 1, 1, 1, 0, 1, 1, 1])
X_embedded = TSNE(n_components=2).fit_transform(X)
X_embedded.shape
```

4 2

METHODS

FITSELF X Y FIT X INTO AN EMBEDDED SPACE

FITTRANSFORM SELF X Y FIT X INTO AN EMBEDDED SPACE AND RETURN THAT TRANS  
FORMED OUTPUT

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELF NCOMPONENTS2 PERPLEXITY300 EARLYEXAGGERATION120 LEARNINGRATE2000

NITER1000 NITERWITHOUTPROGRESS300 MINGRADNORM1E07 METRIC'EUCLIDEAN'

INIT'RANDOM' VERBOSE0 RANDOMSTATENONE METHOD'BARNESHUT' ANGLE05

FITSELFXYNONE

FIT X INTO AN EMBEDDED SPACE

PARAMETERS

XARRAY SHAPE NSAMPLES NFEATURES OR NSAMPLES NSAMPLES IF THE METRIC IS 'PRE  
COMPUTED' X MUST BE A SQUARE DISTANCE MATRIX OTHERWISE IT CONTAINS A SAMPLE PER ROW IF  
THE METHOD IS 'EXACT' X MAY BE A SPARSE MATRIX OF TYPE 'CSR' 'CSC' OR 'COO'

YIGNORED

FITTRANSFORM SELFXYNONE

FIT X INTO AN EMBEDDED SPACE AND RETURN THAT TRANSFORMED OUTPUT

1978 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

PARAMETERS

XARRAY SHAPE NSAMPLES NFEATURES OR NSAMPLES NSAMPLES IF THE METRIC IS 'PRE  
COMPUTED' X MUST BE A SQUARE DISTANCE MATRIX OTHERWISE IT CONTAINS A SAMPLE PER ROW

YIGNORED

RETURNS

XNEW ARRAY SHAPE NSAMPLES NCOMPONENTS EMBEDDING OF THE TRAINING DATA IN LOW  
DIMENSIONAL SPACE

GETPARAMS SELFDEEPTTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE  
PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT  
OF A NESTED OBJECT

RETURNS

SELF

EXAMPLES USING SKLEARNMANIFOLDTSNE

- TSNE THE EFFECT OF VARIOUS PERPLEXITY VALUES ON THE SHAPE
- COMPARISON OF MANIFOLD LEARNING METHODS
- MANIFOLD LEARNING METHODS ON A SEVERED SPHERE
- MANIFOLD LEARNING ON HANDWRITTEN DIGITS LOCALLY LINEAR EMBEDDING ISOMAP

MANIFOLDLOCALLYLINEAREMBEDDING X

PERFORM A LOCALLY LINEAR EMBEDDING ANALYSIS ON THE DATA

MANIFOLDSMACOF DISSIMILARITIES METRIC COMPUTES MULTIDIMENSIONAL SCALING USING THE SMACOF  
ALGORITHM

MANIFOLDSPETRALEMBEDDING ADJACENCY PROJECT THE SAMPLE ON THE FIRST EIGENVECTORS OF THE GRAPH  
LAPLACIAN

6236SKLEARNMANIFOLD LOCALLYLINEAREMBEDDING

SKLEARNMANIFOLD LOCALLYLINEAREMBEDDING XNNEIGHBORS NCOMPONENTS REG0001

EIGENSOLVER'AUTO' TOL1E06 MAXITER100

METHOD'STANDARD' HESSIANTOL00001

MODIFIEDTOL1E12 RANDOMSTATENONE

NJOBSNONE

PERFORM A LOCALLY LINEAR EMBEDDING ANALYSIS ON THE DATA

623SKLEARNMANIFOLD MANIFOLD LEARNING 1979

SCIKITLEARN USER GUIDE RELEASE 0213

READ MORE IN THE USER GUIDE

PARAMETERS

XARRAYLIKE NEARESTNEIGHBORS SAMPLE DATA SHAPE NSAMPLES NFEATURES IN THE FORM OF A NUMPY ARRAY OR A NEARESTNEIGHBORS OBJECT

NNEIGHBORS INTEGER NUMBER OF NEIGHBORS TO CONSIDER FOR EACH POINT

NCOMPONENTS INTEGER NUMBER OF COORDINATES FOR THE MANIFOLD

REGFLOAT REGULARIZATION CONSTANT MULTIPLIES THE TRACE OF THE LOCAL COVARIANCE MATRIX OF THE DISTANCES

EIGENSOLVER STRING ‘AUTO’ ‘ARPACK’ ‘DENSE’ AUTO ALGORITHM WILL ATTEMPT TO CHOOSE THE BEST METHOD FOR INPUT DATA

ARPACK USE ARNOLDI ITERATION IN SHIFTINVERT MODE FOR THIS METHOD M MAY BE A DENSE MATRIX SPARSE MATRIX OR GENERAL LINEAR OPERATOR WARNING ARPACK CAN BE UNSTABLE FOR SOME PROBLEMS IT IS BEST TO TRY SEVERAL RANDOM SEEDS IN ORDER TO CHECK RESULTS

DENSE USE STANDARD DENSE MATRIX OPERATIONS FOR THE EIGENVALUE DECOMPOSITION FOR THIS METHOD M MUST BE AN ARRAY OR MATRIX TYPE THIS METHOD SHOULD BE AVOIDED FOR LARGE PROBLEMS

TOLFLOAT OPTIONAL TOLERANCE FOR ‘ARPACK’ METHOD NOT USED IF EIGENSOLVER‘DENSE’

MAXITER INTEGER MAXIMUM NUMBER OF ITERATIONS FOR THE ARPACK SOLVER

METHOD ‘STANDARD’ ‘HESSIAN’ ‘MODIFIED’ ‘LTSA’

STANDARD USE THE STANDARD LOCALLY LINEAR EMBEDDING ALGORITHM SEE REFERENCE 1

HESSIAN USE THE HESSIAN EIGENMAP METHOD THIS METHOD REQUIRES NNEIGHBORS

NCOMPONENTS 1 NCOMPONENTS 1 2 SEE REFERENCE 2

MODIFIED USE THE MODIFIED LOCALLY LINEAR EMBEDDING ALGORITHM SEE REFERENCE 3

LTSA USE LOCAL TANGENT SPACE ALIGNMENT ALGORITHM SEE REFERENCE 4

HESSIAN\_TOL FLOAT OPTIONAL TOLERANCE FOR HESSIAN EIGENMAPPING METHOD ONLY USED IF METHOD ‘HESSIAN’

MODIFIED\_TOL FLOAT OPTIONAL TOLERANCE FOR MODIFIED LLE METHOD ONLY USED IF METHOD ‘MODIFIED’

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RAN

DOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE

RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NPRANDOM USED WHEN SOLVER ‘ARPACK’

NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF PARALLEL JOBS TO RUN FOR NEIGHBORS

SEARCHNONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS

USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

RETURNS

YARRAYLIKE SHAPE NSAMPLES NCOMPONENTS EMBEDDING VECTORS

SQUAREDERROR FLOAT RECONSTRUCTION ERROR FOR THE EMBEDDING VECTORS EQUIVALENT TO  $\|W(Y - \hat{Y})\|_2^2$  WHERE W ARE THE RECONSTRUCTION WEIGHTS

1980 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

REFERENCES

1234

EXAMPLES USING SKLEARNMANIFOLDLOCALLYLINEAREMBEDDING

•SWISS ROLL REDUCTION WITH LLE

6237SKLEARNMANIFOLD SMACOF

SKLEARNMANIFOLD SMACOFDISSIMILARITIES METRICTRUE NCOMPONENTS2 INITNONE NINIT8

NJOBSNONE MAXITER300 VERBOSE0 EPS0001 RANDOMSTATENONE

RETURNNITERFALSE

COMPUTES MULTIDIMENSIONAL SCALING USING THE SMACOF ALGORITHM

THE SMACOF SCALING BY MAJORIZING A COMPLICATED FUNCTION ALGORITHM IS A MULTIDIMENSIONAL SCALING ALGO

RITHM WHICH MINIMIZES AN OBJECTIVE FUNCTION THE STRESS USING A MAJORIZATION TECHNIQUE STRESS MAJORIZATION

ALSO KNOWN AS THE GUTTMAN TRANSFORM GUARANTEES A MONOTONE CONVERGENCE OF STRESS AND IS MORE POWERFUL THAN

TRADITIONAL TECHNIQUES SUCH AS GRADIENT DESCENT

THE SMACOF ALGORITHM FOR METRIC MDS CAN SUMMARIZED BY THE FOLLOWING STEPS

1 SET AN INITIAL START CONFIGURATION RANDOMLY OR NOT

2 COMPUTE THE STRESS

3 COMPUTE THE GUTTMAN TRANSFORM

4 ITERATE 2 AND 3 UNTIL CONVERGENCE

THE NONMETRIC ALGORITHM ADDS A MONOTONIC REGRESSION STEP BEFORE COMPUTING THE STRESS

PARAMETERS

DISSIMILARITIES NDARRAY SHAPE NSAMPLES NSAMPLES PAIRWISE DISSIMILARITIES BETWEEN THE

POINTS MUST BE SYMMETRIC

METRIC BOOLEAN OPTIONAL DEFAULT TRUE COMPUTE METRIC OR NONMETRIC SMACOF ALGORITHM

NCOMPONENTS INT OPTIONAL DEFAULT 2 NUMBER OF DIMENSIONS IN WHICH TO IMMERSE THE DIS

SIMILARITIES IF AN INIT ARRAY IS PROVIDED THIS OPTION IS OVERRIDDEN AND THE SHAPE OF INIT

IS USED TO DETERMINE THE DIMENSIONALITY OF THE EMBEDDING SPACE

INIT NDARRAY SHAPE NSAMPLES NCOMPONENTS OPTIONAL DEFAULT NONE STARTING CONFIGURA

TION OF THE EMBEDDING TO INITIALIZE THE ALGORITHM BY DEFAULT THE ALGORITHM IS INITIALIZED

WITH A RANDOMLY CHOSEN ARRAY

NINIT INT OPTIONAL DEFAULT 8 NUMBER OF TIMES THE SMACOF ALGORITHM WILL BE RUN WITH

DIFFERENT INITIALIZATIONS THE FINAL RESULTS WILL BE THE BEST OUTPUT OF THE RUNS DETERMINED BY

THE RUN WITH THE SMALLEST FINAL STRESS IF INIT IS PROVIDED THIS OPTION IS OVERRIDDEN AND A

SINGLE RUN IS PERFORMED

NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF JOBS TO USE FOR THE COMPUTATION

IF MULTIPLE INITIALIZATIONS ARE USED NINIT EACH RUN OF THE ALGORITHM IS COMPUTED IN

PARALLEL

NONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL

PROCESSORS SEE GLOSSARY FOR MORE DETAILS

623SKLEARNMANIFOLD MANIFOLD LEARNING 1981

SCIKITLEARN USER GUIDE RELEASE 0213

MAXITER INT OPTIONAL DEFAULT 300 MAXIMUM NUMBER OF ITERATIONS OF THE SMACOF ALGORITHM FOR A SINGLE RUN

VERBOSE INT OPTIONAL DEFAULT 0 LEVEL OF VERBOSITY

EPS FLOAT OPTIONAL DEFAULT 1E3 RELATIVE TOLERANCE WITH RESPECT TO STRESS AT WHICH TO DECLARE CONVERGENCE

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULT NONE THE GENERATOR USED TO INITIALIZE THE CENTERS IF INT RANDOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NPRANDOM

RETURNNITER BOOL OPTIONAL DEFAULT FALSE WHETHER OR NOT TO RETURN THE NUMBER OF ITERATIONS RETURNS

XNDARRAY SHAPE NSAMPLES NCOMPONENTS COORDINATES OF THE POINTS IN A NCOMPONENTS SPACE

STRESS FLOAT THE FINAL VALUE OF THE STRESS SUM OF SQUARED DISTANCE OF THE DISPARITIES AND THE DISTANCES FOR ALL CONSTRAINED POINTS

NITER INT THE NUMBER OF ITERATIONS CORRESPONDING TO THE BEST STRESS RETURNED ONLY IF RETURNNITER IS SET TOTRUE

NOTES

“MODERN MULTIDIMENSIONAL SCALING THEORY AND APPLICATIONS” BORG I GROENEN P SPRINGER SERIES IN STATISTICS 1997

“NONMETRIC MULTIDIMENSIONAL SCALING A NUMERICAL METHOD” KRUSKAL J PSYCHOMETRIKA 29 1964

“MULTIDIMENSIONAL SCALING BY OPTIMIZING GOODNESS OF FIT TO A NONMETRIC HYPOTHESIS” KRUSKAL J PSYCHOMETRIKA 29 1964

6238SKLEARNMANIFOLD SPECTRALEMBEDDING

SKLEARNMANIFOLD SPECTRALEMBEDDING ADJACENCY NCOMPONENTS8 EIGENSOLVERNONE RAN

DOMSTATENONE EIGENTOL00 NORMLAPLACIANTRUE

DROPFIRSTTRUE

PROJECT THE SAMPLE ON THE FIRST EIGENVECTORS OF THE GRAPH LAPLACIAN

THE ADJACENCY MATRIX IS USED TO COMPUTE A NORMALIZED GRAPH LAPLACIAN WHOSE SPECTRUM ESPECIALLY THE EIGENVECTORS ASSOCIATED TO THE SMALLEST EIGENVALUES HAS AN INTERPRETATION IN TERMS OF MINIMAL NUMBER OF CUTS NECESSARY TO SPLIT THE GRAPH INTO COMPARABLY SIZED COMPONENTS

THIS EMBEDDING CAN ALSO ‘WORK’ EVEN IF THE ADJACENCY VARIABLE IS NOT STRICTLY THE ADJACENCY MATRIX OF A GRAPH BUT MORE GENERALLY AN AFFINITY OR SIMILARITY MATRIX BETWEEN SAMPLES FOR INSTANCE THE HEAT KERNEL OF A EUCLIDEAN DISTANCE MATRIX OR A KNN MATRIX

HOWEVER CARE MUST TAKEN TO ALWAYS MAKE THE AFFINITY MATRIX SYMMETRIC SO THAT THE EIGENVECTOR DECOMPOSITION WORKS AS EXPECTED

NOTE LAPLACIAN EIGENMAPS IS THE ACTUAL ALGORITHM IMPLEMENTED HERE

READ MORE IN THE USER GUIDE

PARAMETERS

1982 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

ADJACENCY ARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NSAMPLES THE ADJACENCY MATRIX  
OF THE GRAPH TO EMBED

NCOMPONENTS INTEGER OPTIONAL DEFAULT 8 THE DIMENSION OF THE PROJECTION SUBSPACE

EIGENSOLVER NONE 'ARPACK' 'LOBPCG' OR 'AMG' DEFAULT NONE THE EIGENVALUE DECOMPO  
SITION STRATEGY TO USE AMG REQUIRES PYAMG TO BE INSTALLED IT CAN BE FASTER ON VERY LARGE  
SPARSE PROBLEMS BUT MAY ALSO LEAD TO INSTABILITIES

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULT NONE A PSEUDO RANDOM  
NUMBER GENERATOR USED FOR THE INITIALIZATION OF THE LOBPCG EIGENVECTORS DECOMPOSITION IF INT  
RANDOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE  
RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE  
RANDOMSTATE INSTANCE USED BY NPRANDOM USED WHEN SOLVER 'AMG'

EIGENTOL FLOAT OPTIONAL DEFAULT00 STOPPING CRITERION FOR EIGENDECOMPOSITION OF THE LAPLA  
CIAN MATRIX WHEN USING ARPACK EIGENSOLVER

NORMLAPLACIAN BOOL OPTIONAL DEFAULTTRUE IF TRUE THEN COMPUTE NORMALIZED LAPLACIAN

DROPPFIRST BOOL OPTIONAL DEFAULTTRUE WHETHER TO DROP THE FIRST EIGENVECTOR FOR SPECTRAL EM  
BEDDING THIS SHOULD BE TRUE AS THE FIRST EIGENVECTOR SHOULD BE CONSTANT VECTOR FOR CONNECTED  
GRAPH BUT FOR SPECTRAL CLUSTERING THIS SHOULD BE KEPT AS FALSE TO RETAIN THE FIRST EIGENVECTOR

RETURNS  
EMBEDDING ARRAY SHAPENSAMPLES NCOMPONENTS THE REDUCED SAMPLES

NOTES  
SPECTRAL EMBEDDING LAPLACIAN EIGENMAPS IS MOST USEFUL WHEN THE GRAPH HAS ONE CONNECTED COMPONENT IF  
THERE GRAPH HAS MANY COMPONENTS THE FIRST FEW EIGENVECTORS WILL SIMPLY UNCOVER THE CONNECTED COMPONENTS OF  
THE GRAPH

REFERENCES

- [HTTPSENWIKIPEDIAORGWIKILOBPCG](https://en.wikipedia.org/wiki/LOBPCG)
- TOWARD THE OPTIMAL PRECONDITIONED EIGENSOLVER LOCALLY OPTIMAL BLOCK PRECONDITIONED CONJUGATE GRA  
DIENT METHOD ANDREW V KNYAZEY [HTTPSDOIORG1011372FS1064827500366124](https://doi.org/10.1137/2FS1064827500366124)

624SKLEARNMETRICS METRICS

SEE THE MODEL EVALUATION QUANTIFYING THE QUALITY OF PREDICTIONS SECTION AND THE PAIRWISE METRICS AFFINITIES AND  
KERNELS SECTION OF THE USER GUIDE FOR FURTHER DETAILS THE SKLEARNMETRICS MODULE INCLUDES SCORE FUNCTIONS  
PERFORMANCE METRICS AND PAIRWISE METRICS AND DISTANCE COMPUTATIONS

6241 MODEL SELECTION INTERFACE

SEE THE THE SCORING PARAMETER DEFINING MODEL EVALUATION RULES SECTION OF THE USER GUIDE FOR FURTHER DETAILS

METRICSCHECKSCORING ESTIMATOR SCORING DETERMINE SCORER FROM USER OPTIONS

METRICSGETSCORER SCORING GET A SCORER FROM STRING

CONTINUED ON NEXT PAGE

624SKLEARNMETRICS METRICS 1983

SCIKITLEARN USER GUIDE RELEASE 0213

TABLE 6178 – CONTINUED FROM PREVIOUS PAGE

METRICSMAKESCORER SCOREFUNC MAKE A SCORER FROM A PERFORMANCE METRIC OR LOSS FUNCTION

SKLEARNMETRICS CHECKSCORING

SKLEARNMETRICS CHECKSCORING ESTIMATOR SCORINGNONE ALLOWNONEFALSE

DETERMINE SCORER FROM USER OPTIONS

A TYPEERROR WILL BE THROWN IF THE ESTIMATOR CANNOT BE SCORED

PARAMETERS

ESTIMATOR ESTIMATOR OBJECT IMPLEMENTING ‘FIT’ THE OBJECT TO USE TO FIT THE DATA

SCORING STRING CALLABLE OR NONE OPTIONAL DEFAULT NONE A STRING SEE MODEL EVALUATION DOCUMENTATION OR A SCORER CALLABLE OBJECT FUNCTION WITH SIGNATURE SCORERESTIMATOR

X Y

ALLOWNONE BOOLEAN OPTIONAL DEFAULT FALSE IF NO SCORING IS SPECIFIED AND THE ESTIMATOR HAS NO SCORE FUNCTION WE CAN EITHER RETURN NONE OR RAISE AN EXCEPTION

RETURNS

SCORING CALLABLE A SCORER CALLABLE OBJECT FUNCTION WITH SIGNATURE SCORERESTIMATOR

X Y

SKLEARNMETRICS GETSCORER

SKLEARNMETRICS GETSCORER SCORING

GET A SCORER FROM STRING

PARAMETERS

SCORING STR CALLABLE SCORING METHOD AS STRING IF CALLABLE IT IS RETURNED AS IS

RETURNS

SCORER CALLABLE THE SCORER

SKLEARNMETRICS MAKESCORER

SKLEARNMETRICS MAKESCORER SCOREFUNC GREATERISBETTERTRUE NEEDSPROBAFALSE

NEEDSTHRESHOLDFALSE KWARGS

MAKE A SCORER FROM A PERFORMANCE METRIC OR LOSS FUNCTION

THIS FACTORY FUNCTION WRAPS SCORING FUNCTIONS FOR USE IN GRIDSEARCHCV AND CROSSVALSCORE IT TAKES A SCORE FUNCTION SUCH AS ACCURACYScore MEANSQUAREDERROR ADJUSTEDRANDINDEX OR AVERAGEPRECISION AND RETURNS A CALLABLE THAT SCORES AN ESTIMATOR’S OUTPUT

READ MORE IN THE USER GUIDE

PARAMETERS

SCOREFUNC CALLABLE SCORE FUNCTION OR LOSS FUNCTION WITH SIGNATURE SCOREFUNC

YPRED KWARGS

GREATERISBETTER BOOLEAN DEFAULTTRUE WHETHER SCOREFUNC IS A SCORE FUNCTION DEFAULT MEANING HIGH IS GOOD OR A LOSS FUNCTION MEANING LOW IS GOOD IN THE LATTER CASE THE SCORER OBJECT WILL SIGNFLIP THE OUTCOME OF THE SCOREFUNC

1984 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

NEEDSPROBA BOOLEAN DEFAULTFALSE WHETHER SCOREFUNC REQUIRES PREDICTPROBA TO GET PROBABILITY ESTIMATES OUT OF A CLASSIFIER

IF TRUE FOR BINARY YTRUE THE SCORE FUNCTION IS SUPPOSED TO ACCEPT A 1D YPRED IE PROBABILITY OF THE POSITIVE CLASS SHAPE NSAMPLES

NEEDSTHRESHOLD BOOLEAN DEFAULTFALSE WHETHER SCOREFUNC TAKES A CONTINUOUS DECISION CERTAINTY THIS ONLY WORKS FOR BINARY CLASSIFICATION USING ESTIMATORS THAT HAVE EITHER A DECISIONFUNCTION OR PREDICTPROBA METHOD

IF TRUE FOR BINARY YTRUE THE SCORE FUNCTION IS SUPPOSED TO ACCEPT A 1D YPRED IE PROBABILITY OF THE POSITIVE CLASS OR THE DECISION FUNCTION SHAPE NSAMPLES

FOR EXAMPLE AVERAGEPRECISION OR THE AREA UNDER THE ROC CURVE CAN NOT BE COMPUTED USING DISCRETE PREDICTIONS ALONE

KWARGS ADDITIONAL ARGUMENTS ADDITIONAL PARAMETERS TO BE PASSED TO SCOREFUNC

RETURNS

SCORER CALLABLE CALLABLE OBJECT THAT RETURNS A SCALAR SCORE GREATER IS BETTER

EXAMPLES

```
FROM SKLEARNMETRICS IMPORT FBETASCORE
MAKESCORER
FTWOSCORER MAKESCORERFBETASCORE BETA2
FTWOSCORER
MAKESCORERFBETASCORE BETA2
FROM SKLEARNMODELSELECTION IMPORT GRIDSEARCHCV
FROM SKLEARN SVM IMPORT LINEARSVC
GRID = GRIDSEARCHCVLINEARSVC PARAMGRIDC 1 10
SCORINGFTWOSCORER
```

EXAMPLES USING SKLEARNMETRICSMAKESCORER

- DEMONSTRATION OF MULTIMETRIC EVALUATION ON CROSSVALSCORE AND GRIDSEARCHCV

6242 CLASSIFICATION METRICS

SEE THE CLASSIFICATION METRICS SECTION OF THE USER GUIDE FOR FURTHER DETAILS

METRICSACCURACYScore YTRUE YPRED ACCURACY CLASSIFICATION SCORE

METRICSauc X Y REORDER COMPUTE AREA UNDER THE CURVE AUC USING THE TRAPEZOIDAL RULE

METRICSaverageprecisionscore YTRUE

YSCORECOMPUTE AVERAGE PRECISION AP FROM PREDICTION SCORES

METRICSbalancedaccuracyscore YTRUE

YPREDCOMPUTE THE BALANCED ACCURACY

METRICSbrierscoreloss YTRUE YPROB COMPUTE THE BRIER SCORE

METRICSclassificationreport YTRUE

YPREDBUILD A TEXT REPORT SHOWING THE MAIN CLASSIFICATION METRICS

METRICScohenkappascor Y1 Y2 LABELS COHEN'S KAPPA A STATISTIC THAT MEASURES INTERANNOTATOR AGREEMENT

CONTINUED ON NEXT PAGE

624SKLEARNMETRICS METRICS 1985

SCIKITLEARN USER GUIDE RELEASE 0213

TABLE 6179 – CONTINUED FROM PREVIOUS PAGE

METRICSCONFUSIONMATRIX YTRUE YPRED    COMPUTE CONFUSION MATRIX TO EVALUATE THE ACCURACY OF A CLASSIFICATION

METRICSF1SCORE YTRUE YPRED LABELS    COMPUTE THE F1 SCORE ALSO KNOWN AS BALANCED FSCORE OR FMEASURE

METRICSF1BETAScore YTRUE YPRED BETA    COMPUTE THE FBETA SCORE

METRICSHAMMINGLOSS YTRUE YPRED    COMPUTE THE AVERAGE HAMMING LOSS

METRICSHINGELOSS YTRUE PREDDECISION    AVERAGE HINGE LOSS NONREGULARIZED

METRICSJACCARDScore YTRUE YPRED    JACCARD SIMILARITY COEFFICIENT SCORE

METRICSLGLOSS YTRUE YPRED EPS    LOG LOSS AKA LOGISTIC LOSS OR CROSSENTROPY LOSS

METRICSMATTHEWSCORRcoef YTRUE YPRED  
  COMPUTE THE MATTHEWS CORRELATION COEFFICIENT MCC

METRICSMULTILABELCONFUSIONMATRIX YTRUE  
  COMPUTE A CONFUSION MATRIX FOR EACH CLASS OR SAMPLE

METRICSPRECISSIONRECALLCURVE YTRUE    COMPUTE PRECISIONRECALL PAIRS FOR DIFFERENT PROBABILITY THRESHOLDS

METRICSPRECISSIONRECALLFSCORESUPPORT    COMPUTE PRECISION RECALL FMEASURE AND SUPPORT FOR EACH CLASS

METRICSPRECISSIONSCORE YTRUE YPRED    COMPUTE THE PRECISION

METRICSPRECALLSCORE YTRUE YPRED    COMPUTE THE RECALL

METRICSPROCAUCSCORE YTRUE YSCORE    COMPUTE AREA UNDER THE RECEIVER OPERATING CHARACTERISTIC CURVE ROC AUC FROM PREDICTION SCORES

METRICSPROCCURVE YTRUE YSCORE    COMPUTE RECEIVER OPERATING CHARACTERISTIC ROC

METRICSZEROONELOSS YTRUE YPRED    ZEROONE CLASSIFICATION LOSS

SKLEARNMETRICS ACCURACYScore  
SKLEARNMETRICS ACCURACYScore YTRUE YPRED NORMALIZETRUE SAMPLEWEIGHTNONE  
ACCURACY CLASSIFICATION SCORE

IN MULTILABEL CLASSIFICATION THIS FUNCTION COMPUTES SUBSET ACCURACY THE SET OF LABELS PREDICTED FOR A SAMPLE MUST EXACTLY MATCH THE CORRESPONDING SET OF LABELS IN YTRUE

READ MORE IN THE USER GUIDE

PARAMETERS

YTRUE 1D ARRAYLIKE OR LABEL INDICATOR ARRAY    SPARSE MATRIX GROUND TRUTH CORRECT LABELS

YPRED 1D ARRAYLIKE OR LABEL INDICATOR ARRAY    SPARSE MATRIX PREDICTED LABELS AS RETURNED BY A CLASSIFIER

NORMALIZE BOOL OPTIONAL DEFAULTTRUE IF FALSE    RETURN THE NUMBER OF CORRECTLY CLASSIFIED SAMPLES OTHERWISE RETURN THE FRACTION OF CORRECTLY CLASSIFIED SAMPLES

SAMPLEWEIGHT ARRAYLIKE OF SHAPE    NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT IFNORMALIZE TRUE    RETURN THE FRACTION OF CORRECTLY CLASSIFIED SAMPLES

FLOAT ELSE RETURNS THE NUMBER OF CORRECTLY CLASSIFIED SAMPLES INT

THE BEST PERFORMANCE IS 1 WITH NORMALIZE TRUE AND THE NUMBER OF SAMPLES WITH

NORMALIZE FALSE

SEE ALSO

JACCARDScore HAMMINGLOSS ZEROONELOSS

1986 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

NOTES

IN BINARY AND MULTICLASS CLASSIFICATION THIS FUNCTION IS EQUAL TO THE JACCARDSCORE FUNCTION

EXAMPLES

```
FROM SKLEARNMETRICS IMPORT ACCURACYScore
YPRED 0 2 1 3
YTRUE 0 1 2 3
ACCURACYScoreYTRUE YPRED
05
ACCURACYScoreYTRUE YPRED NORMALIZE FALSE
2
IN THE MULTILABEL CASE WITH BINARY LABEL INDICATORS
IMPORT NUMPY AS NP
ACCURACYScoreNPARRAY0 1 1 1 NPONES2 2
05
EXAMPLES USING SKLEARNMETRICSACCURACYScore
•PLOT CLASSIFICATION PROBABILITY
•MULTICLASS ADABOOSTED DECISION TREES
•PROBABILISTIC PREDICTIONS WITH GAUSSIAN PROCESS CLASSIFICATION GPC
•DEMONSTRATION OF MULTIMETRIC EVALUATION ON CROSSVALScore AND GRIDSEARCHCV
•IMPORTANCE OF FEATURE SCALING
•CLASSIFICATION OF TEXT DOCUMENTS USING SPARSE FEATURES
```

SKLEARNMETRICS AUC

SKLEARNMETRICS AUCXYREORDER'DEPRECATED'

COMPUTE AREA UNDER THE CURVE AUC USING THE TRAPEZOIDAL RULE

THIS IS A GENERAL FUNCTION GIVEN POINTS ON A CURVE FOR COMPUTING THE AREA UNDER THE ROC CURVE SEE ROCAUCScore FOR AN ALTERNATIVE WAY TO SUMMARIZE A PRECISIONRECALL CURVE SEE AVERAGEPRECISIONScore

PARAMETERS

XARRAY SHAPE N X COORDINATES THESE MUST BE EITHER MONOTONIC INCREASING OR MONOTONIC DECREASING

YARRAY SHAPE N Y COORDINATES

REORDER BOOLEAN OPTIONAL DEFAULT'DEPRECATED' WHETHER TO SORT X BEFORE COMPUTING IF FALSE ASSUME THAT X MUST BE EITHER MONOTONIC INCREASING OR MONOTONIC DECREASING IF TRUE Y IS USED TO BREAK TIES WHEN SORTING X MAKE SURE THAT Y HAS A MONOTONIC RELATION TO X WHEN SETTING REORDER TO TRUE

DEPRECATED SINCE VERSION 020 PARAMETER REORDER HAS BEEN DEPRECATED IN VERSION 020 AND WILL BE REMOVED IN 022 IT'S INTRODUCED FOR ROCAUCScore NOT FOR GENERAL USE AND IS 624SKLEARNMETRICS METRICS 1987

SCIKITLEARN USER GUIDE RELEASE 0213

NO LONGER USED THERE WHAT’S MORE THE RESULT FROM AUC WILL BE SIGNIFICANTLY INFLUENCED IF X IS SORTED UNEXPECTEDLY DUE TO SLIGHT FLOATING POINT ERROR SEE ISSUE 9786 FUTURE AND DEFAULT BEHAVIOR IS EQUIVALENT TO REORDERFALSE

RETURNS  
AUC FLOAT  
SEE ALSO

ROCAUCSCORE COMPUTE THE AREA UNDER THE ROC CURVE  
AVERAGEPRECISIONSCORE COMPUTE AVERAGE PRECISION FROM PREDICTION SCORES  
PRECISIONRECALLCURVE COMPUTE PRECISIONRECALL PAIRS FOR DIFFERENT PROBABILITY THRESHOLDS  
EXAMPLES

```
import numpy as np
from sklearn import metrics
y = np.array(1 1 2 2)
pred = np.array(0 1 0 4 0 3 5 0 8)
fpr, tpr, thresholds = metrics.roc_curve(y, pred, poslabel=2)
metrics.auc(fpr, tpr)
```

075  
EXAMPLES USING SKLEARNMETRICS  
•SPECIES DISTRIBUTION MODELING  
•RECEIVER OPERATING CHARACTERISTIC ROC WITH CROSS VALIDATION  
•RECEIVER OPERATING CHARACTERISTIC ROC  
SKLEARNMETRICS AVERAGEPRECISIONSCORE  
SKLEARNMETRICS AVERAGEPRECISIONSCORE YTRUE YSCORE AVERAGE’MACRO’ POSLABEL1  
SAMPLEWEIGHTNONE  
COMPUTE AVERAGE PRECISION AP FROM PREDICTION SCORES

AP SUMMARIZES A PRECISIONRECALL CURVE AS THE WEIGHTED MEAN OF PRECISIONS ACHIEVED AT EACH THRESHOLD WITH THE INCREASE IN RECALL FROM THE PREVIOUS THRESHOLD USED AS THE WEIGHT  
APΣ

$\frac{1}{n} \sum_{i=1}^n p_i r_i$   
WHERE  $p_i$  AND  $r_i$  ARE THE PRECISION AND RECALL AT THE NTH THRESHOLD 1 THIS IMPLEMENTATION IS NOT INTERPOLATED AND IS DIFFERENT FROM COMPUTING THE AREA UNDER THE PRECISIONRECALL CURVE WITH THE TRAPEZOIDAL RULE WHICH USES LINEAR INTERPOLATION AND CAN BE TOO OPTIMISTIC  
NOTE THIS IMPLEMENTATION IS RESTRICTED TO THE BINARY CLASSIFICATION TASK OR MULTILABEL CLASSIFICATION TASK  
READ MORE IN THE USER GUIDE

PARAMETERS  
YTRUE ARRAY SHAPE (NSAMPLES) OR NSAMPLES NCLASSES TRUE BINARY LABELS OR BINARY LABEL INDICATORS  
1988 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

YSCORE ARRAY SHAPE NSAMPLES OR NSAMPLES NCLASSES TARGET SCORES CAN EITHER BE PROBABILITY ESTIMATES OF THE POSITIVE CLASS CONFIDENCE VALUES OR NONTHRESHOLDED MEASURE OF DECISIONS AS RETURNED BY “DECISIONFUNCTION” ON SOME CLASSIFIERS

AVERAGE STRING NONE ‘MICRO’ ‘MACRO’ DEFAULT ‘SAMPLES’ ‘WEIGHTED’ IF NONE THE SCORES FOR EACH CLASS ARE RETURNED OTHERWISE THIS DETERMINES THE TYPE OF AVERAGING PERFORMED ON THE DATA

MICRO CALCULATE METRICS GLOBALLY BY CONSIDERING EACH ELEMENT OF THE LABEL INDICATOR MATRIX AS A LABEL

MACRO CALCULATE METRICS FOR EACH LABEL AND FIND THEIR UNWEIGHTED MEAN THIS DOES NOT TAKE LABEL IMBALANCE INTO ACCOUNT

WEIGHTED CALCULATE METRICS FOR EACH LABEL AND FIND THEIR AVERAGE WEIGHTED BY SUPPORT THE NUMBER OF TRUE INSTANCES FOR EACH LABEL

SAMPLES CALCULATE METRICS FOR EACH INSTANCE AND FIND THEIR AVERAGE

WILL BE IGNORED WHEN YTRUE IS BINARY

POSLABEL INT OR STR DEFAULT1 THE LABEL OF THE POSITIVE CLASS ONLY APPLIED TO BINARY YTRUE FOR MULTILABELINDICATOR YTRUE POSLABEL IS FIXED TO 1

SAMPLEWEIGHT ARRAYLIKE OF SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

AVERAGEPRECISION FLOAT

SEE ALSO

ROCAUCSCORE COMPUTE THE AREA UNDER THE ROC CURVE

PRECISIONRECALLCURVE COMPUTE PRECISIONRECALL PAIRS FOR DIFFERENT PROBABILITY THRESHOLDS

NOTES

CHANGED IN VERSION 019 INSTEAD OF LINEARLY INTERPOLATING BETWEEN OPERATING POINTS PRECISIONS ARE WEIGHTED BY THE CHANGE IN RECALL SINCE THE LAST OPERATING POINT

REFERENCES

1

EXAMPLES

```
import numpy as np
from sklearn.metrics import average_precision_score
ytrue = np.array([0, 1, 1])
yscores = np.array([0.1, 0.4, 0.35, 0.8])
average_precision_score(ytrue, yscores)
```

083

624SKLEARNMETRICS METRICS 1989

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES USING SKLEARNMETRICS

SAVERAGEPRECISIONSCORE

- PRECISIONRECALL

SKLEARNMETRICS BALANCEDACCURACYSORE

SKLEARNMETRICS BALANCEDACCURACYSORE YTRUE YPRED SAMPLEWEIGHTNONE ADJUSTEDFALSE

COMPUTE THE BALANCED ACCURACY

THE BALANCED ACCURACY IN BINARY AND MULTICLASS CLASSIFICATION PROBLEMS TO DEAL WITH IMBALANCED DATASETS IT IS DEFINED AS THE AVERAGE OF RECALL OBTAINED ON EACH CLASS

THE BEST VALUE IS 1 AND THE WORST VALUE IS 0 WHEN ADJUSTEDFALSE

READ MORE IN THE USER GUIDE

PARAMETERS

YTRUE 1D ARRAYLIKE GROUND TRUTH CORRECT TARGET VALUES

YPRED 1D ARRAYLIKE ESTIMATED TARGETS AS RETURNED BY A CLASSIFIER

SAMPLEWEIGHT ARRAYLIKE OF SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

ADJUSTED BOOL DEFAULTFALSE WHEN TRUE THE RESULT IS ADJUSTED FOR CHANCE SO THAT RANDOM PERFORMANCE WOULD SCORE 0 AND PERFECT PERFORMANCE SCORES 1

RETURNS

BALANCEDACCURACY FLOAT

SEE ALSO

RECALLSCORE ROCAUCSCORE

NOTES

SOME LITERATURE PROMOTES ALTERNATIVE DEFINITIONS OF BALANCED ACCURACY OUR DEFINITION IS EQUIVALENT TO ACCURACYSORE WITH CLASSBALANCED SAMPLE WEIGHTS AND SHARES DESIRABLE PROPERTIES WITH THE BINARY CASE

SEE THE USER GUIDE

REFERENCES

12

EXAMPLES

```
FROM SKLEARNMETRICS IMPORT BALANCEDACCURACYSORE
YTRUE 0 1 0 0 1 0
YPRED 0 1 0 0 0 1
BALANCEDACCURACYSOREYTRUE YPRED
```

0625

1990 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

SKLEARNMETRICS BRIERSCORELOSS

SKLEARNMETRICS BRIERSCORELOSS YTRUE YPROB SAMPLEWEIGHTNONE POSLABELNONE

COMPUTE THE BRIER SCORE THE SMALLER THE BRIER SCORE THE BETTER HENCE THE NAMING WITH “LOSS” ACROSS ALL ITEMS IN A SET N PREDICTIONS THE BRIER SCORE MEASURES THE MEAN SQUARED DIFFERENCE BETWEEN 1 THE PREDICTED PROBABILITY ASSIGNED TO THE POSSIBLE OUTCOMES FOR ITEM I AND 2 THE ACTUAL OUTCOME THEREFORE THE LOWER THE BRIER SCORE IS FOR A SET OF PREDICTIONS THE BETTER THE PREDICTIONS ARE CALIBRATED NOTE THAT THE BRIER SCORE ALWAYS TAKES ON A VALUE BETWEEN ZERO AND ONE SINCE THIS IS THE LARGEST POSSIBLE DIFFERENCE BETWEEN A PREDICTED PROBABILITY WHICH MUST BE BETWEEN ZERO AND ONE AND THE ACTUAL OUTCOME WHICH CAN TAKE ON VALUES OF ONLY 0 AND 1 THE BRIER LOSS IS COMPOSED OF REFINEMENT LOSS AND CALIBRATION LOSS THE BRIER SCORE IS APPROPRIATE FOR BINARY AND CATEGORICAL OUTCOMES THAT CAN BE STRUCTURED AS TRUE OR FALSE BUT IS INAPPROPRIATE FOR ORDINAL VARIABLES WHICH CAN TAKE ON THREE OR MORE VALUES THIS IS BECAUSE THE BRIER SCORE ASSUMES THAT ALL POSSIBLE OUTCOMES ARE EQUIVALENTLY “DISTANT” FROM ONE ANOTHER WHICH LABEL IS CONSIDERED TO BE THE POSITIVE LABEL IS CONTROLLED VIA THE PARAMETER POSLABEL WHICH DEFAULTS TO 1 READ MORE IN THE USER GUIDE

PARAMETERS

YTRUE ARRAY SHAPE NSAMPLES TRUE TARGETS

YPROB ARRAY SHAPE NSAMPLES PROBABILITIES OF THE POSITIVE CLASS

SAMPLEWEIGHT ARRAYLIKE OF SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

POSLABEL INT OR STR DEFAULTNONE LABEL OF THE POSITIVE CLASS DEFAULTS TO THE GREATER LABEL

UNLESS YTRUE IS ALL 0 OR ALL 1 IN WHICH CASE POSLABEL DEFAULTS TO 1

RETURNS

SCORE FLOAT BRIER SCORE

REFERENCES

1

EXAMPLES

```
import numpy as np
from sklearn.metrics import brier_score_loss
y_true = np.array([1, 1, 0])
y_prob_categorical = np.array(['spam', 'ham', 'ham', 'spam'])
y_prob = np.array([0.1, 0.9, 0.8, 0.3])
brier_score_loss(y_true, y_prob)
0.037
brier_score_loss(y_true, y_prob, pos_label=0)
0.037
brier_score_loss(y_true_categorical, y_prob, pos='spam')
0.037
brier_score_loss(y_true, y_prob, pos='spam')
0.037
```

EXAMPLES USING SKLEARNMETRICSBRIERSCORELOSS

- PROBABILITY CALIBRATION CURVES

624SKLEARNMETRICS METRICS 1991

SCIKITLEARN USER GUIDE RELEASE 0213

- PROBABILITY CALIBRATION OF CLASSIFIERS

SKLEARNMETRICS CLASSIFICATIONREPORT

SKLEARNMETRICS CLASSIFICATIONREPORT YTRUE YPRED LABELSNONE TARGETNAMESNONE

SAMPLEWEIGHTNONE DIGITS2 OUTPUTDICTFALSE

BUILD A TEXT REPORT SHOWING THE MAIN CLASSIFICATION METRICS

READ MORE IN THE USER GUIDE

PARAMETERS

YTRUE 1D ARRAYLIKE OR LABEL INDICATOR ARRAY SPARSE MATRIX GROUND TRUTH CORRECT TARGET VALUES

YPRED 1D ARRAYLIKE OR LABEL INDICATOR ARRAY SPARSE MATRIX ESTIMATED TARGETS AS RETURNED BY A CLASSIFIER

LABELS ARRAY SHAPE NLABELS OPTIONAL LIST OF LABEL INDICES TO INCLUDE IN THE REPORT

TARGETNAMES LIST OF STRINGS OPTIONAL DISPLAY NAMES MATCHING THE LABELS SAME ORDER

SAMPLEWEIGHT ARRAYLIKE OF SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

DIGITS INT NUMBER OF DIGITS FOR FORMATTING OUTPUT FLOATING POINT VALUES WHEN OUTPUTDICT

ISTRUE THIS WILL BE IGNORED AND THE RETURNED VALUES WILL NOT BE ROUNDED

OUTPUTDICT BOOL DEFAULT FALSE IF TRUE RETURN OUTPUT AS DICT

RETURNS

REPORT STRING DICT TEXT SUMMARY OF THE PRECISION RECALL F1 SCORE FOR EACH CLASS DICTIONARY

RETURNED IF OUTPUTDICT IS TRUE DICTIONARY HAS THE FOLLOWING STRUCTURE

LABEL 1 PRECISION05

RECALL10

F1SCORE067

SUPPORT1

LABEL 2

THE REPORTED AVERAGES INCLUDE MACRO AVERAGE AVERAGING THE UNWEIGHTED MEAN PER LABEL WEIGHTED AVERAGE AVERAGING THE SUPPORTWEIGHTED MEAN PER LABEL SAMPLE AVERAGE ONLY FOR MULTILABEL CLASSIFICATION AND MICRO AVERAGE AVERAGING THE TOTAL TRUE POSITIVES FALSE NEGATIVES AND FALSE POSITIVES IT IS ONLY SHOWN FOR MULTILABEL OR MULTICLASS WITH A SUBSET OF CLASSES BECAUSE IT IS ACCURACY OTHERWISE SEE ALSOFUNCPRECISIONRECALLFSCORESUPPORT FOR MORE DETAILS ON AVERAGES

NOTE THAT IN BINARY CLASSIFICATION RECALL OF THE POSITIVE CLASS IS ALSO KNOWN AS “SENSITIVITY”

RECALL OF THE NEGATIVE CLASS IS “SPECIFICITY”

SEE ALSO

PRECISIONRECALLFSCORESUPPORT CONFUSIONMATRIX

MULTILABELCONFUSIONMATRIX

1992 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES

FROM SKLEARNMETRICS IMPORT CLASSIFICATIONREPORT

YTRUE 0 1 2 2 2

YPRED 0 0 2 2 1

TARGETNAMES CLASS 0 CLASS 1 CLASS 2

PRINTCLASSIFICATIONREPORTYTRUE YPRED TARGETNAMESTARGETNAMES

PRECISION RECALL F1SCORE SUPPORT

CLASS 0 050 100 067 1

CLASS 1 000 000 000 1

CLASS 2 100 067 080 3

ACCURACY 060 5

MACRO AVG 050 056 049 5

WEIGHTED AVG 070 060 061 5

YPRED 1 1 0

YTRUE 1 1 1

PRINTCLASSIFICATIONREPORTYTRUE YPRED LABELS1 2 3

PRECISION RECALL F1SCORE SUPPORT

1 100 067 080 3

2 000 000 000 0

3 000 000 000 0

MICRO AVG 100 067 080 3

MACRO AVG 033 022 027 3

WEIGHTED AVG 100 067 080 3

EXAMPLES USING SKLEARNMETRICSCCLASSIFICATIONREPORT

- FACES RECOGNITION EXAMPLE USING EIGENFACES AND SVMs
- RECOGNIZING HANDWRITTEN DIGITS
- COLUMN TRANSFORMER WITH HETEROGENEOUS DATA SOURCES
- PIPELINE ANOVA SVM
- PARAMETER ESTIMATION USING GRID SEARCH WITH CROSSVALIDATION
- RESTRICTED BOLTZMANN MACHINE FEATURES FOR DIGIT CLASSIFICATION
- LABEL PROPAGATION DIGITS DEMONSTRATING PERFORMANCE
- LABEL PROPAGATION DIGITS ACTIVE LEARNING
- CLASSIFICATION OF TEXT DOCUMENTS USING SPARSE FEATURES

SKLEARNMETRICS COHENKAPPASCORE

SKLEARNMETRICS COHENKAPPASCORE Y1 Y2 LABELSNONE WEIGHTSNONE SAM

PLEWEIGHTNONE

COHEN'S KAPPA A STATISTIC THAT MEASURES INTERANNOTATOR AGREEMENT

624SKLEARNMETRICS METRICS 1993

SCIKITLEARN USER GUIDE RELEASE 0213

THIS FUNCTION COMPUTES COHEN’S KAPPA 1 A SCORE THAT EXPRESSES THE LEVEL OF AGREEMENT BETWEEN TWO ANNOTATORS ON A CLASSIFICATION PROBLEM IT IS DEFINED AS

$$\frac{p_{11} - p_1 p_2}{1 - p_1 p_2}$$

WHERE  $p_{11}$  IS THE EMPIRICAL PROBABILITY OF AGREEMENT ON THE LABEL ASSIGNED TO ANY SAMPLE THE OBSERVED AGREEMENT RATIO AND  $p_1$  IS THE EXPECTED AGREEMENT WHEN BOTH ANNOTATORS ASSIGN LABELS RANDOMLY  $p_2$  IS ESTIMATED USING A PERANNOTATOR EMPIRICAL PRIOR OVER THE CLASS LABELS 2

READ MORE IN THE USER GUIDE

PARAMETERS

Y1ARRAY SHAPE NSAMPLES LABELS ASSIGNED BY THE FIRST ANNOTATOR  
Y2ARRAY SHAPE NSAMPLES LABELS ASSIGNED BY THE SECOND ANNOTATOR THE KAPPA STATISTIC IS SYMMETRIC SO SWAPPING Y1ANDY2DOESN’T CHANGE THE VALUE  
LABELS ARRAY SHAPE NCLASSES OPTIONAL LIST OF LABELS TO INDEX THE MATRIX THIS MAY BE USED TO SELECT A SUBSET OF LABELS IF NONE ALL LABELS THAT APPEAR AT LEAST ONCE IN Y1ORY2ARE USED  
WEIGHTS STR OPTIONAL LIST OF WEIGHTING TYPE TO CALCULATE THE SCORE NONE MEANS NO WEIGHTED “LINEAR” MEANS LINEAR WEIGHTED “QUADRATIC” MEANS QUADRATIC WEIGHTED  
SAMPLEWEIGHT ARRAYLIKE OF SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

KAPPA FLOAT THE KAPPA STATISTIC WHICH IS A NUMBER BETWEEN 1 AND 1 THE MAXIMUM VALUE MEANS COMPLETE AGREEMENT ZERO OR LOWER MEANS CHANCE AGREEMENT

REFERENCES

123

SKLEARNMETRICS CONFUSIONMATRIX

SKLEARNMETRICS CONFUSIONMATRIX YTRUE YPRED LABELSNONE SAMPLEWEIGHTNONE

COMPUTE CONFUSION MATRIX TO EVALUATE THE ACCURACY OF A CLASSIFICATION

BY DEFINITION A CONFUSION MATRIX  $C$  IS SUCH THAT  $C_{ij}$  IS EQUAL TO THE NUMBER OF OBSERVATIONS KNOWN TO BE IN GROUP  $j$  BUT PREDICTED TO BE IN GROUP  $i$

THUS IN BINARY CLASSIFICATION THE COUNT OF TRUE NEGATIVES IS  $C_{00}$  FALSE NEGATIVES IS  $C_{10}$  TRUE POSITIVES IS  $C_{11}$  AND FALSE POSITIVES IS  $C_{01}$

READ MORE IN THE USER GUIDE

PARAMETERS

YTRUE ARRAY SHAPE NSAMPLES GROUND TRUTH CORRECT TARGET VALUES  
YPRED ARRAY SHAPE NSAMPLES ESTIMATED TARGETS AS RETURNED BY A CLASSIFIER  
LABELS ARRAY SHAPE NCLASSES OPTIONAL LIST OF LABELS TO INDEX THE MATRIX THIS MAY BE USED TO REORDER OR SELECT A SUBSET OF LABELS IF NONE IS GIVEN THOSE THAT APPEAR AT LEAST ONCE IN YTRUE ORYPRED ARE USED IN SORTED ORDER  
SAMPLEWEIGHT ARRAYLIKE OF SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

CARRAY SHAPE NCLASSES NCLASSES CONFUSION MATRIX



SCIKITLEARN USER GUIDE RELEASE 0213

REFERENCES

1

EXAMPLES

FROM SKLEARNMETRICS IMPORT CONFUSIONMATRIX

YTRUE 2 0 2 2 0 1

YPRED 0 0 2 2 0 2

CONFUSIONMATRIXYTRUE YPRED

ARRAY2 0 0

0 0 1

1 0 2

YTRUE CAT ANT CAT CAT ANT BIRD

YPRED ANT ANT CAT CAT ANT CAT

CONFUSIONMATRIXYTRUE YPRED LABELSANT BIRD CAT

ARRAY2 0 0

0 0 1

1 0 2

IN THE BINARY CASE WE CAN EXTRACT TRUE POSITIVES ETC AS FOLLOWS

TN FP FN TP CONFUSIONMATRIX0 1 0 1 1 1 1 0RAVEL

TN FP FN TP

0 2 1 1

EXAMPLES USING SKLEARNMETRICSCONFUSIONMATRIX

- FACES RECOGNITION EXAMPLE USING EIGENFACES AND SVMs
- RECOGNIZING HANDWRITTEN DIGITS
- CONFUSION MATRIX
- LABEL PROPAGATION DIGITS DEMONSTRATING PERFORMANCE
- LABEL PROPAGATION DIGITS ACTIVE LEARNING
- CLASSIFICATION OF TEXT DOCUMENTS USING SPARSE FEATURES

SKLEARNMETRICS F1SCORE

SKLEARNMETRICS F1SCORE YTRUE YPRED LABELSNONE POSLABEL1 AVERAGE' BINARY' SAM

PLEWEIGHTNONE

COMPUTE THE F1 SCORE ALSO KNOWN AS BALANCED FSCORE OR FMEASURE

THE F1 SCORE CAN BE INTERPRETED AS A WEIGHTED AVERAGE OF THE PRECISION AND RECALL WHERE AN F1 SCORE REACHES ITS BEST VALUE AT 1 AND WORST SCORE AT 0 THE RELATIVE CONTRIBUTION OF PRECISION AND RECALL TO THE F1 SCORE ARE EQUAL

THE FORMULA FOR THE F1 SCORE IS

F1 2PRECISION RECALL PRECISION RECALL

624SKLEARNMETRICS METRICS 1995

SCIKITLEARN USER GUIDE RELEASE 0213

IN THE MULTICLASS AND MULTILABEL CASE THIS IS THE AVERAGE OF THE F1 SCORE OF EACH CLASS WITH WEIGHTING DEPENDING ON THEAVERAGE PARAMETER

READ MORE IN THE USER GUIDE

PARAMETERS

YTRUE 1D ARRAYLIKE OR LABEL INDICATOR ARRAY SPARSE MATRIX GROUND TRUTH CORRECT TARGET VALUES

YPRED 1D ARRAYLIKE OR LABEL INDICATOR ARRAY SPARSE MATRIX ESTIMATED TARGETS AS RETURNED BY A CLASSIFIER

LABELS LIST OPTIONAL THE SET OF LABELS TO INCLUDE WHEN AVERAGE BINARY AND THEIR ORDER IFAVERAGE IS NONE LABELS PRESENT IN THE DATA CAN BE EXCLUDED FOR EXAMPLE TO CALCULATE A MULTICLASS AVERAGE IGNORING A MAJORITY NEGATIVE CLASS WHILE LABELS NOT PRESENT IN THE DATA WILL RESULT IN 0 COMPONENTS IN A MACRO AVERAGE FOR MULTILABEL TARGETS LABELS ARE COLUMN INDICES BY DEFAULT ALL LABELS IN YTRUE ANDYPRED ARE USED IN SORTED ORDER

CHANGED IN VERSION 017 PARAMETER LABELS IMPROVED FOR MULTICLASS PROBLEM

POSLABEL STR OR INT 1 BY DEFAULT THE CLASS TO REPORT IF AVERAGEBINARY AND THE DATA IS BINARY IF THE DATA ARE MULTICLASS OR MULTILABEL THIS WILL BE IGNORED SETTING LABELSPOSLABEL ANDAVERAGE BINARY WILL REPORT SCORES FOR THAT LABEL ONLY

AVERAGE STRING NONE 'BINARY' DEFAULT 'MICRO' 'MACRO' 'SAMPLES' 'WEIGHTED' THIS PARAMETER IS REQUIRED FOR MULTICLASSMULTILABEL TARGETS IF NONE THE SCORES FOR EACH CLASS ARE RETURNED OTHERWISE THIS DETERMINES THE TYPE OF AVERAGING PERFORMED ON THE DATA

BINARY ONLY REPORT RESULTS FOR THE CLASS SPECIFIED BY POSLABEL THIS IS APPLICABLE ONLY IF TARGETS YTRUEPRED ARE BINARY

MICRO CALCULATE METRICS GLOBALLY BY COUNTING THE TOTAL TRUE POSITIVES FALSE NEGATIVES AND FALSE POSITIVES

MACRO CALCULATE METRICS FOR EACH LABEL AND FIND THEIR UNWEIGHTED MEAN THIS DOES NOT TAKE LABEL IMBALANCE INTO ACCOUNT

WEIGHTED CALCULATE METRICS FOR EACH LABEL AND FIND THEIR AVERAGE WEIGHTED BY SUPPORT THE NUMBER OF TRUE INSTANCES FOR EACH LABEL THIS ALTERS 'MACRO' TO ACCOUNT FOR LABEL IMBALANCE IT CAN RESULT IN AN FSCORE THAT IS NOT BETWEEN PRECISION AND RECALL

SAMPLES CALCULATE METRICS FOR EACH INSTANCE AND FIND THEIR AVERAGE ONLY MEANINGFUL FOR MULTILABEL CLASSIFICATION WHERE THIS DIFFERS FROM ACCURACYScore

SAMPLEWEIGHT ARRAYLIKE OF SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

F1SCORE FLOAT OR ARRAY OF FLOAT SHAPE NUNIQUELABELS F1 SCORE OF THE POSITIVE CLASS IN BINARY CLASSIFICATION OR WEIGHTED AVERAGE OF THE F1 SCORES OF EACH CLASS FOR THE MULTICLASS TASK

SEE ALSO

FBETAScore PRECISIONRECALLFScoresSUPPORT JACCARDScore

MULTILABELCONFUSIONMATRIX

1996 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

NOTES  
WHENTRUE POSITIVE FALSE POSITIVE 0 ORTRUE POSITIVE FALSE NEGATIVE  
0 FSCORE RETURNS 0 AND RAISES UNDEFINEDMETRICWARNING

REFERENCES

1  
EXAMPLES  
FROM SKLEARNMETRICS IMPORT F1SCORE  
YTRUE 0 1 2 0 1 2  
YPRED 0 2 1 0 0 1  
F1SCOREYTRUE YPRED AVERAGEMACRO

026  
F1SCOREYTRUE YPRED AVERAGEMICRO  
033  
F1SCOREYTRUE YPRED AVERAGEWEIGHTED  
026  
F1SCOREYTRUE YPRED AVERAGE NONE  
ARRAY08 0 0

EXAMPLES USING SKLEARNMETRICSF1SCORE  
•PROBABILITY CALIBRATION CURVES  
SKLEARNMETRICS FBETAScore  
SKLEARNMETRICS FBETAScore YTRUE YPRED BETALABELSNONE POSLABEL1 AVERAGE'BINARY'  
SAMPLEWEIGHTNONE  
COMPUTE THE FBETA SCORE

THE FBETA SCORE IS THE WEIGHTED HARMONIC MEAN OF PRECISION AND RECALL REACHING ITS OPTIMAL VALUE AT 1 AND ITS  
WORST VALUE AT 0

THEBETA PARAMETER DETERMINES THE WEIGHT OF RECALL IN THE COMBINED SCORE BETA 1 LENDS MORE WEIGHT TO  
PRECISION WHILE BETA 1 FAVORS RECALL BETA 0 CONSIDERS ONLY PRECISION BETA INF ONLY RECALL  
READ MORE IN THE USER GUIDE

PARAMETERS  
YTRUE 1D ARRAYLIKE OR LABEL INDICATOR ARRAY SPARSE MATRIX GROUND TRUTH CORRECT TARGET  
VALUES  
YPRED 1D ARRAYLIKE OR LABEL INDICATOR ARRAY SPARSE MATRIX ESTIMATED TARGETS AS RETURNED BY  
A CLASSIFIER

BETA FLOAT WEIGHT OF PRECISION IN HARMONIC MEAN  
LABELS LIST OPTIONAL THE SET OF LABELS TO INCLUDE WHEN AVERAGE BINARY AND THEIR  
ORDER IFAVERAGE IS NONE LABELS PRESENT IN THE DATA CAN BE EXCLUDED FOR EXAMPLE TO  
CALCULATE A MULTICLASS AVERAGE IGNORING A MAJORITY NEGATIVE CLASS WHILE LABELS NOT PRESENT IN  
624SKLEARNMETRICS METRICS 1997

SCIKITLEARN USER GUIDE RELEASE 0213

THE DATA WILL RESULT IN 0 COMPONENTS IN A MACRO AVERAGE FOR MULTILABEL TARGETS LABELS ARE COLUMN INDICES BY DEFAULT ALL LABELS IN YTRUE ANDYPRED ARE USED IN SORTED ORDER

CHANGED IN VERSION 017 PARAMETER LABELS IMPROVED FOR MULTICLASS PROBLEM

POSLABEL STR OR INT 1 BY DEFAULT THE CLASS TO REPORT IF AVERAGEBINARY AND THE DATA IS BINARY IF THE DATA ARE MULTICLASS OR MULTILABEL THIS WILL BE IGNORED SETTING LABELSPOSLABEL ANDAVERAGE BINARY WILL REPORT SCORES FOR THAT LABEL ONLY

AVERAGE STRING NONE 'BINARY' DEFAULT 'MICRO' 'MACRO' 'SAMPLES' 'WEIGHTED' THIS PARAMETER IS REQUIRED FOR MULTICLASSMULTILABEL TARGETS IF NONE THE SCORES FOR EACH CLASS ARE RETURNED OTHERWISE THIS DETERMINES THE TYPE OF AVERAGING PERFORMED ON THE DATA

BINARY ONLY REPORT RESULTS FOR THE CLASS SPECIFIED BY POSLABEL THIS IS APPLICABLE ONLY IF TARGETS YTRUEPRED ARE BINARY

MICRO CALCULATE METRICS GLOBALLY BY COUNTING THE TOTAL TRUE POSITIVES FALSE NEGATIVES AND FALSE POSITIVES

MACRO CALCULATE METRICS FOR EACH LABEL AND FIND THEIR UNWEIGHTED MEAN THIS DOES NOT TAKE LABEL IMBALANCE INTO ACCOUNT

WEIGHTED CALCULATE METRICS FOR EACH LABEL AND FIND THEIR AVERAGE WEIGHTED BY SUPPORT THE NUMBER OF TRUE INSTANCES FOR EACH LABEL THIS ALTERS 'MACRO' TO ACCOUNT FOR LABEL IMBALANCE IT CAN RESULT IN AN FSCORE THAT IS NOT BETWEEN PRECISION AND RECALL

SAMPLES CALCULATE METRICS FOR EACH INSTANCE AND FIND THEIR AVERAGE ONLY MEANINGFUL FOR MULTILABEL CLASSIFICATION WHERE THIS DIFFERS FROM ACCURACYSAMPLEWEIGHT ARRAYLIKE OF SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

FBETAScore FLOAT IF AVERAGE IS NOT NONE OR ARRAY OF FLOAT SHAPE NUNIQUELABELS FBETAScore OF THE POSITIVE CLASS IN BINARY CLASSIFICATION OR WEIGHTED AVERAGE OF THE FBETAScore OF EACH CLASS FOR THE MULTICLASS TASK

SEE ALSO

PRECISIONRECALLFSCORESUPPORT MULTILABELCONFUSIONMATRIX

NOTES

WHENTRUE POSITIVE FALSE POSITIVE 0 ORTRUE POSITIVE FALSE NEGATIVE 0 FSCORE RETURNS 0 AND RAISES UNDEFINEDMETRICWARNING

REFERENCES

12

EXAMPLES

```
FROM SKLEARNMETRICS IMPORT FBETAScore
YTRUE 0 1 2 0 1 2
YPRED 0 2 1 0 0 1
FBETAScoreYTRUE YPRED AVERAGEMACRO BETA05
```

1998 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

023  
FBETASCOREYTRUE YPRED AVERAGEMICRO BETA05

033  
FBETASCOREYTRUE YPRED AVERAGEWEIGHTED BETA05

023  
FBETASCOREYTRUE YPRED AVERAGE NONE BETA05

ARRAY071 0 0  
SKLEARNMETRICS HAMMINGLOSS  
SKLEARNMETRICS HAMMINGLOSS YTRUE YPRED LABELSNONE SAMPLEWEIGHTNONE  
COMPUTE THE AVERAGE HAMMING LOSS  
THE HAMMING LOSS IS THE FRACTION OF LABELS THAT ARE INCORRECTLY PREDICTED  
READ MORE IN THE USER GUIDE  
PARAMETERS  
YTRUE 1D ARRAYLIKE OR LABEL INDICATOR ARRAY SPARSE MATRIX GROUND TRUTH CORRECT LABELS  
YPRED 1D ARRAYLIKE OR LABEL INDICATOR ARRAY SPARSE MATRIX PREDICTED LABELS AS RETURNED BY  
A CLASSIFIER  
LABELS ARRAY SHAPE NLABELS OPTIONAL DEFAULT'DEPRECATED' INTEGER ARRAY OF LABELS IF NOT  
PROVIDED LABELS WILL BE INFERRED FROM YTRUE AND YPRED  
NEW IN VERSION 018  
DEPRECATED SINCE VERSION 021 THIS PARAMETER LABELS IS DEPRECATED IN VERSION 021 AND  
WILL BE REMOVED IN VERSION 023 HAMMING LOSS USES YTRUESHAPE1 FOR THE NUMBER  
OF LABELS WHEN YTRUE IS BINARY LABEL INDICATORS SO IT IS UNNECESSARY FOR THE USER TO SPECIFY  
SAMPLEWEIGHT ARRAYLIKE OF SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS  
NEW IN VERSION 018  
RETURNS  
LOSS FLOAT OR INT RETURN THE AVERAGE HAMMING LOSS BETWEEN ELEMENT OF YTRUE ANDYPRED  
SEE ALSO  
ACCURACYScore JACCARDScore ZEROONELOSS  
NOTES  
IN MULTICLASS CLASSIFICATION THE HAMMING LOSS CORRESPONDS TO THE HAMMING DISTANCE BETWEEN YTRUE AND  
YPRED WHICH IS EQUIVALENT TO THE SUBSET ZEROONELOSS FUNCTION WHEN NORMALIZE PARAMETER IS SET TO  
TRUE  
IN MULTILABEL CLASSIFICATION THE HAMMING LOSS IS DIFFERENT FROM THE SUBSET ZEROONE LOSS THE ZEROONE LOSS  
CONSIDERS THE ENTIRE SET OF LABELS FOR A GIVEN SAMPLE INCORRECT IF IT DOES NOT ENTIRELY MATCH THE TRUE SET OF LABELS  
HAMMING LOSS IS MORE FORGIVING IN THAT IT PENALIZES ONLY THE INDIVIDUAL LABELS  
624SKLEARNMETRICS METRICS 1999

SCIKITLEARN USER GUIDE RELEASE 0213  
THE HAMMING LOSS IS UPPERBOUNDED BY THE SUBSET ZEROONE LOSS WHEN NORMALIZE PARAMETER IS SET TO TRUE IT  
IS ALWAYS BETWEEN 0 AND 1 LOWER BEING BETTER  
REFERENCES

12  
EXAMPLES  
FROM SKLEARNMETRICS IMPORT HAMMINGLOSS  
YPRED 1 2 3 4  
YTRUE 2 2 3 4  
HAMMINGLOSSYTRUE YPRED  
025

IN THE MULTILABEL CASE WITH BINARY LABEL INDICATORS  
IMPORT NUMPY AS NP  
HAMMINGLOSSNPARRAY0 1 1 1 NPZEROS2 2  
075

EXAMPLES USING SKLEARNMETRICSHAMMINGLOSS  
•MODEL COMPLEXITY INFLUENCE  
SKLEARNMETRICS HINGELOSS  
SKLEARNMETRICS HINGELOSS YTRUE PREDDECISION LABELSNONE SAMPLEWEIGHTNONE  
AVERAGE HINGE LOSS NONREGULARIZED  
IN BINARY CLASS CASE ASSUMING LABELS IN YTRUE ARE ENCODED WITH 1 AND -1 WHEN A PREDICTION MISTAKE IS  
MADEMARGIN YTRUE PREDDECISION IS ALWAYS NEGATIVE SINCE THE SIGNS DISAGREE IMPLYING  
1 MARGIN IS ALWAYS GREATER THAN 1 THE CUMULATED HINGE LOSS IS THEREFORE AN UPPER BOUND OF THE NUMBER OF  
MISTAKES MADE BY THE CLASSIFIER  
IN MULTICLASS CASE THE FUNCTION EXPECTS THAT EITHER ALL THE LABELS ARE INCLUDED IN YTRUE OR AN OPTIONAL LABELS  
ARGUMENT IS PROVIDED WHICH CONTAINS ALL THE LABELS THE MULTILABEL MARGIN IS CALCULATED ACCORDING TO CRAMMER  
SINGER’S METHOD AS IN THE BINARY CASE THE CUMULATED HINGE LOSS IS AN UPPER BOUND OF THE NUMBER OF MISTAKES  
MADE BY THE CLASSIFIER  
READ MORE IN THE USER GUIDE  
PARAMETERS  
YTRUE ARRAY SHAPE NSAMPLES TRUE TARGET CONSISTING OF INTEGERS OF TWO VALUES THE  
POSITIVE LABEL MUST BE GREATER THAN THE NEGATIVE LABEL  
PREDDECISION ARRAY SHAPE NSAMPLES OR NSAMPLES NCLASSES PREDICTED DECISIONS AS  
OUTPUT BY DECISIONFUNCTION FLOATS  
LABELS ARRAY OPTIONAL DEFAULT NONE CONTAINS ALL THE LABELS FOR THE PROBLEM USED IN MULTICLASS  
HINGE LOSS  
SAMPLEWEIGHT ARRAYLIKE OF SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS  
2000 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213  
RETURNS  
LOSS FLOAT  
REFERENCES  
123  
EXAMPLES  
FROM SKLEARN IMPORT SVM  
FROM SKLEARNMETRICS IMPORT HINGELOSS  
X 0 1  
Y 1 1  
EST SVMLINEARSVCRANDOMSTATE0  
ESTFITX Y  
LINEARSVCC10 CLASSWEIGHTNONE DUALTRUE FITINTERCEPTTRUE  
INTERCEPTSCALING1 LOSSSQUAREDHINGE MAXITER1000  
MULTICLASOVR PENALTYL2 RANDOMSTATE0 TOL00001  
VERBOSE0  
PREDDCISION ESTDECISIONFUNCTION2 3 05  
PREDDCISION  
ARRAY218 236 009  
HINGELOSS1 1 1 PREDDCISION  
030  
IN THE MULTICLASS CASE  
IMPORT NUMPY AS NP  
X NPARRAY0 1 2 3  
Y NPARRAY0 1 2 3  
LABELS NPARRAY0 1 2 3  
EST SVMLINEAR SVC  
ESTFITX Y  
LINEARSVCC10 CLASSWEIGHTNONE DUALTRUE FITINTERCEPTTRUE  
INTERCEPTSCALING1 LOSSSQUAREDHINGE MAXITER1000  
MULTICLASOVR PENALTYL2 RANDOMSTATENONE TOL00001  
VERBOSE0  
PREDDCISION ESTDECISIONFUNCTION1 2 3  
YTRUE 0 2 3  
HINGELOSSYTRUE PREDDCISION LABELS  
056  
SKLEARNMETRICS JACCARDSCORE  
SKLEARNMETRICS JACCARDSCORE YTRUE YPRED LABELSNONE POSLABEL1 AVERAGE' BINARY'  
SAMPLEWEIGHTNONE  
JACCARD SIMILARITY COEFFICIENT SCORE  
THE JACCARD INDEX 1 OR JACCARD SIMILARITY COEFFICIENT DEFINED AS THE SIZE OF THE INTERSECTION DIVIDED BY THE SIZE  
OF THE UNION OF TWO LABEL SETS IS USED TO COMPARE SET OF PREDICTED LABELS FOR A SAMPLE TO THE CORRESPONDING SET OF  
LABELS INYTRUE  
READ MORE IN THE USER GUIDE  
624SKLEARNMETRICS METRICS 2001

SCIKITLEARN USER GUIDE RELEASE 0213

PARAMETERS

YTRUE 1D ARRAYLIKE OR LABEL INDICATOR ARRAY SPARSE MATRIX GROUND TRUTH CORRECT LABELS  
YPRED 1D ARRAYLIKE OR LABEL INDICATOR ARRAY SPARSE MATRIX PREDICTED LABELS AS RETURNED BY  
A CLASSIFIER  
LABELS LIST OPTIONAL THE SET OF LABELS TO INCLUDE WHEN AVERAGE BINARY AND THEIR  
ORDER IFAVERAGE IS NONE LABELS PRESENT IN THE DATA CAN BE EXCLUDED FOR EXAMPLE TO  
CALCULATE A MULTICLASS AVERAGE IGNORING A MAJORITY NEGATIVE CLASS WHILE LABELS NOT PRESENT IN  
THE DATA WILL RESULT IN 0 COMPONENTS IN A MACRO AVERAGE FOR MULTILABEL TARGETS LABELS ARE  
COLUMN INDICES BY DEFAULT ALL LABELS IN YTRUE ANDYPRED ARE USED IN SORTED ORDER  
POSLABEL STR OR INT 1 BY DEFAULT THE CLASS TO REPORT IF AVERAGEBINARY AND THE  
DATA IS BINARY IF THE DATA ARE MULTICLASS OR MULTILABEL THIS WILL BE IGNORED SETTING  
LABELSPOSLABEL ANDAVERAGE BINARY WILL REPORT SCORES FOR THAT LABEL ONLY  
AVERAGE STRING NONE 'BINARY' DEFAULT 'MICRO' 'MACRO' 'SAMPLES' 'WEIGHTED' IF NONE  
THE SCORES FOR EACH CLASS ARE RETURNED OTHERWISE THIS DETERMINES THE TYPE OF AVERAGING  
PERFORMED ON THE DATA  
BINARY ONLY REPORT RESULTS FOR THE CLASS SPECIFIED BY POSLABEL THIS IS APPLICABLE  
ONLY IF TARGETS YTRUEPRED ARE BINARY  
MICRO CALCULATE METRICS GLOBALLY BY COUNTING THE TOTAL TRUE POSITIVES FALSE NEGATIVES  
AND FALSE POSITIVES  
MACRO CALCULATE METRICS FOR EACH LABEL AND FIND THEIR UNWEIGHTED MEAN THIS DOES NOT  
TAKE LABEL IMBALANCE INTO ACCOUNT  
WEIGHTED CALCULATE METRICS FOR EACH LABEL AND FIND THEIR AVERAGE WEIGHTED BY SUP  
PORT THE NUMBER OF TRUE INSTANCES FOR EACH LABEL THIS ALTERS 'MACRO' TO ACCOUNT FOR LABEL  
IMBALANCE  
SAMPLES CALCULATE METRICS FOR EACH INSTANCE AND FIND THEIR AVERAGE ONLY MEANINGFUL  
FOR MULTILABEL CLASSIFICATION  
SAMPLEWEIGHT ARRAYLIKE OF SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS  
RETURNS  
SCORE FLOAT IF AVERAGE IS NOT NONE OR ARRAY OF FLOATS SHAPE NUNIQUELABELS  
SEE ALSO  
ACCURACYScore FScore MULTILABELCONFUSIONMATRIX  
NOTES  
JACCARDScore MAY BE A POOR METRIC IF THERE ARE NO POSITIVES FOR SOME SAMPLES OR CLASSES JACCARD IS  
UNDEFINED IF THERE ARE NO TRUE OR PREDICTED LABELS AND OUR IMPLEMENTATION WILL RETURN A SCORE OF 0 WITH A WARNING  
REFERENCES



SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES

```
import numpy as np
from sklearn.metrics import jaccard_score
y_true = np.array([1, 1, 0])
y_pred = np.array([1, 1, 0])
```

IN THE BINARY CASE

```
jaccard_score(y_true, y_pred)
0.6666
```

IN THE MULTILABEL CASE

```
jaccard_score(y_true, y_pred, average='samples')
0.5833
jaccard_score(y_true, y_pred, average='macro')
0.6666
```

```
jaccard_score(y_true, y_pred, average='none')
```

```
array([0.5, 1.])
```

IN THE MULTICLASS CASE

```
y_pred = [0, 2, 1, 2]
y_true = [0, 1, 2, 2]
jaccard_score(y_true, y_pred, average='none')
```

```
array([1., 0., 0.33])
```

EXAMPLES USING SKLEARNMETRICSJACCARDSCORE

•CLASSIFIER CHAIN

```
sklearn.metrics.log_loss
sklearn.metrics.log_loss(y_true, y_pred, eps=1e-15, normalize=True, sample_weight=None, labels=None)
```

LOG LOSS AKA LOGISTIC LOSS OR CROSSENTROPY LOSS

THIS IS THE LOSS FUNCTION USED IN MULTINOMIAL LOGISTIC REGRESSION AND EXTENSIONS OF IT SUCH AS NEURAL NETWORKS. DEFINED AS THE NEGATIVE LOG-LIKELIHOOD OF THE TRUE LABELS GIVEN A PROBABILISTIC CLASSIFIER'S PREDICTIONS, THE LOG LOSS IS ONLY DEFINED FOR TWO OR MORE LABELS. FOR A SINGLE SAMPLE WITH TRUE LABEL  $y_t$  IN  $\{0, 1\}$  AND ESTIMATED PROBABILITY

$y_p$  THAT  $y_t = 1$  THE LOG LOSS IS

$$-\log(p_{y_t}) = -\log(y_t \log y_t + (1 - y_t) \log(1 - y_t))$$

READ MORE IN THE USER GUIDE

PARAMETERS

**y\_true** ARRAY-LIKE OR LABEL INDICATOR MATRIX: GROUND TRUTH CORRECT LABELS FOR **n** SAMPLES  
**samples** PLES  
6245SKLEARNMETRICS METRICS 2003

SCIKITLEARN USER GUIDE RELEASE 0213

YPRED ARRAYLIKE OF FLOAT SHAPE NSAMPLES NCLASSES OR NSAMPLES PREDICTED PROBABILITIES AS RETURNED BY A CLASSIFIER'S PREDICTPROBA METHOD IF YPREDSHAPE NSAMPLES THE PROBABILITIES PROVIDED ARE ASSUMED TO BE THAT OF THE POSITIVE CLASS THE LABELS INYPRED ARE ASSUMED TO BE ORDERED ALPHABETICALLY AS DONE BY PREPROCESSING LABELBINARIZER

EPS FLOAT LOG LOSS IS UNDEFINED FOR P0 OR P1 SO PROBABILITIES ARE CLIPPED TO MAXEPS MIN1 EPS P

NORMALIZE BOOL OPTIONAL DEFAULTTRUE IF TRUE RETURN THE MEAN LOSS PER SAMPLE OTHERWISE RETURN THE SUM OF THE PERSAMPLE LOSSES

SAMPLEWEIGHT ARRAYLIKE OF SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

LABELS ARRAYLIKE OPTIONAL DEFAULTNONE IF NOT PROVIDED LABELS WILL BE INFERRED FROM YTRUE IFLABELS ISNONE ANDYPRED HAS SHAPE NSAMPLES THE LABELS ARE ASSUMED TO BE BINARY AND ARE INFERRED FROM YTRUE VERSIONADDED 018

RETURNS

LOSS FLOAT

NOTES

THE LOGARITHM USED IS THE NATURAL LOGARITHM BASEE

REFERENCES

CM BISHOP 2006 PATTERN RECOGNITION AND MACHINE LEARNING SPRINGER P 209

EXAMPLES

FROM SKLEARNMETRICS IMPORT LOGLOSS

LOGLOSSSPAM HAM HAM SPAM

1 9 9 1 8 2 35 65

021616

EXAMPLES USING SKLEARNMETRICSLGLOSS

- PROBABILITY CALIBRATION FOR 3CLASS CLASSIFICATION
- PROBABILISTIC PREDICTIONS WITH GAUSSIAN PROCESS CLASSIFICATION GPC

SKLEARNMETRICS MATTHEWSCORRCOE

SKLEARNMETRICS MATTHEWSCORRCOE YTRUE YPRED SAMPLEWEIGHTNONE

COMPUTE THE MATTHEWS CORRELATION COEFFICIENT MCC

THE MATTHEWS CORRELATION COEFFICIENT IS USED IN MACHINE LEARNING AS A MEASURE OF THE QUALITY OF BINARY AND MULTICLASS CLASSIFICATIONS IT TAKES INTO ACCOUNT TRUE AND FALSE POSITIVES AND NEGATIVES AND IS GENERALLY REGARDED AS A BALANCED MEASURE WHICH CAN BE USED EVEN IF THE CLASSES ARE OF VERY DIFFERENT SIZES THE MCC IS IN ESSENCE A CORRELATION COEFFICIENT VALUE BETWEEN 1 AND 1 A COEFFICIENT OF 1 REPRESENTS A PERFECT PREDICTION 0 AN AVERAGE

2004 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

RANDOM PREDICTION AND 1 AN INVERSE PREDICTION THE STATISTIC IS ALSO KNOWN AS THE PHI COEFFICIENT SOURCE WIKIPEDIA

BINARY AND MULTICLASS LABELS ARE SUPPORTED ONLY IN THE BINARY CASE DOES THIS RELATE TO INFORMATION ABOUT TRUE AND FALSE POSITIVES AND NEGATIVES SEE REFERENCES BELOW

READ MORE IN THE USER GUIDE

PARAMETERS

YTRUE ARRAY SHAPE NSAMPLES GROUND TRUTH CORRECT TARGET VALUES

YPRED ARRAY SHAPE NSAMPLES ESTIMATED TARGETS AS RETURNED BY A CLASSIFIER

SAMPLEWEIGHT ARRAYLIKE OF SHAPE NSAMPLES DEFAULT NONE SAMPLE WEIGHTS

RETURNS

MCC FLOAT THE MATTHEWS CORRELATION COEFFICIENT 1 REPRESENTS A PERFECT PREDICTION 0 AN AVERAGE RANDOM PREDICTION AND 1 AND INVERSE PREDICTION

REFERENCES

1234

EXAMPLES

FROM SKLEARNMETRICS IMPORT MATTHEWSCORRCOEFF

YTRUE 1 1 1 1

YPRED 1 1 1 1

MATTHEWSCORRCOEFFYTRUE YPRED

033

SKLEARNMETRICS MULTILABELCONFUSIONMATRIX

SKLEARNMETRICS MULTILABELCONFUSIONMATRIX YTRUE YPRED SAMPLEWEIGHTNONE LABELS NONE SAMPLEWISEFALSE

COMPUTE A CONFUSION MATRIX FOR EACH CLASS OR SAMPLE

NEW IN VERSION 021

COMPUTE CLASSWISE DEFAULT OR SAMPLEWISE SAMPLEWISETRUE MULTILABEL CONFUSION MATRIX TO EVALUATE THE ACCURACY OF A CLASSIFICATION AND OUTPUT CONFUSION MATRICES FOR EACH CLASS OR SAMPLE

IN MULTILABEL CONFUSION MATRIX  $\begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}$  THE COUNT OF TRUE NEGATIVES IS  $\begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}$  00 FALSE NEGATIVES IS  $\begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}$  10 TRUE POSITIVES IS  $\begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}$  11AND FALSE POSITIVES IS  $\begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}$  01

MULTICLASS DATA WILL BE TREATED AS IF BINARIZED UNDER A ONEVSREST TRANSFORMATION RETURNED CONFUSION MATRICES WILL BE IN THE ORDER OF SORTED UNIQUE LABELS IN THE UNION OF YTRUE YPRED

READ MORE IN THE USER GUIDE

PARAMETERS

YTRUE 1D ARRAYLIKE OR LABEL INDICATOR ARRAY SPARSE MATRIX OF SHAPE NSAMPLES NOUTPUTS OR NSAMPLES GROUND TRUTH CORRECT TARGET VALUES

YPRED 1D ARRAYLIKE OR LABEL INDICATOR ARRAY SPARSE MATRIX OF SHAPE NSAMPLES NOUTPUTS OR NSAMPLES ESTIMATED TARGETS AS RETURNED BY A CLASSIFIER

624SKLEARNMETRICS METRICS 2005

SCIKITLEARN USER GUIDE RELEASE 0213

SAMPLEWEIGHT ARRAYLIKE OF SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

LABELS ARRAYLIKE A LIST OF CLASSES OR COLUMN INDICES TO SELECT SOME OR TO FORCE INCLUSION OF CLASSES ABSENT FROM THE DATA

SAMPLEWISE BOOL DEFAULTFALSE IN THE MULTILABEL CASE THIS CALCULATES A CONFUSION MATRIX PER SAMPLE

RETURNS

MULTICONFUSION ARRAY SHAPE NOUTPUTS 2 2 A 2X2 CONFUSION MATRIX CORRESPONDING TO EACH OUTPUT IN THE INPUT WHEN CALCULATING CLASSWISE MULTICONFUSION DEFAULT THEN NOUTPUTS NLABELS WHEN CALCULATING SAMPLEWISE MULTICONFUSION SAMPLEWISETRUE NOUTPUTS NSAMPLES IF LABELS IS DEFINED THE RESULTS WILL BE RETURNED IN THE ORDER SPECIFIED IN LABELS OTHERWISE THE RESULTS WILL BE RETURNED IN SORTED ORDER BY DEFAULT

SEE ALSO

CONFUSIONMATRIX

NOTES

THE MULTILABELCONFUSIONMATRIX CALCULATES CLASSWISE OR SAMPLEWISE MULTILABEL CONFUSION MATRICES AND IN MULTICLASS TASKS LABELS ARE BINARIZED UNDER A ONEVSREST WAY WHILE CONFUSIONMATRIX CALCULATES ONE CONFUSION MATRIX FOR CONFUSION BETWEEN EVERY TWO CLASSES

EXAMPLES

MULTILABELINDICATOR CASE

```
import numpy as np
from sklearn.metrics import multilabel_confusion_matrix

y_true = np.array([0, 1, 0, 1])
y_pred = np.array([0, 0, 1, 1])

multilabel_confusion_matrix(y_true, y_pred)
```

MULTICLASS CASE

```
y_true = ['cat', 'ant', 'cat', 'cat', 'ant', 'bird']
y_pred = ['ant', 'ant', 'cat', 'cat', 'ant', 'cat']

multilabel_confusion_matrix(y_true, y_pred, labels=['ant', 'bird', 'cat'])
```

2006 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

2 1  
1 2

SKLEARNMETRICS PRECISIONRECALLCURVE  
SKLEARNMETRICS PRECISIONRECALLCURVE YTRUE PROBASPRED POSLABELNONE SAM  
PLEWEIGHTNONE

COMPUTE PRECISIONRECALL PAIRS FOR DIFFERENT PROBABILITY THRESHOLDS  
NOTE THIS IMPLEMENTATION IS RESTRICTED TO THE BINARY CLASSIFICATION TASK  
THE PRECISION IS THE RATIO  $TP / (TP + FP)$  WHERE  $TP$  IS THE NUMBER OF TRUE POSITIVES AND  $FP$  THE NUMBER OF FALSE POSITIVES THE PRECISION IS INTUITIVELY THE ABILITY OF THE CLASSIFIER NOT TO LABEL AS POSITIVE A SAMPLE THAT IS NEGATIVE  
THE RECALL IS THE RATIO  $TP / (TP + FN)$  WHERE  $TP$  IS THE NUMBER OF TRUE POSITIVES AND  $FN$  THE NUMBER OF FALSE NEGATIVES THE RECALL IS INTUITIVELY THE ABILITY OF THE CLASSIFIER TO FIND ALL THE POSITIVE SAMPLES  
THE LAST PRECISION AND RECALL VALUES ARE 1 AND 0 RESPECTIVELY AND DO NOT HAVE A CORRESPONDING THRESHOLD THIS ENSURES THAT THE GRAPH STARTS ON THE Y AXIS  
READ MORE IN THE USER GUIDE

PARAMETERS  
YTRUE ARRAY SHAPE NSAMPLES TRUE BINARY LABELS IF LABELS ARE NOT EITHER 1 1 OR 0 1 THEN POSLABEL SHOULD BE EXPLICITLY GIVEN  
PROBASPRED ARRAY SHAPE NSAMPLES ESTIMATED PROBABILITIES OR DECISION FUNCTION  
POSLABEL INT OR STR DEFAULTNONE THE LABEL OF THE POSITIVE CLASS WHEN POSLABELNONE  
IF YTRUE IS IN 1 1 OR 0 1 POSLABEL IS SET TO 1 OTHERWISE AN ERROR WILL BE RAISED  
SAMPLEWEIGHT ARRAYLIKE OF SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS  
RETURNS  
PRECISION ARRAY SHAPE NTHRESHOLDS 1 PRECISION VALUES SUCH THAT ELEMENT I IS THE PRECISION OF PREDICTIONS WITH SCORE THRESHOLDSI AND THE LAST ELEMENT IS 1  
RECALL ARRAY SHAPE NTHRESHOLDS 1 DECREASING RECALL VALUES SUCH THAT ELEMENT I IS THE RECALL OF PREDICTIONS WITH SCORE THRESHOLDSI AND THE LAST ELEMENT IS 0  
THRESHOLDS ARRAY SHAPE NTHRESHOLDS LENNPUNIQUEPROBASPRED INCREASING THRESHOLDS ON THE DECISION FUNCTION USED TO COMPUTE PRECISION AND RECALL  
SEE ALSO  
AVERAGEPRECISIONSCORE COMPUTE AVERAGE PRECISION FROM PREDICTION SCORES  
ROCCURVE COMPUTE RECEIVER OPERATING CHARACTERISTIC ROC CURVE

EXAMPLES  
IMPORT NUMPY AS NP  
FROM SKLEARNMETRICS IMPORT PRECISIONRECALLCURVE  
YTRUE NPARRAY0 1 1  
YSCORES NPARRAY01 04 035 08  
624SKLEARNMETRICS METRICS 2007

SCIKITLEARN USER GUIDE RELEASE 0213

PRECISION RECALL THRESHOLDS PRECISIONRECALLCURVE

YTRUE YSCORES

PRECISION

ARRAY066666667 05 1 1

RECALL

ARRAY1 05 05 0

THRESHOLDS

ARRAY035 04 08

EXAMPLES USING SKLEARNMETRICSPRECISIONRECALLCURVE

- PRECISIONRECALL

SKLEARNMETRICS PRECISIONRECALLFSCORESUPPORT

SKLEARNMETRICS PRECISIONRECALLFSCORESUPPORT YTRUE YPRED BETA10 LA

BELSNONE POSLABEL1 AVER

AGENONE WARNFOR'PRECISION'

'RECALL' 'FSCORE' SAM

PLEWEIGHTNONE

COMPUTE PRECISION RECALL FMEASURE AND SUPPORT FOR EACH CLASS

THE PRECISION IS THE RATIO TP TP FP WHERE TP IS THE NUMBER OF TRUE POSITIVES AND FP THE NUMBER OF FALSE POSITIVES THE PRECISION IS INTUITIVELY THE ABILITY OF THE CLASSIFIER NOT TO LABEL AS POSITIVE A SAMPLE THAT IS NEGATIVE

THE RECALL IS THE RATIO TP TP FN WHERE TP IS THE NUMBER OF TRUE POSITIVES AND FN THE NUMBER OF FALSE NEGATIVES THE RECALL IS INTUITIVELY THE ABILITY OF THE CLASSIFIER TO FIND ALL THE POSITIVE SAMPLES

THE FBETA SCORE CAN BE INTERPRETED AS A WEIGHTED HARMONIC MEAN OF THE PRECISION AND RECALL WHERE AN FBETA SCORE REACHES ITS BEST VALUE AT 1 AND WORST SCORE AT 0

THE FBETA SCORE WEIGHTS RECALL MORE THAN PRECISION BY A FACTOR OF BETA BETA 10 MEANS RECALL AND PRECISION ARE EQUALLY IMPORTANT

THE SUPPORT IS THE NUMBER OF OCCURRENCES OF EACH CLASS IN YTRUE

IF POSLABEL IS NONE AND IN BINARY CLASSIFICATION THIS FUNCTION RETURNS THE AVERAGE PRECISION RECALL AND FMEASURE IF AVERAGE IS ONE OF MICRO MACRO WEIGHTED OR SAMPLES

READ MORE IN THE USER GUIDE

PARAMETERS

YTRUE 1D ARRAYLIKE OR LABEL INDICATOR ARRAY SPARSE MATRIX GROUND TRUTH CORRECT TARGET VALUES

YPRED 1D ARRAYLIKE OR LABEL INDICATOR ARRAY SPARSE MATRIX ESTIMATED TARGETS AS RETURNED BY A CLASSIFIER

BETA FLOAT 10 BY DEFAULT THE STRENGTH OF RECALL VERSUS PRECISION IN THE FSCORE

LABELS LIST OPTIONAL THE SET OF LABELS TO INCLUDE WHEN AVERAGE BINARY AND THEIR ORDER IF AVERAGE IS NONE LABELS PRESENT IN THE DATA CAN BE EXCLUDED FOR EXAMPLE TO CALCULATE A MULTICLASS AVERAGE IGNORING A MAJORITY NEGATIVE CLASS WHILE LABELS NOT PRESENT IN THE DATA WILL RESULT IN 0 COMPONENTS IN A MACRO AVERAGE FOR MULTILABEL TARGETS LABELS ARE COLUMN INDICES BY DEFAULT ALL LABELS IN YTRUE AND YPRED ARE USED IN SORTED ORDER

2008 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

POSLABEL STR OR INT 1 BY DEFAULT THE CLASS TO REPORT IF AVERAGEBINARY AND THE DATA IS BINARY IF THE DATA ARE MULTICLASS OR MULTILABEL THIS WILL BE IGNORED SETTING LABELSPOSLABEL ANDAVERAGE BINARY WILL REPORT SCORES FOR THAT LABEL ONLY

AVERAGE STRING NONE DEFAULT 'BINARY' 'MICRO' 'MACRO' 'SAMPLES' 'WEIGHTED' IF NONE THE SCORES FOR EACH CLASS ARE RETURNED OTHERWISE THIS DETERMINES THE TYPE OF AVERAGING PERFORMED ON THE DATA

BINARY ONLY REPORT RESULTS FOR THE CLASS SPECIFIED BY POSLABEL THIS IS APPLICABLE ONLY IF TARGETS YTRUEPRED ARE BINARY

MICRO CALCULATE METRICS GLOBALLY BY COUNTING THE TOTAL TRUE POSITIVES FALSE NEGATIVES AND FALSE POSITIVES

MACRO CALCULATE METRICS FOR EACH LABEL AND FIND THEIR UNWEIGHTED MEAN THIS DOES NOT TAKE LABEL IMBALANCE INTO ACCOUNT

WEIGHTED CALCULATE METRICS FOR EACH LABEL AND FIND THEIR AVERAGE WEIGHTED BY SUPPORT THE NUMBER OF TRUE INSTANCES FOR EACH LABEL THIS ALTERS 'MACRO' TO ACCOUNT FOR LABEL IMBALANCE IT CAN RESULT IN AN FSCORE THAT IS NOT BETWEEN PRECISION AND RECALL

SAMPLES CALCULATE METRICS FOR EACH INSTANCE AND FIND THEIR AVERAGE ONLY MEANINGFUL FOR MULTILABEL CLASSIFICATION WHERE THIS DIFFERS FROM ACCURACYScore

WARNFOR TUPLE OR SET FOR INTERNAL USE THIS DETERMINES WHICH WARNINGS WILL BE MADE IN THE CASE THAT THIS FUNCTION IS BEING USED TO RETURN ONLY ONE OF ITS METRICS

SAMPLEWEIGHT ARRAYLIKE OF SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

PRECISION FLOAT IF AVERAGE IS NOT NONE OR ARRAY OF FLOAT SHAPE NUNIQUELABELS

RECALL FLOAT IF AVERAGE IS NOT NONE OR ARRAY OF FLOAT SHAPE NUNIQUELABELS

FBETAScore FLOAT IF AVERAGE IS NOT NONE OR ARRAY OF FLOAT SHAPE NUNIQUELABELS

SUPPORT INT IF AVERAGE IS NOT NONE OR ARRAY OF INT SHAPE NUNIQUELABELS THE NUMBER OF OCCURRENCES OF EACH LABEL IN YTRUE

NOTES

WHENTRUE POSITIVE FALSE POSITIVE 0 PRECISION IS UNDEFINED WHEN TRUE POSITIVE FALSE NEGATIVE 0 RECALL IS UNDEFINED IN SUCH CASES THE METRIC WILL BE SET TO 0 AS WILL FSCORE ANDUNDEFINEDMETRICWARNING WILL BE RAISED

REFERENCES

123

EXAMPLES

```
import numpy as np
from sklearn.metrics import precision_recall_fscore_support

ytrue = np.array([cat, dog, pig, cat, dog, pig])
ypred = np.array([cat, pig, dog, cat, cat, dog])

precision, recall, fscore, support = precision_recall_fscore_support(ytrue, ypred, labels=[cat, dog, pig])
```

624SKLEARNMETRICS METRICS 2009

SCIKITLEARN USER GUIDE RELEASE 0213  
PRECISIONRECALLFSCORESUPPORTYTRUE YPRED AVERAGEMACRO

022 033 026 NONE  
PRECISIONRECALLFSCORESUPPORTYTRUE YPRED AVERAGEMICRO

033 033 033 NONE  
PRECISIONRECALLFSCORESUPPORTYTRUE YPRED AVERAGEWEIGHTED

022 033 026 NONE  
IT IS POSSIBLE TO COMPUTE PERLABEL PRECISIONS RECALLS F1SCORES AND SUPPORTS INSTEAD OF AVERAGING  
PRECISIONRECALLFSCORESUPPORTYTRUE YPRED AVERAGE NONE  
LABELSPIG DOG CAT

ARRAY0 0 066  
ARRAY0 0 1 ARRAY0 0 08  
ARRAY2 2 2  
SKLEARNMETRICS PRECISIONSCORE  
SKLEARNMETRICS PRECISIONSCORE YTRUE YPRED LABELSNONE POSLABEL1 AVERAGE'BINARY'  
SAMPLEWEIGHTNONE  
COMPUTE THE PRECISION  
THE PRECISION IS THE RATIO TP TP FP WHERE TP IS THE NUMBER OF TRUE POSITIVES AND FP THE NUMBER OF  
FALSE POSITIVES THE PRECISION IS INTUITIVELY THE ABILITY OF THE CLASSIFIER NOT TO LABEL AS POSITIVE A SAMPLE THAT IS  
NEGATIVE  
THE BEST VALUE IS 1 AND THE WORST VALUE IS 0  
READ MORE IN THE USER GUIDE  
PARAMETERS  
YTRUE 1D ARRAYLIKE OR LABEL INDICATOR ARRAY SPARSE MATRIX GROUND TRUTH CORRECT TARGET  
VALUES  
YPRED 1D ARRAYLIKE OR LABEL INDICATOR ARRAY SPARSE MATRIX ESTIMATED TARGETS AS RETURNED BY  
A CLASSIFIER  
LABELS LIST OPTIONAL THE SET OF LABELS TO INCLUDE WHEN AVERAGE BINARY AND THEIR  
ORDER IF AVERAGE IS NONE LABELS PRESENT IN THE DATA CAN BE EXCLUDED FOR EXAMPLE TO  
CALCULATE A MULTICLASS AVERAGE IGNORING A MAJORITY NEGATIVE CLASS WHILE LABELS NOT PRESENT IN  
THE DATA WILL RESULT IN 0 COMPONENTS IN A MACRO AVERAGE FOR MULTILABEL TARGETS LABELS ARE  
COLUMN INDICES BY DEFAULT ALL LABELS IN YTRUE AND YPRED ARE USED IN SORTED ORDER  
CHANGED IN VERSION 017 PARAMETER LABELS IMPROVED FOR MULTICLASS PROBLEM  
POSLABEL STR OR INT 1 BY DEFAULT THE CLASS TO REPORT IF AVERAGE BINARY AND THE  
DATA IS BINARY IF THE DATA ARE MULTICLASS OR MULTILABEL THIS WILL BE IGNORED SETTING  
LABELS POS LABEL AND AVERAGE BINARY WILL REPORT SCORES FOR THAT LABEL ONLY  
AVERAGE STRING NONE 'BINARY' DEFAULT 'MICRO' 'MACRO' 'SAMPLES' 'WEIGHTED' THIS PARAMETER  
IS REQUIRED FOR MULTICLASS MULTILABEL TARGETS IF NONE THE SCORES FOR EACH CLASS ARE  
RETURNED OTHERWISE THIS DETERMINES THE TYPE OF AVERAGING PERFORMED ON THE DATA  
2010 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

BINARY ONLY REPORT RESULTS FOR THE CLASS SPECIFIED BY POSLABEL THIS IS APPLICABLE ONLY IF TARGETS YTRUEYPRED ARE BINARY

MICRO CALCULATE METRICS GLOBALLY BY COUNTING THE TOTAL TRUE POSITIVES FALSE NEGATIVES AND FALSE POSITIVES

MACRO CALCULATE METRICS FOR EACH LABEL AND FIND THEIR UNWEIGHTED MEAN THIS DOES NOT TAKE LABEL IMBALANCE INTO ACCOUNT

WEIGHTED CALCULATE METRICS FOR EACH LABEL AND FIND THEIR AVERAGE WEIGHTED BY SUPPORT THE NUMBER OF TRUE INSTANCES FOR EACH LABEL THIS ALTERS 'MACRO' TO ACCOUNT FOR LABEL IMBALANCE IT CAN RESULT IN AN FSCORE THAT IS NOT BETWEEN PRECISION AND RECALL

SAMPLES CALCULATE METRICS FOR EACH INSTANCE AND FIND THEIR AVERAGE ONLY MEANINGFUL FOR MULTILABEL CLASSIFICATION WHERE THIS DIFFERS FROM ACCURACYSAMPLEWEIGHT ARRAYLIKE OF SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

PRECISION FLOAT IF AVERAGE IS NOT NONE OR ARRAY OF FLOAT SHAPE NUNIQUELABELS PRECISION OF THE POSITIVE CLASS IN BINARY CLASSIFICATION OR WEIGHTED AVERAGE OF THE PRECISION OF EACH CLASS FOR THE MULTICLASS TASK

SEE ALSO

PRECISIONRECALLFSCORESUPPORT MULTILABELCONFUSIONMATRIX

NOTES

WHENTRUE POSITIVE FALSE POSITIVE 0 PRECISION RETURNS 0 AND RAISES UNDEFINEDMETRICWARNING

EXAMPLES

```
FROM SKLEARNMETRICS IMPORT PRECISIONSCORE
YTRUE 0 1 2 0 1 2
YPRED 0 2 1 0 0 1
PRECISIONSCOREYTRUE YPRED AVERAGEMACRO
022
PRECISIONSCOREYTRUE YPRED AVERAGEMICRO
033
PRECISIONSCOREYTRUE YPRED AVERAGEWEIGHTED
022
PRECISIONSCOREYTRUE YPRED AVERAGE NONE
ARRAY066 0 0
EXAMPLES USING SKLEARNMETRICSPRECISIONSCORE
•PROBABILITY CALIBRATION CURVES
624SKLEARNMETRICS METRICS 2011
```

SCIKITLEARN USER GUIDE RELEASE 0213

SKLEARNMETRICS RECALLSCORE

SKLEARNMETRICS RECALLSCORE YTRUE YPRED LABELSNONE POSLABEL1 AVERAGE' BINARY' SAM  
PLEWEIGHTNONE

COMPUTE THE RECALL

THE RECALL IS THE RATIO  $TP / (TP + FN)$  WHERE  $TP$  IS THE NUMBER OF TRUE POSITIVES AND  $FN$  THE NUMBER OF FALSE  
NEGATIVES THE RECALL IS INTUITIVELY THE ABILITY OF THE CLASSIFIER TO FIND ALL THE POSITIVE SAMPLES  
THE BEST VALUE IS 1 AND THE WORST VALUE IS 0

READ MORE IN THE USER GUIDE

PARAMETERS

YTRUE 1D ARRAYLIKE OR LABEL INDICATOR ARRAY SPARSE MATRIX GROUND TRUTH CORRECT TARGET  
VALUES

YPRED 1D ARRAYLIKE OR LABEL INDICATOR ARRAY SPARSE MATRIX ESTIMATED TARGETS AS RETURNED BY  
A CLASSIFIER

LABELS LIST OPTIONAL THE SET OF LABELS TO INCLUDE WHEN AVERAGE BINARY AND THEIR  
ORDER IF AVERAGE IS NONE LABELS PRESENT IN THE DATA CAN BE EXCLUDED FOR EXAMPLE TO  
CALCULATE A MULTICLASS AVERAGE IGNORING A MAJORITY NEGATIVE CLASS WHILE LABELS NOT PRESENT IN  
THE DATA WILL RESULT IN 0 COMPONENTS IN A MACRO AVERAGE FOR MULTILABEL TARGETS LABELS ARE  
COLUMN INDICES BY DEFAULT ALL LABELS IN YTRUE AND YPRED ARE USED IN SORTED ORDER  
CHANGED IN VERSION 017 PARAMETER LABELS IMPROVED FOR MULTICLASS PROBLEM

POSLABEL STR OR INT 1 BY DEFAULT THE CLASS TO REPORT IF AVERAGE BINARY AND THE  
DATA IS BINARY IF THE DATA ARE MULTICLASS OR MULTILABEL THIS WILL BE IGNORED SETTING  
LABELS POSLABEL AND AVERAGE BINARY WILL REPORT SCORES FOR THAT LABEL ONLY

AVERAGE STRING NONE 'BINARY' DEFAULT 'MICRO' 'MACRO' 'SAMPLES' 'WEIGHTED' THIS PA  
RAMETER IS REQUIRED FOR MULTICLASS MULTILABEL TARGETS IF NONE THE SCORES FOR EACH CLASS ARE  
RETURNED OTHERWISE THIS DETERMINES THE TYPE OF AVERAGING PERFORMED ON THE DATA

BINARY ONLY REPORT RESULTS FOR THE CLASS SPECIFIED BY POSLABEL THIS IS APPLICABLE  
ONLY IF TARGETS YTRUE YPRED ARE BINARY

MICRO CALCULATE METRICS GLOBALLY BY COUNTING THE TOTAL TRUE POSITIVES FALSE NEGATIVES  
AND FALSE POSITIVES

MACRO CALCULATE METRICS FOR EACH LABEL AND FIND THEIR UNWEIGHTED MEAN THIS DOES NOT  
TAKE LABEL IMBALANCE INTO ACCOUNT

WEIGHTED CALCULATE METRICS FOR EACH LABEL AND FIND THEIR AVERAGE WEIGHTED BY SUPPORT  
THE NUMBER OF TRUE INSTANCES FOR EACH LABEL THIS ALTERS 'MACRO' TO ACCOUNT FOR LABEL  
IMBALANCE IT CAN RESULT IN AN FSCORE THAT IS NOT BETWEEN PRECISION AND RECALL

SAMPLES CALCULATE METRICS FOR EACH INSTANCE AND FIND THEIR AVERAGE ONLY MEANINGFUL  
FOR MULTILABEL CLASSIFICATION WHERE THIS DIFFERS FROM ACCURACYScore

SAMPLEWEIGHT ARRAYLIKE OF SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

RECALL FLOAT IF AVERAGE IS NOT NONE OR ARRAY OF FLOAT SHAPE NUNIQUE LABELS RECALL OF THE  
POSITIVE CLASS IN BINARY CLASSIFICATION OR WEIGHTED AVERAGE OF THE RECALL OF EACH CLASS FOR THE  
MULTICLASS TASK

SEE ALSO

2012 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

PRECISIONRECALLFScoresSUPPORT BALANCEDACCURACYScore

MULTILABELCONFUSIONMATRIX

NOTES

WHENTRUE POSITIVE FALSE NEGATIVE 0 RECALL RETURNS 0 AND RAISES

UNDEFINEDMETRICWARNING

EXAMPLES

FROM SKLEARNMETRICS IMPORT RECALLScore

YTRUE 0 1 2 0 1 2

YPRED 0 2 1 0 0 1

RECALLScoreYTRUE YPRED AVERAGEMACRO

033

RECALLScoreYTRUE YPRED AVERAGEMICRO

033

RECALLScoreYTRUE YPRED AVERAGEWEIGHTED

033

RECALLScoreYTRUE YPRED AVERAGE NONE

ARRAY1 0 0

EXAMPLES USING SKLEARNMETRICSRECALLScore

- PROBABILITY CALIBRATION CURVES

SKLEARNMETRICS ROCAUCScore

SKLEARNMETRICS ROCAUCScore YTRUE YScore AVERAGE'MACRO' SAMPLEWEIGHTNONE

MAXFPRNONE

COMPUTE AREA UNDER THE RECEIVER OPERATING CHARACTERISTIC CURVE ROC AUC FROM PREDICTION SCORES

NOTE THIS IMPLEMENTATION IS RESTRICTED TO THE BINARY CLASSIFICATION TASK OR MULTILABEL CLASSIFICATION TASK IN LABEL

INDICATOR FORMAT

READ MORE IN THE USER GUIDE

PARAMETERS

YTRUE ARRAY SHAPE NSAMPLES OR NSAMPLES NCLASSES TRUE BINARY LABELS OR BINARY LABEL

INDICATORS

YScore ARRAY SHAPE NSAMPLES OR NSAMPLES NCLASSES TARGET SCORES CAN EITHER BE

PROBABILITY ESTIMATES OF THE POSITIVE CLASS CONFIDENCE VALUES OR NONTHRESHOLDED MEASURE

OF DECISIONS AS RETURNED BY "DECISIONFUNCTION" ON SOME CLASSIFIERS FOR BINARY YTRUE

YScore IS SUPPOSED TO BE THE SCORE OF THE CLASS WITH GREATER LABEL

AVERAGE STRING NONE 'MICRO' 'MACRO' DEFAULT 'SAMPLES' 'WEIGHTED' IF NONE THE SCORES

FOR EACH CLASS ARE RETURNED OTHERWISE THIS DETERMINES THE TYPE OF AVERAGING PERFORMED ON

THE DATA

MICRO CALCULATE METRICS GLOBALLY BY CONSIDERING EACH ELEMENT OF THE LABEL INDICATOR

MATRIX AS A LABEL

624SKLEARNMETRICS METRICS 2013

SCIKITLEARN USER GUIDE RELEASE 0213

MACRO CALCULATE METRICS FOR EACH LABEL AND FIND THEIR UNWEIGHTED MEAN THIS DOES NOT TAKE LABEL IMBALANCE INTO ACCOUNT

WEIGHTED CALCULATE METRICS FOR EACH LABEL AND FIND THEIR AVERAGE WEIGHTED BY SUPPORT THE NUMBER OF TRUE INSTANCES FOR EACH LABEL

SAMPLES CALCULATE METRICS FOR EACH INSTANCE AND FIND THEIR AVERAGE WILL BE IGNORED WHEN YTRUE IS BINARY

SAMPLEWEIGHT ARRAYLIKE OF SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

MAXFPR FLOAT 0 AND 1 OPTIONAL IF NOT NONE THE STANDARDIZED PARTIAL AUC 3OVER THE RANGE 0 MAXFPR IS RETURNED

RETURNS

AUC FLOAT

SEE ALSO

AVERAGEPRECISIONSCORE AREA UNDER THE PRECISIONRECALL CURVE

ROCCURVE COMPUTE RECEIVER OPERATING CHARACTERISTIC ROC CURVE

REFERENCES

123

EXAMPLES

IMPORT NUMPY AS NP

FROM SKLEARNMETRICS IMPORT ROCAUCSCORE

YTRUE NPARRAY0 1 1

YSCORES NPARRAY01 04 035 08

ROCAUCSCOREYTRUE YSCORES

075

SKLEARNMETRICS ROCCURVE

SKLEARNMETRICS ROCCURVE YTRUE YSCORE POSLABELNONE SAMPLEWEIGHTNONE

DROPINTERMEDIATETRUE

COMPUTE RECEIVER OPERATING CHARACTERISTIC ROC

NOTE THIS IMPLEMENTATION IS RESTRICTED TO THE BINARY CLASSIFICATION TASK

READ MORE IN THE USER GUIDE

PARAMETERS

YTRUE ARRAY SHAPE NSAMPLES TRUE BINARY LABELS IF LABELS ARE NOT EITHER 1 1 OR 0 1 THEN POSLABEL SHOULD BE EXPLICITLY GIVEN

YSCORE ARRAY SHAPE NSAMPLES TARGET SCORES CAN EITHER BE PROBABILITY ESTIMATES OF THE POSITIVE CLASS CONFIDENCE VALUES OR NONTHRESHOLDED MEASURE OF DECISIONS AS RETURNED BY "DECISIONFUNCTION" ON SOME CLASSIFIERS

POSLABEL INT OR STR DEFAULTNONE THE LABEL OF THE POSITIVE CLASS WHEN POSLABELNONE

IF YTRUE IS IN 1 1 OR 0 1 POSLABEL IS SET TO 1 OTHERWISE AN ERROR WILL BE RAISED

2014 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

SAMPLEWEIGHT ARRAYLIKE OF SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS  
DROPINTERMEDIATE BOOLEAN OPTIONAL DEFAULTTRUE WHETHER TO DROP SOME SUBOPTIMAL  
THRESHOLDS WHICH WOULD NOT APPEAR ON A PLOTTED ROC CURVE THIS IS USEFUL IN ORDER TO  
CREATE LIGHTER ROC CURVES

NEW IN VERSION 017 PARAMETER DROPINTERMEDIATE

RETURNS  
FPRARRAY SHAPE 2 INCREASING FALSE POSITIVE RATES SUCH THAT ELEMENT I IS THE FALSE POSITIVE  
RATE OF PREDICTIONS WITH SCORE THRESHOLDSI  
TPRARRAY SHAPE 2 INCREASING TRUE POSITIVE RATES SUCH THAT ELEMENT I IS THE TRUE POSITIVE  
RATE OF PREDICTIONS WITH SCORE THRESHOLDSI  
THRESHOLDS ARRAY SHAPE NTHRESHOLDS DECREASING THRESHOLDS ON THE DECISION FUNCTION USED  
TO COMPUTE FPR AND TPR THRESHOLDS0 REPRESENTS NO INSTANCES BEING PREDICTED AND IS  
ARBITRARILY SET TO MAXYSORE 1

SEE ALSO  
ROCAUCSCORE COMPUTE THE AREA UNDER THE ROC CURVE

NOTES  
SINCE THE THRESHOLDS ARE SORTED FROM LOW TO HIGH VALUES THEY ARE REVERSED UPON RETURNING THEM TO ENSURE THEY  
CORRESPOND TO BOTH FPR ANDTPR WHICH ARE SORTED IN REVERSED ORDER DURING THEIR CALCULATION

REFERENCES

12

EXAMPLES

```
import numpy as np
from sklearn import metrics
Y = np.array([1, 1, 2, 2])
scores = np.array([0.1, 0.4, 0.35, 0.8])
fpr, tpr, thresholds = metrics.roc_curve(scores, poslabel=2)
```

```
fpr
array([ 0.  0.5  0.5  1.])
tpr
array([ 0.5  0.5  1.  1.])
```

```
thresholds
array([ 0.8  0.4  0.35  0.1])
```

EXAMPLES USING SKLEARNMETRICSROCCURVE

- SPECIES DISTRIBUTION MODELING
  - FEATURE TRANSFORMATIONS WITH ENSEMBLES OF TREES
  - RECEIVER OPERATING CHARACTERISTIC ROC WITH CROSS VALIDATION
- 624SKLEARNMETRICS METRICS 2015

SCIKITLEARN USER GUIDE RELEASE 0213

- RECEIVER OPERATING CHARACTERISTIC ROC

SKLEARNMETRICS ZEROONELOSS

SKLEARNMETRICS ZEROONELOSS YTRUE YPRED NORMALIZETRUE SAMPLEWEIGHTNONE

ZEROONE CLASSIFICATION LOSS

IF NORMALIZE IS TRUE RETURN THE FRACTION OF MISCLASSIFICATIONS FLOAT ELSE IT RETURNS THE NUMBER OF MISCLASSIFICATIONS INT THE BEST PERFORMANCE IS 0

READ MORE IN THE USER GUIDE

PARAMETERS

YTRUE 1D ARRAYLIKE OR LABEL INDICATOR ARRAY SPARSE MATRIX GROUND TRUTH CORRECT LABELS

YPRED 1D ARRAYLIKE OR LABEL INDICATOR ARRAY SPARSE MATRIX PREDICTED LABELS AS RETURNED BY A CLASSIFIER

NORMALIZE BOOL OPTIONAL DEFAULTTRUE IF FALSE RETURN THE NUMBER OF MISCLASSIFICATIONS OTHERWISE RETURN THE FRACTION OF MISCLASSIFICATIONS

SAMPLEWEIGHT ARRAYLIKE OF SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

LOSS FLOAT OR INT IF NORMALIZE TRUE RETURN THE FRACTION OF MISCLASSIFICATIONS FLOAT ELSE IT RETURNS THE NUMBER OF MISCLASSIFICATIONS INT

SEE ALSO

ACCURACYScore HAMMINGLOSS JACCARDScore

NOTES

IN MULTILABEL CLASSIFICATION THE ZEROONELOSS FUNCTION CORRESPONDS TO THE SUBSET ZEROONE LOSS FOR EACH SAMPLE THE ENTIRE SET OF LABELS MUST BE CORRECTLY PREDICTED OTHERWISE THE LOSS FOR THAT SAMPLE IS EQUAL TO ONE

EXAMPLES

```
from sklearn.metrics import zeroone_loss
ypred = [1, 2, 3, 4]
ytrue = [2, 2, 3, 4]
zeroone_loss(ytrue, ypred)
0.25
zeroone_loss(ytrue, ypred, normalize=False)
1
```

IN THE MULTILABEL CASE WITH BINARY LABEL INDICATORS

```
import numpy as np
zeroone_loss(np.array([0, 1, 1, 1]), np.ones(2, 2))
0.5
```

2016 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES USING SKLEARNMETRICSZEROONELOSS

- DISCRETE VERSUS REAL ADABOOST

6243 REGRESSION METRICS

SEE THE REGRESSION METRICS SECTION OF THE USER GUIDE FOR FURTHER DETAILS

METRICSEXPLAINEDVARIANCESCORE YTRUE

YPREDEXPLAINED VARIANCE REGRESSION SCORE FUNCTION

METRICSMAXERROR YTRUE YPRED MAXERROR METRIC CALCULATES THE MAXIMUM RESIDUAL ERROR

METRICSMEANABSOLUTEERROR YTRUE YPRED MEAN ABSOLUTE ERROR REGRESSION LOSS

METRICSMEANSQUAREDERROR YTRUE YPRED

MEAN SQUARED ERROR REGRESSION LOSS

METRICSMEANSQUAREDLOGERROR YTRUE

YPREDMEAN SQUARED LOGARITHMIC ERROR REGRESSION LOSS

METRICSMEDIANABSOLUTEERROR YTRUE

YPREDMEDIAN ABSOLUTE ERROR REGRESSION LOSS

METRICSR2SCORE YTRUE YPRED R2 COEFFICIENT OF DETERMINATION REGRESSION SCORE FUNCTION

SKLEARNMETRICS EXPLAINEDVARIANCESCORE

SKLEARNMETRICS EXPLAINEDVARIANCESCORE YTRUE YPRED SAMPLEWEIGHTNONE MULTIOUTPUT

PUT'UNIFORMAVERAGE'

EXPLAINED VARIANCE REGRESSION SCORE FUNCTION

BEST POSSIBLE SCORE IS 10 LOWER VALUES ARE WORSE

READ MORE IN THE USER GUIDE

PARAMETERS

YTRUE ARRAYLIKE OF SHAPE NSAMPLES OR NSAMPLES NOUTPUTS GROUND TRUTH CORRECT TARGET VALUES

YPRED ARRAYLIKE OF SHAPE NSAMPLES OR NSAMPLES NOUTPUTS ESTIMATED TARGET VALUES

SAMPLEWEIGHT ARRAYLIKE OF SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

MULTIOUTPUT STRING IN 'RAWVALUES' 'UNIFORMAVERAGE' 'VARIANCEWEIGHTED' OR ARRAYLIKE OF SHAPE NOUTPUTS DEFINES AGGREGATING OF MULTIPLE OUTPUT SCORES ARRAYLIKE VALUE DEFINES WEIGHTS USED TO AVERAGE SCORES

'RAWVALUES' RETURNS A FULL SET OF SCORES IN CASE OF MULTIOUTPUT INPUT

'UNIFORMAVERAGE' SCORES OF ALL OUTPUTS ARE AVERAGED WITH UNIFORM WEIGHT

'VARIANCEWEIGHTED' SCORES OF ALL OUTPUTS ARE AVERAGED WEIGHTED BY THE VARIANCES OF EACH INDIVIDUAL OUTPUT

RETURNS

SCORE FLOAT OR NDARRAY OF FLOATS THE EXPLAINED VARIANCE OR NDARRAY IF 'MULTIOUTPUT' IS 'RAWVALUES'

624SKLEARNMETRICS METRICS 2017

SCIKITLEARN USER GUIDE RELEASE 0213

NOTES

THIS IS NOT A SYMMETRIC FUNCTION

EXAMPLES

FROM SKLEARNMETRICS IMPORT EXPLAINEDVARIANCESCORE

YTRUE 3 05 2 7

YPRED 25 00 2 8

EXPLAINEDVARIANCESCOREYTRUE YPRED

0957

YTRUE 05 1 1 1 7 6

YPRED 0 2 1 2 8 5

EXPLAINEDVARIANCESCOREYTRUE YPRED MULTIOUTPUTUNIFORMAVERAGE

0983

SKLEARNMETRICS MAXERROR

SKLEARNMETRICS MAXERROR YTRUE YPRED

MAXERROR METRIC CALCULATES THE MAXIMUM RESIDUAL ERROR

READ MORE IN THE USER GUIDE

PARAMETERS

YTRUE ARRAYLIKE OF SHAPE NSAMPLES GROUND TRUTH CORRECT TARGET VALUES

YPRED ARRAYLIKE OF SHAPE NSAMPLES ESTIMATED TARGET VALUES

RETURNS

MAXERROR FLOAT A POSITIVE FLOATING POINT VALUE THE BEST VALUE IS 00

EXAMPLES

FROM SKLEARNMETRICS IMPORT MAXERROR

YTRUE 3 2 7 1

YPRED 4 2 7 1

MAXERRORYTRUE YPRED

1

SKLEARNMETRICS MEANABSOLUTEERROR

SKLEARNMETRICS MEANABSOLUTEERROR YTRUE YPRED SAMPLEWEIGHTNONE MULTIOUT

PUT'UNIFORMAVERAGE'

MEAN ABSOLUTE ERROR REGRESSION LOSS

READ MORE IN THE USER GUIDE

PARAMETERS

YTRUE ARRAYLIKE OF SHAPE NSAMPLES OR NSAMPLES NOUTPUTS GROUND TRUTH CORRECT

TARGET VALUES

2018 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

YPRED ARRAYLIKE OF SHAPE NSAMPLES OR NSAMPLES NOUTPUTS ESTIMATED TARGET VALUES  
SAMPLEWEIGHT ARRAYLIKE OF SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS  
MULTIOUTPUT STRING IN 'RAWVALUES' 'UNIFORMAVERAGE' OR ARRAYLIKE OF SHAPE NOUTPUTS  
DEFINES AGGREGATING OF MULTIPLE OUTPUT VALUES ARRAYLIKE VALUE DEFINES WEIGHTS USED TO  
AVERAGE ERRORS  
'RAWVALUES' RETURNS A FULL SET OF ERRORS IN CASE OF MULTIOUTPUT INPUT  
'UNIFORMAVERAGE' ERRORS OF ALL OUTPUTS ARE AVERAGED WITH UNIFORM WEIGHT  
RETURNS  
LOSS FLOAT OR NDARRAY OF FLOATS IF MULTIOUTPUT IS 'RAWVALUES' THEN MEAN ABSOLUTE ERROR IS RE  
TURNED FOR EACH OUTPUT SEPARATELY IF MULTIOUTPUT IS 'UNIFORMAVERAGE' OR AN NDARRAY OF  
WEIGHTS THEN THE WEIGHTED AVERAGE OF ALL OUTPUT ERRORS IS RETURNED  
MAE OUTPUT IS NONNEGATIVE FLOATING POINT THE BEST VALUE IS 00

EXAMPLES  
FROM SKLEARNMETRICS IMPORT MEANABSOLUTEERROR  
YTRUE 3 05 2 7  
YPRED 25 00 2 8  
MEANABSOLUTEERRORYTRUE YPRED  
05  
YTRUE 05 1 1 1 7 6  
YPRED 0 2 1 2 8 5  
MEANABSOLUTEERRORYTRUE YPRED  
075  
MEANABSOLUTEERRORYTRUE YPRED MULTIOUTPUTRAWVALUES  
ARRAY05 1  
MEANABSOLUTEERRORYTRUE YPRED MULTIOUTPUT03 07

085  
SKLEARNMETRICS MEANSQUAREDERROR  
SKLEARNMETRICS MEANSQUAREDERROR YTRUE YPRED SAMPLEWEIGHTNONE MULTIOUT  
PUT'UNIFORMAVERAGE'  
MEAN SQUARED ERROR REGRESSION LOSS  
READ MORE IN THE USER GUIDE

PARAMETERS  
YTRUE ARRAYLIKE OF SHAPE NSAMPLES OR NSAMPLES NOUTPUTS GROUND TRUTH CORRECT  
TARGET VALUES  
YPRED ARRAYLIKE OF SHAPE NSAMPLES OR NSAMPLES NOUTPUTS ESTIMATED TARGET VALUES  
SAMPLEWEIGHT ARRAYLIKE OF SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS  
MULTIOUTPUT STRING IN 'RAWVALUES' 'UNIFORMAVERAGE' OR ARRAYLIKE OF SHAPE NOUTPUTS  
DEFINES AGGREGATING OF MULTIPLE OUTPUT VALUES ARRAYLIKE VALUE DEFINES WEIGHTS USED TO  
AVERAGE ERRORS  
'RAWVALUES' RETURNS A FULL SET OF ERRORS IN CASE OF MULTIOUTPUT INPUT  
624SKLEARNMETRICS METRICS 2019

SCIKITLEARN USER GUIDE RELEASE 0213

‘UNIFORMAVERAGE’ ERRORS OF ALL OUTPUTS ARE AVERAGED WITH UNIFORM WEIGHT

RETURNS

LOSS FLOAT OR NDARRAY OF FLOATS A NONNEGATIVE FLOATING POINT VALUE THE BEST VALUE IS 00 OR AN  
ARRAY OF FLOATING POINT VALUES ONE FOR EACH INDIVIDUAL TARGET

EXAMPLES

```
FROM SKLEARNMETRICS IMPORT MEANSQUAREDERROR
YTRUE 3 05 2 7
YPRED 25 00 2 8
MEANSQUAREDERRORYTRUE YPRED
0375
YTRUE 05 11 17 6
YPRED 0 21 28 5
MEANSQUAREDERRORYTRUE YPRED
0708
MEANSQUAREDERRORYTRUE YPRED MULTIOUTPUTRAWVALUES
```

ARRAY041666667 1

MEANSQUAREDERRORYTRUE YPRED MULTIOUTPUT03 07

0825

EXAMPLES USING SKLEARNMETRICSMEANSQUAREDERROR

- MODEL COMPLEXITY INFLUENCE
- GRADIENT BOOSTING REGRESSION
- PLOT RIDGE COEFFICIENTS AS A FUNCTION OF THE L2 REGULARIZATION
- LINEAR REGRESSION EXAMPLE
- ROBUST LINEAR ESTIMATOR FITTING

SKLEARNMETRICS MEANSQUAREDLOGERROR

SKLEARNMETRICS MEANSQUAREDLOGERROR YTRUE YPRED SAMPLEWEIGHTNONE MULTIOUT  
PUT‘UNIFORMAVERAGE’

MEAN SQUARED LOGARITHMIC ERROR REGRESSION LOSS

READ MORE IN THE USER GUIDE

PARAMETERS

YTRUE ARRAYLIKE OF SHAPE NSAMPLES OR NSAMPLES NOUTPUTS GROUND TRUTH CORRECT  
TARGET VALUES

YPRED ARRAYLIKE OF SHAPE NSAMPLES OR NSAMPLES NOUTPUTS ESTIMATED TARGET VALUES

SAMPLEWEIGHT ARRAYLIKE OF SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

MULTIOUTPUT STRING IN ‘RAWVALUES’ ‘UNIFORMAVERAGE’ OR ARRAYLIKE OF SHAPE NOUTPUTS  
DEFINES AGGREGATING OF MULTIPLE OUTPUT VALUES ARRAYLIKE VALUE DEFINES WEIGHTS USED TO  
AVERAGE ERRORS

2020 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

‘RAWVALUES’ RETURNS A FULL SET OF ERRORS WHEN THE INPUT IS OF MULTIOUTPUT FORMAT

‘UNIFORMAVERAGE’ ERRORS OF ALL OUTPUTS ARE AVERAGED WITH UNIFORM WEIGHT

RETURNS

LOSS FLOAT OR NDARRAY OF FLOATS A NONNEGATIVE FLOATING POINT VALUE THE BEST VALUE IS 00 OR AN

ARRAY OF FLOATING POINT VALUES ONE FOR EACH INDIVIDUAL TARGET

EXAMPLES

FROM SKLEARNMETRICS IMPORT MEANSQUAREDLOGERROR

YTRUE 3 5 25 7

YPRED 25 5 4 8

MEANSQUAREDLOGERRORYTRUE YPRED

0039

YTRUE 05 1 1 2 7 6

YPRED 05 2 1 25 8 8

MEANSQUAREDLOGERRORYTRUE YPRED

0044

MEANSQUAREDLOGERRORYTRUE YPRED MULTIOUTPUTRAWVALUES

ARRAY000462428 008377444

MEANSQUAREDLOGERRORYTRUE YPRED MULTIOUTPUT03 07

0060

SKLEARNMETRICS MEDIANABSOLUTEERROR

SKLEARNMETRICS MEDIANABSOLUTEERROR YTRUE YPRED

MEDIAN ABSOLUTE ERROR REGRESSION LOSS

READ MORE IN THE USER GUIDE

PARAMETERS

YTRUE ARRAYLIKE OF SHAPE NSAMPLES GROUND TRUTH CORRECT TARGET VALUES

YPRED ARRAYLIKE OF SHAPE NSAMPLES ESTIMATED TARGET VALUES

RETURNS

LOSS FLOAT A POSITIVE FLOATING POINT VALUE THE BEST VALUE IS 00

EXAMPLES

FROM SKLEARNMETRICS IMPORT MEDIANABSOLUTEERROR

YTRUE 3 05 2 7

YPRED 25 00 2 8

MEDIANABSOLUTEERRORYTRUE YPRED

05

EXAMPLES USING SKLEARNMETRICSMEDIANABSOLUTEERROR

•EFFECT OF TRANSFORMING THE TARGETS IN REGRESSION MODEL

624SKLEARNMETRICS METRICS 2021

SCIKITLEARN USER GUIDE RELEASE 0213

SKLEARNMETRICS R2SCORE

SKLEARNMETRICS R2SCORE YTRUE YPRED SAMPLEWEIGHTNONE MULTIOUTPUT'UNIFORMAVERAGE'

R2 COEFFICIENT OF DETERMINATION REGRESSION SCORE FUNCTION

BEST POSSIBLE SCORE IS 10 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 00

READ MORE IN THE USER GUIDE

PARAMETERS

YTRUE ARRAYLIKE OF SHAPE NSAMPLES OR NSAMPLES NOUTPUTS GROUND TRUTH CORRECT

TARGET VALUES

YPRED ARRAYLIKE OF SHAPE NSAMPLES OR NSAMPLES NOUTPUTS ESTIMATED TARGET VALUES

SAMPLEWEIGHT ARRAYLIKE OF SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

MULTIOUTPUT STRING IN 'RAWVALUES' 'UNIFORMAVERAGE' 'VARIANCEWEIGHTED' OR NONE OR

ARRAYLIKE OF SHAPE NOUTPUTS DEFINES AGGREGATING OF MULTIPLE OUTPUT SCORES ARRAYLIKE

VALUE DEFINES WEIGHTS USED TO AVERAGE SCORES DEFAULT IS "UNIFORMAVERAGE"

'RAWVALUES' RETURNS A FULL SET OF SCORES IN CASE OF MULTIOUTPUT INPUT

'UNIFORMAVERAGE' SCORES OF ALL OUTPUTS ARE AVERAGED WITH UNIFORM WEIGHT

'VARIANCEWEIGHTED' SCORES OF ALL OUTPUTS ARE AVERAGED WEIGHTED BY THE VARIANCES OF

EACH INDIVIDUAL OUTPUT

CHANGED IN VERSION 019 DEFAULT VALUE OF MULTIOUTPUT IS 'UNIFORMAVERAGE'

RETURNS

ZFLOAT OR NDARRAY OF FLOATS THE R2 SCORE OR NDARRAY OF SCORES IF 'MULTIOUTPUT' IS 'RAWVALUES'

NOTES

THIS IS NOT A SYMMETRIC FUNCTION

UNLIKE MOST OTHER SCORES R2 SCORE MAY BE NEGATIVE IT NEED NOT ACTUALLY BE THE SQUARE OF A QUANTITY R

THIS METRIC IS NOT WELLDEFINED FOR SINGLE SAMPLES AND WILL RETURN A NAN VALUE IF NSAMPLES IS LESS THAN TWO

REFERENCES

1

EXAMPLES

FROM SKLEARNMETRICS IMPORT R2SCORE

YTRUE 3 05 2 7

YPRED 25 00 2 8

R2SCOREYTRUE YPRED

0948

YTRUE 05 1 1 1 7 6

YPRED 0 2 1 2 8 5

R2SCOREYTRUE YPRED

MULTIOUTPUTVARIANCEWEIGHTED

2022 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

0938

YTRUE 1 2 3

YPRED 1 2 3

R2SCOREYTRUE YPRED

10

YTRUE 1 2 3

YPRED 2 2 2

R2SCOREYTRUE YPRED

00

YTRUE 1 2 3

YPRED 3 2 1

R2SCOREYTRUE YPRED

30

EXAMPLES USING SKLEARNMETRICSR2SCORE

•EFFECT OF TRANSFORMING THE TARGETS IN REGRESSION MODEL

•LINEAR REGRESSION EXAMPLE

•LASSO AND ELASTIC NET FOR SPARSE SIGNALS

6244 MULTILABEL RANKING METRICS

SEE THE MULTILABEL RANKING METRICS SECTION OF THE USER GUIDE FOR FURTHER DETAILS

METRICSCOVERAGEERROR YTRUE YSCORE COVERAGE ERROR MEASURE

METRICSLABELRANKINGAVERAGEPRECISIONSCORE COMPUTE RANKINGBASED AVERAGE PRECISION

METRICSLABELRANKINGLOSS YTRUE YSCORE COMPUTE RANKING LOSS MEASURE

SKLEARNMETRICS COVERAGEERROR

SKLEARNMETRICS COVERAGEERROR YTRUE YSCORE SAMPLEWEIGHTNONE

COVERAGE ERROR MEASURE

COMPUTE HOW FAR WE NEED TO GO THROUGH THE RANKED SCORES TO COVER ALL TRUE LABELS THE BEST VALUE IS EQUAL TO THE

AVERAGE NUMBER OF LABELS IN YTRUE PER SAMPLE

TIES INSCORES ARE BROKEN BY GIVING MAXIMAL RANK THAT WOULD HAVE BEEN ASSIGNED TO ALL TIED VALUES

NOTE OUR IMPLEMENTATION'S SCORE IS 1 GREATER THAN THE ONE GIVEN IN TSOUMAKAS ET AL 2010 THIS EXTENDS IT TO

HANDLE THE DEGENERATE CASE IN WHICH AN INSTANCE HAS 0 TRUE LABELS

READ MORE IN THE USER GUIDE

PARAMETERS

YTRUE ARRAY SHAPE NSAMPLES NLABELS TRUE BINARY LABELS IN BINARY INDICATOR FORMAT

YSCORE ARRAY SHAPE NSAMPLES NLABELS TARGET SCORES CAN EITHER BE PROBABILITY ESTI

MATES OF THE POSITIVE CLASS CONFIDENCE VALUES OR NONTHRESHOLDED MEASURE OF DECISIONS AS

RETURNED BY "DECISIONFUNCTION" ON SOME CLASSIFIERS

SAMPLEWEIGHT ARRAYLIKE OF SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

6245SKLEARNMETRICS METRICS 2023

SKLEARNMETRICS LABELRANKINGAVERAGEPRECISIONSCORE

SKLEARNMETRICS LABELRANKINGAVERAGEPRECISIONSCORE YTRUE YSCORE SAM

PLEWEIGHTNONE

COMPUTE RANKINGBASED AVERAGE PRECISION

LABEL RANKING AVERAGE PRECISION LRAP IS THE AVERAGE OVER EACH GROUND TRUTH LABEL ASSIGNED TO EACH SAMPLE OF THE RATIO OF TRUE VS TOTAL LABELS WITH LOWER SCORE

THIS METRIC IS USED IN MULTILABEL RANKING PROBLEM WHERE THE GOAL IS TO GIVE BETTER RANK TO THE LABELS ASSOCIATED TO EACH SAMPLE

THE OBTAINED SCORE IS ALWAYS STRICTLY GREATER THAN 0 AND THE BEST VALUE IS 1

READ MORE IN THE USER GUIDE

PARAMETERS

YTRUE ARRAY OR SPARSE MATRIX SHAPE NSAMPLES NLABELS TRUE BINARY LABELS IN BINARY

INDICATOR FORMAT

YSCORE ARRAY SHAPE NSAMPLES NLABELS TARGET SCORES CAN EITHER BE PROBABILITY ESTI

MATES OF THE POSITIVE CLASS CONFIDENCE VALUES OR NONTHRESHOLDED MEASURE OF DECISIONS AS RETURNED BY “DECISIONFUNCTION” ON SOME CLASSIFIERS

SAMPLEWEIGHT ARRAYLIKE OF SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT

EXAMPLES

IMPORT NUMPY AS NP

FROM SKLEARNMETRICS IMPORT LABELRANKINGAVERAGEPRECISIONSCORE

YTRUE NPARRAY1 0 0 0 0 1

YSCORE NPARRAY075 05 1 1 02 01

LABELRANKINGAVERAGEPRECISIONSCOREYTRUE YSCORE

0416

SKLEARNMETRICS LABELRANKINGLOSS

SKLEARNMETRICS LABELRANKINGLOSS YTRUE YSCORE SAMPLEWEIGHTNONE

COMPUTE RANKING LOSS MEASURE

COMPUTE THE AVERAGE NUMBER OF LABEL PAIRS THAT ARE INCORRECTLY ORDERED GIVEN YSCORE WEIGHTED BY THE SIZE OF THE LABEL SET AND THE NUMBER OF LABELS NOT IN THE LABEL SET

THIS IS SIMILAR TO THE ERROR SET SIZE BUT WEIGHTED BY THE NUMBER OF RELEVANT AND IRRELEVANT LABELS THE BEST PERFORMANCE IS ACHIEVED WITH A RANKING LOSS OF ZERO

SCIKITLEARN USER GUIDE RELEASE 0213

READ MORE IN THE USER GUIDE

NEW IN VERSION 017 A FUNCTION LABELRANKINGLOSS

PARAMETERS

YTRUE ARRAY OR SPARSE MATRIX SHAPE NSAMPLES NLABELS TRUE BINARY LABELS IN BINARY

INDICATOR FORMAT

YSCORE ARRAY SHAPE NSAMPLES NLABELS TARGET SCORES CAN EITHER BE PROBABILITY ESTI

MATES OF THE POSITIVE CLASS CONFIDENCE VALUES OR NONTHRESHOLDED MEASURE OF DECISIONS AS

RETURNED BY “DECISIONFUNCTION” ON SOME CLASSIFIERS

SAMPLEWEIGHT ARRAYLIKE OF SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

LOSS FLOAT

REFERENCES

1

6245 CLUSTERING METRICS

SEE THE CLUSTERING PERFORMANCE EVALUATION SECTION OF THE USER GUIDE FOR FURTHER DETAILS THE SKLEARNMETRICS

CLUSTER SUBMODULE CONTAINS EVALUATION METRICS FOR CLUSTER ANALYSIS RESULTS THERE ARE TWO FORMS OF EVALUATION

- SUPERVISED WHICH USES A GROUND TRUTH CLASS VALUES FOR EACH SAMPLE
- UNSUPERVISED WHICH DOES NOT AND MEASURES THE ‘QUALITY’ OF THE MODEL ITSELF

METRICSDJUSTEDMUTUALINFOSCORE

ADJUSTED MUTUAL INFORMATION BETWEEN TWO CLUSTERINGS

METRICSDJUSTEDRANDSCORE LABELSTRUE RAND INDEX ADJUSTED FOR CHANCE

METRICSCALINSKIHARABASZSCORE X LABELS COMPUTE THE CALINSKI AND HARABASZ SCORE

METRICSDAVIESBOULDINSCORE X LABELS COMPUTES THE DAVIESBOULDIN SCORE

METRICSCOMPLETENESSSCORE LABELSTRUE COMPLETENESS METRIC OF A CLUSTER LABELING GIVEN A GROUND TRUTH

METRICSCUSTERCONTINGENCYMATRIX

BUILD A CONTINGENCY MATRIX DESCRIBING THE RELATIONSHIP BETWEEN LABELS

METRICSFOWLKESMALLOWSSCORE LABELSTRUE

MEASURE THE SIMILARITY OF TWO CLUSTERINGS OF A SET OF POINTS

METRICSHOMOGENEITYCOMPLETENESSVMEASURE COMPUTE THE HOMOGENEITY AND COMPLETENESS AND V MEASURE SCORES AT ONCE

METRICSHOMOGENEITYSCORE LABELSTRUE HOMOGENEITY METRIC OF A CLUSTER LABELING GIVEN A GROUND TRUTH

METRICSMUTUALINFOSCORE LABELSTRUE MUTUAL INFORMATION BETWEEN TWO CLUSTERINGS

METRICSNORMALIZEDMUTUALINFOSCORE

NORMALIZED MUTUAL INFORMATION BETWEEN TWO CLUSTERINGS

METRICSSILHOUETTESCORE X LABELS COMPUTE THE MEAN SILHOUETTE COEFFICIENT OF ALL SAMPLES

METRICSSILHOUETTESAMPLES X LABELS METRIC COMPUTE THE SILHOUETTE COEFFICIENT FOR EACH SAMPLE

METRICSVMEASURESORE LABELSTRUE LA

BELSPREDVMEASURE CLUSTER LABELING GIVEN A GROUND TRUTH

6245SKLEARNMETRICS METRICS 2025

SCIKITLEARN USER GUIDE RELEASE 0213

SKLEARNMETRICS ADJUSTEDMUTUALINFOSCORE

SKLEARNMETRICS ADJUSTEDMUTUALINFOSCORE LABELSTRUE LABELSPRED AVER  
AGEMETHOD'WARN'

ADJUSTED MUTUAL INFORMATION BETWEEN TWO CLUSTERINGS

ADJUSTED MUTUAL INFORMATION AMI IS AN ADJUSTMENT OF THE MUTUAL INFORMATION MI SCORE TO ACCOUNT FOR  
CHANCE IT ACCOUNTS FOR THE FACT THAT THE MI IS GENERALLY HIGHER FOR TWO CLUSTERINGS WITH A LARGER NUMBER OF  
CLUSTERS REGARDLESS OF WHETHER THERE IS ACTUALLY MORE INFORMATION SHARED FOR TWO CLUSTERINGS [AND] THE AMI  
IS GIVEN AS

AMIU V MIU V EMIU V AVGHU HV EMIU V

THIS METRIC IS INDEPENDENT OF THE ABSOLUTE VALUES OF THE LABELS A PERMUTATION OF THE CLASS OR CLUSTER LABEL VALUES  
WON'T CHANGE THE SCORE VALUE IN ANY WAY

THIS METRIC IS FURTHERMORE SYMMETRIC SWITCHING LABELTRUE WITHLABELPRED WILL RETURN THE SAME SCORE  
VALUE THIS CAN BE USEFUL TO MEASURE THE AGREEMENT OF TWO INDEPENDENT LABEL ASSIGNMENTS STRATEGIES ON THE SAME  
DATASET WHEN THE REAL GROUND TRUTH IS NOT KNOWN

BE MINDFUL THAT THIS FUNCTION IS AN ORDER OF MAGNITUDE SLOWER THAN OTHER METRICS SUCH AS THE ADJUSTED RAND  
INDEX

READ MORE IN THE USER GUIDE

PARAMETERS

LABELSTRUE INT ARRAY SHAPE NSAMPLES A CLUSTERING OF THE DATA INTO DISJOINT SUBSETS

LABELSPRED ARRAY SHAPE NSAMPLES A CLUSTERING OF THE DATA INTO DISJOINT SUBSETS

AVERAGEMETHOD STRING OPTIONAL DEFAULT 'WARN' HOW TO COMPUTE THE NORMALIZER IN THE  
DENOMINATOR POSSIBLE OPTIONS ARE 'MIN' 'GEOMETRIC' 'ARITHMETIC' AND 'MAX' IF 'WARN'  
'MAX' WILL BE USED THE DEFAULT WILL CHANGE TO 'ARITHMETIC' IN VERSION 022

NEW IN VERSION 020

RETURNS

AMI FLOAT UPPERLIMITED BY 10 THE AMI RETURNS A VALUE OF 1 WHEN THE TWO PARTITIONS ARE  
IDENTICAL IE PERFECTLY MATCHED RANDOM PARTITIONS INDEPENDENT LABELLINGS HAVE AN EX  
PECTED AMI AROUND 0 ON AVERAGE HENCE CAN BE NEGATIVE

SEE ALSO

ADJUSTEDRANDSCORE ADJUSTED RAND INDEX

MUTUALINFOSCORE MUTUAL INFORMATION NOT ADJUSTED FOR CHANCE

REFERENCES

12

EXAMPLES

PERFECT LABELINGS ARE BOTH HOMOGENEOUS AND COMPLETE HENCE HAVE SCORE 10

2026 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213  
FROM SKLEARNMETRICSCUSTER IMPORT ADJUSTEDMUTUALINFOSCORE  
ADJUSTEDMUTUALINFOSCORE0 0 1 1 0 0 1 1

10  
ADJUSTEDMUTUALINFOSCORE0 0 1 1 1 1 0 0

10  
IF CLASSES MEMBERS ARE COMPLETELY SPLIT ACROSS DIFFERENT CLUSTERS THE ASSIGNMENT IS TOTALLY INCOMPLETE HENCE  
THE AMI IS NULL  
ADJUSTEDMUTUALINFOSCORE0 0 0 0 0 1 2 3

00  
EXAMPLES USING SKLEARNMETRICSDJUSTEDMUTUALINFOSCORE  
•DEMO OF AFFINITY PROPAGATION CLUSTERING ALGORITHM  
•DEMO OF DBSCAN CLUSTERING ALGORITHM  
•ADJUSTMENT FOR CHANCE IN CLUSTERING PERFORMANCE EVALUATION  
•A DEMO OF KMEANS CLUSTERING ON THE HANDWRITTEN DIGITS DATA

SKLEARNMETRICS ADJUSTEDRANDSCORE  
SKLEARNMETRICS ADJUSTEDRANDSCORE LABELSTRUE LABELSPRED  
RAND INDEX ADJUSTED FOR CHANCE

THE RAND INDEX COMPUTES A SIMILARITY MEASURE BETWEEN TWO CLUSTERINGS BY CONSIDERING ALL PAIRS OF SAMPLES AND  
COUNTING PAIRS THAT ARE ASSIGNED IN THE SAME OR DIFFERENT CLUSTERS IN THE PREDICTED AND TRUE CLUSTERINGS  
THE RAW RI SCORE IS THEN “ADJUSTED FOR CHANCE” INTO THE ARI SCORE USING THE FOLLOWING SCHEME

ARI RI EXPECTEDRI MAXRI EXPECTEDRI  
THE ADJUSTED RAND INDEX IS THUS ENSURED TO HAVE A VALUE CLOSE TO 00 FOR RANDOM LABELING INDEPENDENTLY OF THE  
NUMBER OF CLUSTERS AND SAMPLES AND EXACTLY 10 WHEN THE CLUSTERINGS ARE IDENTICAL UP TO A PERMUTATION

ARI IS A SYMMETRIC MEASURE  
ADJUSTEDRANDSCOREA B ADJUSTEDRANDSCOREB A  
READ MORE IN THE USER GUIDE

PARAMETERS  
LABELSTRUE INT ARRAY SHAPE NSAMPLES GROUND TRUTH CLASS LABELS TO BE USED AS A REFERENCE  
LABELSPRED ARRAY SHAPE NSAMPLES CLUSTER LABELS TO EVALUATE

RETURNS  
ARIFLOAT SIMILARITY SCORE BETWEEN 10 AND 10 RANDOM LABELINGS HAVE AN ARI CLOSE TO 00  
10 STANDS FOR PERFECT MATCH  
SEE ALSO  
624SKLEARNMETRICS METRICS 2027

SCIKITLEARN USER GUIDE RELEASE 0213

ADJUSTEDMUTUALINFOSCORE ADJUSTED MUTUAL INFORMATION

REFERENCES

HUBERT1985 WK

EXAMPLES

PERFECTLY MATCHING LABELINGS HAVE A SCORE OF 1 EVEN

FROM SKLEARNMETRICSCluster import AdjustedRandScore

AdjustedRandScore0 0 1 1 0 0 1 1

10

AdjustedRandScore0 0 1 1 1 1 0 0

10

LABELINGS THAT ASSIGN ALL CLASSES MEMBERS TO THE SAME CLUSTERS ARE COMPLETE BE NOT ALWAYS PURE HENCE PENALIZED

AdjustedRandScore0 0 1 2 0 0 1 1

057

ARI IS SYMMETRIC SO LABELINGS THAT HAVE PURE CLUSTERS WITH MEMBERS COMING FROM THE SAME CLASSES BUT UNNECESSARY SPLITS ARE PENALIZED

AdjustedRandScore0 0 1 1 0 0 1 2

057

IF CLASSES MEMBERS ARE COMPLETELY SPLIT ACROSS DIFFERENT CLUSTERS THE ASSIGNMENT IS TOTALLY INCOMPLETE HENCE THE ARI IS VERY LOW

AdjustedRandScore0 0 0 0 0 1 2 3

00

EXAMPLES USING SKLEARNMETRICSAAdjustedRandScore

- DEMO OF AFFINITY PROPAGATION CLUSTERING ALGORITHM
- DEMO OF DBSCAN CLUSTERING ALGORITHM
- ADJUSTMENT FOR CHANCE IN CLUSTERING PERFORMANCE EVALUATION
- A DEMO OF KMEANS CLUSTERING ON THE HANDWRITTEN DIGITS DATA
- CLUSTERING TEXT DOCUMENTS USING KMEANS

SKLEARNMETRICS CALINSKI HARABASZ SCORE

SKLEARNMETRICS CALINSKI HARABASZ SCORE XLABELS

COMPUTE THE CALINSKI AND HARABASZ SCORE

IT IS ALSO KNOWN AS THE VARIANCE RATIO CRITERION

THE SCORE IS DEFINED AS RATIO BETWEEN THE WITHINCLUSTER DISPERSION AND THE BETWEENCLUSTER DISPERSION

READ MORE IN THE USER GUIDE

2028 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES LIST OFNFEATURES DIMENSIONAL DATA

POINTS EACH ROW CORRESPONDS TO A SINGLE DATA POINT

LABELS ARRAYLIKE SHAPE NSAMPLES PREDICTED LABELS FOR EACH SAMPLE

RETURNS

SCORE FLOAT THE RESULTING CALINSKI HARABASZ SCORE

REFERENCES

1

SKLEARNMETRICS DAVIESBOULDIN SCORE

SKLEARNMETRICS DAVIESBOULDIN SCORE XLABELS

COMPUTES THE DAVIESBOULDIN SCORE

THE SCORE IS DEFINED AS THE AVERAGE SIMILARITY MEASURE OF EACH CLUSTER WITH ITS MOST SIMILAR CLUSTER WHERE SIMILARITY IS THE RATIO OF WITHINCLUSTER DISTANCES TO BETWEENCLUSTER DISTANCES THUS CLUSTERS WHICH ARE FARTHER APART AND LESS DISPERSED WILL RESULT IN A BETTER SCORE

THE MINIMUM SCORE IS ZERO WITH LOWER VALUES INDICATING BETTER CLUSTERING

READ MORE IN THE USER GUIDE

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES LIST OFNFEATURES DIMENSIONAL DATA

POINTS EACH ROW CORRESPONDS TO A SINGLE DATA POINT

LABELS ARRAYLIKE SHAPE NSAMPLES PREDICTED LABELS FOR EACH SAMPLE

RETURNS

SCORE FLOAT THE RESULTING DAVIESBOULDIN SCORE

REFERENCES

1

SKLEARNMETRICS COMPLETENESS SCORE

SKLEARNMETRICS COMPLETENESS SCORE LABEL TRUE LABELS PRED

COMPLETENESS METRIC OF A CLUSTER LABELING GIVEN A GROUND TRUTH

A CLUSTERING RESULT SATISFIES COMPLETENESS IF ALL THE DATA POINTS THAT ARE MEMBERS OF A GIVEN CLASS ARE ELEMENTS OF THE SAME CLUSTER THIS METRIC IS INDEPENDENT OF THE ABSOLUTE VALUES OF THE LABELS A PERMUTATION OF THE CLASS OR CLUSTER LABEL VALUES WON'T CHANGE THE SCORE VALUE IN ANY WAY

THIS METRIC IS NOT SYMMETRIC SWITCHING LABEL TRUE WITH LABEL PRED WILL RETURN THE HOMOGENEITY SCORE WHICH WILL BE DIFFERENT IN GENERAL

READ MORE IN THE USER GUIDE

624 SKLEARNMETRICS METRICS 2029

SCIKITLEARN USER GUIDE RELEASE 0213

PARAMETERS

LABELSTRUE INT ARRAY SHAPE NSAMPLES GROUND TRUTH CLASS LABELS TO BE USED AS A REFERENCE

LABELSPRED ARRAY SHAPE NSAMPLES CLUSTER LABELS TO EVALUATE

RETURNS

COMPLETENESS FLOAT SCORE BETWEEN 00 AND 10 10 STANDS FOR PERFECTLY COMPLETE LABELING

SEE ALSO

HOMOGENEITYSCORE

VMEASURESCORE

REFERENCES

1

EXAMPLES

PERFECT LABELINGS ARE COMPLETE

FROM SKLEARNMETRICSCUSTER IMPORT COMPLETENESSSCORE

COMPLETENESSSCORE0 0 1 1 1 1 0 0

10

NONPERFECT LABELINGS THAT ASSIGN ALL CLASSES MEMBERS TO THE SAME CLUSTERS ARE STILL COMPLETE

PRINTCOMPLETENESSSCORE0 0 1 1 0 0 0 0

10

PRINTCOMPLETENESSSCORE0 1 2 3 0 0 1 1

0999

IF CLASSES MEMBERS ARE SPLIT ACROSS DIFFERENT CLUSTERS THE ASSIGNMENT CANNOT BE COMPLETE

PRINTCOMPLETENESSSCORE0 0 1 1 0 1 0 1

00

PRINTCOMPLETENESSSCORE0 0 0 0 0 1 2 3

00

EXAMPLES USING SKLEARNMETRICSCOMPLETENESSSCORE

- DEMO OF AFFINITY PROPAGATION CLUSTERING ALGORITHM
- DEMO OF DBSCAN CLUSTERING ALGORITHM
- A DEMO OF KMEANS CLUSTERING ON THE HANDWRITTEN DIGITS DATA
- CLUSTERING TEXT DOCUMENTS USING KMEANS

2030 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

SKLEARNMETRICSCCLUSTER CONTINGENCYMATRIX

SKLEARNMETRICSCCLUSTER CONTINGENCYMATRIX LABELSTRUE LABELSPRED EPSNONE

SPARSEFALSE

BUILD A CONTINGENCY MATRIX DESCRIBING THE RELATIONSHIP BETWEEN LABELS

PARAMETERS

LABELSTRUE INT ARRAY SHAPE NSAMPLES GROUND TRUTH CLASS LABELS TO BE USED AS A REFERENCE

LABELSPRED ARRAY SHAPE NSAMPLES CLUSTER LABELS TO EVALUATE

EPS NONE OR FLOAT OPTIONAL IF A FLOAT THAT VALUE IS ADDED TO ALL VALUES IN THE CONTINGENCY

MATRIX THIS HELPS TO STOP NAN PROPAGATION IF NONE NOTHING IS ADJUSTED

SPARSE BOOLEAN OPTIONAL IF TRUE RETURN A SPARSE CSR CONTINENCY MATRIX IF EPS IS NOT

NONE ANDSPARSE IS TRUE WILL THROW VALUEERROR

NEW IN VERSION 018

RETURNS

CONTINGENCY ARRAYLIKE SPARSE SHAPENCLASSESTRUE NCLASSESPRED MATRIX  $\frac{1}{n}$  SUCH THAT

$\frac{1}{n}$  IS THE NUMBER OF SAMPLES IN TRUE CLASS  $\frac{1}{n}$  AND IN PREDICTED CLASS  $\frac{1}{n}$  IFEPS IS NONE

THE DTYPE OF THIS ARRAY WILL BE INTEGER IF EPS IS GIVEN THE DTYPE WILL BE FLOAT WILL BE A

SCIPYSPARSECSRMATRIX IFSPARSETRUE

SKLEARNMETRICS FOWLKESMALLOWSSCORE

SKLEARNMETRICS FOWLKESMALLOWSSCORE LABELSTRUE LABELSPRED SPARSEFALSE

MEASURE THE SIMILARITY OF TWO CLUSTERINGS OF A SET OF POINTS

THE FOWLKESMALLOWS INDEX FMI IS DEFINED AS THE GEOMETRIC MEAN BETWEEN OF THE PRECISION AND RECALL

$FMI = \frac{TP}{\sqrt{TP * FP}}$

WHERE TP IS THE NUMBER OF TRUE POSITIVE IE THE NUMBER OF PAIR OF POINTS THAT BELONGS IN THE SAME CLUSTERS

IN BOTH LABELSTRUE AND LABELSPRED FP IS THE NUMBER OF FALSE POSITIVE IE THE NUMBER OF PAIR OF

POINTS THAT BELONGS IN THE SAME CLUSTERS IN LABELSTRUE AND NOT IN LABELSPRED AND FN IS THE NUMBER

OF FALSE NEGATIVE IE THE NUMBER OF PAIR OF POINTS THAT BELONGS IN THE SAME CLUSTERS IN LABELSPRED AND NOT

IN LABELSTRUE

THE SCORE RANGES FROM 0 TO 1 A HIGH VALUE INDICATES A GOOD SIMILARITY BETWEEN TWO CLUSTERS

READ MORE IN THE USER GUIDE

PARAMETERS

LABELSTRUE INT ARRAY SHAPE NSAMPLES A CLUSTERING OF THE DATA INTO DISJOINT SUBSETS

LABELSPRED ARRAY SHAPE NSAMPLES A CLUSTERING OF THE DATA INTO DISJOINT SUBSETS

SPARSE BOOL COMPUTE CONTINGENCY MATRIX INTERNALLY WITH SPARSE MATRIX

RETURNS

SCORE FLOAT THE RESULTING FOWLKESMALLOWS SCORE

REFERENCES

12

624 SKLEARNMETRICS METRICS 2031

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES

PERFECT LABELINGS ARE BOTH HOMOGENEOUS AND COMPLETE HENCE HAVE SCORE 10

FROM SKLEARNMETRICSCCLUSTER IMPORT FOWLKESMALLOWSSCORE

FOWLKESMALLOWSSCORE0 0 1 1 0 0 1 1

10

FOWLKESMALLOWSSCORE0 0 1 1 1 1 0 0

10

IF CLASSES MEMBERS ARE COMPLETELY SPLIT ACROSS DIFFERENT CLUSTERS THE ASSIGNMENT IS TOTALLY RANDOM HENCE THE

FMI IS NULL

FOWLKESMALLOWSSCORE0 0 0 0 0 1 2 3

00

SKLEARNMETRICS HOMOGENEITYCOMPLETENESSVMEASURE

SKLEARNMETRICS HOMOGENEITYCOMPLETENESSVMEASURE LABELSTRUE LABELSPRED

BETA10

COMPUTE THE HOMOGENEITY AND COMPLETENESS AND VMEASURE SCORES AT ONCE

THOSE METRICS ARE BASED ON NORMALIZED CONDITIONAL ENTROPY MEASURES OF THE CLUSTERING LABELING TO EVALUATE GIVEN

THE KNOWLEDGE OF A GROUND TRUTH CLASS LABELS OF THE SAME SAMPLES

A CLUSTERING RESULT SATISFIES HOMOGENEITY IF ALL OF ITS CLUSTERS CONTAIN ONLY DATA POINTS WHICH ARE MEMBERS OF A SINGLE CLASS

A CLUSTERING RESULT SATISFIES COMPLETENESS IF ALL THE DATA POINTS THAT ARE MEMBERS OF A GIVEN CLASS ARE ELEMENTS OF THE SAME CLUSTER

BOTH SCORES HAVE POSITIVE VALUES BETWEEN 00 AND 10 LARGER VALUES BEING DESIRABLE

THOSE 3 METRICS ARE INDEPENDENT OF THE ABSOLUTE VALUES OF THE LABELS A PERMUTATION OF THE CLASS OR CLUSTER LABEL

VALUES WON'T CHANGE THE SCORE VALUES IN ANY WAY

VMEASURE IS FURTHERMORE SYMMETRIC SWAPPING LABELSTRUE ANDLABELPRED WILL GIVE THE

SAME SCORE THIS DOES NOT HOLD FOR HOMOGENEITY AND COMPLETENESS VMEASURE IS IDENTICAL TO

NORMALIZEDMUTUALINFOSCORE WITH THE ARITHMETIC AVERAGING METHOD

READ MORE IN THE USER GUIDE

PARAMETERS

LABELSTRUE INT ARRAY SHAPE NSAMPLES GROUND TRUTH CLASS LABELS TO BE USED AS A REFERENCE

LABELSPRED ARRAY SHAPE NSAMPLES CLUSTER LABELS TO EVALUATE

BETA FLOAT RATIO OF WEIGHT ATTRIBUTED TO HOMOGENEITY VSCompleteness IFBETA IS

GREATER THAN 1 COMPLETENESS IS WEIGHTED MORE STRONGLY IN THE CALCULATION IF BETA IS

LESS THAN 1 HOMOGENEITY IS WEIGHTED MORE STRONGLY

RETURNS

HOMOGENEITY FLOAT SCORE BETWEEN 00 AND 10 10 STANDS FOR PERFECTLY HOMOGENEOUS LABELING

Completeness FLOAT SCORE BETWEEN 00 AND 10 10 STANDS FOR PERFECTLY COMPLETE LABELING

VMEASURE FLOAT HARMONIC MEAN OF THE FIRST TWO

SEE ALSO

2032 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

HOMOGENEITYSCORE

COMPLETENESSSCORE

VMEASURESORE

SKLEARNMETRICS HOMOGENEITYSCORE

SKLEARNMETRICS HOMOGENEITYSCORE LABELTRUE LABELSPRED

HOMOGENEITY METRIC OF A CLUSTER LABELING GIVEN A GROUND TRUTH

A CLUSTERING RESULT SATISFIES HOMOGENEITY IF ALL OF ITS CLUSTERS CONTAIN ONLY DATA POINTS WHICH ARE MEMBERS OF A SINGLE CLASS

THIS METRIC IS INDEPENDENT OF THE ABSOLUTE VALUES OF THE LABELS A PERMUTATION OF THE CLASS OR CLUSTER LABEL VALUES WON'T CHANGE THE SCORE VALUE IN ANY WAY

THIS METRIC IS NOT SYMMETRIC SWITCHING LABELTRUE WITHLABELPRED WILL RETURN THE COMPLETENESSSCORE WHICH WILL BE DIFFERENT IN GENERAL

READ MORE IN THE USER GUIDE

PARAMETERS

LABELTRUE INT ARRAY SHAPE NSAMPLES GROUND TRUTH CLASS LABELS TO BE USED AS A REFERENCE

LABELSPRED ARRAY SHAPE NSAMPLES CLUSTER LABELS TO EVALUATE

RETURNS

HOMOGENEITY FLOAT SCORE BETWEEN 00 AND 10 10 STANDS FOR PERFECTLY HOMOGENEOUS LABELING

SEE ALSO

COMPLETENESSSCORE

VMEASURESORE

REFERENCES

1

EXAMPLES

PERFECT LABELINGS ARE HOMOGENEOUS

FROM SKLEARNMETRICSCUSTER IMPORT HOMOGENEITYSCORE

HOMOGENEITYSCORE0 0 1 1 1 1 0 0

10

NONPERFECT LABELINGS THAT FURTHER SPLIT CLASSES INTO MORE CLUSTERS CAN BE PERFECTLY HOMOGENEOUS

PRINT6F HOMOGENEITYSCORE0 0 1 1 0 0 1 2

1000000

PRINT6F HOMOGENEITYSCORE0 0 1 1 0 1 2 3

1000000

624SKLEARNMETRICS METRICS 2033

SCIKITLEARN USER GUIDE RELEASE 0213  
CLUSTERS THAT INCLUDE SAMPLES FROM DIFFERENT CLASSES DO NOT MAKE FOR AN HOMOGENEOUS LABELING  
PRINT6F HOMOGENEITYSCORE0 0 1 1 0 1 0 1

00  
PRINT6F HOMOGENEITYSCORE0 0 1 1 0 0 0 0

00  
EXAMPLES USING SKLEARNMETRICSHOMOGENEITYSCORE  
•DEMO OF AFFINITY PROPAGATION CLUSTERING ALGORITHM  
•DEMO OF DBSCAN CLUSTERING ALGORITHM  
•A DEMO OF KMEANS CLUSTERING ON THE HANDWRITTEN DIGITS DATA  
•CLUSTERING TEXT DOCUMENTS USING KMEANS

SKLEARNMETRICS MUTUALINFOSCORE  
SKLEARNMETRICS MUTUALINFOSCORE LABELSTRUE LABELSPRED CONTINGENCYNONE  
MUTUAL INFORMATION BETWEEN TWO CLUSTERINGS  
THE MUTUAL INFORMATION IS A MEASURE OF THE SIMILARITY BETWEEN TWO LABELS OF THE SAME DATA WHERE  $n_i$  IS THE  
NUMBER OF THE SAMPLES IN CLUSTER  $i$  AND  $n_j$  IS THE NUMBER OF THE SAMPLES IN CLUSTER  $j$  THE MUTUAL INFORMATION  
BETWEEN CLUSTERINGS  $i$  AND  $j$  IS GIVEN AS

$$I(i, j) = \frac{1}{n} \sum_{k=1}^n \log \frac{n}{n_i n_j} n_{ij}$$

THIS METRIC IS INDEPENDENT OF THE ABSOLUTE VALUES OF THE LABELS A PERMUTATION OF THE CLASS OR CLUSTER LABEL VALUES  
WON'T CHANGE THE SCORE VALUE IN ANY WAY

THIS METRIC IS FURTHERMORE SYMMETRIC SWITCHING LABELSTRUE WITHLABELPRED WILL RETURN THE SAME SCORE  
VALUE THIS CAN BE USEFUL TO MEASURE THE AGREEMENT OF TWO INDEPENDENT LABEL ASSIGNMENTS STRATEGIES ON THE SAME  
DATASET WHEN THE REAL GROUND TRUTH IS NOT KNOWN

READ MORE IN THE USER GUIDE

PARAMETERS

LABELSTRUE INT ARRAY SHAPE NSAMPLES A CLUSTERING OF THE DATA INTO DISJOINT SUBSETS

LABELSPRED ARRAY SHAPE NSAMPLES A CLUSTERING OF THE DATA INTO DISJOINT SUBSETS

CONTINGENCY NONE ARRAY SPARSE MATRIX SHAPE NCLASSESTRUE NCLASSES PRED A CON  
TINGENCY MATRIX GIVEN BY THE CONTINGENCYMATRIX FUNCTION IF VALUE IS NONE IT WILL  
BE COMPUTED OTHERWISE THE GIVEN VALUE IS USED WITH LABELSTRUE ANDLABELSPRED

IGNORED

RETURNS

MIFLOAT MUTUAL INFORMATION A NONNEGATIVE VALUE

SEE ALSO

ADJUSTEDMUTUALINFOSCORE ADJUSTED AGAINST CHANCE MUTUAL INFORMATION

NORMALIZEDMUTUALINFOSCORE NORMALIZED MUTUAL INFORMATION

2034 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES USING SKLEARNMETRICSMUTUALINFOSCORE

- ADJUSTMENT FOR CHANCE IN CLUSTERING PERFORMANCE EVALUATION

SKLEARNMETRICS NORMALIZEDMUTUALINFOSCORE

SKLEARNMETRICS NORMALIZEDMUTUALINFOSCORE LABELSTRUE LABELSPRED AVER  
AGEMETHOD'WARN'

NORMALIZED MUTUAL INFORMATION BETWEEN TWO CLUSTERINGS

NORMALIZED MUTUAL INFORMATION NMI IS A NORMALIZATION OF THE MUTUAL INFORMATION MI SCORE TO SCALE THE RESULTS BETWEEN 0 NO MUTUAL INFORMATION AND 1 PERFECT CORRELATION IN THIS FUNCTION MUTUAL INFORMATION IS NORMALIZED BY SOME GENERALIZED MEAN OF HLABELSTRUE ANDHLABELSPRED DEFINED BY THE AVERAGEMETHOD

THIS MEASURE IS NOT ADJUSTED FOR CHANCE THEREFORE ADJUSTEDMUTUALINFOSCORE MIGHT BE PREFERRED

THIS METRIC IS INDEPENDENT OF THE ABSOLUTE VALUES OF THE LABELS A PERMUTATION OF THE CLASS OR CLUSTER LABEL VALUES WON'T CHANGE THE SCORE VALUE IN ANY WAY

THIS METRIC IS FURTHERMORE SYMMETRIC SWITCHING LABELTRUE WITHLABELPRED WILL RETURN THE SAME SCORE

VALUE THIS CAN BE USEFUL TO MEASURE THE AGREEMENT OF TWO INDEPENDENT LABEL ASSIGNMENTS STRATEGIES ON THE SAME DATASET WHEN THE REAL GROUND TRUTH IS NOT KNOWN

READ MORE IN THE USER GUIDE

PARAMETERS

LABELSTRUE INT ARRAY SHAPE NSAMPLES A CLUSTERING OF THE DATA INTO DISJOINT SUBSETS

LABELSPRED ARRAY SHAPE NSAMPLES A CLUSTERING OF THE DATA INTO DISJOINT SUBSETS

AVERAGEMETHOD STRING OPTIONAL DEFAULT 'WARN' HOW TO COMPUTE THE NORMALIZER IN THE DENOMINATOR POSSIBLE OPTIONS ARE 'MIN' 'GEOMETRIC' 'ARITHMETIC' AND 'MAX' IF 'WARN' 'GEOMETRIC' WILL BE USED THE DEFAULT WILL CHANGE TO 'ARITHMETIC' IN VERSION 022

NEW IN VERSION 020

RETURNS

NMI FLOAT SCORE BETWEEN 00 AND 10 10 STANDS FOR PERFECTLY COMPLETE LABELING

SEE ALSO

VMEASURESCORE VMEASURE NMI WITH ARITHMETIC MEAN OPTION

ADJUSTEDRANDSCORE ADJUSTED RAND INDEX

ADJUSTEDMUTUALINFOSCORE ADJUSTED MUTUAL INFORMATION ADJUSTED AGAINST CHANCE

EXAMPLES

PERFECT LABELINGS ARE BOTH HOMOGENEOUS AND COMPLETE HENCE HAVE SCORE 10

FROM SKLEARNMETRICSCLUSTR IMPORT NORMALIZEDMUTUALINFOSCORE

NORMALIZEDMUTUALINFOSCORE0 0 1 1 0 0 1 1

10

NORMALIZEDMUTUALINFOSCORE0 0 1 1 1 1 0 0

624SKLEARNMETRICS METRICS 2035

10  
IF CLASSES MEMBERS ARE COMPLETELY SPLIT ACROSS DIFFERENT CLUSTERS THE ASSIGNMENT IS TOTALLY INCOMPLETE HENCE  
THE NMI IS NULL  
NORMALIZEDMUTUALINFOSCORE0 0 0 0 0 1 2 3

00  
SKLEARNMETRICS SILHOUETTESCORE  
SKLEARNMETRICS SILHOUETTESCORE XLABELS METRIC'EUCLIDEAN' SAMPLESIZENONE RAN  
DOMSTATENONE KWDS  
COMPUTE THE MEAN SILHOUETTE COEFFICIENT OF ALL SAMPLES  
THE SILHOUETTE COEFFICIENT IS CALCULATED USING THE MEAN INTRACLUSTER DISTANCE A AND THE MEAN NEARESTCLUSTER  
DISTANCE B FOR EACH SAMPLE THE SILHOUETTE COEFFICIENT FOR A SAMPLE IS  $B - A / \max(A, B)$  TO CLARIFY  
BIS THE DISTANCE BETWEEN A SAMPLE AND THE NEAREST CLUSTER THAT THE SAMPLE IS NOT A PART OF NOTE THAT SILHOUETTE  
COEFFICIENT IS ONLY DEFINED IF NUMBER OF LABELS IS 2 NLABELS NSAMPLES 1  
THIS FUNCTION RETURNS THE MEAN SILHOUETTE COEFFICIENT OVER ALL SAMPLES TO OBTAIN THE VALUES FOR EACH SAMPLE  
USESILHOUETTESAMPLES  
THE BEST VALUE IS 1 AND THE WORST VALUE IS 1 VALUES NEAR 0 INDICATE OVERLAPPING CLUSTERS NEGATIVE VALUES  
GENERALLY INDICATE THAT A SAMPLE HAS BEEN ASSIGNED TO THE WRONG CLUSTER AS A DIFFERENT CLUSTER IS MORE SIMILAR  
READ MORE IN THE USER GUIDE  
PARAMETERS  
XARRAY NSAMPLESA NSAMPLESA IF METRIC "PRECOMPUTED" OR NSAMPLESA  
NFEATURES OTHERWISE ARRAY OF PAIRWISE DISTANCES BETWEEN SAMPLES OR A FEATURE ARRAY  
LABELS ARRAY SHAPE NSAMPLES PREDICTED LABELS FOR EACH SAMPLE  
METRIC STRING OR CALLABLE THE METRIC TO USE WHEN CALCULATING DISTANCE BETWEEN INSTANCES  
IN A FEATURE ARRAY IF METRIC IS A STRING IT MUST BE ONE OF THE OPTIONS ALLOWED BY  
METRICSPAIRWISEPAIRWISEDISTANCES IF X IS THE DISTANCE ARRAY ITSELF USE  
METRICPRECOMPUTED  
SAMPLESIZE INT OR NONE THE SIZE OF THE SAMPLE TO USE WHEN COMPUTING THE SILHOUETTE CO  
EFFICIENT ON A RANDOM SUBSET OF THE DATA IF SAMPLESIZE IS NONE NO SAMPLING IS  
USED  
RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE THE GENERATOR  
USED TO RANDOMLY SELECT A SUBSET OF SAMPLES IF INT RANDOMSTATE IS THE SEED USED BY THE  
RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS THE RANDOM NUMBER  
GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NP  
RANDOM USED WHEN SAMPLESIZE IS NOT NONE  
KWDS OPTIONAL KEYWORD PARAMETERS ANY FURTHER PARAMETERS ARE PASSED DIRECTLY TO THE DIS  
TANCE FUNCTION IF USING A SCIPYSPATIALDISTANCE METRIC THE PARAMETERS ARE STILL METRIC DE  
PENDENT SEE THE SCIPY DOCS FOR USAGE EXAMPLES  
RETURNS  
SILHOUETTE FLOAT MEAN SILHOUETTE COEFFICIENT FOR ALL SAMPLES  
2036 CHAPTER 6 API REFERENCE

REFERENCES

12

EXAMPLES USING SKLEARNMETRICSSILHOUETTESCORE

- DEMO OF AFFINITY PROPAGATION CLUSTERING ALGORITHM
- DEMO OF DBSCAN CLUSTERING ALGORITHM
- A DEMO OF KMEANS CLUSTERING ON THE HANDWRITTEN DIGITS DATA
- SELECTING THE NUMBER OF CLUSTERS WITH SILHOUETTE ANALYSIS ON KMEANS CLUSTERING
- CLUSTERING TEXT DOCUMENTS USING KMEANS

SKLEARNMETRICS SILHOUETTESAMPLES

SKLEARNMETRICS SILHOUETTESAMPLES XLABELS METRIC'EUCLIDEAN' KWDS

COMPUTE THE SILHOUETTE COEFFICIENT FOR EACH SAMPLE

THE SILHOUETTE COEFFICIENT IS A MEASURE OF HOW WELL SAMPLES ARE CLUSTERED WITH SAMPLES THAT ARE SIMILAR TO THEMSELVES CLUSTERING MODELS WITH A HIGH SILHOUETTE COEFFICIENT ARE SAID TO BE DENSE WHERE SAMPLES IN THE SAME CLUSTER ARE SIMILAR TO EACH OTHER AND WELL SEPARATED WHERE SAMPLES IN DIFFERENT CLUSTERS ARE NOT VERY SIMILAR TO EACH OTHER

THE SILHOUETTE COEFFICIENT IS CALCULATED USING THE MEAN INTRACLUSTER DISTANCE  $A$  AND THE MEAN NEARESTCLUSTER

DISTANCE  $B$  FOR EACH SAMPLE THE SILHOUETTE COEFFICIENT FOR A SAMPLE IS  $B - A / \max(A, B)$  NOTE THAT

SILHOUETTE COEFFICIENT IS ONLY DEFINED IF NUMBER OF LABELS IS 2  $N_{LABELS} - N_{SAMPLES} - 1$

THIS FUNCTION RETURNS THE SILHOUETTE COEFFICIENT FOR EACH SAMPLE

THE BEST VALUE IS 1 AND THE WORST VALUE IS 1 VALUES NEAR 0 INDICATE OVERLAPPING CLUSTERS

READ MORE IN THE USER GUIDE

PARAMETERS

XARRAY  $N_{SAMPLES}$  A  $N_{SAMPLES}$  IF METRIC "PRECOMPUTED" OR  $N_{SAMPLES}$  A

$N_{FEATURES}$  OTHERWISE ARRAY OF PAIRWISE DISTANCES BETWEEN SAMPLES OR A FEATURE ARRAY

LABELS ARRAY SHAPE  $N_{SAMPLES}$  LABEL VALUES FOR EACH SAMPLE

METRIC STRING OR CALLABLE THE METRIC TO USE WHEN CALCULATING DISTANCE BETWEEN INSTANCES IN

A FEATURE ARRAY IF METRIC IS A STRING IT MUST BE ONE OF THE OPTIONS ALLOWED BY SKLEARN

METRICSPAIRWISEPAIRWISEDISTANCES IF X IS THE DISTANCE ARRAY ITSELF USE

"PRECOMPUTED" AS THE METRIC

'KWDS' OPTIONAL KEYWORD PARAMETERS ANY FURTHER PARAMETERS ARE PASSED DIRECTLY TO THE DIS

TANCE FUNCTION IF USING A SCIPYSPATIALDISTANCE METRIC THE PARAMETERS ARE STILL

METRIC DEPENDENT SEE THE SCIPY DOCS FOR USAGE EXAMPLES

RETURNS

SILHOUETTE ARRAY SHAPE  $N_{SAMPLES}$  SILHOUETTE COEFFICIENT FOR EACH SAMPLES

REFERENCES

12

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES USING SKLEARNMETRICSSILHOUETTESAMPLES

- SELECTING THE NUMBER OF CLUSTERS WITH SILHOUETTE ANALYSIS ON KMEANS CLUSTERING

SKLEARNMETRICS VMEASURESCORE

SKLEARNMETRICS VMEASURESCORE LABELSTRUE LABELSPRED BETA10

VMEASURE CLUSTER LABELING GIVEN A GROUND TRUTH

THIS SCORE IS IDENTICAL TO NORMALIZEDMUTUALINFOSCORE WITH THEARITHMETIC OPTION FOR AVERAGING

THE VMEASURE IS THE HARMONIC MEAN BETWEEN HOMOGENEITY AND COMPLETENESS

V 1 BETA HOMOGENEITY COMPLETENESS

BETA HOMOGENEITY COMPLETENESS

THIS METRIC IS INDEPENDENT OF THE ABSOLUTE VALUES OF THE LABELS A PERMUTATION OF THE CLASS OR CLUSTER LABEL VALUES WON'T CHANGE THE SCORE VALUE IN ANY WAY

THIS METRIC IS FURTHERMORE SYMMETRIC SWITCHING LABELTRUE WITHLABELPRED WILL RETURN THE SAME SCORE

VALUE THIS CAN BE USEFUL TO MEASURE THE AGREEMENT OF TWO INDEPENDENT LABEL ASSIGNMENTS STRATEGIES ON THE SAME DATASET WHEN THE REAL GROUND TRUTH IS NOT KNOWN

READ MORE IN THE USER GUIDE

PARAMETERS

LABELSTRUE INT ARRAY SHAPE NSAMPLES GROUND TRUTH CLASS LABELS TO BE USED AS A REFERENCE

LABELSPRED ARRAY SHAPE NSAMPLES CLUSTER LABELS TO EVALUATE

BETA FLOAT RATIO OF WEIGHT ATTRIBUTED TO HOMOGENEITY VSCompleteness IFBETA IS GREATER THAN 1 COMPLETENESS IS WEIGHTED MORE STRONGLY IN THE CALCULATION IF BETA IS LESS THAN 1 HOMOGENEITY IS WEIGHTED MORE STRONGLY

RETURNS

VMEASURE FLOAT SCORE BETWEEN 00 AND 10 10 STANDS FOR PERFECTLY COMPLETE LABELING

SEE ALSO

HOMOGENEITYSCORE

COMPLETENESSSCORE

NORMALIZEDMUTUALINFOSCORE

REFERENCES

1

EXAMPLES

PERFECT LABELINGS ARE BOTH HOMOGENEOUS AND COMPLETE HENCE HAVE SCORE 10

2038 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213  
 FROM SKLEARNMETRICSCUSTER IMPORT VMEASURESORE  
 VMEASURESORE0 0 1 1 0 0 1 1  
 10  
 VMEASURESORE0 0 1 1 1 1 0 0  
 10  
 LABELINGS THAT ASSIGN ALL CLASSES MEMBERS TO THE SAME CLUSTERS ARE COMPLETE BE NOT HOMOGENEOUS HENCE PENAL  
 IZED  
 PRINT6F VMEASURESORE0 0 1 2 0 0 1 1  
 08  
 PRINT6F VMEASURESORE0 1 2 3 0 0 1 1  
 066  
 LABELINGS THAT HAVE PURE CLUSTERS WITH MEMBERS COMING FROM THE SAME CLASSES ARE HOMOGENEOUS BUT UN  
 NECESSARY SPLITS HARMS COMPLETENESS AND THUS PENALIZE VMEASURE AS WELL  
 PRINT6F VMEASURESORE0 0 1 1 0 0 1 2  
 08  
 PRINT6F VMEASURESORE0 0 1 1 0 1 2 3  
 066  
 IF CLASSES MEMBERS ARE COMPLETELY SPLIT ACROSS DIFFERENT CLUSTERS THE ASSIGNMENT IS TOTALLY INCOMPLETE HENCE THE  
 VMEASURE IS NULL  
 PRINT6F VMEASURESORE0 0 0 0 0 1 2 3  
 00  
 CLUSTERS THAT INCLUDE SAMPLES FROM TOTALLY DIFFERENT CLASSES TOTALLY DESTROY THE HOMOGENEITY OF THE LABELING  
 HENCE  
 PRINT6F VMEASURESORE0 0 1 1 0 0 0 0  
 00  
 EXAMPLES USING SKLEARNMETRICSVMEASURESORE  
 •BICCLUSTERING DOCUMENTS WITH THE SPECTRAL COCLUSTERING ALGORITHM  
 •DEMO OF AFFINITY PROPAGATION CLUSTERING ALGORITHM  
 •DEMO OF DBSCAN CLUSTERING ALGORITHM  
 •ADJUSTMENT FOR CHANCE IN CLUSTERING PERFORMANCE EVALUATION  
 •A DEMO OF KMEANS CLUSTERING ON THE HANDWRITTEN DIGITS DATA  
 •CLUSTERING TEXT DOCUMENTS USING KMEANS  
 6246 BICCLUSTERING METRICS  
 SEE THE BICCLUSTERING EVALUATION SECTION OF THE USER GUIDE FOR FURTHER DETAILS  
 624SKLEARNMETRICS METRICS 2039

SCIKITLEARN USER GUIDE RELEASE 0213

METRICSCONSSENSUSSCORE A B SIMILARITY THE SIMILARITY OF TWO SETS OF BICLUSTERS

SKLEARNMETRICS CONSENSUSSCORE

SKLEARNMETRICS CONSENSUSSCORE ABSIMILARITY'JACCARD'

THE SIMILARITY OF TWO SETS OF BICLUSTERS

SIMILARITY BETWEEN INDIVIDUAL BICLUSTERS IS COMPUTED THEN THE BEST MATCHING BETWEEN SETS IS FOUND USING THE HUNGARIAN ALGORITHM THE FINAL SCORE IS THE SUM OF SIMILARITIES DIVIDED BY THE SIZE OF THE LARGER SET

READ MORE IN THE USER GUIDE

PARAMETERS

AROWS COLUMNS TUPLE OF ROW AND COLUMN INDICATORS FOR A SET OF BICLUSTERS

BROWS COLUMNS ANOTHER SET OF BICLUSTERS LIKE A

SIMILARITY STRING OR FUNCTION OPTIONAL DEFAULT "JACCARD" MAY BE THE STRING "JACCARD" TO USE THE JACCARD COEFFICIENT OR ANY FUNCTION THAT TAKES FOUR ARGUMENTS EACH OF WHICH IS A 1D INDICATOR VECTOR AROWS ACOLUMNS BROWS BCOLUMNS

REFERENCES

- HOCHREITER BODENHOFER ET AL 2010 FABIA FACTOR ANALYSIS FOR BICLUSTER ACQUISITION

EXAMPLES USING SKLEARNMETRICSCONSSENSUSSCORE

- A DEMO OF THE SPECTRAL COCLUSTERING ALGORITHM
- A DEMO OF THE SPECTRAL BICLUSTERING ALGORITHM

6247 PAIRWISE METRICS

SEE THE PAIRWISE METRICS AFFINITIES AND KERNELS SECTION OF THE USER GUIDE FOR FURTHER DETAILS

METRICSPAIRWISEADDITIVECHI2KERNEL X

YCOMPUTES THE ADDITIVE CHISQUARED KERNEL BETWEEN OBSERVATIONS IN X AND Y

METRICSPAIRWISECHI2KERNEL X Y GAMMA COMPUTES THE EXPONENTIAL CHISQUARED KERNEL X AND Y

METRICSPAIRWISECOSINESIMILARITY X Y

COMPUTE COSINE SIMILARITY BETWEEN SAMPLES IN X AND Y

METRICSPAIRWISECOSINEDISTANCES X Y COMPUTE COSINE DISTANCE BETWEEN SAMPLES IN X AND Y

METRICSPAIRWISEDISTANCEMETRICS VALID METRICS FOR PAIRWISEDISTANCES

METRICSPAIRWISEEUCLIDEANDISTANCES X

Y CONSIDERING THE ROWS OF X AND YX AS VECTORS COMPUTE THE DISTANCE MATRIX BETWEEN EACH PAIR OF VECTORS

METRICSPAIRWISEHAVERSINEDISTANCES X

YCOMPUTE THE HAVERSINE DISTANCE BETWEEN SAMPLES IN X AND Y

METRICSPAIRWISEKERNELMETRICS VALID METRICS FOR PAIRWISEKERNELS

METRICSPAIRWISELAPLACIANKERNEL X Y

GAMMA COMPUTE THE LAPLACIAN KERNEL BETWEEN X AND Y

METRICSPAIRWISELINEARKERNEL X Y COMPUTE THE LINEAR KERNEL BETWEEN X AND Y

CONTINUED ON NEXT PAGE

2040 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

TABLE 6184 – CONTINUED FROM PREVIOUS PAGE

METRICSPAIRWISEMANHATTANDISTANCES X  
Y COMPUTE THE L1 DISTANCES BETWEEN THE VECTORS IN X AND Y

METRICSPAIRWISEPAIRWISEKERNELS X Y  
COMPUTE THE KERNEL BETWEEN ARRAYS X AND OPTIONAL ARRAY Y

METRICSPAIRWISEPOLYNOMIALKERNEL X Y  
COMPUTE THE POLYNOMIAL KERNEL BETWEEN X AND Y

METRICSPAIRWISERBFKERNEL X Y GAMMA COMPUTE THE RBF GAUSSIAN KERNEL BETWEEN X AND Y

METRICSPAIRWISESIGMOIDKERNEL X Y COMPUTE THE SIGMOID KERNEL BETWEEN X AND Y

METRICSPAIRWISEPAIREDDEUCLIDEANDISTANCES X  
YCOMPUTES THE PAIRED EUCLIDEAN DISTANCES BETWEEN X AND Y

METRICSPAIRWISEPAIREDMANHATTANDISTANCES X  
YCOMPUTE THE L1 DISTANCES BETWEEN THE VECTORS IN X AND Y

METRICSPAIRWISEPAIREDCOSINEDISTANCES X  
YCOMPUTES THE PAIRED COSINE DISTANCES BETWEEN X AND Y

METRICSPAIRWISEPAIREDDISTANCES X Y  
METRICCOMPUTES THE PAIRED DISTANCES BETWEEN X AND Y

METRICSPAIRWISEDISTANCES X Y METRIC COMPUTE THE DISTANCE MATRIX FROM A VECTOR ARRAY X AND OP  
TIONAL Y

METRICSPAIRWISEDISTANCESARGMIN X Y  
COMPUTE MINIMUM DISTANCES BETWEEN ONE POINT AND A SET  
OF POINTS

METRICSPAIRWISEDISTANCESARGMINMIN X  
YCOMPUTE MINIMUM DISTANCES BETWEEN ONE POINT AND A SET  
OF POINTS

METRICSPAIRWISEDISTANCESCHUNKED X Y  
GENERATE A DISTANCE MATRIX CHUNK BY CHUNK WITH OPTIONAL  
REDUCTION

SKLEARNMETRICSPAIRWISE ADDITIVECHI2KERNEL  
SKLEARNMETRICSPAIRWISE ADDITIVECHI2KERNEL XYNONE  
COMPUTES THE ADDITIVE CHISQUARED KERNEL BETWEEN OBSERVATIONS IN X AND Y  
THE CHISQUARED KERNEL IS COMPUTED BETWEEN EACH PAIR OF ROWS IN X AND Y X AND Y HAVE TO BE NONNEGATIVE  
THIS KERNEL IS MOST COMMONLY APPLIED TO HISTOGRAMS  
THE CHISQUARED KERNEL IS GIVEN BY  
$$K_{X,Y} = \sum (X - Y)^2$$
  
IT CAN BE INTERPRETED AS A WEIGHTED DIFFERENCE PER ENTRY  
READ MORE IN THE USER GUIDE

PARAMETERS  
XARRAYLIKE OF SHAPE NSAMPLESX NFEATURES  
YARRAY OF SHAPE NSAMPLESY NFEATURES  
RETURNS  
KERNELMATRIX ARRAY OF SHAPE NSAMPLESX NSAMPLESY

SEE ALSO  
CHI2KERNEL THE EXPONENTIATED VERSION OF THE KERNEL WHICH IS USUALLY PREFERABLE  
SKLEARNKERNELAPPROXIMATIONADDITIVECHI2SAMPLER A FOURIER APPROXIMATION TO THIS KER  
NEL

624SKLEARNMETRICS METRICS 2041

SCIKITLEARN USER GUIDE RELEASE 0213

NOTES

AS THE NEGATIVE OF A DISTANCE THIS KERNEL IS ONLY CONDITIONALLY POSITIVE DEFINITE

REFERENCES

• ZHANG J AND MARSZALEK M AND LAZEBNIK S AND SCHMID C LOCAL FEATURES AND KERNELS FOR CLASSIFICATION OF TEXTURE AND OBJECT CATEGORIES A COMPREHENSIVE STUDY INTERNATIONAL JOURNAL OF COMPUTER VISION 2007 [HTTPSRESEARCHMICROSOFTCOMENUSUMPEOPLEMANIKPROJECTSTRADEOFFPAPERSZHANGIJCV06PDF](https://research.microsoft.com/en-us/projects/stradeoff/papers/zhangijcv06.pdf)

SKLEARNMETRICSPAIRWISE CHI2KERNEL

SKLEARNMETRICSPAIRWISE CHI2KERNEL XYNONE GAMMA10

COMPUTES THE EXPONENTIAL CHISQUARED KERNEL X AND Y

THE CHISQUARED KERNEL IS COMPUTED BETWEEN EACH PAIR OF ROWS IN X AND Y X AND Y HAVE TO BE NONNEGATIVE

THIS KERNEL IS MOST COMMONLY APPLIED TO HISTOGRAMS

THE CHISQUARED KERNEL IS GIVEN BY

$$K_{X,Y} = \exp(\gamma \sum_i x_i^2 - x_i y_i)$$

IT CAN BE INTERPRETED AS A WEIGHTED DIFFERENCE PER ENTRY

READ MORE IN THE USER GUIDE

PARAMETERS

XARRAYLIKE OF SHAPE NSAMPLESX NFEATURES

YARRAY OF SHAPE NSAMPLESY NFEATURES

GAMMA FLOAT DEFAULT1 SCALING PARAMETER OF THE CHI2 KERNEL

RETURNS

KERNELMATRIX ARRAY OF SHAPE NSAMPLESX NSAMPLESY

SEE ALSO

ADDITIVECHI2KERNEL THE ADDITIVE VERSION OF THIS KERNEL

SKLEARNKERNELAPPROXIMATIONADDITIVECHI2SAMPLER A FOURIER APPROXIMATION TO THE ADDI

TIVE VERSION OF THIS KERNEL

REFERENCES

• ZHANG J AND MARSZALEK M AND LAZEBNIK S AND SCHMID C LOCAL FEATURES AND KERNELS FOR CLASSIFICATION OF TEXTURE AND OBJECT CATEGORIES A COMPREHENSIVE STUDY INTERNATIONAL JOURNAL OF COMPUTER VISION 2007 [HTTPSRESEARCHMICROSOFTCOMENUSUMPEOPLEMANIKPROJECTSTRADEOFFPAPERSZHANGIJCV06PDF](https://research.microsoft.com/en-us/projects/stradeoff/papers/zhangijcv06.pdf)

SKLEARNMETRICSPAIRWISE COSINESIMILARITY

SKLEARNMETRICSPAIRWISE COSINESIMILARITY XYNONE DENSEOUTPUTTRUE

COMPUTE COSINE SIMILARITY BETWEEN SAMPLES IN X AND Y

COSINE SIMILARITY OR THE COSINE KERNEL COMPUTES SIMILARITY AS THE NORMALIZED DOT PRODUCT OF X AND Y

2042 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

KX Y X Y XY

ON L2NORMALIZED DATA THIS FUNCTION IS EQUIVALENT TO LINEARKERNEL

READ MORE IN THE USER GUIDE

PARAMETERS

XNDARRAY OR SPARSE ARRAY SHAPE NSAMPLESX NFEATURES INPUT DATA

YNDARRAY OR SPARSE ARRAY SHAPE NSAMPLESY NFEATURES INPUT DATA IF NONE THE OUTPUT

WILL BE THE PAIRWISE SIMILARITIES BETWEEN ALL SAMPLES IN X

DENSEOUTPUT BOOLEAN OPTIONAL DEFAULT TRUE WHETHER TO RETURN DENSE OUTPUT EVEN WHEN THE

INPUT IS SPARSE IF FALSE THE OUTPUT IS SPARSE IF BOTH INPUT ARRAYS ARE SPARSE

NEW IN VERSION 017 PARAMETER DENSEOUTPUT FOR DENSE OUTPUT

RETURNS

KERNEL MATRIX ARRAY AN ARRAY WITH SHAPE NSAMPLESX NSAMPLESY

SKLEARNMETRICSPAIRWISE COSINEDISTANCES

SKLEARNMETRICSPAIRWISE COSINEDISTANCES XYNONE

COMPUTE COSINE DISTANCE BETWEEN SAMPLES IN X AND Y

COSINE DISTANCE IS DEFINED AS 10 MINUS THE COSINE SIMILARITY

READ MORE IN THE USER GUIDE

PARAMETERS

XARRAYLIKE SPARSE MATRIX WITH SHAPE NSAMPLESX NFEATURES

YARRAYLIKE SPARSE MATRIX OPTIONAL WITH SHAPE NSAMPLESY NFEATURES

RETURNS

DISTANCE MATRIX ARRAY AN ARRAY WITH SHAPE NSAMPLESX NSAMPLESY

SEE ALSO

SKLEARNMETRICSPAIRWISECOSINESIMILARITY

SCIPYSPATIALDISTANCECOSINE DENSE MATRICES ONLY

SKLEARNMETRICSPAIRWISE DISTANCEMETRICS

SKLEARNMETRICSPAIRWISE DISTANCEMETRICS

VALID METRICS FOR PAIRWISEDISTANCES

THIS FUNCTION SIMPLY RETURNS THE VALID PAIRWISE DISTANCE METRICS IT EXISTS TO ALLOW FOR A DESCRIPTION OF THE

MAPPING FOR EACH OF THE VALID STRINGS

THE VALID DISTANCE METRICS AND THE FUNCTION THEY MAP TO ARE

624SKLEARNMETRICS METRICS 2043

SCIKITLEARN USER GUIDE RELEASE 0213

METRIC FUNCTION

‘CITYBLOCK’ METRICSPAIRWISEMANHATTANDISTANCES

‘COSINE’ METRICSPAIRWISECOSINEDISTANCES

‘EUCLIDEAN’ METRICSPAIRWISEEUCLIDEANDISTANCES

‘HAVERSINE’ METRICSPAIRWISEHAVERSINEDISTANCES

‘L1’ METRICSPAIRWISEMANHATTANDISTANCES

‘L2’ METRICSPAIRWISEEUCLIDEANDISTANCES

‘MANHATTAN’ METRICSPAIRWISEMANHATTANDISTANCES

READ MORE IN THE USER GUIDE

SKLEARNMETRICSPAIRWISE EUCLIDEANDISTANCES

SKLEARNMETRICSPAIRWISE EUCLIDEANDISTANCES XYNONE YNORMSQUAREDNONE

SQUAREDFALSE XNORMSQUAREDNONE

CONSIDERING THE ROWS OF X AND YX AS VECTORS COMPUTE THE DISTANCE MATRIX BETWEEN EACH PAIR OF VECTORS

FOR EFFICIENCY REASONS THE EUCLIDEAN DISTANCE BETWEEN A PAIR OF ROW VECTOR X AND Y IS COMPUTED AS

$\sqrt{\sum (x_i - y_i)^2}$

THIS FORMULATION HAS TWO ADVANTAGES OVER OTHER WAYS OF COMPUTING DISTANCES FIRST IT IS COMPUTATIONALLY EF

FICIENT WHEN DEALING WITH SPARSE DATA SECOND IF ONE ARGUMENT VARIES BUT THE OTHER REMAINS UNCHANGED THEN

$\sum (x_i - y_i)^2$  CAN BE PRECOMPUTED

HOWEVER THIS IS NOT THE MOST PRECISE WAY OF DOING THIS COMPUTATION AND THE DISTANCE MATRIX RETURNED BY THIS

FUNCTION MAY NOT BE EXACTLY SYMMETRIC AS REQUIRED BY EG SCIPYSPATIALDISTANCE FUNCTIONS

READ MORE IN THE USER GUIDE

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES1 NFEATURES

YARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES2 NFEATURES

YNORMSQUARED ARRAYLIKE SHAPE NSAMPLES2 OPTIONAL PRECOMPUTED DOTPRODUCTS OF

VECTORS IN Y EG  $\sum y_i^2$  MAY BE IGNORED IN SOME CASES SEE THE NOTE

BELOW

SQUARED BOOLEAN OPTIONAL RETURN SQUARED EUCLIDEAN DISTANCES

XNORMSQUARED ARRAYLIKE SHAPE NSAMPLES1 OPTIONAL PRECOMPUTED DOTPRODUCTS OF

VECTORS IN X EG  $\sum x_i^2$  MAY BE IGNORED IN SOME CASES SEE THE NOTE

BELOW

RETURNS

DISTANCES ARRAY SHAPE NSAMPLES1 NSAMPLES2

SEE ALSO

PAIREDDISTANCES DISTANCES BETWEEN PAIRS OF ELEMENTS OF X AND Y

NOTES

TO ACHIEVE BETTER ACCURACY XNORMSQUARED ANDYNORMSQUARED MAY BE UNUSED IF THEY ARE PASSED AS

FLOAT32

2044 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES

FROM SKLEARNMETRICSPAIRWISE IMPORT EUCLIDEANDISTANCES

X 0 1 1 1

DISTANCE BETWEEN ROWS OF X

EUCLIDEANDISTANCESX X

ARRAY0 1

1 0

GET DISTANCE TO ORIGIN

EUCLIDEANDISTANCESX 0 0

ARRAY1

141421356

SKLEARNMETRICSPAIRWISE HAVERSINEDISTANCES

SKLEARNMETRICSPAIRWISE HAVERSINEDISTANCES XYNONE

COMPUTE THE HAVERSINE DISTANCE BETWEEN SAMPLES IN X AND Y

THE HAVERSINE OR GREAT CIRCLE DISTANCE IS THE ANGULAR DISTANCE BETWEEN TWO POINTS ON THE SURFACE OF A SPHERE

THE FIRST DISTANCE OF EACH POINT IS ASSUMED TO BE THE LATITUDE THE SECOND IS THE LONGITUDE GIVEN IN RADIANS THE

DIMENSION OF THE DATA MUST BE 2

$$\sqrt{2 \arcsin \sqrt{\sin^2 \frac{\phi_1 - \phi_2}{2} + \cos \phi_1 \cos \phi_2 \sin^2 \frac{\lambda_1 - \lambda_2}{2}}}$$

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES1 2

YARRAYLIKE SHAPE NSAMPLES2 2 OPTIONAL

RETURNS

DISTANCE ARRAY SHAPE NSAMPLES1 NSAMPLES2

NOTES

AS THE EARTH IS NEARLY SPHERICAL THE HAVERSINE FORMULA PROVIDES A GOOD APPROXIMATION OF THE DISTANCE BETWEEN

TWO POINTS OF THE EARTH SURFACE WITH A LESS THAN 1 ERROR ON AVERAGE

EXAMPLES

WE WANT TO CALCULATE THE DISTANCE BETWEEN THE EZEIZA AIRPORT BUENOS AIRES ARGENTINA AND THE CHARLES DE

GAULLE AIRPORT PARIS FRANCE

FROM SKLEARNMETRICSPAIRWISE IMPORT HAVERSINEDISTANCES

BSAS 34833333 585166646

PARIS 490083899664 253844117956

RESULT HAVERSINEDISTANCESBSAS PARIS

RESULT63710001000 MULTIPLY BY EARTH RADIUS TO GET KILOMETERS

ARRAY 0 1127945379464

1127945379464 0

624SKLEARNMETRICS METRICS 2045

SCIKITLEARN USER GUIDE RELEASE 0213  
SKLEARNMETRICSPAIRWISE KERNELMETRICS  
SKLEARNMETRICSPAIRWISE KERNELMETRICS  
VALID METRICS FOR PAIRWISEKERNELS  
THIS FUNCTION SIMPLY RETURNS THE VALID PAIRWISE DISTANCE METRICS IT EXISTS HOWEVER TO ALLOW FOR A VERBOSE DESCRIPTION OF THE MAPPING FOR EACH OF THE VALID STRINGS  
THE VALID DISTANCE METRICS AND THE FUNCTION THEY MAP TO ARE  
METRIC FUNCTION  
'ADDITIVECHI2' SKLEARNPAIRWISEADDITIVECHI2KERNEL  
'CHI2' SKLEARNPAIRWISECHI2KERNEL  
'LINEAR' SKLEARNPAIRWISELINEARKERNEL  
'POLY' SKLEARNPAIRWISEPOLYNOMIALKERNEL  
'POLYNOMIAL' SKLEARNPAIRWISEPOLYNOMIALKERNEL  
'RBF' SKLEARNPAIRWISERBFBFKERNEL  
'LAPLACIAN' SKLEARNPAIRWISELAPLACIANKERNEL  
'SIGMOID' SKLEARNPAIRWISESIGMOIDKERNEL  
'COSINE' SKLEARNPAIRWISECOSINESIMILARITY  
READ MORE IN THE USER GUIDE  
SKLEARNMETRICSPAIRWISE LAPLACIANKERNEL  
SKLEARNMETRICSPAIRWISE LAPLACIANKERNEL XYNONE GAMMANONE  
COMPUTE THE LAPLACIAN KERNEL BETWEEN X AND Y  
THE LAPLACIAN KERNEL IS DEFINED AS  
$$K(X, Y) = \exp(-\gamma \|X - Y\|_1)$$
  
FOR EACH PAIR OF ROWS X IN X AND Y IN Y READ MORE IN THE USER GUIDE  
NEW IN VERSION 017  
PARAMETERS  
XARRAY OF SHAPE NSAMPLESX NFEATURES  
YARRAY OF SHAPE NSAMPLESY NFEATURES  
GAMMA FLOAT DEFAULT NONE IF NONE DEFAULTS TO 10 NFEATURES  
RETURNS  
KERNELMATRIX ARRAY OF SHAPE NSAMPLESX NSAMPLESY  
SKLEARNMETRICSPAIRWISE LINEARKERNEL  
SKLEARNMETRICSPAIRWISE LINEARKERNEL XYNONE DENSEOUTPUTTRUE  
COMPUTE THE LINEAR KERNEL BETWEEN X AND Y  
READ MORE IN THE USER GUIDE  
PARAMETERS  
XARRAY OF SHAPE NSAMPLES1 NFEATURES

2046 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

YARRAY OF SHAPE NSAMPLES2 NFEATURES

DENSEOUTPUT BOOLEAN OPTIONAL DEFAULT TRUE WHETHER TO RETURN DENSE OUTPUT EVEN WHEN THE INPUT IS SPARSE IF FALSE THE OUTPUT IS SPARSE IF BOTH INPUT ARRAYS ARE SPARSE

NEW IN VERSION 020

RETURNS

GRAM MATRIX ARRAY OF SHAPE NSAMPLES1 NSAMPLES2

SKLEARNMETRICSPAIRWISE MANHATTANDISTANCES

SKLEARNMETRICSPAIRWISE MANHATTANDISTANCES XYNONE SUMOVERFEATURESTRUE

COMPUTE THE L1 DISTANCES BETWEEN THE VECTORS IN X AND Y

WITH SUMOVERFEATURES EQUAL TO FALSE IT RETURNS THE COMPONENTWISE DISTANCES

READ MORE IN THE USER GUIDE

PARAMETERS

XARRAYLIKE AN ARRAY WITH SHAPE NSAMPLESX NFEATURES

YARRAYLIKE OPTIONAL AN ARRAY WITH SHAPE NSAMPLESY NFEATURES

SUMOVERFEATURES BOOL DEFAULTTRUE IF TRUE THE FUNCTION RETURNS THE PAIRWISE DISTANCE MATRIX ELSE IT RETURNS THE COMPONENTWISE L1 PAIRWISEDISTANCES NOT SUPPORTED FOR SPARSE MATRIX INPUTS

RETURNS

DARRAY IF SUMOVERFEATURES IS FALSE SHAPE IS NSAMPLESX NSAMPLESY NFEATURES AND D CONTAINS THE COMPONENTWISE L1 PAIRWISEDISTANCES IE ABSOLUTE DIFFERENCE ELSE SHAPE IS NSAMPLESX NSAMPLESY AND D CONTAINS THE PAIRWISE L1 DISTANCES

EXAMPLES

FROM SKLEARNMETRICSPAIRWISE IMPORT MANHATTANDISTANCES

MANHATTANDISTANCES3 3

ARRAY0

MANHATTANDISTANCES3 2

ARRAY1

MANHATTANDISTANCES2 3

ARRAY1

MANHATTANDISTANCES1 2 3 4 1 2 0 3

ARRAY0 2

4 4

IMPORT NUMPY AS NP

X NPONES1 2

Y NPFULL2 2 2

MANHATTANDISTANCESX Y SUMOVERFEATURES FALSE

ARRAY1 1

1 1

624SKLEARNMETRICS METRICS 2047

SCIKITLEARN USER GUIDE RELEASE 0213

SKLEARNMETRICSPAIRWISE PAIRWISEKERNELS

SKLEARNMETRICSPAIRWISE PAIRWISEKERNELS X YNONE METRIC'LINEAR' FIL

TERPARAMSFALSE NJOBSNONE KWDS

COMPUTE THE KERNEL BETWEEN ARRAYS X AND OPTIONAL ARRAY Y

THIS METHOD TAKES EITHER A VECTOR ARRAY OR A KERNEL MATRIX AND RETURNS A KERNEL MATRIX IF THE INPUT IS A VECTOR

ARRAY THE KERNELS ARE COMPUTED IF THE INPUT IS A KERNEL MATRIX IT IS RETURNED INSTEAD

THIS METHOD PROVIDES A SAFE WAY TO TAKE A KERNEL MATRIX AS INPUT WHILE PRESERVING COMPATIBILITY WITH MANY

OTHER ALGORITHMS THAT TAKE A VECTOR ARRAY

IF Y IS GIVEN DEFAULT IS NONE THEN THE RETURNED MATRIX IS THE PAIRWISE KERNEL BETWEEN THE ARRAYS FROM BOTH X

AND Y

VALID VALUES FOR METRIC ARE

'ADDITIVECHI2' 'CHI2' 'LINEAR' 'POLY' 'POLYNOMIAL' 'RBF' 'LAPLACIAN' 'SIGMOID' 'COSINE'

READ MORE IN THE USER GUIDE

PARAMETERS

XARRAY NSAMPLESA NSAMPLESA IF METRIC "PRECOMPUTED" OR NSAMPLESA

NFEATURES OTHERWISE ARRAY OF PAIRWISE KERNELS BETWEEN SAMPLES OR A FEATURE ARRAY

YARRAY NSAMPLESB NFEATURES A SECOND FEATURE ARRAY ONLY IF X HAS SHAPE NSAMPLESA

NFEATURES

METRIC STRING OR CALLABLE THE METRIC TO USE WHEN CALCULATING KERNEL BETWEEN INSTANCES

IN A FEATURE ARRAY IF METRIC IS A STRING IT MUST BE ONE OF THE METRICS IN PAIR

WISEPAIRWISEKERNELFUNCTIONS IF METRIC IS "PRECOMPUTED" X IS ASSUMED TO

BE A KERNEL MATRIX ALTERNATIVELY IF METRIC IS A CALLABLE FUNCTION IT IS CALLED ON EACH PAIR OF

INSTANCES ROWS AND THE RESULTING VALUE RECORDED THE CALLABLE SHOULD TAKE TWO ARRAYS FROM

X AS INPUT AND RETURN A VALUE INDICATING THE DISTANCE BETWEEN THEM

FILTERPARAMS BOOLEAN WHETHER TO FILTER INVALID PARAMETERS OR NOT

NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF JOBS TO USE FOR THE COMPUTATION

THIS WORKS BY BREAKING DOWN THE PAIRWISE MATRIX INTO NJOBS EVEN SLICES AND COMPUTING

THEM IN PARALLEL

NONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL

PROCESSORS SEE GLOSSARY FOR MORE DETAILS

KWDS OPTIONAL KEYWORD PARAMETERS ANY FURTHER PARAMETERS ARE PASSED DIRECTLY TO THE KERNEL

FUNCTION

RETURNS

KARRAY NSAMPLESA NSAMPLESA OR NSAMPLESA NSAMPLESB A KERNEL MATRIX K SUCH

THAT KI J IS THE KERNEL BETWEEN THE ITH AND JTH VECTORS OF THE GIVEN MATRIX X IF Y IS NONE

IF Y IS NOT NONE THEN KI J IS THE KERNEL BETWEEN THE ITH ARRAY FROM X AND THE JTH ARRAY

FROM Y

NOTES

IF METRIC IS 'PRECOMPUTED' Y IS IGNORED AND X IS RETURNED

2048 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213  
SKLEARNMETRICSPAIRWISE POLYNOMIALKERNEL  
SKLEARNMETRICSPAIRWISE POLYNOMIALKERNEL XYNONE DEGREE3 GAMMANONE  
COEF01  
COMPUTE THE POLYNOMIAL KERNEL BETWEEN X AND Y  
KX Y GAMMA X Y COEF0DEGREE  
READ MORE IN THE USER GUIDE  
PARAMETERS  
XNDARRAY OF SHAPE NSAMPLES1 NFEATURES  
YNDARRAY OF SHAPE NSAMPLES2 NFEATURES  
DEGREE INT DEFAULT 3  
GAMMA FLOAT DEFAULT NONE IF NONE DEFAULTS TO 10 NFEATURES  
COEF0 FLOAT DEFAULT 1  
RETURNS  
GRAM MATRIX ARRAY OF SHAPE NSAMPLES1 NSAMPLES2  
SKLEARNMETRICSPAIRWISE RBFKERNEL  
SKLEARNMETRICSPAIRWISE RBFKERNEL XYNONE GAMMANONE  
COMPUTE THE RBF GAUSSIAN KERNEL BETWEEN X AND Y  
KX Y EXPGAMMA XY2  
FOR EACH PAIR OF ROWS X IN X AND Y IN Y  
READ MORE IN THE USER GUIDE  
PARAMETERS  
XARRAY OF SHAPE NSAMPLESX NFEATURES  
YARRAY OF SHAPE NSAMPLESY NFEATURES  
GAMMA FLOAT DEFAULT NONE IF NONE DEFAULTS TO 10 NFEATURES  
RETURNS  
KERNELMATRIX ARRAY OF SHAPE NSAMPLESX NSAMPLESY  
SKLEARNMETRICSPAIRWISE SIGMOIDKERNEL  
SKLEARNMETRICSPAIRWISE SIGMOIDKERNEL XYNONE GAMMANONE COEF01  
COMPUTE THE SIGMOID KERNEL BETWEEN X AND Y  
KX Y TANHGAMMA X Y COEF0  
READ MORE IN THE USER GUIDE  
PARAMETERS  
XNDARRAY OF SHAPE NSAMPLES1 NFEATURES  
YNDARRAY OF SHAPE NSAMPLES2 NFEATURES  
624SKLEARNMETRICS METRICS 2049

SCIKITLEARN USER GUIDE RELEASE 0213  
GAMMA FLOAT DEFAULT NONE IF NONE DEFAULTS TO 10 NFEATURES  
COEF0 FLOAT DEFAULT 1  
RETURNS  
GRAM MATRIX ARRAY OF SHAPE NSAMPLES1 NSAMPLES2  
SKLEARNMETRICSPAIRWISE PAIREDEUCLIDEANDISTANCES  
SKLEARNMETRICSPAIRWISE PAIREDEUCLIDEANDISTANCES XY  
COMPUTES THE PAIRED EUCLIDEAN DISTANCES BETWEEN X AND Y  
READ MORE IN THE USER GUIDE  
PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES  
YARRAYLIKE SHAPE NSAMPLES NFEATURES  
RETURNS  
DISTANCES NDARRAY NSAMPLES  
SKLEARNMETRICSPAIRWISE PAIREDMANHATTANDISTANCES  
SKLEARNMETRICSPAIRWISE PAIREDMANHATTANDISTANCES XY  
COMPUTE THE L1 DISTANCES BETWEEN THE VECTORS IN X AND Y  
READ MORE IN THE USER GUIDE  
PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES  
YARRAYLIKE SHAPE NSAMPLES NFEATURES  
RETURNS  
DISTANCES NDARRAY NSAMPLES  
SKLEARNMETRICSPAIRWISE PAIREDCOSINEDISTANCES  
SKLEARNMETRICSPAIRWISE PAIREDCOSINEDISTANCES XY  
COMPUTES THE PAIRED COSINE DISTANCES BETWEEN X AND Y  
READ MORE IN THE USER GUIDE  
PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES  
YARRAYLIKE SHAPE NSAMPLES NFEATURES  
RETURNS  
DISTANCES NDARRAY SHAPE NSAMPLES  
2050 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

NOTES

THE COSINE DISTANCE IS EQUIVALENT TO THE HALF THE SQUARED EUCLIDEAN DISTANCE IF EACH SAMPLE IS NORMALIZED TO UNIT NORM

SKLEARNMETRICSPAIRWISE PAIREDDISTANCES

SKLEARNMETRICSPAIRWISE PAIREDDISTANCES XYMETRIC'EUCLIDEAN' KWDS

COMPUTES THE PAIRED DISTANCES BETWEEN X AND Y

COMPUTES THE DISTANCES BETWEEN X0 Y0 X1 Y1 ETC

READ MORE IN THE USER GUIDE

PARAMETERS

XNDARRAY NSAMPLES NFEATURES ARRAY 1 FOR DISTANCE COMPUTATION

YNDARRAY NSAMPLES NFEATURES ARRAY 2 FOR DISTANCE COMPUTATION

METRIC STRING OR CALLABLE THE METRIC TO USE WHEN CALCULATING DISTANCE BETWEEN INSTANCES

IN A FEATURE ARRAY IF METRIC IS A STRING IT MUST BE ONE OF THE OPTIONS SPECIFIED IN

PAIREDDISTANCES INCLUDING "EUCLIDEAN" "MANHATTAN" OR "COSINE" ALTERNATIVELY IF

METRIC IS A CALLABLE FUNCTION IT IS CALLED ON EACH PAIR OF INSTANCES ROWS AND THE RESULTING

VALUE RECORDED THE CALLABLE SHOULD TAKE TWO ARRAYS FROM X AS INPUT AND RETURN A VALUE

INDICATING THE DISTANCE BETWEEN THEM

RETURNS

DISTANCES NDARRAY NSAMPLES

SEE ALSO

PAIRWISEDISTANCES COMPUTES THE DISTANCE BETWEEN EVERY PAIR OF SAMPLES

EXAMPLES

FROM SKLEARNMETRICSPAIRWISE IMPORT PAIREDDISTANCES

X 0 1 1 1

Y 0 1 2 1

PAIREDDISTANCESX Y

ARRAY0 1

SKLEARNMETRICS PAIRWISEDISTANCES

SKLEARNMETRICS PAIRWISEDISTANCES XYNONE METRIC'EUCLIDEAN' NJOBSNONE KWDS

COMPUTE THE DISTANCE MATRIX FROM A VECTOR ARRAY X AND OPTIONAL Y

THIS METHOD TAKES EITHER A VECTOR ARRAY OR A DISTANCE MATRIX AND RETURNS A DISTANCE MATRIX IF THE INPUT IS A VECTOR

ARRAY THE DISTANCES ARE COMPUTED IF THE INPUT IS A DISTANCES MATRIX IT IS RETURNED INSTEAD

THIS METHOD PROVIDES A SAFE WAY TO TAKE A DISTANCE MATRIX AS INPUT WHILE PRESERVING COMPATIBILITY WITH MANY

OTHER ALGORITHMS THAT TAKE A VECTOR ARRAY

IF Y IS GIVEN DEFAULT IS NONE THEN THE RETURNED MATRIX IS THE PAIRWISE DISTANCE BETWEEN THE ARRAYS FROM BOTH X

AND Y

624SKLEARNMETRICS METRICS 2051

SCIKITLEARN USER GUIDE RELEASE 0213

VALID VALUES FOR METRIC ARE

- FROM SCIKITLEARN ‘CITYBLOCK’ ‘COSINE’ ‘EUCLIDEAN’ ‘L1’ ‘L2’ ‘MANHATTAN’ THESE METRICS SUPPORT SPARSE MATRIX INPUTS
- FROM SCIPYSPATIALDISTANCE ‘BRAYCURTIS’ ‘CANBERRA’ ‘CHEBYSHEV’ ‘CORRELATION’ ‘DICE’ ‘HAMMING’ ‘JACCARD’ ‘KULSINSKI’ ‘MAHALANOBIS’ ‘MINKOWSKI’ ‘ROGERSTANIMOTO’ ‘RUSSELLRAO’ ‘SEUCLIDEAN’ ‘SOKALMICHENER’ ‘SOKALSNEATH’ ‘SQEUCLIDEAN’ ‘YULE’ SEE THE DOCUMENTATION FOR SCIPYSPATIALDISTANCE FOR DETAILS ON THESE METRICS THESE METRICS DO NOT SUPPORT SPARSE MATRIX INPUTS

NOTE THAT IN THE CASE OF ‘CITYBLOCK’ ‘COSINE’ AND ‘EUCLIDEAN’ WHICH ARE VALID SCIPYSPATIALDISTANCE METRICS THE SCIKITLEARN IMPLEMENTATION WILL BE USED WHICH IS FASTER AND HAS SUPPORT FOR SPARSE MATRICES EXCEPT FOR ‘CITYBLOCK’ FOR A VERBOSE DESCRIPTION OF THE METRICS FROM SCIKITLEARN SEE THE DOC OF THE SKLEARNPAIRWISEDISTANCEMETRICS FUNCTION

READ MORE IN THE USER GUIDE

PARAMETERS

XARRAY NSAMPLESA NSAMPLESA IF METRIC “PRECOMPUTED” OR NSAMPLESA

NFEATURES OTHERWISE ARRAY OF PAIRWISE DISTANCES BETWEEN SAMPLES OR A FEATURE ARRAY

YARRAY NSAMPLESB NFEATURES OPTIONAL AN OPTIONAL SECOND FEATURE ARRAY ONLY ALLOWED

IF METRIC “PRECOMPUTED”

METRIC STRING OR CALLABLE THE METRIC TO USE WHEN CALCULATING DISTANCE BETWEEN IN

STANCES IN A FEATURE ARRAY IF METRIC IS A STRING IT MUST BE ONE OF THE OPTIONS AL

LOWED BY SCIPYSPATIALDISTANCEPDIST FOR ITS METRIC PARAMETER OR A METRIC LISTED IN PAIR

WISEPAIRWISEDISTANCEFUNCTIONS IF METRIC IS “PRECOMPUTED” X IS ASSUMED TO

BE A DISTANCE MATRIX ALTERNATIVELY IF METRIC IS A CALLABLE FUNCTION IT IS CALLED ON EACH PAIR

OF INSTANCES ROWS AND THE RESULTING VALUE RECORDED THE CALLABLE SHOULD TAKE TWO ARRAYS

FROM X AS INPUT AND RETURN A VALUE INDICATING THE DISTANCE BETWEEN THEM

NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF JOBS TO USE FOR THE COMPUTATION

THIS WORKS BY BREAKING DOWN THE PAIRWISE MATRIX INTO NJOBS EVEN SLICES AND COMPUTING

THEM IN PARALLEL

NONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL

PROCESSORS SEE GLOSSARY FOR MORE DETAILS

KWDS OPTIONAL KEYWORD PARAMETERS ANY FURTHER PARAMETERS ARE PASSED DIRECTLY TO THE DIS

TANCE FUNCTION IF USING A SCIPYSPATIALDISTANCE METRIC THE PARAMETERS ARE STILL METRIC DE

PENDENT SEE THE SCIPY DOCS FOR USAGE EXAMPLES

RETURNS

DARRAY NSAMPLESA NSAMPLESA OR NSAMPLESA NSAMPLESB A DISTANCE MATRIX D

SUCH THAT DI J IS THE DISTANCE BETWEEN THE ITH AND JTH VECTORS OF THE GIVEN MATRIX X IF Y

IS NONE IF Y IS NOT NONE THEN DI J IS THE DISTANCE BETWEEN THE ITH ARRAY FROM X AND THE

JTH ARRAY FROM Y

SEE ALSO

PAIRWISEDISTANCESCHUNKED PERFORMS THE SAME CALCULATION AS THIS FUNCTION BUT RETURNS A GENERATOR

OF CHUNKS OF THE DISTANCE MATRIX IN ORDER TO LIMIT MEMORY USAGE

PAIREDDISTANCES COMPUTES THE DISTANCES BETWEEN CORRESPONDING ELEMENTS OF TWO ARRAYS

2052 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES USING SKLEARNMETRICSPAIRWISEDISTANCES

- AGGLOMERATIVE CLUSTERING WITH DIFFERENT METRICS

SKLEARNMETRICS PAIRWISEDISTANCESARGMIN

SKLEARNMETRICS PAIRWISEDISTANCESARGMIN X Y AXIS1 METRIC'EUCLIDEAN'

BATCHSIZENONE METRICKWARGSNONE

COMPUTE MINIMUM DISTANCES BETWEEN ONE POINT AND A SET OF POINTS

THIS FUNCTION COMPUTES FOR EACH ROW IN X THE INDEX OF THE ROW OF Y WHICH IS CLOSEST ACCORDING TO THE SPECIFIED DISTANCE

THIS IS MOSTLY EQUIVALENT TO CALLING

PAIRWISEDISTANCESX YY METRICMETRICARGMINAXISAXIS

BUT USES MUCH LESS MEMORY AND IS FASTER FOR LARGE ARRAYS

THIS FUNCTION WORKS WITH DENSE 2D ARRAYS ONLY

PARAMETERS

XARRAYLIKE ARRAYS CONTAINING POINTS RESPECTIVE SHAPES NSAMPLES1 NFEATURES AND

NSAMPLES2 NFEATURES

YARRAYLIKE ARRAYS CONTAINING POINTS RESPECTIVE SHAPES NSAMPLES1 NFEATURES AND

NSAMPLES2 NFEATURES

AXIS INT OPTIONAL DEFAULT 1 AXIS ALONG WHICH THE ARGMIN AND DISTANCES ARE TO BE COMPUTED

METRIC STRING OR CALLABLE METRIC TO USE FOR DISTANCE COMPUTATION ANY METRIC FROM SCIKITLEARN

OR SCIPYSPATIALDISTANCE CAN BE USED

IF METRIC IS A CALLABLE FUNCTION IT IS CALLED ON EACH PAIR OF INSTANCES ROWS AND THE RESULTING

VALUE RECORDED THE CALLABLE SHOULD TAKE TWO ARRAYS AS INPUT AND RETURN ONE VALUE INDICATING

THE DISTANCE BETWEEN THEM THIS WORKS FOR SCIPY'S METRICS BUT IS LESS EFFICIENT THAN PASSING

THE METRIC NAME AS A STRING

DISTANCE MATRICES ARE NOT SUPPORTED

VALID VALUES FOR METRIC ARE

- FROM SCIKITLEARN 'CITYBLOCK' 'COSINE' 'EUCLIDEAN' 'L1' 'L2' 'MANHATTAN'
- FROM SCIPYSPATIALDISTANCE 'BRAYCURTIS' 'CANBERRA' 'CHEBYSHEV' 'CORRELATION' 'DICE'

'HAMMING' 'JACCARD' 'KULSINSKI' 'MAHALANOBIS' 'MINKOWSKI' 'ROGERSTANIMOTO' 'RUS

SELLRAO' 'SEUCLIDEAN' 'SOKALMICHENER' 'SOKALSNEATH' 'SQUCLIDEAN' 'YULE'

SEE THE DOCUMENTATION FOR SCIPYSPATIALDISTANCE FOR DETAILS ON THESE METRICS

BATCHSIZE INTEGER DEPRECATED SINCE VERSION 020 DEPRECATED FOR REMOVAL IN 022 USE

SKLEARNSETCONFIGWORKINGMEMORY INSTEAD

METRICKWARGS DICT KEYWORD ARGUMENTS TO PASS TO SPECIFIED METRIC FUNCTION

RETURNS

ARGMIN NUMPYNDARRAY YARGMINI IS THE ROW IN Y THAT IS CLOSEST TO XI

SEE ALSO

SKLEARNMETRICSPAIRWISEDISTANCES

624SKLEARNMETRICS METRICS 2053

SCIKITLEARN USER GUIDE RELEASE 0213

SKLEARNMETRICSPAIRWISEDISTANCESARGMINMIN

EXAMPLES USING SKLEARNMETRICSPAIRWISEDISTANCESARGMIN

- COLOR QUANTIZATION USING KMEANS
- COMPARISON OF THE KMEANS AND MINIBATCHKMEANS CLUSTERING ALGORITHMS

SKLEARNMETRICS PAIRWISEDISTANCESARGMINMIN

SKLEARNMETRICS PAIRWISEDISTANCESARGMINMIN XYAXIS1 METRIC'EUCLIDEAN'

BATCHSIZENONE METRICKWARGSNONE

COMPUTE MINIMUM DISTANCES BETWEEN ONE POINT AND A SET OF POINTS

THIS FUNCTION COMPUTES FOR EACH ROW IN X THE INDEX OF THE ROW OF Y WHICH IS CLOSEST ACCORDING TO THE SPECIFIED DISTANCE THE MINIMAL DISTANCES ARE ALSO RETURNED

THIS IS MOSTLY EQUIVALENT TO CALLING

PAIRWISEDISTANCESX YY METRICMETRICARGMINAXISAXIS PAIRWISEDISTANCESX YY

METRICMETRICMINAXISAXIS

BUT USES MUCH LESS MEMORY AND IS FASTER FOR LARGE ARRAYS

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES1 NFEATURES ARRAY CONTAINING POINTS

YARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES2 NFEATURES ARRAYS CONTAINING POINTS

AXIS INT OPTIONAL DEFAULT 1 AXIS ALONG WHICH THE ARGMIN AND DISTANCES ARE TO BE COMPUTED

METRIC STRING OR CALLABLE DEFAULT 'EUCLIDEAN' METRIC TO USE FOR DISTANCE COMPUTATION ANY METRIC FROM SCIKITLEARN OR SCIPYSPATIALDISTANCE CAN BE USED

IF METRIC IS A CALLABLE FUNCTION IT IS CALLED ON EACH PAIR OF INSTANCES ROWS AND THE RESULTING VALUE RECORDED THE CALLABLE SHOULD TAKE TWO ARRAYS AS INPUT AND RETURN ONE VALUE INDICATING THE DISTANCE BETWEEN THEM THIS WORKS FOR SCIPY'S METRICS BUT IS LESS EFFICIENT THAN PASSING THE METRIC NAME AS A STRING

DISTANCE MATRICES ARE NOT SUPPORTED

VALID VALUES FOR METRIC ARE

- FROM SCIKITLEARN 'CITYBLOCK' 'COSINE' 'EUCLIDEAN' 'L1' 'L2' 'MANHATTAN'
- FROM SCIPYSPATIALDISTANCE 'BRAYCURTIS' 'CANBERRA' 'CHEBYSHEV' 'CORRELATION' 'DICE' 'HAMMING' 'JACCARD' 'KULSINSKI' 'MAHALANOBIS' 'MINKOWSKI' 'ROGERSTANIMOTO' 'RUSSELLRAO' 'SEUCLIDEAN' 'SOKALMICHENER' 'SOKALSNEATH' 'SQUCLIDEAN' 'YULE'

SEE THE DOCUMENTATION FOR SCIPYSPATIALDISTANCE FOR DETAILS ON THESE METRICS

BATCHSIZE INTEGER DEPRECATED SINCE VERSION 020 DEPRECATED FOR REMOVAL IN 022 USE SKLEARNSETCONFIGWORKINGMEMORY INSTEAD

METRICKWARGS DICT OPTIONAL KEYWORD ARGUMENTS TO PASS TO SPECIFIED METRIC FUNCTION

RETURNS

ARGMIN NUMPYNDARRAY YARGMINI IS THE ROW IN Y THAT IS CLOSEST TO XI

DISTANCES NUMPYNDARRAY DISTANCESI IS THE DISTANCE BETWEEN THE ITH ROW IN X AND THE ARGMINITH ROW IN Y

2054 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

SEE ALSO

SKLEARNMETRICSPAIRWISEDISTANCES

SKLEARNMETRICSPAIRWISEDISTANCESARGMIN

SKLEARNMETRICS PAIRWISEDISTANCESCHUNKED

SKLEARNMETRICS PAIRWISEDISTANCESCHUNKED XYNONE REDUCEFUNCNONE METRIC'EUCLIDEAN' NJOBSNONE WORKINGMEMORYNONE KWDS

GENERATE A DISTANCE MATRIX CHUNK BY CHUNK WITH OPTIONAL REDUCTION

IN CASES WHERE NOT ALL OF A PAIRWISE DISTANCE MATRIX NEEDS TO BE STORED AT ONCE THIS IS USED TO CALCULATE PAIRWISE DISTANCES IN WORKINGMEMORY SIZED CHUNKS IF REDUCEFUNC IS GIVEN IT IS RUN ON EACH CHUNK AND ITS RETURN VALUES ARE CONCATENATED INTO LISTS ARRAYS OR SPARSE MATRICES

PARAMETERS

XARRAY NSAMPLESA NSAMPLESA IF METRIC "PRECOMPUTED" OR NSAMPLESA NFEATURES OTHERWISE ARRAY OF PAIRWISE DISTANCES BETWEEN SAMPLES OR A FEATURE ARRAY

YARRAY NSAMPLESB NFEATURES OPTIONAL AN OPTIONAL SECOND FEATURE ARRAY ONLY ALLOWED IF METRIC "PRECOMPUTED"

REDUCEFUNC CALLABLE OPTIONAL THE FUNCTION WHICH IS APPLIED ON EACH CHUNK OF THE DISTANCE MATRIX REDUCING IT TO NEEDED VALUES REDUCEFUNCDCHUNK START IS CALLED REPEATEDLY WHERE DCHUNK IS A CONTIGUOUS VERTICAL SLICE OF THE PAIRWISE DISTANCE MATRIX STARTING AT ROW START IT SHOULD RETURN AN ARRAY A LIST OR A SPARSE MATRIX OF LENGTH DCHUNKSHAPE0 OR A TUPLE OF SUCH OBJECTS

IF NONE PAIRWISEDISTANCESCHUNKED RETURNS A GENERATOR OF VERTICAL CHUNKS OF THE DISTANCE MATRIX

METRIC STRING OR CALLABLE THE METRIC TO USE WHEN CALCULATING DISTANCE BETWEEN INSTANCES IN A FEATURE ARRAY IF METRIC IS A STRING IT MUST BE ONE OF THE OPTIONS ALLOWED BY SCIPYSPATIALDISTANCEPDIST FOR ITS METRIC PARAMETER OR A METRIC LISTED IN PAIRWISEPAIRWISEDISTANCEFUNCTIONS IF METRIC IS "PRECOMPUTED" X IS ASSUMED TO BE A DISTANCE MATRIX ALTERNATIVELY IF METRIC IS A CALLABLE FUNCTION IT IS CALLED ON EACH PAIR OF INSTANCES ROWS AND THE RESULTING VALUE RECORDED THE CALLABLE SHOULD TAKE TWO ARRAYS FROM X AS INPUT AND RETURN A VALUE INDICATING THE DISTANCE BETWEEN THEM

NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF JOBS TO USE FOR THE COMPUTATION THIS WORKS BY BREAKING DOWN THE PAIRWISE MATRIX INTO NJOBS EVEN SLICES AND COMPUTING THEM IN PARALLEL

NONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

WORKINGMEMORY INT OPTIONAL THE SOUGHT MAXIMUM MEMORY FOR TEMPORARY DISTANCE MATRIX CHUNKS WHEN NONE DEFAULT THE VALUE OF SKLEARN GETCONFIGWORKINGMEMORY IS USED

'KWDS' OPTIONAL KEYWORD PARAMETERS ANY FURTHER PARAMETERS ARE PASSED DIRECTLY TO THE DISTANCE FUNCTION IF USING A SCIPYSPATIALDISTANCE METRIC THE PARAMETERS ARE STILL METRIC DEPENDENT SEE THE SCIPY DOCS FOR USAGE EXAMPLES

YIELDS

624SKLEARNMETRICS METRICS 2055

SCIKITLEARN USER GUIDE RELEASE 0213

DCHUNK ARRAY OR SPARSE MATRIX A CONTIGUOUS SLICE OF DISTANCE MATRIX OPTIONALLY PROCESSED BYREDUCEFUNC

EXAMPLES

WITHOUT REDUCEFUNC

IMPORT NUMPY AS NP

FROM SKLEARNMETRICS IMPORT PAIRWISEDISTANCESCHUNKED

X NPRANDOMRANDOMSTATE0RAND5 3

DCHUNK NEXTPAIRWISEDISTANCESCHUNKEDX

DCHUNK

ARRAY0 029 041 019 057

029 0 057 041 076

041 057 0 044 090

019 041 044 0 051

057 076 090 051 0

RETRIEVE ALL NEIGHBORS AND AVERAGE DISTANCE WITHIN RADIUS R

R 2

DEF REDUCEFUNCDCHUNK START

NEIGH NPFLATNONZEROD R FORDINDCHUNK

AVGDIST DCHUNK DCHUNK RMEANAXIS1

RETURN NEIGH AVGDIST

GEN PAIRWISEDISTANCESCHUNKEDX REDUCEFUNCREDUCEFUNC

NEIGH AVGDIST NEXTGEN

NEIGH

ARRAY0 3 ARRAY1 ARRAY2 ARRAY0 3 ARRAY4

AVGDIST

ARRAY0039 0 0 0039 0

WHERE R IS DEFINED PER SAMPLE WE NEED TO MAKE USE OF START

R 2 4 4 3 1

DEF REDUCEFUNCDCHUNK START

NEIGH NPFLATNONZEROD RI

FOR I DINENUMERATEDCHUNK START

RETURN NEIGH

NEIGH NEXTPAIRWISEDISTANCESCHUNKEDX REDUCEFUNCREDUCEFUNC

NEIGH

ARRAY0 3 ARRAY0 1 ARRAY2 ARRAY0 3 ARRAY4

FORCE ROWBYROW GENERATION BY REDUCING WORKINGMEMORY

GEN PAIRWISEDISTANCESCHUNKEDX REDUCEFUNCREDUCEFUNC

WORKINGMEMORY0

NEXTGEN

ARRAY0 3

NEXTGEN

ARRAY0 1

625SKLEARNMIXTURE GAUSSIAN MIXTURE MODELS

THESKLEARNMIXTURE MODULE IMPLEMENTS MIXTURE MODELING ALGORITHMS

2056 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

USER GUIDE SEE THE GAUSSIAN MIXTURE MODELS SECTION FOR FURTHER DETAILS  
MIXTUREBAYESIANGAUSSIANMIXTURE VARIATIONAL BAYESIAN ESTIMATION OF A GAUSSIAN MIXTURE  
MIXTUREGAUSSIANMIXTURE NCOMPONENTS GAUSSIAN MIXTURE  
6251SKLEARNMIXTURE BAYESIANGAUSSIANMIXTURE  
CLASSSKLEARNMIXTURE BAYESIANGAUSSIANMIXTURE NCOMPONENTS1 COVARIANCETYPE'FULL'

TOL0001 REGCOVAR1E06 MAXITER100  
NINIT1 INITPARAMS'KMEANS'  
WEIGHTCONCENTRATIONPRIORTYPE'DIRICHLETPROCESS'  
WEIGHTCONCENTRATIONPRIORNONE  
MEANPRECISIONPRIORNONE

MEANPRIORNONE DE  
GREESOFFREEDOMPRIORNONE COVARI  
ANCEPRIORNONE RANDOMSTATENONE  
WARMSTARTFALSE VERBOSE0 VER  
BOSEINTERVAL10

VARIATIONAL BAYESIAN ESTIMATION OF A GAUSSIAN MIXTURE

THIS CLASS ALLOWS TO INFER AN APPROXIMATE POSTERIOR DISTRIBUTION OVER THE PARAMETERS OF A GAUSSIAN MIXTURE  
DISTRIBUTION THE EFFECTIVE NUMBER OF COMPONENTS CAN BE INFERRED FROM THE DATA  
THIS CLASS IMPLEMENTS TWO TYPES OF PRIOR FOR THE WEIGHTS DISTRIBUTION A FINITE MIXTURE MODEL WITH DIRICHLET  
DISTRIBUTION AND AN INFINITE MIXTURE MODEL WITH THE DIRICHLET PROCESS IN PRACTICE DIRICHLET PROCESS INFERENCE  
ALGORITHM IS APPROXIMATED AND USES A TRUNCATED DISTRIBUTION WITH A FIXED MAXIMUM NUMBER OF COMPONENTS  
CALLED THE STICKBREAKING REPRESENTATION THE NUMBER OF COMPONENTS ACTUALLY USED ALMOST ALWAYS DEPENDS ON  
THE DATA

NEW IN VERSION 018  
READ MORE IN THE USER GUIDE  
PARAMETERS

NCOMPONENTS INT DEFAULTS TO 1 THE NUMBER OF MIXTURE COMPONENTS DEPENDING ON THE DATA  
AND THE VALUE OF THE WEIGHTCONCENTRATIONPRIOR THE MODEL CAN DECIDE TO NOT USE  
ALL THE COMPONENTS BY SETTING SOME COMPONENT WEIGHTS TO VALUES VERY CLOSE TO ZERO  
THE NUMBER OF EFFECTIVE COMPONENTS IS THEREFORE SMALLER THAN NCOMPONENTS  
COVARIANCETYPE 'FULL' 'TIED' 'DIAG' 'SPHERICAL' DEFAULTS TO 'FULL' STRING DESCRIBING THE  
TYPE OF COVARIANCE PARAMETERS TO USE MUST BE ONE OF  
FULL EACH COMPONENT HAS ITS OWN GENERAL COVARIANCE MATRIX  
TIED ALL COMPONENTS SHARE THE SAME GENERAL COVARIANCE MATRIX  
DIAG EACH COMPONENT HAS ITS OWN DIAGONAL COVARIANCE MATRIX  
SPHERICAL EACH COMPONENT HAS ITS OWN SINGLE VARIANCE  
TOLFLOAT DEFAULTS TO 1E3 THE CONVERGENCE THRESHOLD EM ITERATIONS WILL STOP WHEN THE LOWER  
BOUND AVERAGE GAIN ON THE LIKELIHOOD OF THE TRAINING DATA WITH RESPECT TO THE MODEL IS BELOW  
THIS THRESHOLD  
REGCOVAR FLOAT DEFAULTS TO 1E6 NONNEGATIVE REGULARIZATION ADDED TO THE DIAGONAL OF CO  
VARIANCE ALLOWS TO ASSURE THAT THE COVARIANCE MATRICES ARE ALL POSITIVE  
MAXITER INT DEFAULTS TO 100 THE NUMBER OF EM ITERATIONS TO PERFORM  
625SKLEARNMIXTURE GAUSSIAN MIXTURE MODELS 2057

SCIKITLEARN USER GUIDE RELEASE 0213

NINIT INT DEFAULTS TO 1 THE NUMBER OF INITIALIZATIONS TO PERFORM THE RESULT WITH THE HIGHEST LOWER BOUND VALUE ON THE LIKELIHOOD IS KEPT

INITPARAMS ‘KMEANS’ ‘RANDOM’ DEFAULTS TO ‘KMEANS’ THE METHOD USED TO INITIALIZE THE WEIGHTS THE MEANS AND THE COVARIANCES MUST BE ONE OF

KMEANS RESPONSIBILITIES ARE INITIALIZED USING KMEANS

RANDOM RESPONSIBILITIES ARE INITIALIZED RANDOMLY

WEIGHTCONCENTRATIONPRIORTYPE STR DEFAULTS TO ‘DIRICHLETPROCESS’ STRING DESCRIBING THE TYPE OF THE WEIGHT CONCENTRATION PRIOR MUST BE ONE OF

DIRICHLETPROCESS USING THE STICKBREAKING REPRESENTATION

DIRICHLETDISTRIBUTION CAN FAVOR MORE UNIFORM WEIGHTS

WEIGHTCONCENTRATIONPRIOR FLOAT NONE OPTIONAL THE DIRICHLET CONCENTRATION OF EACH COMPONENT ON THE WEIGHT DISTRIBUTION DIRICHLET THIS IS COMMONLY CALLED GAMMA IN THE LITERATURE THE HIGHER CONCENTRATION PUTS MORE MASS IN THE CENTER AND WILL LEAD TO MORE COMPONENTS BEING ACTIVE WHILE A LOWER CONCENTRATION PARAMETER WILL LEAD TO MORE MASS AT THE EDGE OF THE MIXTURE WEIGHTS SIMPLEX THE VALUE OF THE PARAMETER MUST BE GREATER THAN 0 IF IT IS NONE IT’S SET TO 1 NCOMPONENTS

MEANPRECISIONPRIOR FLOAT NONE OPTIONAL THE PRECISION PRIOR ON THE MEAN DISTRIBUTION GAUSSIAN CONTROLS THE EXTEND TO WHERE MEANS CAN BE PLACED LARGER VALUES CONCENTRATE THE MEANS OF EACH CLUSTERS AROUND MEANPRIOR THE VALUE OF THE PARAMETER MUST BE GREATER THAN 0 IF IT IS NONE IT’S SET TO 1

MEANPRIOR ARRAYLIKE SHAPE NFEATURES OPTIONAL THE PRIOR ON THE MEAN DISTRIBUTION GAUSSIAN IF IT IS NONE IT’S SET TO THE MEAN OF X

DEGREESOFFREEDOMPRIOR FLOAT NONE OPTIONAL THE PRIOR OF THE NUMBER OF DEGREES OF FREEDOM ON THE COVARIANCE DISTRIBUTIONS WISHART IF IT IS NONE IT’S SET TO NFEATURES

COVARIANCEPRIOR FLOAT OR ARRAYLIKE OPTIONAL THE PRIOR ON THE COVARIANCE DISTRIBUTION WISHART IF IT IS NONE THE EMIPRICAL COVARIANCE PRIOR IS INITIALIZED USING THE COVARIANCE OF X THE SHAPE DEPENDS ON COVARIANCETYPE

NFEATURES NFEATURES IFFULL

NFEATURES NFEATURES IFTIED

NFEATURES IFDIAG

FLOAT IFSPHERICAL

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RANDOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NRANDOM

WARMSTART BOOL DEFAULT TO FALSE IF ‘WARMSTART’ IS TRUE THE SOLUTION OF THE LAST FITTING IS USED AS INITIALIZATION FOR THE NEXT CALL OF FIT THIS CAN SPEED UP CONVERGENCE WHEN FIT IS CALLED SEVERAL TIMES ON SIMILAR PROBLEMS SEE THE GLOSSARY

VERBOSE INT DEFAULT TO 0 ENABLE VERBOSE OUTPUT IF 1 THEN IT PRINTS THE CURRENT INITIALIZATION AND EACH ITERATION STEP IF GREATER THAN 1 THEN IT PRINTS ALSO THE LOG PROBABILITY AND THE TIME NEEDED FOR EACH STEP

VERBOSEINTERVAL INT DEFAULT TO 10 NUMBER OF ITERATION DONE BEFORE THE NEXT PRINT

ATTRIBUTES

WEIGHTS ARRAYLIKE SHAPE NCOMPONENTS THE WEIGHTS OF EACH MIXTURE COMPONENTS

2058 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

MEANS ARRAYLIKE SHAPE NCOMPONENTS NFEATURES THE MEAN OF EACH MIXTURE COMPONENT  
COVARIANCES ARRAYLIKE THE COVARIANCE OF EACH MIXTURE COMPONENT THE SHAPE DEPENDS ON  
COVARIANCETYPE

NCOMPONENTS IFSPHERICAL  
NFEATURES NFEATURES IFTIED  
NCOMPONENTS NFEATURES IFDIAG

NCOMPONENTS NFEATURES NFEATURES IFFULL  
PRECISIONS ARRAYLIKE THE PRECISION MATRICES FOR EACH COMPONENT IN THE MIXTURE A PRECI  
SION MATRIX IS THE INVERSE OF A COVARIANCE MATRIX A COVARIANCE MATRIX IS SYMMETRIC POSI  
TIVE DEFINITE SO THE MIXTURE OF GAUSSIAN CAN BE EQUIVALENTLY PARAMETERIZED BY THE PRECISION  
MATRICES STORING THE PRECISION MATRICES INSTEAD OF THE COVARIANCE MATRICES MAKES IT MORE  
EFFICIENT TO COMPUTE THE LOGLIKELIHOOD OF NEW SAMPLES AT TEST TIME THE SHAPE DEPENDS ON  
COVARIANCETYPE

NCOMPONENTS IFSPHERICAL  
NFEATURES NFEATURES IFTIED  
NCOMPONENTS NFEATURES IFDIAG  
NCOMPONENTS NFEATURES NFEATURES IFFULL

PRECISIONSCHOLESKY ARRAYLIKE THE CHOLESKY DECOMPOSITION OF THE PRECISION MATRICES OF  
EACH MIXTURE COMPONENT A PRECISION MATRIX IS THE INVERSE OF A COVARIANCE MATRIX A COVARI  
ANCE MATRIX IS SYMMETRIC POSITIVE DEFINITE SO THE MIXTURE OF GAUSSIAN CAN BE EQUIVALENTLY  
PARAMETERIZED BY THE PRECISION MATRICES STORING THE PRECISION MATRICES INSTEAD OF THE CO  
VARIANCE MATRICES MAKES IT MORE EFFICIENT TO COMPUTE THE LOGLIKELIHOOD OF NEW SAMPLES AT  
TEST TIME THE SHAPE DEPENDS ON COVARIANCETYPE

NCOMPONENTS IFSPHERICAL  
NFEATURES NFEATURES IFTIED  
NCOMPONENTS NFEATURES IFDIAG  
NCOMPONENTS NFEATURES NFEATURES IFFULL

CONVERGED BOOL TRUE WHEN CONVERGENCE WAS REACHED IN FIT FALSE OTHERWISE  
NITER INT NUMBER OF STEP USED BY THE BEST FIT OF INFERENCE TO REACH THE CONVERGENCE  
LOWERBOUND FLOAT LOWER BOUND VALUE ON THE LIKELIHOOD OF THE TRAINING DATA WITH RESPECT TO  
THE MODEL OF THE BEST FIT OF INFERENCE

WEIGHTCONCENTRATIONPRIOR TUPLE OR FLOAT THE DIRICHLET CONCENTRATION OF EACH  
COMPONENT ON THE WEIGHT DISTRIBUTION DIRICHLET THE TYPE DEPENDS ON  
WEIGHTCONCENTRATIONPRIORTYPE

FLOAT FLOAT IFDIRICHLETPROCESS BETA PARAMETERS  
FLOAT IFDIRICHLETDISTRIBUTION DIRICHLET PARAMETERS

THE HIGHER CONCENTRATION PUTS MORE MASS IN THE CENTER AND WILL LEAD TO MORE COMPONENTS  
BEING ACTIVE WHILE A LOWER CONCENTRATION PARAMETER WILL LEAD TO MORE MASS AT THE EDGE OF  
THE SIMPLEX

WEIGHTCONCENTRATION ARRAYLIKE SHAPE NCOMPONENTS THE DIRICHLET CONCENTRATION OF  
EACH COMPONENT ON THE WEIGHT DISTRIBUTION DIRICHLET  
MEANPRECISIONPRIOR FLOAT THE PRECISION PRIOR ON THE MEAN DISTRIBUTION GAUSSIAN CON  
TROLS THE EXTEND TO WHERE MEANS CAN BE PLACED LARGER VALUES CONCENTRATE THE MEANS OF EACH  
CLUSTERS AROUND MEANPRIOR

SCIKITLEARN USER GUIDE RELEASE 0213

MEANPRECISION ARRAYLIKE SHAPE NCOMPONENTS THE PRECISION OF EACH COMPONENTS ON THE MEAN DISTRIBUTION GAUSSIAN

MEANPRIOR ARRAYLIKE SHAPE NFEATURES THE PRIOR ON THE MEAN DISTRIBUTION GAUSSIAN

DEGREESOFFREEDOMPRIOR FLOAT THE PRIOR OF THE NUMBER OF DEGREES OF FREEDOM ON THE CO VARIANCE DISTRIBUTIONS WISHART

DEGREESOFFREEDOM ARRAYLIKE SHAPE NCOMPONENTS THE NUMBER OF DEGREES OF FREEDOM OF EACH COMPONENTS IN THE MODEL

COVARIANCEPRIOR FLOAT OR ARRAYLIKE THE PRIOR ON THE COVARIANCE DISTRIBUTION WISHART THE SHAPE DEPENDS ON COVARIANCETYPE

NFEATURES NFEATURES IFFULL

NFEATURES NFEATURES IFTIED

NFEATURES IFDIAG

FLOAT IFSPHERICAL

SEE ALSO

GAUSSIANMIXTURE FINITE GAUSSIAN MIXTURE FIT WITH EM

REFERENCES

R16529824BFF21 R16529824BFF22 R16529824BFF23

METHODS

FITSELF X Y ESTIMATE MODEL PARAMETERS WITH THE EM ALGORITHM

FITPREDICT SELF X Y ESTIMATE MODEL PARAMETERS USING X AND PREDICT THE LA BELS FOR X

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PREDICT SELF X PREDICT THE LABELS FOR THE DATA SAMPLES IN X USING TRAINED MODEL

PREDICTPROBA SELF X PREDICT POSTERIOR PROBABILITY OF EACH COMPONENT GIVEN THE DATA

SAMPLE SELF NSAMPLES GENERATE RANDOM SAMPLES FROM THE FITTED GAUSSIAN DIS TRIBUTION

SCORE SELF X Y COMPUTE THE PERSAMPLE AVERAGE LOGLIKELIHOOD OF THE GIVEN DATA X

SCORESAMPLES SELF X COMPUTE THE WEIGHTED LOG PROBABILITIES FOR EACH SAM PLE

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELFNCOMPONENTS1 COVARIANCETYPE'FULL' TOL0001 REGCOVAR1E06 MAXITER100

NINIT1 INITPARAMS'KMEANS' WEIGHTCONCENTRATIONPRIORTYPE'DIRICHLETPROCESS'

WEIGHTCONCENTRATIONPRIORNONE MEANPRECISIONPRIORNONE MEANPRIORNONE

DEGREESOFFREEDOMPRIORNONE COVARIANCEPRIORNONE RANDOMSTATENONE

WARMSTARTFALSE VERBOSE0 VERBOSEINTERVAL10

FITSELFXYNONE

ESTIMATE MODEL PARAMETERS WITH THE EM ALGORITHM

2060 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

THE METHOD FITS THE MODEL NINIT TIMES AND SETS THE PARAMETERS WITH WHICH THE MODEL HAS THE LARGEST LIKELIHOOD OR LOWER BOUND WITHIN EACH TRIAL THE METHOD ITERATES BETWEEN ESTEP AND MSTEP FORMAXITER TIMES UNTIL THE CHANGE OF LIKELIHOOD OR LOWER BOUND IS LESS THAN TOL OTHERWISE A CONVERGENCEWARNING IS RAISED IF WARMSTART ISTRUE THENNINIT IS IGNORED AND A SINGLE INITIALIZATION IS PERFORMED UPON THE FIRST CALL UPON CONSECUTIVE CALLS TRAINING STARTS WHERE IT LEFT OFF

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES LIST OF NFEATURESDIMENSIONAL DATA POINTS

EACH ROW CORRESPONDS TO A SINGLE DATA POINT

RETURNS

SELF

FITPREDICT SELFXYNONE

ESTIMATE MODEL PARAMETERS USING X AND PREDICT THE LABELS FOR X

THE METHOD FITS THE MODEL NINIT TIMES AND SETS THE PARAMETERS WITH WHICH THE MODEL HAS THE LARGEST LIKELIHOOD OR LOWER BOUND WITHIN EACH TRIAL THE METHOD ITERATES BETWEEN ESTEP AND MSTEP FOR MAXITER TIMES UNTIL THE CHANGE OF LIKELIHOOD OR LOWER BOUND IS LESS THAN TOL OTHERWISE A CONVERGENCEWARNING IS RAISED AFTER FITTING IT PREDICTS THE MOST PROBABLE LABEL FOR THE INPUT DATA POINTS

NEW IN VERSION 020

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES LIST OF NFEATURESDIMENSIONAL DATA POINTS

EACH ROW CORRESPONDS TO A SINGLE DATA POINT

RETURNS

LABELS ARRAY SHAPE NSAMPLES COMPONENT LABELS

GETPARAMS SELFDEEPTTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELF

PREDICT THE LABELS FOR THE DATA SAMPLES IN X USING TRAINED MODEL

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES LIST OF NFEATURESDIMENSIONAL DATA POINTS

EACH ROW CORRESPONDS TO A SINGLE DATA POINT

RETURNS

LABELS ARRAY SHAPE NSAMPLES COMPONENT LABELS

PREDICTPROBA SELF

PREDICT POSTERIOR PROBABILITY OF EACH COMPONENT GIVEN THE DATA

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES LIST OF NFEATURESDIMENSIONAL DATA POINTS

EACH ROW CORRESPONDS TO A SINGLE DATA POINT

625SKLEARNMIXTURE GAUSSIAN MIXTURE MODELS 2061

SCIKITLEARN USER GUIDE RELEASE 0213

RETURNS

RESP ARRAY SHAPE NSAMPLES NCOMPONENTS RETURNS THE PROBABILITY EACH GAUSSIAN STATE IN THE MODEL GIVEN EACH SAMPLE

SAMPLESELFNSAMPLES1

GENERATE RANDOM SAMPLES FROM THE FITTED GAUSSIAN DISTRIBUTION

PARAMETERS

NSAMPLES INT OPTIONAL NUMBER OF SAMPLES TO GENERATE DEFAULTS TO 1

RETURNS

XARRAY SHAPE NSAMPLES NFEATURES RANDOMLY GENERATED SAMPLE

YARRAY SHAPE NSAMPLES COMPONENT LABELS

SCORESELFXYNONE

COMPUTE THE PERSAMPLE AVERAGE LOGLIKELIHOOD OF THE GIVEN DATA X

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NDIMENSIONS LIST OF NFEATURESDIMENSIONAL DATA POINTS EACH ROW CORRESPONDS TO A SINGLE DATA POINT

RETURNS

LOGLIKELIHOOD FLOAT LOG LIKELIHOOD OF THE GAUSSIAN MIXTURE GIVEN X

SCORESAMPLES SELF

COMPUTE THE WEIGHTED LOG PROBABILITIES FOR EACH SAMPLE

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES LIST OF NFEATURESDIMENSIONAL DATA POINTS EACH ROW CORRESPONDS TO A SINGLE DATA POINT

RETURNS

LOGPROB ARRAY SHAPE NSAMPLES LOG PROBABILITIES OF EACH DATA POINT IN X

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

EXAMPLES USING SKLEARNMIXTUREBAYESIANGAUSSIANMIXTURE

- GAUSSIAN MIXTURE MODEL ELLIPSOIDS
- GAUSSIAN MIXTURE MODEL SINE CURVE
- CONCENTRATION PRIOR TYPE ANALYSIS OF VARIATION BAYESIAN GAUSSIAN MIXTURE

2062 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213  
6252SKLEARNMIXTURE GAUSSIANMIXTURE  
CLASSSSKLEARNMIXTURE GAUSSIANMIXTURE NCOMPONENTS1 COVARIANCETYPE'FULL' TOL0001  
REGCOVAR1E06 MAXITER100 NINIT1  
INITPARAMS'KMEANS' WEIGHTSINITNONE  
MEANSINITNONE PRECISIONSINITNONE RAN  
DOMSTATENONE WARMSTARTFALSE VERBOSE0  
VERBOSEINTERVAL10  
GAUSSIAN MIXTURE  
REPRESENTATION OF A GAUSSIAN MIXTURE MODEL PROBABILITY DISTRIBUTION THIS CLASS ALLOWS TO ESTIMATE THE PARAME  
TERS OF A GAUSSIAN MIXTURE DISTRIBUTION  
READ MORE IN THE USER GUIDE  
NEW IN VERSION 018  
PARAMETERS  
NCOMPONENTS INT DEFAULTS TO 1 THE NUMBER OF MIXTURE COMPONENTS  
COVARIANCETYPE 'FULL' DEFAULT 'TIED' 'DIAG' 'SPHERICAL' STRING DESCRIBING THE TYPE OF  
COVARIANCE PARAMETERS TO USE MUST BE ONE OF  
'FULL' EACH COMPONENT HAS ITS OWN GENERAL COVARIANCE MATRIX  
'TIED' ALL COMPONENTS SHARE THE SAME GENERAL COVARIANCE MATRIX  
'DIAG' EACH COMPONENT HAS ITS OWN DIAGONAL COVARIANCE MATRIX  
'SPHERICAL' EACH COMPONENT HAS ITS OWN SINGLE VARIANCE  
TOLFLOAT DEFAULTS TO 1E3 THE CONVERGENCE THRESHOLD EM ITERATIONS WILL STOP WHEN THE LOWER  
BOUND AVERAGE GAIN IS BELOW THIS THRESHOLD  
REGCOVAR FLOAT DEFAULTS TO 1E6 NONNEGATIVE REGULARIZATION ADDED TO THE DIAGONAL OF CO  
VARIANCE ALLOWS TO ASSURE THAT THE COVARIANCE MATRICES ARE ALL POSITIVE  
MAXITER INT DEFAULTS TO 100 THE NUMBER OF EM ITERATIONS TO PERFORM  
NINIT INT DEFAULTS TO 1 THE NUMBER OF INITIALIZATIONS TO PERFORM THE BEST RESULTS ARE KEPT  
INITPARAMS 'KMEANS' 'RANDOM' DEFAULTS TO 'KMEANS' THE METHOD USED TO INITIALIZE THE  
WEIGHTS THE MEANS AND THE PRECISIONS MUST BE ONE OF  
KMEANS RESPONSIBILITIES ARE INITIALIZED USING KMEANS  
RANDOM RESPONSIBILITIES ARE INITIALIZED RANDOMLY  
WEIGHTSINIT ARRAYLIKE SHAPE NCOMPONENTS OPTIONAL THE USERPROVIDED INITIAL WEIGHTS  
DEFAULTS TO NONE IF IT NONE WEIGHTS ARE INITIALIZED USING THE INITPARAMS METHOD  
MEANSINIT ARRAYLIKE SHAPE NCOMPONENTS NFEATURES OPTIONAL THE USERPROVIDED INITIAL  
MEANS DEFAULTS TO NONE IF IT NONE MEANS ARE INITIALIZED USING THE INITPARAMS METHOD  
PRECISIONSINIT ARRAYLIKE OPTIONAL THE USERPROVIDED INITIAL PRECISIONS INVERSE OF THE CO  
VARIANCE MATRICES DEFAULTS TO NONE IF IT NONE PRECISIONS ARE INITIALIZED USING THE  
'INITPARAMS' METHOD THE SHAPE DEPENDS ON 'COVARIANCETYPE'  
NCOMPONENTS IFSPHERICAL  
NFEATURES NFEATURES IFTIED  
NCOMPONENTS NFEATURES IFDIAG  
NCOMPONENTS NFEATURES NFEATURES IFFULL  
625SKLEARNMIXTURE GAUSSIAN MIXTURE MODELS 2063

SCIKITLEARN USER GUIDE RELEASE 0213

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RAN  
DOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE  
RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS  
THE RANDOMSTATE INSTANCE USED BY NPRANDOM  
WARMSTART BOOL DEFAULT TO FALSE IF 'WARMSTART' IS TRUE THE SOLUTION OF THE LAST FITTING IS  
USED AS INITIALIZATION FOR THE NEXT CALL OF FIT THIS CAN SPEED UP CONVERGENCE WHEN FIT IS  
CALLED SEVERAL TIMES ON SIMILAR PROBLEMS IN THAT CASE 'NINIT' IS IGNORED AND ONLY A SINGLE  
INITIALIZATION OCCURS UPON THE FIRST CALL SEE THE GLOSSARY  
VERBOSE INT DEFAULT TO 0 ENABLE VERBOSE OUTPUT IF 1 THEN IT PRINTS THE CURRENT INITIALIZATION  
AND EACH ITERATION STEP IF GREATER THAN 1 THEN IT PRINTS ALSO THE LOG PROBABILITY AND THE TIME  
NEEDED FOR EACH STEP  
VERBOSEINTERVAL INT DEFAULT TO 10 NUMBER OF ITERATION DONE BEFORE THE NEXT PRINT  
ATTRIBUTES  
WEIGHTS ARRAYLIKE SHAPE NCOMPONENTS THE WEIGHTS OF EACH MIXTURE COMPONENTS  
MEANS ARRAYLIKE SHAPE NCOMPONENTS NFEATURES THE MEAN OF EACH MIXTURE COMPONENT  
COVARIANCES ARRAYLIKE THE COVARIANCE OF EACH MIXTURE COMPONENT THE SHAPE DEPENDS ON  
COVARIANCETYPE  
NCOMPONENTS IFSPHERICAL  
NFEATURES NFEATURES IFTIED  
NCOMPONENTS NFEATURES IFDIAG  
NCOMPONENTS NFEATURES NFEATURES IFFULL  
PRECISIONS ARRAYLIKE THE PRECISION MATRICES FOR EACH COMPONENT IN THE MIXTURE A PRECI  
SION MATRIX IS THE INVERSE OF A COVARIANCE MATRIX A COVARIANCE MATRIX IS SYMMETRIC POSI  
TIVE DEFINITE SO THE MIXTURE OF GAUSSIAN CAN BE EQUIVALENTLY PARAMETERIZED BY THE PRECISION  
MATRICES STORING THE PRECISION MATRICES INSTEAD OF THE COVARIANCE MATRICES MAKES IT MORE  
EFFICIENT TO COMPUTE THE LOGLIKELIHOOD OF NEW SAMPLES AT TEST TIME THE SHAPE DEPENDS ON  
COVARIANCETYPE  
NCOMPONENTS IFSPHERICAL  
NFEATURES NFEATURES IFTIED  
NCOMPONENTS NFEATURES IFDIAG  
NCOMPONENTS NFEATURES NFEATURES IFFULL  
PRECISIONSCHOLESKY ARRAYLIKE THE CHOLESKY DECOMPOSITION OF THE PRECISION MATRICES OF  
EACH MIXTURE COMPONENT A PRECISION MATRIX IS THE INVERSE OF A COVARIANCE MATRIX A COVARI  
ANCE MATRIX IS SYMMETRIC POSITIVE DEFINITE SO THE MIXTURE OF GAUSSIAN CAN BE EQUIVALENTLY  
PARAMETERIZED BY THE PRECISION MATRICES STORING THE PRECISION MATRICES INSTEAD OF THE CO  
VARIANCE MATRICES MAKES IT MORE EFFICIENT TO COMPUTE THE LOGLIKELIHOOD OF NEW SAMPLES AT  
TEST TIME THE SHAPE DEPENDS ON COVARIANCETYPE  
NCOMPONENTS IFSPHERICAL  
NFEATURES NFEATURES IFTIED  
NCOMPONENTS NFEATURES IFDIAG  
NCOMPONENTS NFEATURES NFEATURES IFFULL  
CONVERGED BOOL TRUE WHEN CONVERGENCE WAS REACHED IN FIT FALSE OTHERWISE  
NITER INT NUMBER OF STEP USED BY THE BEST FIT OF EM TO REACH THE CONVERGENCE  
LOWERBOUND FLOAT LOWER BOUND VALUE ON THE LOGLIKELIHOOD OF THE TRAINING DATA WITH RE  
SPECT TO THE MODEL OF THE BEST FIT OF EM  
2064 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

SEE ALSO

BAYESIANGAUSSIANMIXTURE GAUSSIAN MIXTURE MODEL FIT WITH A VARIATIONAL INFERENCE METHODS

AICSELF X AKAIKE INFORMATION CRITERION FOR THE CURRENT MODEL ON THE INPUT X

BICSELF X BAYESIAN INFORMATION CRITERION FOR THE CURRENT MODEL ON THE INPUT X

FITSELF X Y ESTIMATE MODEL PARAMETERS WITH THE EM ALGORITHM

FITPREDICT SELF X Y ESTIMATE MODEL PARAMETERS USING X AND PREDICT THE LABELS FOR X

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PREDICT SELF X PREDICT THE LABELS FOR THE DATA SAMPLES IN X USING TRAINED MODEL

PREDICTPROBA SELF X PREDICT POSTERIOR PROBABILITY OF EACH COMPONENT GIVEN THE DATA

SAMPLE SELF NSAMPLES GENERATE RANDOM SAMPLES FROM THE FITTED GAUSSIAN DISTRIBUTION

SCORE SELF X Y COMPUTE THE PERSAMPLE AVERAGE LOGLIKELIHOOD OF THE GIVEN DATA X

SCORESAMPLES SELF X COMPUTE THE WEIGHTED LOG PROBABILITIES FOR EACH SAMPLE

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELF NCOMPONENTS1 COVARIANCETYPE'FULL' TOL0001 REGCOVAR1E06 MAXITER100 NINIT1 INITPARAMS'KMEANS' WEIGHTSINITNONE MEANSINITNONE PRECISIONSINITNONE RANDOMSTATENONE WARMSTARTFALSE VERBOSE0 VERBOSEINTERVAL10

AICSELF X

AKAIKE INFORMATION CRITERION FOR THE CURRENT MODEL ON THE INPUT X

PARAMETERS

XARRAY OF SHAPE NSAMPLES NDIMENSIONS

RETURNS

AICFLOAT THE LOWER THE BETTER

BICSELF X

BAYESIAN INFORMATION CRITERION FOR THE CURRENT MODEL ON THE INPUT X

PARAMETERS

XARRAY OF SHAPE NSAMPLES NDIMENSIONS

RETURNS

BICFLOAT THE LOWER THE BETTER

FITSELFXYNONE

ESTIMATE MODEL PARAMETERS WITH THE EM ALGORITHM

THE METHOD FITS THE MODEL NINIT TIMES AND SETS THE PARAMETERS WITH WHICH THE MODEL HAS THE LARGEST LIKELIHOOD OR LOWER BOUND WITHIN EACH TRIAL THE METHOD ITERATES BETWEEN ESTEP AND MSTEP

625SKLEARNMIXTURE GAUSSIAN MIXTURE MODELS 2065

SCIKITLEARN USER GUIDE RELEASE 0213

FORMAXITER TIMES UNTIL THE CHANGE OF LIKELIHOOD OR LOWER BOUND IS LESS THAN TOL OTHERWISE A CONVERGENCEWARNING IS RAISED IF WARMSTART ISTRUE THENNINIT IS IGNORED AND A SINGLE INITIALIZATION IS PERFORMED UPON THE FIRST CALL UPON CONSECUTIVE CALLS TRAINING STARTS WHERE IT LEFT OFF

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES LIST OF NFEATURESDIMENSIONAL DATA POINTS

EACH ROW CORRESPONDS TO A SINGLE DATA POINT

RETURNS

SELF

FITPREDICT SELFXYNONE

ESTIMATE MODEL PARAMETERS USING X AND PREDICT THE LABELS FOR X

THE METHOD FITS THE MODEL NINIT TIMES AND SETS THE PARAMETERS WITH WHICH THE MODEL HAS THE LARGEST LIKELIHOOD OR LOWER BOUND WITHIN EACH TRIAL THE METHOD ITERATES BETWEEN ESTEP AND MSTEP FOR MAXITER TIMES UNTIL THE CHANGE OF LIKELIHOOD OR LOWER BOUND IS LESS THAN TOL OTHERWISE A CONVERGENCEWARNING IS RAISED AFTER FITTING IT PREDICTS THE MOST PROBABLE LABEL FOR THE INPUT DATA POINTS

NEW IN VERSION 020

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES LIST OF NFEATURESDIMENSIONAL DATA POINTS

EACH ROW CORRESPONDS TO A SINGLE DATA POINT

RETURNS

LABELS ARRAY SHAPE NSAMPLES COMPONENT LABELS

GETPARAMS SELFDEEPTTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELF

PREDICT THE LABELS FOR THE DATA SAMPLES IN X USING TRAINED MODEL

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES LIST OF NFEATURESDIMENSIONAL DATA POINTS

EACH ROW CORRESPONDS TO A SINGLE DATA POINT

RETURNS

LABELS ARRAY SHAPE NSAMPLES COMPONENT LABELS

PREDICTPROBA SELF

PREDICT POSTERIOR PROBABILITY OF EACH COMPONENT GIVEN THE DATA

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES LIST OF NFEATURESDIMENSIONAL DATA POINTS

EACH ROW CORRESPONDS TO A SINGLE DATA POINT

RETURNS

2066 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

RESP ARRAY SHAPE NSAMPLES NCOMPONENTS RETURNS THE PROBABILITY EACH GAUSSIAN STATE IN THE MODEL GIVEN EACH SAMPLE

SAMPLESELFNSAMPLES1

GENERATE RANDOM SAMPLES FROM THE FITTED GAUSSIAN DISTRIBUTION

PARAMETERS

NSAMPLES INT OPTIONAL NUMBER OF SAMPLES TO GENERATE DEFAULTS TO 1

RETURNS

XARRAY SHAPE NSAMPLES NFEATURES RANDOMLY GENERATED SAMPLE

YARRAY SHAPE NSAMPLES COMPONENT LABELS

SCORESELFXYNONE

COMPUTE THE PERSAMPLE AVERAGE LOGLIKELIHOOD OF THE GIVEN DATA X

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NDIMENSIONS LIST OF NFEATURESDIMENSIONAL DATA POINTS

EACH ROW CORRESPONDS TO A SINGLE DATA POINT

RETURNS

LOGLIKELIHOOD FLOAT LOG LIKELIHOOD OF THE GAUSSIAN MIXTURE GIVEN X

SCORESAMPLES SELF X

COMPUTE THE WEIGHTED LOG PROBABILITIES FOR EACH SAMPLE

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES LIST OF NFEATURESDIMENSIONAL DATA POINTS

EACH ROW CORRESPONDS TO A SINGLE DATA POINT

RETURNS

LOGPROB ARRAY SHAPE NSAMPLES LOG PROBABILITIES OF EACH DATA POINT IN X

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

EXAMPLES USING SKLEARNMIXTUREGAUSSIANMIXTURE

- COMPARING DIFFERENT CLUSTERING ALGORITHMS ON TOY DATASETS
- DENSITY ESTIMATION FOR A GAUSSIAN MIXTURE
- GAUSSIAN MIXTURE MODEL ELLIPSOIDS
- GAUSSIAN MIXTURE MODEL SELECTION
- GMM COVARIANCES
- GAUSSIAN MIXTURE MODEL SINE CURVE

625SKLEARNMIXTURE GAUSSIAN MIXTURE MODELS 2067

SCIKITLEARN USER GUIDE RELEASE 0213

626SKLEARNMODELSELECTION MODEL SELECTION  
USER GUIDE SEE THE CROSSVALIDATION EVALUATING ESTIMATOR PERFORMANCE TUNING THE HYPERPARAMETERS OF AN ESTIMATOR AND LEARNING CURVE SECTIONS FOR FURTHER DETAILS

6261 SPLITTER CLASSES

MODELSELECTIONGROUPKFOLD NSPLITS KFOLD ITERATOR VARIANT WITH NONOVERLAPPING GROUPS

MODELSELECTIONGROUPSHUFFLESPLIT SHUFFLEGROUPSOUT CROSSVALIDATION ITERATOR

MODELSELECTIONKFOLD NSPLITS SHUFFLE KFOLDS CROSSVALIDATOR

MODELSELECTIONLEAVEONEGROUPOUT LEAVE ONE GROUP OUT CROSSVALIDATOR

MODELSELECTIONLEAVEPGROUPSOUT NGROUPS LEAVE P GROUPS OUT CROSSVALIDATOR

MODELSELECTIONLEAVEONEOUT LEAVEONEOUT CROSSVALIDATOR

MODELSELECTIONLEAVEPOUT P LEAVEPOUT CROSSVALIDATOR

MODELSELECTIONPREDEFINEDSPLIT TESTFOLD PREDEFINED SPLIT CROSSVALIDATOR

MODELSELECTIONREPEATEDKFOLD NSPLITS  
REPEATED KFOLD CROSS VALIDATOR

MODELSELECTIONREPEATEDSTRATIFIEDKFOLD REPEATED STRATIFIED KFOLD CROSS VALIDATOR

MODELSELECTIONSHUFFLESPLIT NSPLITS RANDOM PERMUTATION CROSSVALIDATOR

MODELSELECTIONSTRATIFIEDKFOLD NSPLITS  
STRATIFIED KFOLDS CROSSVALIDATOR

MODELSELECTIONSTRATIFIEDSHUFFLESPLIT STRATIFIED SHUFFLESPLIT CROSSVALIDATOR

MODELSELECTIONTIMESERIESSPLIT NSPLITS  
TIME SERIES CROSSVALIDATOR

SKLEARNMODELSELECTION GROUPKFOLD

CLASSSSKLEARNMODELSELECTION GROUPKFOLD NSPLITS'WARN'

KFOLD ITERATOR VARIANT WITH NONOVERLAPPING GROUPS

THE SAME GROUP WILL NOT APPEAR IN TWO DIFFERENT FOLDS THE NUMBER OF DISTINCT GROUPS HAS TO BE AT LEAST EQUAL TO THE NUMBER OF FOLDS

THE FOLDS ARE APPROXIMATELY BALANCED IN THE SENSE THAT THE NUMBER OF DISTINCT GROUPS IS APPROXIMATELY THE SAME IN EACH FOLD

PARAMETERS

NSPLITS INT DEFAULT3 NUMBER OF FOLDS MUST BE AT LEAST 2

CHANGED IN VERSION 020 NSPLITS DEFAULT VALUE WILL CHANGE FROM 3 TO 5 IN V022

SEE ALSO

LEAVEONEGROUPOUT FOR SPLITTING THE DATA ACCORDING TO EXPLICIT DOMAINSPECIFIC STRATIFICATION OF THE DATASET

EXAMPLES

```
import numpy as np
from sklearn.model_selection import GroupKFold
X = np.array(1 2 3 4 5 6 7 8)
Y = np.array(1 2 3 4)
groups = np.array(0 2 2)
```

2068 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213  
GROUPKFOLD GROUPKFOLDNSPLITS2  
GROUPKFOLDGETNSPLITSX Y GROUPS  
2  
PRINTGROUPKFOLD  
GROUPKFOLDNSPLITS2  
FOR TRAININDEX TESTINDEX INGROUPKFOLDSPPLITX Y GROUPS  
PRINTTRAIN TRAININDEX TEST TESTINDEX  
XTRAIN XTEST XTRAININDEX XTESTINDEX  
YTRAIN YTEST YTRAININDEX YTESTINDEX  
PRINTXTRAIN XTEST YTRAIN YTEST  
  
TRAIN 0 1 TEST 2 3  
1 2  
3 4 5 6  
7 8 1 2 3 4  
TRAIN 2 3 TEST 0 1  
5 6  
7 8 1 2  
3 4 3 4 1 2  
METHODS  
GETNSPLITS SELF X Y GROUPS RETURNS THE NUMBER OF SPLITTING ITERATIONS IN THE CROSS  
VALIDATOR  
SPLIT SELF X Y GROUPS GENERATE INDICES TO SPLIT DATA INTO TRAINING AND TEST SET  
INIT SELFNSPLITS'WARN'  
GETNSPLITS SELFXNONE YNONE GROUPSNONE  
RETURNS THE NUMBER OF SPLITTING ITERATIONS IN THE CROSSVALIDATOR  
PARAMETERS  
XOBJECT ALWAYS IGNORED EXISTS FOR COMPATIBILITY  
YOBJECT ALWAYS IGNORED EXISTS FOR COMPATIBILITY  
GROUPS OBJECT ALWAYS IGNORED EXISTS FOR COMPATIBILITY  
RETURNS  
NSPLITS INT RETURNS THE NUMBER OF SPLITTING ITERATIONS IN THE CROSSVALIDATOR  
SPLITSELFXYNONE GROUPSNONE  
GENERATE INDICES TO SPLIT DATA INTO TRAINING AND TEST SET  
PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA WHERE NSAMPLES IS THE NUMBER  
OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES  
YARRAYLIKE SHAPE NSAMPLES OPTIONAL THE TARGET VARIABLE FOR SUPERVISED LEARNING PROB  
LEMS  
GROUPS ARRAYLIKE WITH SHAPE NSAMPLES GROUP LABELS FOR THE SAMPLES USED WHILE SPLIT  
TING THE DATASET INTO TRRAINTEST SET  
YIELDS  
626SKLEARNMODELSELECTION MODEL SELECTION 2069

SCIKITLEARN USER GUIDE RELEASE 0213

TRAIN NDARRAY THE TRAINING SET INDICES FOR THAT SPLIT

TEST NDARRAY THE TESTING SET INDICES FOR THAT SPLIT

EXAMPLES USING SKLEARNMODELSELECTIONGROUPKFOLD

- VISUALIZING CROSSVALIDATION BEHAVIOR IN SCIKITLEARN

SKLEARNMODELSELECTION GROUPSHUFFLESPLIT

CLASSSSKLEARNMODELSELECTION GROUPSHUFFLESPLIT NSPLITS5 TESTSIZENONE

TRAINSIZENONE RANDOMSTATENONE

SHUFFLEGROUPSOUT CROSSVALIDATION ITERATOR

PROVIDES RANDOMIZED TRAINTEST INDICES TO SPLIT DATA ACCORDING TO A THIRDPARTY PROVIDED GROUP THIS GROUP INFORMATION CAN BE USED TO ENCODE ARBITRARY DOMAIN SPECIFIC STRATIFICATIONS OF THE SAMPLES AS INTEGERS

FOR INSTANCE THE GROUPS COULD BE THE YEAR OF COLLECTION OF THE SAMPLES AND THUS ALLOW FOR CROSSVALIDATION AGAINST TIMEBASED SPLITS

THE DIFFERENCE BETWEEN LEAVEPGROUPSOUT AND GROUPSHUFFLESPLIT IS THAT THE FORMER GENERATES SPLITS USING ALL SUBSETS OF SIZE PUNIQUE GROUPS WHEREAS GROUPSHUFFLESPLIT GENERATES A USERDETERMINED NUMBER OF RANDOM TEST SPLITS EACH WITH A USERDETERMINED FRACTION OF UNIQUE GROUPS

FOR EXAMPLE A LESS COMPUTATIONALLY INTENSIVE ALTERNATIVE TO LEAVEPGROUPSOUTP10 WOULD BE

GROUPSHUFFLESPLITTESTSIZE10 NSPLITS100

NOTE THE PARAMETERS TESTSIZE ANDTRAINSIZES REFER TO GROUPS AND NOT TO SAMPLES AS IN SHUFFLESPLIT

PARAMETERS

NSPLITS INT DEFAULT 5 NUMBER OF RESHUFFLING SPLITTING ITERATIONS

TESTSIZE FLOAT INT NONE OPTIONAL DEFAULTNONE IF FLOAT SHOULD BE BETWEEN 00 AND 10

AND REPRESENT THE PROPORTION OF THE DATASET TO INCLUDE IN THE TEST SPLIT IF INT REPRESENTS THE ABSOLUTE NUMBER OF TEST GROUPS IF NONE THE VALUE IS SET TO THE COMPLEMENT OF THE TRAIN SIZE

IFTRAINSIZES IS ALSO NONE IT WILL BE SET TO 02

TRAINSIZES FLOAT INT OR NONE DEFAULT IS NONE IF FLOAT SHOULD BE BETWEEN 00 AND 10 AND REPRESENT THE PROPORTION OF THE GROUPS TO INCLUDE IN THE TRAIN SPLIT IF INT REPRESENTS THE ABSOLUTE NUMBER OF TRAIN GROUPS IF NONE THE VALUE IS AUTOMATICALLY SET TO THE COMPLEMENT OF THE TEST SIZE

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RAN

DOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE

RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS

THE RANDOMSTATE INSTANCE USED BY NPRANDOM

METHODS

GETNSPLITS SELF X Y GROUPS RETURNS THE NUMBER OF SPLITTING ITERATIONS IN THE CROSS

VALIDATOR

SPLIT SELF X Y GROUPS GENERATE INDICES TO SPLIT DATA INTO TRAINING AND TEST SET

INIT SELFNSPLITS5 TESTSIZENONE TRAINSIZENONE RANDOMSTATENONE

2070 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

GETNSPLITS SELFYNONE YNONE GROUPSNONE

RETURNS THE NUMBER OF SPLITTING ITERATIONS IN THE CROSSVALIDATOR

PARAMETERS

XOBJECT ALWAYS IGNORED EXISTS FOR COMPATIBILITY

YOBJECT ALWAYS IGNORED EXISTS FOR COMPATIBILITY

GROUPS OBJECT ALWAYS IGNORED EXISTS FOR COMPATIBILITY

RETURNS

NSPLITS INT RETURNS THE NUMBER OF SPLITTING ITERATIONS IN THE CROSSVALIDATOR

SPLITSELFYNONE GROUPSNONE

GENERATE INDICES TO SPLIT DATA INTO TRAINING AND TEST SET

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

YARRAYLIKE SHAPE NSAMPLES OPTIONAL THE TARGET VARIABLE FOR SUPERVISED LEARNING PROBLEMS

GROUPS ARRAYLIKE WITH SHAPE NSAMPLES GROUP LABELS FOR THE SAMPLES USED WHILE SPLITTING THE DATASET INTO TRAINTEST SET

YIELDS

TRAIN NDARRAY THE TRAINING SET INDICES FOR THAT SPLIT

TEST NDARRAY THE TESTING SET INDICES FOR THAT SPLIT

NOTES

RANDOMIZED CV SPLITTERS MAY RETURN DIFFERENT RESULTS FOR EACH CALL OF SPLIT YOU CAN MAKE THE RESULTS IDENTICAL BY SETTING RANDOMSTATE TO AN INTEGER

EXAMPLES USING SKLEARNMODELSELECTIONGROUPSHUFFLESPLIT

- VISUALIZING CROSSVALIDATION BEHAVIOR IN SCIKITLEARN

SKLEARNMODELSELECTION KFOLD

CLASSSKLEARNMODELSELECTION KFOLDNSPLITS'WARN' SHUFFLEFALSE RANDOMSTATENONE

KFOLDS CROSSVALIDATOR

PROVIDES TRAINTEST INDICES TO SPLIT DATA IN TRAINTEST SETS SPLIT DATASET INTO K CONSECUTIVE FOLDS WITHOUT SHUFFLING BY DEFAULT

EACH FOLD IS THEN USED ONCE AS A VALIDATION WHILE THE K - 1 REMAINING FOLDS FORM THE TRAINING SET

READ MORE IN THE USER GUIDE

PARAMETERS

NSPLITS INT DEFAULT3 NUMBER OF FOLDS MUST BE AT LEAST 2

CHANGED IN VERSION 020 NSPLITS DEFAULT VALUE WILL CHANGE FROM 3 TO 5 IN V022

626SKLEARNMODELSELECTION MODEL SELECTION 2071

SCIKITLEARN USER GUIDE RELEASE 0213

SHUFFLE BOOLEAN OPTIONAL WHETHER TO SHUFFLE THE DATA BEFORE SPLITTING INTO BATCHES

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RAN

DOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE

RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS

THE RANDOMSTATE INSTANCE USED BY NPRANDOM USED WHEN SHUFFLE TRUE

SEE ALSO

STRATIFIEDKFOLD TAKES GROUP INFORMATION INTO ACCOUNT TO AVOID BUILDING FOLDS WITH IMBALANCED CLASS

DISTRIBUTIONS FOR BINARY OR MULTICLASS CLASSIFICATION TASKS

GROUPKFOLD KFOLD ITERATOR VARIANT WITH NONOVERLAPPING GROUPS

REPEATEDKFOLD REPEATS KFOLD N TIMES

NOTES

THE FIRSTNSAMPLES NSPLITS FOLDS HAVE SIZE NSAMPLES NSPLITS 1 OTHER FOLDS HAVE

SIZENSAMPLES NSPLITS WHERE NSAMPLES IS THE NUMBER OF SAMPLES

RANDOMIZED CV SPLITTERS MAY RETURN DIFFERENT RESULTS FOR EACH CALL OF SPLIT YOU CAN MAKE THE RESULTS IDENTICAL

BY SETTINGRANDOMSTATE TO AN INTEGER

EXAMPLES

```
import numpy as np
from sklearn.model_selection import kfold
X = np.array([1, 2, 3, 4, 1, 2, 3, 4])
Y = np.array([1, 2, 3, 4])
kf = kfold(n_splits=2)
for train_index, test_index in kf.split(X):
    print("TRAIN: %s TEST: %s" % (train_index, test_index))
    X_train, X_test, Y_train, Y_test = X[train_index], X[test_index], Y[train_index], Y[test_index]
```

TRAIN 2 3 TEST 0 1

TRAIN 0 1 TEST 2 3

METHODS

GETNSPLITS SELF X Y GROUPS RETURNS THE NUMBER OF SPLITTING ITERATIONS IN THE CROSS

VALIDATOR

SPLIT SELF X Y GROUPS GENERATE INDICES TO SPLIT DATA INTO TRAINING AND TEST SET

INIT SELFNSPLITS'WARN' SHUFFLEFALSE RANDOMSTATENONE

GETNSPLITS SELF X Y GROUPS RETURNS THE NUMBER OF SPLITTING ITERATIONS IN THE CROSSVALIDATOR

PARAMETERS

2072 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

XOBJECT ALWAYS IGNORED EXISTS FOR COMPATIBILITY

YOBJECT ALWAYS IGNORED EXISTS FOR COMPATIBILITY

GROUPS OBJECT ALWAYS IGNORED EXISTS FOR COMPATIBILITY

RETURNS

NSPLITS INT RETURNS THE NUMBER OF SPLITTING ITERATIONS IN THE CROSSVALIDATOR

SPLITSELFXYNONE GROUPSNONE

GENERATE INDICES TO SPLIT DATA INTO TRAINING AND TEST SET

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

YARRAYLIKE SHAPE NSAMPLES THE TARGET VARIABLE FOR SUPERVISED LEARNING PROBLEMS

GROUPS ARRAYLIKE WITH SHAPE NSAMPLES OPTIONAL GROUP LABELS FOR THE SAMPLES USED WHILE SPLITTING THE DATASET INTO TRAI NTEST SET

YIELDS

TRAIN NDARRAY THE TRAINING SET INDICES FOR THAT SPLIT

TEST NDARRAY THE TESTING SET INDICES FOR THAT SPLIT

EXAMPLES USING SKLEARNMODELSELECTIONKFOLD

- FEATURE AGGLOMERATION VS UNIVARIATE SELECTION
- GRADIENT BOOSTING OUTFBAG ESTIMATES
- CROSSVALIDATION ON DIABETES DATASET EXERCISE
- NESTED VERSUS NONNESTED CROSSVALIDATION
- VISUALIZING CROSSVALIDATION BEHAVIOR IN SCIKITLEARN

SKLEARNMODELSELECTION LEAVEONEGROUPOUT

CLASSSSKLEARNMODELSELECTION LEAVEONEGROUPOUT

LEAVE ONE GROUP OUT CROSSVALIDATOR

PROVIDES TRAI NTEST INDICES TO SPLIT DATA ACCORDING TO A THIRDPARTY PROVIDED GROUP THIS GROUP INFORMATION CAN BE USED TO ENCODE ARBITRARY DOMAIN SPECIFIC STRATIFICATIONS OF THE SAMPLES AS INTEGERS

FOR INSTANCE THE GROUPS COULD BE THE YEAR OF COLLECTION OF THE SAMPLES AND THUS ALLOW FOR CROSSVALIDATION AGAINST TIMEBASED SPLITS

READ MORE IN THE USER GUIDE

EXAMPLES

```
import numpy as np
from sklearn.model_selection import LeaveOneGroupOut
X = np.array(1 2 3 4 5 6 7 8)
Y = np.array(2 1 2)
groups = np.array(1 1 2 2)
626sklearn.model_selection MODEL SELECTION 2073
```

SCIKITLEARN USER GUIDE RELEASE 0213  
LOGO LEAVEONEGROUPOUT  
LOGOGETNSPLITSX Y GROUPS  
2  
LOGOGETNSPLITSGROUPSGROUPS GROUPS IS ALWAYS REQUIRED  
2  
PRINTLOGO  
LEAVEONEGROUPOUT  
FOR TRAININDEX TESTINDEX INLOGOSPLITX Y GROUPS  
PRINTTRAIN TRAININDEX TEST TESTINDEX  
XTRAIN XTEST XTRAININDEX XTESTINDEX  
YTRAIN YTEST YTRAININDEX YTESTINDEX  
PRINTXTRAIN XTEST YTRAIN YTEST  
TRAIN 2 3 TEST 0 1  
5 6  
7 8 1 2  
3 4 1 2 1 2  
TRAIN 0 1 TEST 2 3  
1 2  
3 4 5 6  
7 8 1 2 1 2  
METHODS  
GETNSPLITS SELF X Y GROUPS RETURNS THE NUMBER OF SPLITTING ITERATIONS IN THE CROSS  
VALIDATOR  
SPLIT SELF X Y GROUPS GENERATE INDICES TO SPLIT DATA INTO TRAINING AND TEST SET  
INIT SELFARGS KWARGS  
INITIALIZE SELF SEE HELPTYPESELF FOR ACCURATE SIGNATURE  
GETNSPLITS SELFXNONE YNONE GROUPSNONE  
RETURNS THE NUMBER OF SPLITTING ITERATIONS IN THE CROSSVALIDATOR  
PARAMETERS  
XOBJECT ALWAYS IGNORED EXISTS FOR COMPATIBILITY  
YOBJECT ALWAYS IGNORED EXISTS FOR COMPATIBILITY  
GROUPS ARRAYLIKE WITH SHAPE NSAMPLES GROUP LABELS FOR THE SAMPLES USED WHILE SPLIT  
TING THE DATASET INTO TRAINTEST SET THIS 'GROUPS' PARAMETER MUST ALWAYS BE SPECIFIED TO  
CALCULATE THE NUMBER OF SPLITS THOUGH THE OTHER PARAMETERS CAN BE OMITTED  
RETURNS  
NSPLITS INT RETURNS THE NUMBER OF SPLITTING ITERATIONS IN THE CROSSVALIDATOR  
SPLITSELFXYNONE GROUPSNONE  
GENERATE INDICES TO SPLIT DATA INTO TRAINING AND TEST SET  
PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA WHERE NSAMPLES IS THE NUMBER  
OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES  
YARRAYLIKE OF LENGTH NSAMPLES OPTIONAL THE TARGET VARIABLE FOR SUPERVISED LEARNING  
PROBLEMS  
2074 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

GROUPS ARRAYLIKE WITH SHAPE NSAMPLES GROUP LABELS FOR THE SAMPLES USED WHILE SPLITTING THE DATASET INTO TRAINTEST SET YIELDS

TRAIN NDARRAY THE TRAINING SET INDICES FOR THAT SPLIT

TEST NDARRAY THE TESTING SET INDICES FOR THAT SPLIT

SKLEARNMODELSELECTION LEAVEPGROUPSOUT

CLASSSSKLEARNMODELSELECTION LEAVEPGROUPSOUT NGROUPS

LEAVE P GROUPS OUT CROSSVALIDATOR

PROVIDES TRAINTEST INDICES TO SPLIT DATA ACCORDING TO A THIRDPARTY PROVIDED GROUP THIS GROUP INFORMATION CAN BE USED TO ENCODE ARBITRARY DOMAIN SPECIFIC STRATIFICATIONS OF THE SAMPLES AS INTEGERS

FOR INSTANCE THE GROUPS COULD BE THE YEAR OF COLLECTION OF THE SAMPLES AND THUS ALLOW FOR CROSSVALIDATION AGAINST TIMEBASED SPLITS

THE DIFFERENCE BETWEEN LEAVEPGROUPSOUT AND LEAVEONEGROUPOUT IS THAT THE FORMER BUILDS THE TEST SETS WITH ALL THE SAMPLES ASSIGNED TO PDIFFERENT VALUES OF THE GROUPS WHILE THE LATTER USES SAMPLES ALL ASSIGNED THE SAME GROUPS

READ MORE IN THE USER GUIDE

PARAMETERS

NGROUPS INT NUMBER OF GROUPS P TO LEAVE OUT IN THE TEST SPLIT

SEE ALSO

GROUPKFOLD KFOLD ITERATOR VARIANT WITH NONOVERLAPPING GROUPS

EXAMPLES

```
import numpy as np
from sklearn.model_selection import LeavePGroupsOut
X = np.array(1 2 3 4 5 6)
Y = np.array(1 2 1)
groups = np.array(1 2 3)
lpgo = LeavePGroupsOut(n_groups=2)
lpgo.get_n_splits(X, Y, groups)
3
lpgo.get_n_splits(X, groups, groups) groups is always required
3
print(lpgo)
LeavePGroupsOut(n_groups=2)
for train_index, test_index in lpgo.split(X, Y, groups):
    print(train_index, test_index)
X_train, X_test, X_train_index, X_test_index
Y_train, Y_test, Y_train_index, Y_test_index
print(X_train, X_test, Y_train, Y_test)
train 2 test 0 1
5 6 1 2
3 4 1 1 2
train 1 test 0 2
3 4 1 2
5 6 2 1 1
626sklearn.model_selection model selection 2075
```

SCIKITLEARN USER GUIDE RELEASE 0213

TRAIN 0 TEST 1 2

1 2 3 4

5 6 1 2 1

METHODS

GETNSPLITS SELF X Y GROUPS RETURNS THE NUMBER OF SPLITTING ITERATIONS IN THE CROSS VALIDATOR

SPLIT SELF X Y GROUPS GENERATE INDICES TO SPLIT DATA INTO TRAINING AND TEST SET

INIT SELFNGROUPS

GETNSPLITS SELF X NONE Y NONE GROUPS NONE

RETURNS THE NUMBER OF SPLITTING ITERATIONS IN THE CROSSVALIDATOR

PARAMETERS

X OBJECT ALWAYS IGNORED EXISTS FOR COMPATIBILITY

Y OBJECT ALWAYS IGNORED EXISTS FOR COMPATIBILITY

GROUPS ARRAYLIKE WITH SHAPE NSAMPLES GROUP LABELS FOR THE SAMPLES USED WHILE SPLITTING THE DATASET INTO TRAI NTEST SET THIS 'GROUPS' PARAMETER MUST ALWAYS BE SPECIFIED TO CALCULATE THE NUMBER OF SPLITS THOUGH THE OTHER PARAMETERS CAN BE OMITTED

RETURNS

NSPLITS INT RETURNS THE NUMBER OF SPLITTING ITERATIONS IN THE CROSSVALIDATOR

SPLITSELF X Y NONE GROUPS NONE

GENERATE INDICES TO SPLIT DATA INTO TRAINING AND TEST SET

PARAMETERS

X ARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

Y ARRAYLIKE OF LENGTH NSAMPLES OPTIONAL THE TARGET VARIABLE FOR SUPERVISED LEARNING PROBLEMS

GROUPS ARRAYLIKE WITH SHAPE NSAMPLES GROUP LABELS FOR THE SAMPLES USED WHILE SPLITTING THE DATASET INTO TRAI NTEST SET

YIELDS

TRAIN NDARRAY THE TRAINING SET INDICES FOR THAT SPLIT

TEST NDARRAY THE TESTING SET INDICES FOR THAT SPLIT

SKLEARNMODELSELECTION LEAVEONEOUT

CLASS SKLEARNMODELSELECTION LEAVEONEOUT

LEAVEONEOUT CROSSVALIDATOR

PROVIDES TRAI NTEST INDICES TO SPLIT DATA IN TRAI NTEST SETS EACH SAMPLE IS USED ONCE AS A TEST SET SINGLETON WHILE THE REMAINING SAMPLES FORM THE TRAINING SET

NOTE LEAVEONEOUT IS EQUIVALENT TO KFOLD NSPLITS N AND LEAVEPOUT P1 WHERE N IS THE NUMBER OF SAMPLES

2076 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

DUE TO THE HIGH NUMBER OF TEST SETS WHICH IS THE SAME AS THE NUMBER OF SAMPLES THIS CROSSVALIDATION METHOD CAN BE VERY COSTLY FOR LARGE DATASETS ONE SHOULD FAVOR KFOLD SHUFFLESPLIT ORSTRATIFIEDKFOLD

READ MORE IN THE USER GUIDE

SEE ALSO

LEAVEONEGROUPOUT FOR SPLITTING THE DATA ACCORDING TO EXPLICIT DOMAINSPECIFIC STRATIFICATION OF THE DATASET

GROUPKFOLD KFOLD ITERATOR VARIANT WITH NONOVERLAPPING GROUPS

EXAMPLES

```
import numpy as np
from sklearn.model_selection import LeaveOneOut
X = np.array(1 2 3 4)
Y = np.array(1 2)
loo = LeaveOneOut()
loo.get_n_splits(X)
2
print(loo)
LeaveOneOut()
for train_index, test_index in loo.split(X):
    print('TRAIN:', train_index, 'TEST:', test_index)
X_train, X_test = X[train_index], X[test_index]
Y_train, Y_test = Y[train_index], Y[test_index]
print(X_train, X_test, Y_train, Y_test)
TRAIN 1 TEST 0
3 4 1 2 2 1
TRAIN 0 TEST 1
1 2 3 4 1 2
```

METHODS

GETNSPLITS SELF X Y GROUPS RETURNS THE NUMBER OF SPLITTING ITERATIONS IN THE CROSS VALIDATOR

SPLIT SELF X Y GROUPS GENERATE INDICES TO SPLIT DATA INTO TRAINING AND TEST SET

INIT SELF FARGS KWARGS

INITIALIZE SELF SEE HELPTYPESELF FOR ACCURATE SIGNATURE

GETNSPLITS SELF X Y NONE GROUPS NONE

RETURNS THE NUMBER OF SPLITTING ITERATIONS IN THE CROSSVALIDATOR

PARAMETERS

X ARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

Y OBJECT ALWAYS IGNORED EXISTS FOR COMPATIBILITY

GROUPS OBJECT ALWAYS IGNORED EXISTS FOR COMPATIBILITY

RETURNS

NSPLITS INT RETURNS THE NUMBER OF SPLITTING ITERATIONS IN THE CROSSVALIDATOR

626SKLEARNMODELSELECTION MODEL SELECTION 2077

SCIKITLEARN USER GUIDE RELEASE 0213

SPLITSELFYNONE GROUPSNONE

GENERATE INDICES TO SPLIT DATA INTO TRAINING AND TEST SET

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

YARRAYLIKE OF LENGTH NSAMPLES THE TARGET VARIABLE FOR SUPERVISED LEARNING PROBLEMS

GROUPS ARRAYLIKE WITH SHAPE NSAMPLES OPTIONAL GROUP LABELS FOR THE SAMPLES USED WHILE SPLITTING THE DATASET INTO TRAINTEST SET

YIELDS

TRAIN NDARRAY THE TRAINING SET INDICES FOR THAT SPLIT

TEST NDARRAY THE TESTING SET INDICES FOR THAT SPLIT

SKLEARNMODELSELECTION LEAVEPOUT

CLASSSSKLEARNMODELSELECTION LEAVEPOUT P

LEAVEPOUT CROSSVALIDATOR

PROVIDES TRAINTEST INDICES TO SPLIT DATA IN TRAINTEST SETS THIS RESULTS IN TESTING ON ALL DISTINCT SAMPLES OF SIZE P WHILE THE REMAINING N - P SAMPLES FORM THE TRAINING SET IN EACH ITERATION

NOTELEAVEPOUTP IS NOT EQUIVALENT TO KFOLDNSPLITSNSAMPLES - P WHICH CREATES NON OVERLAPPING TEST SETS

DUE TO THE HIGH NUMBER OF ITERATIONS WHICH GROWS COMBINATORICALLY WITH THE NUMBER OF SAMPLES THIS CROSS VALIDATION METHOD CAN BE VERY COSTLY FOR LARGE DATASETS ONE SHOULD FAVOR KFOLD STRATIFIEDKFOLD OR SHUFFLESPLIT

READ MORE IN THE USER GUIDE

PARAMETERS

PINT SIZE OF THE TEST SETS MUST BE STRICTLY GREATER THAN THE NUMBER OF SAMPLES

EXAMPLES

```
import numpy as np
from sklearn.model_selection import LeavePout
X = np.array(1 2 3 4 5 6 7 8)
Y = np.array(1 2 3 4)
lpo = LeavePout(2)
lpo.get_n_splits(X)
6
print(lpo)
LeavePout(2)
for train_index, test_index in lpo.split(X):
    print('TRAIN: %s, TEST: %s' % (train_index, test_index))
X_train, X_test = X[train_index], X[test_index]
Y_train, Y_test = Y[train_index], Y[test_index]
TRAIN 2 3 TEST 0 1
TRAIN 1 3 TEST 0 2
TRAIN 1 2 TEST 0 3
TRAIN 0 3 TEST 1 2
```

2078 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

TRAIN 0 2 TEST 1 3

TRAIN 0 1 TEST 2 3

METHODS

GETNSPLITS SELF X Y GROUPS RETURNS THE NUMBER OF SPLITTING ITERATIONS IN THE CROSS VALIDATOR

SPLIT SELF X Y GROUPS GENERATE INDICES TO SPLIT DATA INTO TRAINING AND TEST SET

INIT SELF

GETNSPLITS SELFXYNONE GROUPSNONE

RETURNS THE NUMBER OF SPLITTING ITERATIONS IN THE CROSSVALIDATOR

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

YOBJECT ALWAYS IGNORED EXISTS FOR COMPATIBILITY

GROUPS OBJECT ALWAYS IGNORED EXISTS FOR COMPATIBILITY

SPLITSELFXYNONE GROUPSNONE

GENERATE INDICES TO SPLIT DATA INTO TRAINING AND TEST SET

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

YARRAYLIKE OF LENGTH NSAMPLES THE TARGET VARIABLE FOR SUPERVISED LEARNING PROBLEMS

GROUPS ARRAYLIKE WITH SHAPE NSAMPLES OPTIONAL GROUP LABELS FOR THE SAMPLES USED WHILE SPLITTING THE DATASET INTO TRAINTEST SET

YIELDS

TRAIN NDARRAY THE TRAINING SET INDICES FOR THAT SPLIT

TEST NDARRAY THE TESTING SET INDICES FOR THAT SPLIT

SKLEARNMODELSELECTION PREDEFINEDSPLIT

CLASSSKLEARNMODELSELECTION PREDEFINEDSPLIT TESTFOLD

PREDEFINED SPLIT CROSSVALIDATOR

PROVIDES TRAINTEST INDICES TO SPLIT DATA INTO TRAINTEST SETS USING A PREDEFINED SCHEME SPECIFIED BY THE USER WITH THE TESTFOLD PARAMETER

READ MORE IN THE USER GUIDE

PARAMETERS

TESTFOLD ARRAYLIKE SHAPE NSAMPLES THE ENTRY TESTFOLDI REPRESENTS THE INDEX OF THE TEST SET THAT SAMPLE I BELONGS TO IT IS POSSIBLE TO EXCLUDE SAMPLE I FROM ANY TEST SET IE INCLUDE SAMPLE I IN EVERY TRAINING SET BY SETTING TESTFOLDI EQUAL TO 1

626SKLEARNMODELSELECTION MODEL SELECTION 2079

```
SCIKITLEARN USER GUIDE RELEASE 0213
EXAMPLES
  IMPORT NUMPY AS NP
  FROM SKLEARNMODELSELECTION IMPORT PREDEFINEDSPLIT
X  NPARRAY1 2 3 4 1 2 3 4
Y  NPARRAY0 0 1 1
TESTFOLD  0 1 1 1
PS  PREDEFINEDSPLITTESTFOLD
PSGETNSPLITS
2
PRINTPS
PREDEFINEDSPLITTESTFOLDARRAY 0 1 1 1
  FOR TRAININDEX TESTINDEX INPSSPLIT
PRINTTRAIN TRAININDEX TEST TESTINDEX
XTRAIN XTEST  XTRAININDEX XTESTINDEX
YTRAIN YTEST  YTRAININDEX YTESTINDEX
TRAIN 1 2 3 TEST 0
TRAIN 0 2 TEST 1 3
METHODS
GETNSPLITS SELF X Y GROUPS RETURNS THE NUMBER OF SPLITTING ITERATIONS IN THE CROSS
VALIDATOR
SPLIT SELF X Y GROUPS GENERATE INDICES TO SPLIT DATA INTO TRAINING AND TEST SET
INIT SELFTESTFOLD
GETNSPLITS SELFXNONE YNONE GROUPSNONE
RETURNS THE NUMBER OF SPLITTING ITERATIONS IN THE CROSSVALIDATOR
PARAMETERS
XOBJECT ALWAYS IGNORED EXISTS FOR COMPATIBILITY
YOBJECT ALWAYS IGNORED EXISTS FOR COMPATIBILITY
GROUPS OBJECT ALWAYS IGNORED EXISTS FOR COMPATIBILITY
RETURNS
NSPLITS INT RETURNS THE NUMBER OF SPLITTING ITERATIONS IN THE CROSSVALIDATOR
SPLITSELFXNONE YNONE GROUPSNONE
GENERATE INDICES TO SPLIT DATA INTO TRAINING AND TEST SET
PARAMETERS
XOBJECT ALWAYS IGNORED EXISTS FOR COMPATIBILITY
YOBJECT ALWAYS IGNORED EXISTS FOR COMPATIBILITY
GROUPS OBJECT ALWAYS IGNORED EXISTS FOR COMPATIBILITY
YIELDS
TRAIN NDARRAY THE TRAINING SET INDICES FOR THAT SPLIT
TEST NDARRAY THE TESTING SET INDICES FOR THAT SPLIT
2080 CHAPTER 6 API REFERENCE
```

SCIKITLEARN USER GUIDE RELEASE 0213  
SKLEARNMODELSELECTION REPEATEDKFOLD  
CLASSSSKLEARNMODELSELECTION REPEATEDKFOLD NSPLITS5 NREPEATS10 RAN  
DOMSTATENONE  
REPEATED KFOLD CROSS VALIDATOR  
REPEATS KFOLD N TIMES WITH DIFFERENT RANDOMIZATION IN EACH REPETITION  
READ MORE IN THE USER GUIDE  
PARAMETERS  
NSPLITS INT DEFAULT5 NUMBER OF FOLDS MUST BE AT LEAST 2  
NREPEATS INT DEFAULT10 NUMBER OF TIMES CROSSVALIDATOR NEEDS TO BE REPEATED  
RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RAN  
DOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE  
RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS  
THE RANDOMSTATE INSTANCE USED BY NPRANDOM  
SEE ALSO  
REPEATEDSTRATIFIEDKFOLD REPEATS STRATIFIED KFOLD N TIMES  
NOTES  
RANDOMIZED CV SPLITTERS MAY RETURN DIFFERENT RESULTS FOR EACH CALL OF SPLIT YOU CAN MAKE THE RESULTS IDENTICAL  
BY SETTINGRANDOMSTATE TO AN INTEGER  
EXAMPLES  
IMPORT NUMPY AS NP  
FROM SKLEARNMODELSELECTION IMPORT REPEATEDKFOLD  
X NPARRAY1 2 3 4 1 2 3 4  
Y NPARRAY0 0 1 1  
RKF REPEATEDKFOLDNSPLITS2 NREPEATS2 RANDOMSTATE2652124  
FOR TRAININDEX TESTINDEX INRKFSPPLITX  
PRINTTRAIN TRAININDEX TEST TESTINDEX  
XTRAIN XTEST XTRAININDEX XTESTINDEX  
YTRAIN YTEST YTRAININDEX YTESTINDEX  
  
TRAIN 0 1 TEST 2 3  
TRAIN 2 3 TEST 0 1  
TRAIN 1 2 TEST 0 3  
TRAIN 0 3 TEST 1 2  
METHODS  
GETNSPLITS SELF X Y GROUPS RETURNS THE NUMBER OF SPLITTING ITERATIONS IN THE CROSS  
VALIDATOR  
SPLIT SELF X Y GROUPS GENERATES INDICES TO SPLIT DATA INTO TRAINING AND TEST SET  
INIT SELFNSPLITS5 NREPEATS10 RANDOMSTATENONE  
626SKLEARNMODELSELECTION MODEL SELECTION 2081

SCIKITLEARN USER GUIDE RELEASE 0213  
GETNSPLITS SELFYNONE YNONE GROUPSNONE  
RETURNS THE NUMBER OF SPLITTING ITERATIONS IN THE CROSSVALIDATOR  
PARAMETERS  
XOBJECT ALWAYS IGNORED EXISTS FOR COMPATIBILITY NPZEROSNSAMPLES MAY BE  
USED AS A PLACEHOLDER  
YOBJECT ALWAYS IGNORED EXISTS FOR COMPATIBILITY NPZEROSNSAMPLES MAY BE  
USED AS A PLACEHOLDER  
GROUPS ARRAYLIKE WITH SHAPE NSAMPLES OPTIONAL GROUP LABELS FOR THE SAMPLES USED  
WHILE SPLITTING THE DATASET INTO TRAINTEST SET  
RETURNS  
NSPLITS INT RETURNS THE NUMBER OF SPLITTING ITERATIONS IN THE CROSSVALIDATOR  
SPLITSELFYNONE GROUPSNONE  
GENERATES INDICES TO SPLIT DATA INTO TRAINING AND TEST SET  
PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA WHERE NSAMPLES IS THE NUMBER  
OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES  
YARRAYLIKE OF LENGTH NSAMPLES THE TARGET VARIABLE FOR SUPERVISED LEARNING PROBLEMS  
GROUPS ARRAYLIKE WITH SHAPE NSAMPLES OPTIONAL GROUP LABELS FOR THE SAMPLES USED  
WHILE SPLITTING THE DATASET INTO TRAINTEST SET  
YIELDS  
TRAIN NDARRAY THE TRAINING SET INDICES FOR THAT SPLIT  
TEST NDARRAY THE TESTING SET INDICES FOR THAT SPLIT  
SKLEARNMODELSELECTION REPEATEDSTRATIFIEDKFOLD  
CLASSSSKLEARNMODELSELECTION REPEATEDSTRATIFIEDKFOLD NSPLITS5 NREPEATS10 RAN  
DOMSTATENONE  
REPEATED STRATIFIED KFOLD CROSS VALIDATOR  
REPEATS STRATIFIED KFOLD N TIMES WITH DIFFERENT RANDOMIZATION IN EACH REPETITION  
READ MORE IN THE USER GUIDE  
PARAMETERS  
NSPLITS INT DEFAULT5 NUMBER OF FOLDS MUST BE AT LEAST 2  
NREPEATS INT DEFAULT10 NUMBER OF TIMES CROSSVALIDATOR NEEDS TO BE REPEATED  
RANDOMSTATE NONE INT OR RANDOMSTATE DEFAULTNONE RANDOM STATE TO BE USED TO GENERATE  
RANDOM STATE FOR EACH REPETITION  
SEE ALSO  
REPEATEDKFOLD REPEATS KFOLD N TIMES  
2082 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

NOTES

RANDOMIZED CV SPLITTERS MAY RETURN DIFFERENT RESULTS FOR EACH CALL OF SPLIT YOU CAN MAKE THE RESULTS IDENTICAL BY SETTINGRANDOMSTATE TO AN INTEGER

EXAMPLES

```
import numpy as np
from sklearn.model_selection import RepeatedStratifiedKFold
X = np.array([2, 3, 4, 1, 2, 3, 4])
Y = np.array([0, 1, 1])
rskf = RepeatedStratifiedKFold(n_splits=2, n_repeats=2,
                                random_state=36851234)
for train_index, test_index in rskf.split(X):
    print('TRAIN: %s TEST: %s' % (train_index, test_index))
    x_train, x_test, y_train, y_test = X[train_index], X[test_index], Y[train_index], Y[test_index]
```

TRAIN 1 2 TEST 0 3  
TRAIN 0 3 TEST 1 2  
TRAIN 1 3 TEST 0 2  
TRAIN 0 2 TEST 1 3

METHODS  
get\_n\_splits(self, X, Y, groups) returns the number of splitting iterations in the cross validator

split(self, X, Y, groups) generates indices to split data into training and test set

init(self, n\_splits=5, n\_repeats=10, random\_state=None)

get\_n\_splits(self, X=None, Y=None, groups=None)

returns the number of splitting iterations in the cross validator

PARAMETERS

X OBJECT ALWAYS IGNORED EXISTS FOR COMPATIBILITY NPZEROSNSAMPLES MAY BE

USED AS A PLACEHOLDER

Y OBJECT ALWAYS IGNORED EXISTS FOR COMPATIBILITY NPZEROSNSAMPLES MAY BE

USED AS A PLACEHOLDER

GROUPS ARRAYLIKE WITH SHAPE NSAMPLES OPTIONAL GROUP LABELS FOR THE SAMPLES USED

WHILE SPLITTING THE DATASET INTO TRAINTEST SET

RETURNS

nsplits int returns the number of splitting iterations in the cross validator

splitself X Y None groups None

generates indices to split data into training and test set

PARAMETERS

X ARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA WHERE NSAMPLES IS THE NUMBER

OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

SCIKITLEARN USER GUIDE RELEASE 0213  
 YARRAYLIKE OF LENGTH NSAMPLES THE TARGET VARIABLE FOR SUPERVISED LEARNING PROBLEMS  
 GROUPS ARRAYLIKE WITH SHAPE NSAMPLES OPTIONAL GROUP LABELS FOR THE SAMPLES USED  
 WHILE SPLITTING THE DATASET INTO TRAINTEST SET  
 YIELDS  
 TRAIN NDARRAY THE TRAINING SET INDICES FOR THAT SPLIT  
 TEST NDARRAY THE TESTING SET INDICES FOR THAT SPLIT  
 SKLEARNMODELSELECTION SHUFFLESPLIT  
 CLASSSSKLEARNMODELSELECTION SHUFFLESPLIT NSPLITS10 TESTSIZENONE TRAINSIZENONE  
 RANDOMSTATENONE  
 RANDOM PERMUTATION CROSSVALIDATOR  
 YIELDS INDICES TO SPLIT DATA INTO TRAINING AND TEST SETS  
 NOTE CONTRARY TO OTHER CROSSVALIDATION STRATEGIES RANDOM SPLITS DO NOT GUARANTEE THAT ALL FOLDS WILL BE DIFFERENT  
 ALTHOUGH THIS IS STILL VERY LIKELY FOR SIZEABLE DATASETS  
 READ MORE IN THE USER GUIDE  
 PARAMETERS  
 NSPLITS INT DEFAULT 10 NUMBER OF RESHUFFLING SPLITTING ITERATIONS  
 TESTSIZE FLOAT INT NONE DEFAULTNONE IF FLOAT SHOULD BE BETWEEN 00 AND 10 AND REPRE  
 SENT THE PROPORTION OF THE DATASET TO INCLUDE IN THE TEST SPLIT IF INT REPRESENTS THE ABSOLUTE  
 NUMBER OF TEST SAMPLES IF NONE THE VALUE IS SET TO THE COMPLEMENT OF THE TRAIN SIZE IF  
 TRAINSIZE IS ALSO NONE IT WILL BE SET TO 01  
 TRAINSIZE FLOAT INT OR NONE DEFAULTNONE IF FLOAT SHOULD BE BETWEEN 00 AND 10 AND REPRE  
 SENT THE PROPORTION OF THE DATASET TO INCLUDE IN THE TRAIN SPLIT IF INT REPRESENTS THE ABSOLUTE  
 NUMBER OF TRAIN SAMPLES IF NONE THE VALUE IS AUTOMATICALLY SET TO THE COMPLEMENT OF THE  
 TEST SIZE  
 RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RAN  
 DOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE  
 RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS  
 THE RANDOMSTATE INSTANCE USED BY NPRANDOM  
 EXAMPLES  
 IMPORT NUMPY AS NP  
 FROM SKLEARNMODELSELECTION IMPORT SHUFFLESPLIT  
 X NPARRAY1 2 3 4 5 6 7 8 3 4 5 6  
 Y NPARRAY1 2 1 2 1 2  
 RS SHUFFLESPLITNSPLITS5 TESTSIZE25 RANDOMSTATE0  
 RSGETNSPLITSX  
 5  
 PRINTRS  
 SHUFFLESPLITNSPLITS5 RANDOMSTATE0 TESTSIZE025 TRAINSIZENONE  
 FOR TRAININDEX TESTINDEX INRSPLITX  
 PRINTTRAIN TRAININDEX TEST TESTINDEX  
  
 TRAIN 1 3 0 4 TEST 5 2  
 TRAIN 4 0 2 5 TEST 1 3  
 2084 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

TRAIN 1 2 4 0 TEST 3 5

TRAIN 3 4 1 0 TEST 5 2

TRAIN 3 5 1 0 TEST 2 4

RS SHUFFLESPLITNSPLITS5 TRAINSIZE05 TESTSIZE25

RANDOMSTATE0

FOR TRAININDEX TESTINDEX INRSSPLITX

PRINTTRAIN TRAININDEX TEST TESTINDEX

  

TRAIN 1 3 0 TEST 5 2

TRAIN 4 0 2 TEST 1 3

TRAIN 1 2 4 TEST 3 5

TRAIN 3 4 1 TEST 5 2

TRAIN 3 5 1 TEST 2 4

METHODS

GETNSPLITS SELF X Y GROUPS RETURNS THE NUMBER OF SPLITTING ITERATIONS IN THE CROSS VALIDATOR

SPLIT SELF X Y GROUPS GENERATE INDICES TO SPLIT DATA INTO TRAINING AND TEST SET

INIT SELFNSPLITS10 TESTSIZENONE TRAINSIZENONE RANDOMSTATENONE

GETNSPLITS SELFYNONE YNONE GROUPSNONE

RETURNS THE NUMBER OF SPLITTING ITERATIONS IN THE CROSSVALIDATOR

PARAMETERS

XOBJECT ALWAYS IGNORED EXISTS FOR COMPATIBILITY

YOBJECT ALWAYS IGNORED EXISTS FOR COMPATIBILITY

GROUPS OBJECT ALWAYS IGNORED EXISTS FOR COMPATIBILITY

RETURNS

NSPLITS INT RETURNS THE NUMBER OF SPLITTING ITERATIONS IN THE CROSSVALIDATOR

SPLITSELFYNONE GROUPSNONE

GENERATE INDICES TO SPLIT DATA INTO TRAINING AND TEST SET

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

YARRAYLIKE SHAPE NSAMPLES THE TARGET VARIABLE FOR SUPERVISED LEARNING PROBLEMS

GROUPS ARRAYLIKE WITH SHAPE NSAMPLES OPTIONAL GROUP LABELS FOR THE SAMPLES USED WHILE SPLITTING THE DATASET INTO TRAINTEST SET

YIELDS

TRAIN NDARRAY THE TRAINING SET INDICES FOR THAT SPLIT

TEST NDARRAY THE TESTING SET INDICES FOR THAT SPLIT

626SKLEARNMODELSELECTION MODEL SELECTION 2085

SCIKITLEARN USER GUIDE RELEASE 0213

NOTES

RANDOMIZED CV SPLITTERS MAY RETURN DIFFERENT RESULTS FOR EACH CALL OF SPLIT YOU CAN MAKE THE RESULTS IDENTICAL BY SETTING RANDOMSTATE TO AN INTEGER

EXAMPLES USING SKLEARNMODELSELECTIONSHUFFLESPLIT

•VISUALIZING CROSSVALIDATION BEHAVIOR IN SCIKITLEARN

•PLOTING LEARNING CURVES

•SCALING THE REGULARIZATION PARAMETER FOR SVCS

SKLEARNMODELSELECTION STRATIFIEDKFOLD

CLASSSKLEARNMODELSELECTION STRATIFIEDKFOLD NSPLITS'WARN' SHUFFLEFALSE RAN

DOMSTATENONE

STRATIFIED KFOLDS CROSSVALIDATOR

PROVIDES TRAINTEST INDICES TO SPLIT DATA IN TRAINTEST SETS

THIS CROSSVALIDATION OBJECT IS A VARIATION OF KFOLD THAT RETURNS STRATIFIED FOLDS THE FOLDS ARE MADE BY PRESERVING THE PERCENTAGE OF SAMPLES FOR EACH CLASS

READ MORE IN THE USER GUIDE

PARAMETERS

NSPLITS INT DEFAULT3 NUMBER OF FOLDS MUST BE AT LEAST 2

CHANGED IN VERSION 020 NSPLITS DEFAULT VALUE WILL CHANGE FROM 3 TO 5 IN V022

SHUFFLE BOOLEAN OPTIONAL WHETHER TO SHUFFLE EACH CLASS'S SAMPLES BEFORE SPLITTING INTO BATCHES

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RAN

DOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE

RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS

THE RANDOMSTATE INSTANCE USED BY NPRANDOM USED WHEN SHUFFLE TRUE

SEE ALSO

REPEATEDSTRATIFIEDKFOLD REPEATS STRATIFIED KFOLD N TIMES

NOTES

TRAIN AND TEST SIZES MAY BE DIFFERENT IN EACH FOLD WITH A DIFFERENCE OF AT MOST NCLASSES

EXAMPLES

IMPORT NUMPY AS NP

FROM SKLEARNMODELSELECTION IMPORT STRATIFIEDKFOLD

X NPARRAY1 2 3 4 1 2 3 4

Y NPARRAY0 0 1 1

SKF STRATIFIEDKFOLDNSPLITS2

SKFGETNSPLITSX Y

2

2086 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

PRINTSKF

STRATIFIEDKFOLDNSPLITS2 RANDOMSTATENONE SHUFFLEFALSE

FOR TRAININDEX TESTINDEX INSKFSPLITX Y

PRINTTRAIN TRAININDEX TEST TESTINDEX

XTRAIN XTEST XTRAININDEX XTESTINDEX

YTRAIN YTEST YTRAININDEX YTESTINDEX

TRAIN 1 3 TEST 0 2

TRAIN 0 2 TEST 1 3

METHODS

GETNSPLITS SELF X Y GROUPS RETURNS THE NUMBER OF SPLITTING ITERATIONS IN THE CROSS VALIDATOR

SPLIT SELF X Y GROUPS GENERATE INDICES TO SPLIT DATA INTO TRAINING AND TEST SET

INIT SELFNSPLITS'WARN' SHUFFLEFALSE RANDOMSTATENONE

GETNSPLITS SELFXNONE YNONE GROUPSNONE

RETURNS THE NUMBER OF SPLITTING ITERATIONS IN THE CROSSVALIDATOR

PARAMETERS

XOBJECT ALWAYS IGNORED EXISTS FOR COMPATIBILITY

YOBJECT ALWAYS IGNORED EXISTS FOR COMPATIBILITY

GROUPS OBJECT ALWAYS IGNORED EXISTS FOR COMPATIBILITY

RETURNS

NSPLITS INT RETURNS THE NUMBER OF SPLITTING ITERATIONS IN THE CROSSVALIDATOR

SPLITSELFXYGROUPSNONE

GENERATE INDICES TO SPLIT DATA INTO TRAINING AND TEST SET

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

NOTE THAT PROVIDING YIS SUFFICIENT TO GENERATE THE SPLITS AND HENCE NP

ZEROSNSAMPLES MAY BE USED AS A PLACEHOLDER FOR XINSTEAD OF ACTUAL TRAINING DATA

YARRAYLIKE SHAPE NSAMPLES THE TARGET VARIABLE FOR SUPERVISED LEARNING PROBLEMS

STRATIFICATION IS DONE BASED ON THE Y LABELS

GROUPS OBJECT ALWAYS IGNORED EXISTS FOR COMPATIBILITY

YIELDS

TRAIN NDARRAY THE TRAINING SET INDICES FOR THAT SPLIT

TEST NDARRAY THE TESTING SET INDICES FOR THAT SPLIT

NOTES

RANDOMIZED CV SPLITTERS MAY RETURN DIFFERENT RESULTS FOR EACH CALL OF SPLIT YOU CAN MAKE THE RESULTS IDENTICAL BY SETTING RANDOMSTATE TO AN INTEGER

6265KLEARNMODELSELECTION MODEL SELECTION 2087

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES USING SKLEARNMODELSELECTIONSTRATIFIEDKFOLD

- RECURSIVE FEATURE ELIMINATION WITH CROSSVALIDATION
- TEST WITH PERMUTATIONS THE SIGNIFICANCE OF A CLASSIFICATION SCORE
- GMM COVARIANCES
- RECEIVER OPERATING CHARACTERISTIC ROC WITH CROSS VALIDATION
- VISUALIZING CROSSVALIDATION BEHAVIOR IN SCIKITLEARN

SKLEARNMODELSELECTION STRATIFIEDSHUFFLESPLIT

CLASSSSKLEARNMODELSELECTION STRATIFIEDSHUFFLESPLIT NSPLITS10 TESTSIZENONE

TRAINSIZENONE RAN

DOMSTATENONE

STRATIFIED SHUFFLESPLIT CROSSVALIDATOR

PROVIDES TRAINTEST INDICES TO SPLIT DATA IN TRAINTEST SETS

THIS CROSSVALIDATION OBJECT IS A MERGE OF STRATIFIEDKFOLD AND SHUFFLESPLIT WHICH RETURNS STRATIFIED RANDOMIZED FOLDS THE FOLDS ARE MADE BY PRESERVING THE PERCENTAGE OF SAMPLES FOR EACH CLASS

NOTE LIKE THE SHUFFLESPLIT STRATEGY STRATIFIED RANDOM SPLITS DO NOT GUARANTEE THAT ALL FOLDS WILL BE DIFFERENT ALTHOUGH THIS IS STILL VERY LIKELY FOR SIZEABLE DATASETS

READ MORE IN THE USER GUIDE

PARAMETERS

NSPLITS INT DEFAULT 10 NUMBER OF RESHUFFLING SPLITTING ITERATIONS

TESTSIZE FLOAT INT NONE OPTIONAL DEFAULTNONE IF FLOAT SHOULD BE BETWEEN 00 AND 10

AND REPRESENT THE PROPORTION OF THE DATASET TO INCLUDE IN THE TEST SPLIT IF INT REPRESENTS THE ABSOLUTE NUMBER OF TEST SAMPLES IF NONE THE VALUE IS SET TO THE COMPLEMENT OF THE TRAIN SIZE IFTRAINSIZENONE IT WILL BE SET TO 01

TRAINSIZENONE FLOAT INT OR NONE DEFAULT IS NONE IF FLOAT SHOULD BE BETWEEN 00 AND 10 AND REPRESENT THE PROPORTION OF THE DATASET TO INCLUDE IN THE TRAIN SPLIT IF INT REPRESENTS THE ABSOLUTE NUMBER OF TRAIN SAMPLES IF NONE THE VALUE IS AUTOMATICALLY SET TO THE COMPLEMENT OF THE TEST SIZE

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RAN

DOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE

RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS

THE RANDOMSTATE INSTANCE USED BY NRANDOM

EXAMPLES

```
import numpy as np
from sklearn.model_selection import StratifiedShuffleSplit
X = np.array([1, 2, 3, 4, 1, 2, 3, 4, 1, 2, 3, 4])
Y = np.array([0, 0, 0, 1, 1, 1, 1, 1, 0, 0, 0, 0])
sss = StratifiedShuffleSplit(n_splits=5, test_size=0.5, random_state=0)
sss.get_n_splits(X, Y)
5
print(sss)
StratifiedShuffleSplit(n_splits=5, test_size=0.5, random_state=0)
```

2088 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213  
FOR TRAININDEX TESTINDEX INSSSSPLITX Y  
PRINTTRAIN TRAININDEX TEST TESTINDEX  
XTRAIN XTEST XTRAININDEX XTESTINDEX  
YTRAIN YTEST YTRAININDEX YTESTINDEX  
TRAIN 5 2 3 TEST 4 1 0  
TRAIN 5 1 4 TEST 0 2 3  
TRAIN 5 0 2 TEST 4 3 1  
TRAIN 4 1 0 TEST 2 3 5  
TRAIN 0 5 1 TEST 3 4 2  
METHODS  
GETNSPLITS SELF X Y GROUPS RETURNS THE NUMBER OF SPLITTING ITERATIONS IN THE CROSS  
VALIDATOR  
SPLIT SELF X Y GROUPS GENERATE INDICES TO SPLIT DATA INTO TRAINING AND TEST SET  
INIT SELFNSPLITS10 TESTSIZENONE TRAINSIZENONE RANDOMSTATENONE  
GETNSPLITS SELFXXNONE YNONE GROUPSNONE  
RETURNS THE NUMBER OF SPLITTING ITERATIONS IN THE CROSSVALIDATOR  
PARAMETERS  
XOBJECT ALWAYS IGNORED EXISTS FOR COMPATIBILITY  
YOBJECT ALWAYS IGNORED EXISTS FOR COMPATIBILITY  
GROUPS OBJECT ALWAYS IGNORED EXISTS FOR COMPATIBILITY  
RETURNS  
NSPLITS INT RETURNS THE NUMBER OF SPLITTING ITERATIONS IN THE CROSSVALIDATOR  
SPLITSELFXYGROUPSNONE  
GENERATE INDICES TO SPLIT DATA INTO TRAINING AND TEST SET  
PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA WHERE NSAMPLES IS THE NUMBER  
OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES  
NOTE THAT PROVIDING YIS SUFFICIENT TO GENERATE THE SPLITS AND HENCE NP  
ZEROSNSAMPLES MAY BE USED AS A PLACEHOLDER FOR XINSTEAD OF ACTUAL TRAINING DATA  
YARRAYLIKE SHAPE NSAMPLES THE TARGET VARIABLE FOR SUPERVISED LEARNING PROBLEMS  
STRATIFICATION IS DONE BASED ON THE Y LABELS  
GROUPS OBJECT ALWAYS IGNORED EXISTS FOR COMPATIBILITY  
YIELDS  
TRAIN NDARRAY THE TRAINING SET INDICES FOR THAT SPLIT  
TEST NDARRAY THE TESTING SET INDICES FOR THAT SPLIT  
NOTES  
RANDOMIZED CV SPLITTERS MAY RETURN DIFFERENT RESULTS FOR EACH CALL OF SPLIT YOU CAN MAKE THE RESULTS  
IDENTICAL BY SETTING RANDOMSTATE TO AN INTEGER  
6265KLEARNMODELSELECTION MODEL SELECTION 2089

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES USING SKLEARNMODELSELECTIONSTRATIFIEDSHUFFLESPLIT

- VISUALIZING CROSSVALIDATION BEHAVIOR IN SCIKITLEARN
- RBF SVM PARAMETERS

SKLEARNMODELSELECTION TIMESERIESSPLIT

CLASSSSKLEARNMODELSELECTION TIMESERIESSPLIT NSPLITS'WARN' MAXTRAINSIZEONE

TIME SERIES CROSSVALIDATOR

PROVIDES TRAINTEST INDICES TO SPLIT TIME SERIES DATA SAMPLES THAT ARE OBSERVED AT FIXED TIME INTERVALS IN TRAINTEST SETS IN EACH SPLIT TEST INDICES MUST BE HIGHER THAN BEFORE AND THUS SHUFFLING IN CROSS VALIDATOR IS INAPPROPRIATE THIS CROSSVALIDATION OBJECT IS A VARIATION OF KFOLD IN THE KTH SPLIT IT RETURNS FIRST K FOLDS AS TRAIN SET AND THE K1TH FOLD AS TEST SET

NOTE THAT UNLIKE STANDARD CROSSVALIDATION METHODS SUCCESSIVE TRAINING SETS ARE SUPERSETS OF THOSE THAT COME BEFORE THEM

READ MORE IN THE USER GUIDE

PARAMETERS

NSPLITS INT DEFAULT3 NUMBER OF SPLITS MUST BE AT LEAST 2

CHANGED IN VERSION 020 NSPLITS DEFAULT VALUE WILL CHANGE FROM 3 TO 5 IN V022

MAXTRAINSIZE INT OPTIONAL MAXIMUM SIZE FOR A SINGLE TRAINING SET

NOTES

THE TRAINING SET HAS SIZE INSAMPLES NSPLITS 1 NSAMPLES NSPLITS 1 IN THEITH SPLIT WITH A TEST SET OF SIZE NSAMPLESNSPLITS 1

WHERE NSAMPLES IS THE NUMBER OF SAMPLES

EXAMPLES

```
import numpy as np
from sklearn.model_selection import TimeSeriesSplit
X = np.array([1, 2, 3, 4, 1, 2, 3, 4, 1, 2, 3, 4])
Y = np.array([1, 2, 3, 4, 5, 6])
tscv = TimeSeriesSplit(n_splits=5)
print(tscv)
TimeSeriesSplit(max_train_size=None, n_splits=5)
for train_index, test_index in tscv.split(X):
    print('TRAIN:', train_index, 'TEST:', test_index)
X_train, X_test = X[train_index], X[test_index]
Y_train, Y_test = Y[train_index], Y[test_index]
train_0_test_1
train_0_1_test_2
train_0_1_2_test_3
train_0_1_2_3_test_4
train_0_1_2_3_4_test_5
```

2090 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

METHODS

GETNSPLITS SELF X Y GROUPS RETURNS THE NUMBER OF SPLITTING ITERATIONS IN THE CROSS VALIDATOR

SPLIT SELF X Y GROUPS GENERATE INDICES TO SPLIT DATA INTO TRAINING AND TEST SET

INIT SELFNSPLITS'WARN' MAXTRAINSIZENONE

GETNSPLITS SELFXXNONE YNONE GROUPSNONE

RETURNS THE NUMBER OF SPLITTING ITERATIONS IN THE CROSSVALIDATOR

PARAMETERS

XOBJECT ALWAYS IGNORED EXISTS FOR COMPATIBILITY

YOBJECT ALWAYS IGNORED EXISTS FOR COMPATIBILITY

GROUPS OBJECT ALWAYS IGNORED EXISTS FOR COMPATIBILITY

RETURNS

NSPLITS INT RETURNS THE NUMBER OF SPLITTING ITERATIONS IN THE CROSSVALIDATOR

SPLITSELFXYNONE GROUPSNONE

GENERATE INDICES TO SPLIT DATA INTO TRAINING AND TEST SET

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

YARRAYLIKE SHAPE NSAMPLES ALWAYS IGNORED EXISTS FOR COMPATIBILITY

GROUPS ARRAYLIKE WITH SHAPE NSAMPLES ALWAYS IGNORED EXISTS FOR COMPATIBILITY

YIELDS

TRAIN NDARRAY THE TRAINING SET INDICES FOR THAT SPLIT

TEST NDARRAY THE TESTING SET INDICES FOR THAT SPLIT

EXAMPLES USING SKLEARNMODELSELECTIONTIMESERIESSPLIT

•VISUALIZING CROSSVALIDATION BEHAVIOR IN SCIKITLEARN

6262 SPLITTER FUNCTIONS

MODELSELECTIONCHECKCV CV Y CLASSIFIER INPUT CHECKER UTILITY FOR BUILDING A CROSSVALIDATOR

MODELSELECTIONTRAINTESTSPLIT ARRAYS

    SPLIT ARRAYS OR MATRICES INTO RANDOM TRAIN AND TEST SUBSETS

SKLEARNMODELSELECTION CHECKCV

SKLEARNMODELSELECTION CHECKCV CV'WARN' YNONE CLASSIFIERFALSE

INPUT CHECKER UTILITY FOR BUILDING A CROSSVALIDATOR

PARAMETERS

6265SKLEARNMODELSELECTION  MODEL SELECTION 2091

SCIKITLEARN USER GUIDE RELEASE 0213

CVINT CROSSVALIDATION GENERATOR OR AN ITERABLE OPTIONAL DETERMINES THE CROSSVALIDATION SPLITTING STRATEGY POSSIBLE INPUTS FOR CV ARE

- NONE TO USE THE DEFAULT 3FOLD CROSSVALIDATION
- INTEGER TO SPECIFY THE NUMBER OF FOLDS
- CV SPLITTER
- AN ITERABLE YIELDING TRAIN TEST SPLITS AS ARRAYS OF INDICES

FOR INTEGER NONE INPUTS IF CLASSIFIER IS TRUE AND YIS EITHER BINARY OR MULTICLASS STRATIFIEDKFOLD IS USED IN ALL OTHER CASES KFOLD IS USED

REFER USER GUIDE FOR THE VARIOUS CROSSVALIDATION STRATEGIES THAT CAN BE USED HERE

CHANGED IN VERSION 020 CVDEFAULT VALUE WILL CHANGE FROM 3FOLD TO 5FOLD IN V022

YARRAYLIKE OPTIONAL THE TARGET VARIABLE FOR SUPERVISED LEARNING PROBLEMS

CLASSIFIER BOOLEAN OPTIONAL DEFAULT FALSE WHETHER THE TASK IS A CLASSIFICATION TASK IN WHICH CASE STRATIFIED KFOLD WILL BE USED

RETURNS

CHECKEDCV A CROSSVALIDATOR INSTANCE THE RETURN VALUE IS A CROSSVALIDATOR WHICH GENERATES THE TRAINTEST SPLITS VIA THE SPLIT METHOD

SKLEARNMODELSELECTION TRAINTESTSPLIT

SKLEARNMODELSELECTION TRAINTESTSPLIT ARRAYS OPTIONS

SPLIT ARRAYS OR MATRICES INTO RANDOM TRAIN AND TEST SUBSETS

QUICK UTILITY THAT WRAPS INPUT VALIDATION AND NEXTSHUFFLESPLITSPLITX Y AND APPLICATION TO INPUT DATA INTO A SINGLE CALL FOR SPLITTING AND OPTIONALLY SUBSAMPLING DATA IN A ONELINER

READ MORE IN THE USER GUIDE

PARAMETERS

ARRAYS SEQUENCE OF INDEXABLES WITH SAME LENGTH SHAPE0 ALLOWED INPUTS ARE LISTS NUMPY ARRAYS SCIPYSPARSE MATRICES OR PANDAS DATAFRAMES

TESTSIZE FLOAT INT OR NONE OPTIONAL DEFAULTNONE IF FLOAT SHOULD BE BETWEEN 00 AND 10 AND REPRESENT THE PROPORTION OF THE DATASET TO INCLUDE IN THE TEST SPLIT IF INT REPRESENTS THE ABSOLUTE NUMBER OF TEST SAMPLES IF NONE THE VALUE IS SET TO THE COMPLEMENT OF THE TRAIN SIZE IFTRAINSIZ IS ALSO NONE IT WILL BE SET TO 025

TRAINSIZ FLOAT INT OR NONE DEFAULTNONE IF FLOAT SHOULD BE BETWEEN 00 AND 10 AND REPRESENT THE PROPORTION OF THE DATASET TO INCLUDE IN THE TRAIN SPLIT IF INT REPRESENTS THE ABSOLUTE NUMBER OF TRAIN SAMPLES IF NONE THE VALUE IS AUTOMATICALLY SET TO THE COMPLEMENT OF THE TEST SIZE

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RANDOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NRANDOM

SHUFFLE BOOLEAN OPTIONAL DEFAULTTRUE WHETHER OR NOT TO SHUFFLE THE DATA BEFORE SPLITTING IF SHUFFLEFALSE THEN STRATIFY MUST BE NONE

STRATIFY ARRAYLIKE OR NONE DEFAULTNONE IF NOT NONE DATA IS SPLIT IN A STRATIFIED FASHION USING THIS AS THE CLASS LABELS

2092 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

RETURNS

SPLITTING LIST LENGTH2 LENARRAYS LIST CONTAINING TRAINTEST SPLIT OF INPUTS  
NEW IN VERSION 016 IF THE INPUT IS SPARSE THE OUTPUT WILL BE A SCIPYSPARSE  
CSRMATRIX ELSE OUTPUT TYPE IS THE SAME AS THE INPUT TYPE

EXAMPLES

```
import numpy as np
from sklearn.model_selection import train_test_split
X, Y, NPARANGE10, RESHAPE5, 2, RANGE5
```

```
X
array([[0, 1],
       [2, 3],
       [4, 5],
       [6, 7],
       [8, 9]])
list([0, 1, 2, 3, 4])
XTRAIN, XTEST, YTRAIN, YTEST, TRAINTESTSPLIT,
X, Y, TESTSIZE033, RANDOMSTATE42
```

```
XTRAIN
array([[0, 1],
       [6, 7]])
YTRAIN
array([2, 0, 3])
XTEST
array([[8, 9]])
YTEST
array([1, 4])
TRAINTESTSPLIT, SHUFFLE, FALSE
[0, 1, 2, 3, 4]
```

EXAMPLES USING SKLEARNMODELSELECTIONTRAINTESTSPLIT

- FACES RECOGNITION EXAMPLE USING EIGENFACES AND SVMs
  - PREDICTION LATENCY
  - PROBABILITY CALIBRATION CURVES
  - PROBABILITY CALIBRATION OF CLASSIFIERS
  - CLASSIFIER COMPARISON
  - COLUMN TRANSFORMER WITH MIXED TYPES
  - EFFECT OF TRANSFORMING THE TARGETS IN REGRESSION MODEL
  - COMPARING RANDOM FORESTS AND THE MULTIOUTPUT META ESTIMATOR
  - EARLY STOPPING OF GRADIENT BOOSTING
- 626SKLEARNMODELSELECTION MODEL SELECTION 2093

SCIKITLEARN USER GUIDE RELEASE 0213

- FEATURE TRANSFORMATIONS WITH ENSEMBLES OF TREES
- GRADIENT BOOSTING OUTFBAG ESTIMATES
- PIPELINE ANOVA SVM
- COMPARING VARIOUS ONLINE SOLVERS
- MNIST CLASSFICATION USING MULTINOMIAL LOGISTIC L1
- MULTICLASS SPARSE LOGISITIC REGRESSION ON NEWGROUPS20
- EARLY STOPPING OF STOCHASTIC GRADIENT DESCENT
- PARAMETER ESTIMATION USING GRID SEARCH WITH CROSSVALIDATION
- CONFUSION MATRIX
- RECEIVER OPERATING CHARACTERISTIC ROC
- PRECISIONRECALL
- CLASSIFIER CHAIN
- COMPARING NEAREST NEIGHBORS WITH AND WITHOUT NEIGHBORHOOD COMPONENTS ANALYSIS
- DIMENSIONALITY REDUCTION WITH NEIGHBORHOOD COMPONENTS ANALYSIS
- RESTRICTED BOLTZMANN MACHINE FEATURES FOR DIGIT CLASSIFICATION
- VARYING REGULARIZATION IN MULTILAYER PERCEPTRON
- USING FUNCTIONTRANSFORMER TO SELECT COLUMNS
- IMPORTANCE OF FEATURE SCALING
- MAP DATA TO A NORMAL DISTRIBUTION
- FEATURE DISCRETIZATION
- UNDERSTANDING THE DECISION TREE STRUCTURE

6263 HYPERPARAMETER OPTIMIZERS

MODELSELECTIONGRIDSEARCHCV ESTIMATOR EXHAUSTIVE SEARCH OVER SPECIFIED PARAMETER VALUES FOR AN ESTIMATOR

MODELSELECTIONPARAMETERGRID PARAMGRID GRID OF PARAMETERS WITH A DISCRETE NUMBER OF VALUES FOR EACH

MODELSELECTIONPARAMETERSAMPLER GENERATOR ON PARAMETERS SAMPLED FROM GIVEN DISTRIBUTIONS

MODELSELECTIONRANDOMIZEDSEARCHCV

RANDOMIZED SEARCH ON HYPER PARAMETERS

SKLEARNMODELSELECTION GRIDSEARCHCV

CLASSSSKLEARNMODELSELECTION GRIDSEARCHCV ESTIMATOR PARAMGRID SCORINGNONE

NJOBSNONE IID'WARN' REFITTRUE

CV'WARN' VERBOSE0 PREDISPATCH'2NJOBS'

ERRORSCORE'RAISEDEPRECATING' RE

TURNTRAINSCOREFALSE

EXHAUSTIVE SEARCH OVER SPECIFIED PARAMETER VALUES FOR AN ESTIMATOR

IMPORTANT MEMBERS ARE FIT PREDICT

2094 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213  
GRIDSEARCHCV IMPLEMENTS A “FIT” AND A “SCORE” METHOD IT ALSO IMPLEMENTS “PREDICT” “PREDICTPROBA” “DECISIONFUNCTION” “TRANSFORM” AND “INVERSETRANSFORM” IF THEY ARE IMPLEMENTED IN THE ESTIMATOR USED  
THE PARAMETERS OF THE ESTIMATOR USED TO APPLY THESE METHODS ARE OPTIMIZED BY CROSSVALIDATED GRIDSEARCH OVER A PARAMETER GRID

READ MORE IN THE USER GUIDE

PARAMETERS

ESTIMATOR ESTIMATOR OBJECT THIS IS ASSUMED TO IMPLEMENT THE SCIKITLEARN ESTIMATOR INTERFACE  
EITHER ESTIMATOR NEEDS TO PROVIDE A SCORE FUNCTION OR SCORING MUST BE PASSED

PARAMGRID DICT OR LIST OF DICTIONARIES DICTIONARY WITH PARAMETERS NAMES STRING AS KEYS AND  
LISTS OF PARAMETER SETTINGS TO TRY AS VALUES OR A LIST OF SUCH DICTIONARIES IN WHICH CASE THE  
GRIDS SPANNED BY EACH DICTIONARY IN THE LIST ARE EXPLORED THIS ENABLES SEARCHING OVER ANY  
SEQUENCE OF PARAMETER SETTINGS

SCORING STRING CALLABLE LISTTUPLE DICT OR NONE DEFAULT NONE A SINGLE STRING SEE THE SCORING  
PARAMETER DEFINING MODEL EVALUATION RULES OR A CALLABLE SEE DEFINING YOUR SCORING STRATEGY  
FROM METRIC FUNCTIONS TO EVALUATE THE PREDICTIONS ON THE TEST SET

FOR EVALUATING MULTIPLE METRICS EITHER GIVE A LIST OF UNIQUE STRINGS OR A DICT WITH NAMES AS  
KEYS AND CALLABLES AS VALUES

NOTE THAT WHEN USING CUSTOM SCORERS EACH SCORER SHOULD RETURN A SINGLE VALUE METRIC  
FUNCTIONS RETURNING A LISTARRAY OF VALUES CAN BE WRAPPED INTO MULTIPLE SCORERS THAT RETURN  
ONE VALUE EACH

SEESPECIFYING MULTIPLE METRICS FOR EVALUATION FOR AN EXAMPLE

IF NONE THE ESTIMATOR’S SCORE METHOD IS USED

NJOBS INT OR NONE OPTIONAL DEFAULTNONE NUMBER OF JOBS TO RUN IN PARALLEL NONE MEANS 1  
UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL PROCESSORS SEE

GLOSSARY FOR MORE DETAILS

PREDISPATCH INT OR STRING OPTIONAL CONTROLS THE NUMBER OF JOBS THAT GET DISPATCHED DURING  
PARALLEL EXECUTION REDUCING THIS NUMBER CAN BE USEFUL TO AVOID AN EXPLOSION OF MEMORY  
CONSUMPTION WHEN MORE JOBS GET DISPATCHED THAN CPUS CAN PROCESS THIS PARAMETER CAN  
BE

- NONE IN WHICH CASE ALL THE JOBS ARE IMMEDIATELY CREATED AND SPAWNED USE THIS FOR  
LIGHTWEIGHT AND FASTRUNNING JOBS TO AVOID DELAYS DUE TO ONDEMAND SPAWNING OF THE JOBS
- AN INT GIVING THE EXACT NUMBER OF TOTAL JOBS THAT ARE SPAWNED
- A STRING GIVING AN EXPRESSION AS A FUNCTION OF NJOBS AS IN ‘2NJOBS’

IIDBOOLEAN DEFAULT‘WARN’ IF TRUE RETURN THE AVERAGE SCORE ACROSS FOLDS WEIGHTED BY THE  
NUMBER OF SAMPLES IN EACH TEST SET IN THIS CASE THE DATA IS ASSUMED TO BE IDENTICALLY DIS  
TRIBUTED ACROSS THE FOLDS AND THE LOSS MINIMIZED IS THE TOTAL LOSS PER SAMPLE AND NOT THE  
MEAN LOSS ACROSS THE FOLDS IF FALSE RETURN THE AVERAGE SCORE ACROSS FOLDS DEFAULT IS TRUE  
BUT WILL CHANGE TO FALSE IN VERSION 022 TO CORRESPOND TO THE STANDARD DEFINITION OF CROSS  
VALIDATION

CHANGED IN VERSION 020 PARAMETER IID WILL CHANGE FROM TRUE TO FALSE BY DEFAULT IN  
VERSION 022 AND WILL BE REMOVED IN 024

CVINT CROSSVALIDATION GENERATOR OR AN ITERABLE OPTIONAL DETERMINES THE CROSSVALIDATION SPLIT  
TING STRATEGY POSSIBLE INPUTS FOR CV ARE

- NONE TO USE THE DEFAULT 3FOLD CROSS VALIDATION

626SKLEARNMODELSELECTION MODEL SELECTION 2095

SCIKITLEARN USER GUIDE RELEASE 0213

- INTEGER TO SPECIFY THE NUMBER OF FOLDS IN A STRATIFIEDKFOLD
- CV SPLITTER
- AN ITERABLE YIELDING TRAIN TEST SPLITS AS ARRAYS OF INDICES

FOR INTEGERNONE INPUTS IF THE ESTIMATOR IS A CLASSIFIER AND YIS EITHER BINARY OR MULTICLASS STRATIFIEDKFOLD IS USED IN ALL OTHER CASES KFOLD IS USED

REFER USER GUIDE FOR THE VARIOUS CROSSVALIDATION STRATEGIES THAT CAN BE USED HERE  
CHANGED IN VERSION 020 CVDEFAULT VALUE IF NONE WILL CHANGE FROM 3FOLD TO 5FOLD IN V022

REFIT BOOLEAN STRING OR CALLABLE DEFAULTTRUE REFIT AN ESTIMATOR USING THE BEST FOUND PARAM ETERS ON THE WHOLE DATASET

FOR MULTIPLE METRIC EVALUATION THIS NEEDS TO BE A STRING DENOTING THE SCORER THAT WOULD BE USED TO FIND THE BEST PARAMETERS FOR REFITTING THE ESTIMATOR AT THE END

WHERE THERE ARE CONSIDERATIONS OTHER THAN MAXIMUM SCORE IN CHOOSING A BEST ESTIMA TORREFIT CAN BE SET TO A FUNCTION WHICH RETURNS THE SELECTED BESTINDEX GIVEN

CVRESULTS

THE REFITTED ESTIMATOR IS MADE AVAILABLE AT THE BESTESTIMATOR ATTRIBUTE AND PERMITS USINGPREDICT DIRECTLY ON THIS GRIDSEARCHCV INSTANCE

ALSO FOR MULTIPLE METRIC EVALUATION THE ATTRIBUTES BESTINDEX BESTSCORE AND BESTPARAMS WILL ONLY BE AVAILABLE IF REFIT IS SET AND ALL OF THEM WILL BE DETERMINED

WRT THIS SPECIFIC SCORER BESTSCORE IS NOT RETURNED IF REFIT IS CALLABLE SEESCORING PARAMETER TO KNOW MORE ABOUT MULTIPLE METRIC EVALUATION

CHANGED IN VERSION 020 SUPPORT FOR CALLABLE ADDED

VERBOSE INTEGER CONTROLS THE VERBOSITY THE HIGHER THE MORE MESSAGES

ERRORSCORE 'RAISE' OR NUMERIC VALUE TO ASSIGN TO THE SCORE IF AN ERROR OCCURS IN ESTIMATOR FITTING IF SET TO 'RAISE' THE ERROR IS RAISED IF A NUMERIC VALUE IS GIVEN FITFAILEDWARNING

IS RAISED IF THIS PARAMETER DOES NOT AFFECT THE REFIT STEP WHICH WILL ALWAYS RAISE THE ERROR DEFAULT IS 'RAISE' BUT FROM VERSION 022 IT WILL CHANGE TO NPNAN

RETURNTRAINSCORE BOOLEAN DEFAULTFALSE IF FALSE THECVRESULTS ATTRIBUTE WILL NOT

INCLUDE TRAINING SCORES COMPUTING TRAINING SCORES IS USED TO GET INSIGHTS ON HOW DIFFER ENT PARAMETER SETTINGS IMPACT THE OVERFITTINGUNDERFITTING TRADEOFF HOWEVER COMPUTING THE

SCORES ON THE TRAINING SET CAN BE COMPUTATIONALLY EXPENSIVE AND IS NOT STRICTLY REQUIRED TO SELECT THE PARAMETERS THAT YIELD THE BEST GENERALIZATION PERFORMANCE

ATTRIBUTES

CVRESULTS DICT OF NUMPY MASKED NDARRAYS A DICT WITH KEYS AS COLUMN HEADERS AND VALUES AS COLUMNS THAT CAN BE IMPORTED INTO A PANDAS DATAFRAME

FOR INSTANCE THE BELOW GIVEN TABLE

PARAMKERNEL PARAMGAMMA PARAMDEGREE SPLIT0TESTSCORE RANKT

'POLY' - 2 080 2

'POLY' - 3 070 4

'RBF' 01 - 080 3

'RBF' 02 - 093 1

WILL BE REPRESENTED BY A CVRESULTS DICT OF

2096 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

PARAMKERNEL MASKEDARRAYDATA POLY POLY RBF RBF  
MASK FALSE FALSE FALSE FALSE  
PARAMGAMMA MASKEDARRAYDATA 01 02  
MASK TRUE TRUE FALSE FALSE  
PARAMDEGREE MASKEDARRAYDATA 20 30  
MASK FALSE FALSE TRUE TRUE  
SPLIT0TESTSCORE 080 070 080 093  
SPLIT1TESTSCORE 082 050 070 078  
MEANTESTSCORE 081 060 075 085  
STDTESTSCORE 001 010 005 008  
RANKTESTSCORE 2 4 3 1  
SPLIT0TRAINSCORE 080 092 070 093  
SPLIT1TRAINSCORE 082 055 070 087  
MEANTRAINSCORE 081 074 070 090  
STDTRAINSCORE 001 019 000 003  
MEANFITTIME 073 063 043 049  
STDFITTIME 001 002 001 001  
MEANSCORETIME 001 006 004 004  
STDSCORETIME 000 000 000 001  
PARAMS KERNEL POLY DEGREE 2

NOTE  
THE KEYPARAMS IS USED TO STORE A LIST OF PARAMETER SETTINGS DICTS FOR ALL THE PARAMETER CANDIDATES  
THEMEANFITTIME STDFITTIME MEANSCORETIME AND  
STDSCORETIME ARE ALL IN SECONDS  
FOR MULTIMETRIC EVALUATION THE SCORES FOR ALL THE SCORERS ARE AVAILABLE IN THE CVRESULTS DICT AT THE KEYS ENDING WITH THAT SCORER'S NAME SCORERNAME  
INSTEAD OFSCORE SHOWN ABOVE 'SPLIT0TESTPRECISION' 'MEANTRAINPRECISION' ETC  
BESTESTIMATOR ESTIMATOR OR DICT ESTIMATOR THAT WAS CHOSEN BY THE SEARCH IE ESTIMATOR WHICH GAVE HIGHEST SCORE OR SMALLEST LOSS IF SPECIFIED ON THE LEFT OUT DATA NOT AVAILABLE IF  
REFITFALSE  
SEEREFIT PARAMETER FOR MORE INFORMATION ON ALLOWED VALUES  
BESTSCORE FLOAT MEAN CROSSVALIDATED SCORE OF THE BESTESTIMATOR  
FOR MULTIMETRIC EVALUATION THIS IS PRESENT ONLY IF REFIT IS SPECIFIED  
BESTPARAMS DICT PARAMETER SETTING THAT GAVE THE BEST RESULTS ON THE HOLD OUT DATA  
FOR MULTIMETRIC EVALUATION THIS IS PRESENT ONLY IF REFIT IS SPECIFIED  
BESTINDEX INT THE INDEX OF THE CVRESULTS ARRAYS WHICH CORRESPONDS TO THE BEST CANDIDATE PARAMETER SETTING  
THE DICT AT SEARCHCVRESULTSPARAMSSEARCHBESTINDEX GIVES  
THE PARAMETER SETTING FOR THE BEST MODEL THAT GIVES THE HIGHEST MEAN SCORE SEARCH  
BESTSCORE  
FOR MULTIMETRIC EVALUATION THIS IS PRESENT ONLY IF REFIT IS SPECIFIED  
SCORER FUNCTION OR A DICT SCORER FUNCTION USED ON THE HELD OUT DATA TO CHOOSE THE BEST PARAMETERS FOR THE MODEL  
626SKLEARNMODELSELECTION MODEL SELECTION 2097

SCIKITLEARN USER GUIDE RELEASE 0213

FOR MULTIMETRIC EVALUATION THIS ATTRIBUTE HOLDS THE VALIDATED SCORING DICT WHICH MAPS THE SCORER KEY TO THE SCORER CALLABLE

NSPLITS INT THE NUMBER OF CROSSVALIDATION SPLITS FOLDSITERATIONS

REFITTIME FLOAT SECONDS USED FOR REFITTING THE BEST MODEL ON THE WHOLE DATASET

THIS IS PRESENT ONLY IF REFIT IS NOT FALSE

SEE ALSO

PARAMETERGRID GENERATES ALL THE COMBINATIONS OF A HYPERPARAMETER GRID

SKLEARNMODELSELECTIONTRAINTESTSPLIT UTILITY FUNCTION TO SPLIT THE DATA INTO A DEVELOPMENT SET USABLE FOR FITTING A GRIDSEARCHCV INSTANCE AND AN EVALUATION SET FOR ITS FINAL EVALUATION

SKLEARNMETRICSMAKESCORER MAKE A SCORER FROM A PERFORMANCE METRIC OR LOSS FUNCTION

NOTES

THE PARAMETERS SELECTED ARE THOSE THAT MAXIMIZE THE SCORE OF THE LEFT OUT DATA UNLESS AN EXPLICIT SCORE IS PASSED IN WHICH CASE IT IS USED INSTEAD

IFNJOBS WAS SET TO A VALUE HIGHER THAN ONE THE DATA IS COPIED FOR EACH POINT IN THE GRID AND NOT NJOBS TIMES THIS IS DONE FOR EFFICIENCY REASONS IF INDIVIDUAL JOBS TAKE VERY LITTLE TIME BUT MAY RAISE ERRORS IF THE DATASET IS LARGE AND NOT ENOUGH MEMORY IS AVAILABLE A WORKAROUND IN THIS CASE IS TO SET PREDISPATCH THEN THE MEMORY IS COPIED ONLY PREDISPATCH MANY TIMES A REASONABLE VALUE FOR PREDISPATCH IS2

NJOBS

EXAMPLES

FROM SKLEARN IMPORT SVM DATASETS

FROM SKLEARNMODELSELECTION IMPORT GRIDSEARCHCV

IRIS DATASETSLOADIRIS

PARAMETERS KERNELLINEAR RBF C1 10

SVC SVMSCVCGAMMASCALE

CLF GRIDSEARCHCVSVC PARAMETERS CV5

CLFFITIRISDATA IRISTARGET

GRIDSEARCHCVCV5 ERRORSORE

ESTIMATORSVCC10 CACHESIZE CLASSWEIGHT COEF0

DECISIONFUNCTIONSHAPEOVR DEGREE GAMMA

KERNELRBF MAXITER1 PROBABILITYFALSE

RANDOMSTATENONE SHRINKINGTRUE TOL

VERBOSEFALSE

IID NJOBSNONE

PARAMGRID PREDISPATCH REFIT RETURNTRAINSCORE

SCORING VERBOSE

SORTEDCLFCVRESULTSKEYS

MEANFITTIME MEANSORETIME MEANTESTSCORE

PARAMC PARAMKERNEL PARAMS

RANKTESTSCORE SPLIT0TESTSCORE

SPLIT2TESTSCORE

STDFITTIME STDSCORETIME STDTESTSCORE

2098 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

METHODS

DECISIONFUNCTION SELF X CALL DECISIONFUNCTION ON THE ESTIMATOR WITH THE BEST FOUND PARAMETERS

FITSELF X Y GROUPS RUN FIT WITH ALL SETS OF PARAMETERS

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

INVERSETRANSFORM SELF XT CALL INVERSETRANSFORM ON THE ESTIMATOR WITH THE BEST FOUND PARAMS

PREDICT SELF X CALL PREDICT ON THE ESTIMATOR WITH THE BEST FOUND PARAM ETERS

PREDICTLOGPROBA SELF X CALL PREDICTLOGPROBA ON THE ESTIMATOR WITH THE BEST FOUND PARAMETERS

PREDICTPROBA SELF X CALL PREDICTPROBA ON THE ESTIMATOR WITH THE BEST FOUND PARAMETERS

SCORE SELF X Y RETURNS THE SCORE ON THE GIVEN DATA IF THE ESTIMATOR HAS BEEN REFIT

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF X CALL TRANSFORM ON THE ESTIMATOR WITH THE BEST FOUND PA RAMETERS

INIT SELFESTIMATOR PARAMGRID SCORINGNONE NJOBSNONE IID'WARN' REFITTRUE CV'WARN' VERBOSE0 PREDISPATCH'2NJOBS' ERRORSORE'RAISEDEPRECATING' RE TURNTRAINSCOREFALSE

DECISIONFUNCTION SELF X

CALL DECISIONFUNCTION ON THE ESTIMATOR WITH THE BEST FOUND PARAMETERS ONLY AVAILABLE IF REFITTRUE AND THE UNDERLYING ESTIMATOR SUPPORTS DECISIONFUNCTION PARAMETERS

XINDEXABLE LENGTH NSAMPLES MUST FULFILL THE INPUT ASSUMPTIONS OF THE UNDERLYING ESTI MATOR

FITSELFXYNONE GROUPSNONE FITPARAMS

RUN FIT WITH ALL SETS OF PARAMETERS

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING VECTOR WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUT OPTIONAL TARGET RELATIVE TO X FOR CLASSIFICATION OR REGRESSION NONE FOR UNSUPERVISED LEARNING

GROUPS ARRAYLIKE WITH SHAPE NSAMPLES OPTIONAL GROUP LABELS FOR THE SAMPLES USED WHILE SPLITTING THE DATASET INTO TRAINTEST SET ONLY USED IN CONJUNCTION WITH A "GROUP" CV INSTANCE EG GROUPKFOLD

FITPARAMS DICT OF STRING OBJECT PARAMETERS PASSED TO THE FIT METHOD OF THE ESTIMA TOR

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

626SKLEARNMODELSELECTION MODEL SELECTION 2099

SCIKITLEARN USER GUIDE RELEASE 0213

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

INVERSETRANSFORM SELFXT

CALL INVERSETRANSFORM ON THE ESTIMATOR WITH THE BEST FOUND PARAMS

ONLY AVAILABLE IF THE UNDERLYING ESTIMATOR IMPLEMENTS INVERSETRANSFORM ANDREFITTRUE

PARAMETERS

XTINDEXABLE LENGTH NSAMPLES MUST FULFILL THE INPUT ASSUMPTIONS OF THE UNDERLYING ESTI

MATOR

PREDICTSELF

CALL PREDICT ON THE ESTIMATOR WITH THE BEST FOUND PARAMETERS

ONLY AVAILABLE IF REFITTRUE AND THE UNDERLYING ESTIMATOR SUPPORTS PREDICT

PARAMETERS

XINDEXABLE LENGTH NSAMPLES MUST FULFILL THE INPUT ASSUMPTIONS OF THE UNDERLYING ESTI

MATOR

PREDICTLOGPROBA SELF

CALL PREDICTLOGPROBA ON THE ESTIMATOR WITH THE BEST FOUND PARAMETERS

ONLY AVAILABLE IF REFITTRUE AND THE UNDERLYING ESTIMATOR SUPPORTS PREDICTLOGPROBA

PARAMETERS

XINDEXABLE LENGTH NSAMPLES MUST FULFILL THE INPUT ASSUMPTIONS OF THE UNDERLYING ESTI

MATOR

PREDICTPROBA SELF

CALL PREDICTPROBA ON THE ESTIMATOR WITH THE BEST FOUND PARAMETERS

ONLY AVAILABLE IF REFITTRUE AND THE UNDERLYING ESTIMATOR SUPPORTS PREDICTPROBA

PARAMETERS

XINDEXABLE LENGTH NSAMPLES MUST FULFILL THE INPUT ASSUMPTIONS OF THE UNDERLYING ESTI

MATOR

SCORESELFXYNONE

RETURNS THE SCORE ON THE GIVEN DATA IF THE ESTIMATOR HAS BEEN REFIT

THIS USES THE SCORE DEFINED BY SCORING WHERE PROVIDED AND THE BESTESTIMATORSORE METHOD

OTHERWISE

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES INPUT DATA WHERE NSAMPLES IS THE NUMBER

OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUT OPTIONAL TARGET RELATIVE TO X

FOR CLASSIFICATION OR REGRESSION NONE FOR UNSUPERVISED LEARNING

RETURNS

SCORE FLOAT

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

2100 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS  
SELF

TRANSFORM SELF  
CALL TRANSFORM ON THE ESTIMATOR WITH THE BEST FOUND PARAMETERS  
ONLY AVAILABLE IF THE UNDERLYING ESTIMATOR SUPPORTS TRANSFORM ANDREFITTRUE

PARAMETERS  
XINDEXABLE LENGTH NSAMPLES MUST FULFILL THE INPUT ASSUMPTIONS OF THE UNDERLYING ESTI  
MATOR

EXAMPLES USING SKLEARNMODELSELECTIONGRIDSEARCHCV

- COMPARISON OF KERNEL RIDGE REGRESSION AND SVR
- FACES RECOGNITION EXAMPLE USING EIGENFACES AND SVMs
- FEATURE AGGLOMERATION VS UNIVARIATE SELECTION
- CONCATENATING MULTIPLE FEATURE EXTRACTION METHODS
- PIPELINING CHAINING A PCA AND A LOGISTIC REGRESSION
- COLUMN TRANSFORMER WITH MIXED TYPES
- SELECTING DIMENSIONALITY REDUCTION WITH PIPELINE AND GRIDSEARCHCV
- SHRINKAGE COVARIANCE ESTIMATION LEDOITWOLF VS OAS AND MAXLIKELIHOOD
- MODEL SELECTION WITH PROBABILISTIC PCA AND FACTOR ANALYSIS FA
- CROSSVALIDATION ON DIABETES DATASET EXERCISE
- COMPARISON OF KERNEL RIDGE AND GAUSSIAN PROCESS REGRESSION
- PARAMETER ESTIMATION USING GRID SEARCH WITH CROSSVALIDATION
- COMPARING RANDOMIZED SEARCH AND GRID SEARCH FOR HYPERPARAMETER ESTIMATION
- NESTED VERSUS NONNESTED CROSSVALIDATION
- DEMONSTRATION OF MULTIMETRIC EVALUATION ON CROSSVALSCORE AND GRIDSEARCHCV
- BALANCE MODEL COMPLEXITY AND CROSSVALIDATED SCORE
- SAMPLE PIPELINE FOR TEXT FEATURE EXTRACTION AND EVALUATION
- KERNEL DENSITY ESTIMATION
- FEATURE DISCRETIZATION
- SCALING THE REGULARIZATION PARAMETER FOR SVCS
- RBF SVM PARAMETERS

626SKLEARNMODELSELECTION MODEL SELECTION 2101

SCIKITLEARN USER GUIDE RELEASE 0213

SKLEARNMODELSELECTION PARAMETERGRID

CLASSSSKLEARNMODELSELECTION PARAMETERGRID PARAMGRID

GRID OF PARAMETERS WITH A DISCRETE NUMBER OF VALUES FOR EACH

CAN BE USED TO ITERATE OVER PARAMETER VALUE COMBINATIONS WITH THE PYTHON BUILTIN FUNCTION ITER

READ MORE IN THE USER GUIDE

PARAMETERS

PARAMGRID DICT OF STRING TO SEQUENCE OR SEQUENCE OF SUCH THE PARAMETER GRID TO EXPLORE AS

A DICTIONARY MAPPING ESTIMATOR PARAMETERS TO SEQUENCES OF ALLOWED VALUES

AN EMPTY DICT SIGNIFIES DEFAULT PARAMETERS

A SEQUENCE OF DICTS SIGNIFIES A SEQUENCE OF GRIDS TO SEARCH AND IS USEFUL TO AVOID EXPLORING

PARAMETER COMBINATIONS THAT MAKE NO SENSE OR HAVE NO EFFECT SEE THE EXAMPLES BELOW

SEE ALSO

GRIDSEARCHCV USESPARAMETERGRID TO PERFORM A FULL PARALLELIZED PARAMETER SEARCH

EXAMPLES

```
FROM SKLEARNMODELSELECTION IMPORT PARAMETERGRID
PARAMGRID A 1 2 B TRUEFALSE
LISTPARAMETERGRIDPARAMGRID
A 1 B TRUE A 1 B FALSE
A 2 B TRUE A 2 B FALSE
TRUE
GRID KERNEL LINEAR KERNEL RBF GAMMA 1 10
LISTPARAMETERGRIDGRID KERNEL LINEAR
KERNEL RBF GAMMA 1
KERNEL RBF GAMMA 10
TRUE
PARAMETERGRIDGRID1 KERNEL RBF GAMMA 1
TRUE
INIT SELFPARAMGRID
SKLEARNMODELSELECTION PARAMETERSAMPLER
CLASSSSKLEARNMODELSELECTION PARAMETERSAMPLER PARAMDISTRIBUTIONS NITER RAN
DOMSTATENONE
GENERATOR ON PARAMETERS SAMPLED FROM GIVEN DISTRIBUTIONS
NONDETERMINISTIC ITERABLE OVER RANDOM CANDIDATE COMBINATIONS FOR HYPER PARAMETER SEARCH IF ALL PARAMETERS
ARE PRESENTED AS A LIST SAMPLING WITHOUT REPLACEMENT IS PERFORMED IF AT LEAST ONE PARAMETER IS GIVEN AS A
DISTRIBUTION SAMPLING WITH REPLACEMENT IS USED IT IS HIGHLY RECOMMENDED TO USE CONTINUOUS DISTRIBUTIONS FOR
CONTINUOUS PARAMETERS
NOTE THAT BEFORE SCIPY 016 THE SCIPYSTATSDISTRIBUTIONS DO NOT ACCEPT A CUSTOM RNG INSTANCE
AND ALWAYS USE THE SINGLETON RNG FROM NUMPYRANDOM HENCE SETTING RANDOMSTATE WILL NOT GUARANTEE
A DETERMINISTIC ITERATION WHENEVER SCIPYSTATS DISTRIBUTIONS ARE USED TO DEFINE THE PARAMETER SEARCH SPACE
DETERMINISTIC BEHAVIOR IS HOWEVER GUARANTEED FROM SCIPY 016 ONWARDS
2102 CHAPTER 6 API REFERENCE
```

SCIKITLEARN USER GUIDE RELEASE 0213  
READ MORE IN THE USER GUIDE  
PARAMETERS  
PARAMDISTRIBUTIONS DICT DICTIONARY WHERE THE KEYS ARE PARAMETERS AND VALUES ARE DISTRIBUTIONS FROM WHICH A PARAMETER IS TO BE SAMPLED DISTRIBUTIONS EITHER HAVE TO PROVIDE A RVS FUNCTION TO SAMPLE FROM THEM OR CAN BE GIVEN AS A LIST OF VALUES WHERE A UNIFORM DISTRIBUTION IS ASSUMED  
NITER INTEGER NUMBER OF PARAMETER SETTINGS THAT ARE PRODUCED  
RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE PSEUDO RANDOM NUMBER GENERATOR STATE USED FOR RANDOM UNIFORM SAMPLING FROM LISTS OF POSSIBLE VALUES INSTEAD OF SCIPYSTATS DISTRIBUTIONS IF INT RANDOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NPRANDOM  
RETURNS  
PARAMS DICT OF STRING TO ANY YIELDS DICTIONARIES MAPPING EACH ESTIMATOR PARAMETER TO AS SAMPLED VALUE  
EXAMPLES  
FROM SKLEARNMODELSELECTION IMPORT PARAMETERSAMPLER  
FROM SCIPYSTATSDISTRIBUTIONS IMPORT EXPON  
IMPORT NUMPY AS NP  
RNG NPRANDOMRANDOMSTATE0  
PARAMGRID A1 2 B EXPON  
PARAMLIST LISTPARAMETERSAMPLERPARAMGRID NITER4  
RANDOMSTATERNG  
ROUNDEDLIST DICTK ROUNDV 6 FORK VINDITEMS  
FOR DINPARAMLIST  
ROUNDEDLIST B 089856 A 1  
B 0923223 A 1  
B 1878964 A 2  
B 1038159 A 2  
TRUE  
INIT SELFPARAMDISTRIBUTIONS NITER RANDOMSTATENONE  
SKLEARNMODELSELECTION RANDOMIZEDSEARCHCV  
CLASSSSKLEARNMODELSELECTION RANDOMIZEDSEARCHCV ESTIMATOR PARAMDISTRIBUTIONS  
NITER10 SCORINGNONE  
NJOBSNONE IID'WARN' RE  
FITTRUE CV'WARN' VER  
BOSE0 PREDISPATCH'2NJOBS'  
RANDOMSTATENONE  
ERRORSCORE'RAISEDEPRECATING' RE  
TURNTRAINSCOREFALSE  
RANDOMIZED SEARCH ON HYPER PARAMETERS  
RANDOMIZEDSEARCHCV IMPLEMENTS A "FIT" AND A "SCORE" METHOD IT ALSO IMPLEMENTS "PREDICT" "PREDICTPROBA" "DECISIONFUNCTION" "TRANSFORM" AND "INVERSETRANSFORM" IF THEY ARE IMPLEMENTED IN THE ESTIMATOR USED  
626SKLEARNMODELSELECTION MODEL SELECTION 2103

SCIKITLEARN USER GUIDE RELEASE 0213

THE PARAMETERS OF THE ESTIMATOR USED TO APPLY THESE METHODS ARE OPTIMIZED BY CROSSVALIDATED SEARCH OVER PARAMETER SETTINGS

IN CONTRAST TO GRIDSEARCHCV NOT ALL PARAMETER VALUES ARE TRIED OUT BUT RATHER A FIXED NUMBER OF PARAMETER SETTINGS IS SAMPLED FROM THE SPECIFIED DISTRIBUTIONS THE NUMBER OF PARAMETER SETTINGS THAT ARE TRIED IS GIVEN BY NITER

IF ALL PARAMETERS ARE PRESENTED AS A LIST SAMPLING WITHOUT REPLACEMENT IS PERFORMED IF AT LEAST ONE PARAMETER IS GIVEN AS A DISTRIBUTION SAMPLING WITH REPLACEMENT IS USED IT IS HIGHLY RECOMMENDED TO USE CONTINUOUS DISTRIBUTIONS FOR CONTINUOUS PARAMETERS

NOTE THAT BEFORE SCIPY 016 THE SCIPYSTATSDISTRIBUTIONS DO NOT ACCEPT A CUSTOM RNG INSTANCE AND ALWAYS USE THE SINGLETON RNG FROM NUMPYRANDOM HENCE SETTING RANDOMSTATE WILL NOT GUARANTEE A DETERMINISTIC ITERATION WHENEVER SCIPYSTATS DISTRIBUTIONS ARE USED TO DEFINE THE PARAMETER SEARCH SPACE READ MORE IN THE USER GUIDE

PARAMETERS

ESTIMATOR ESTIMATOR OBJECT A OBJECT OF THAT TYPE IS INSTANTIATED FOR EACH GRID POINT THIS IS ASSUMED TO IMPLEMENT THE SCIKITLEARN ESTIMATOR INTERFACE EITHER ESTIMATOR NEEDS TO PROVIDE ASCORE FUNCTION OR SCORING MUST BE PASSED

PARAMDISTRIBUTIONS DICT DICTIONARY WITH PARAMETERS NAMES STRING AS KEYS AND DISTRIBUTIONS OR LISTS OF PARAMETERS TO TRY DISTRIBUTIONS MUST PROVIDE A RVS METHOD FOR SAMPLING SUCH AS THOSE FROM SCIPYSTATSDISTRIBUTIONS IF A LIST IS GIVEN IT IS SAMPLED UNIFORMLY

NITER INT DEFAULT10 NUMBER OF PARAMETER SETTINGS THAT ARE SAMPLED NITER TRADES OFF RUNTIME VS QUALITY OF THE SOLUTION

SCORING STRING CALLABLE LISTTUPLE DICT OR NONE DEFAULT NONE A SINGLE STRING SEE THE SCORING PARAMETER DEFINING MODEL EVALUATION RULES OR A CALLABLE SEE DEFINING YOUR SCORING STRATEGY FROM METRIC FUNCTIONS TO EVALUATE THE PREDICTIONS ON THE TEST SET

FOR EVALUATING MULTIPLE METRICS EITHER GIVE A LIST OF UNIQUE STRINGS OR A DICT WITH NAMES AS KEYS AND CALLABLES AS VALUES

NOTE THAT WHEN USING CUSTOM SCORERS EACH SCORER SHOULD RETURN A SINGLE VALUE METRIC FUNCTIONS RETURNING A LISTARRAY OF VALUES CAN BE WRAPPED INTO MULTIPLE SCORERS THAT RETURN ONE VALUE EACH

SEESPECIFYING MULTIPLE METRICS FOR EVALUATION FOR AN EXAMPLE

IF NONE THE ESTIMATOR’S SCORE METHOD IS USED

NJOBS INT OR NONE OPTIONAL DEFAULTNONE NUMBER OF JOBS TO RUN IN PARALLEL NONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

PREDISPATCH INT OR STRING OPTIONAL CONTROLS THE NUMBER OF JOBS THAT GET DISPATCHED DURING PARALLEL EXECUTION REDUCING THIS NUMBER CAN BE USEFUL TO AVOID AN EXPLOSION OF MEMORY CONSUMPTION WHEN MORE JOBS GET DISPATCHED THAN CPUS CAN PROCESS THIS PARAMETER CAN BE

- NONE IN WHICH CASE ALL THE JOBS ARE IMMEDIATELY CREATED AND SPAWNED USE THIS FOR LIGHTWEIGHT AND FASTRUNNING JOBS TO AVOID DELAYS DUE TO ONDEMAND SPAWNING OF THE JOBS
- AN INT GIVING THE EXACT NUMBER OF TOTAL JOBS THAT ARE SPAWNED
- A STRING GIVING AN EXPRESSION AS A FUNCTION OF NJOBS AS IN ‘2NJOBS’

2104 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

IIDBOOLEAN DEFAULT'WARN' IF TRUE RETURN THE AVERAGE SCORE ACROSS FOLDS WEIGHTED BY THE NUMBER OF SAMPLES IN EACH TEST SET IN THIS CASE THE DATA IS ASSUMED TO BE IDENTICALLY DISTRIBUTED ACROSS THE FOLDS AND THE LOSS MINIMIZED IS THE TOTAL LOSS PER SAMPLE AND NOT THE MEAN LOSS ACROSS THE FOLDS IF FALSE RETURN THE AVERAGE SCORE ACROSS FOLDS DEFAULT IS TRUE BUT WILL CHANGE TO FALSE IN VERSION 022 TO CORRESPOND TO THE STANDARD DEFINITION OF CROSS VALIDATION

CHANGED IN VERSION 020 PARAMETER IID WILL CHANGE FROM TRUE TO FALSE BY DEFAULT IN VERSION 022 AND WILL BE REMOVED IN 024

CVINT CROSSVALIDATION GENERATOR OR AN ITERABLE OPTIONAL DETERMINES THE CROSSVALIDATION SPLITTING STRATEGY POSSIBLE INPUTS FOR CV ARE

- NONE TO USE THE DEFAULT 3FOLD CROSS VALIDATION
- INTEGER TO SPECIFY THE NUMBER OF FOLDS IN A STRATIFIEDKFOLD
- CV SPLITTER
- AN ITERABLE YIELDING TRAIN TEST SPLITS AS ARRAYS OF INDICES

FOR INTEGERNONE INPUTS IF THE ESTIMATOR IS A CLASSIFIER AND YIS EITHER BINARY OR MULTICLASS STRATIFIEDKFOLD IS USED IN ALL OTHER CASES KFOLD IS USED

REFER USER GUIDE FOR THE VARIOUS CROSSVALIDATION STRATEGIES THAT CAN BE USED HERE  
CHANGED IN VERSION 020 CVDEFAULT VALUE IF NONE WILL CHANGE FROM 3FOLD TO 5FOLD IN V022

REFIT BOOLEAN STRING OR CALLABLE DEFAULTTRUE REFIT AN ESTIMATOR USING THE BEST FOUND PARAMETERS ON THE WHOLE DATASET

FOR MULTIPLE METRIC EVALUATION THIS NEEDS TO BE A STRING DENOTING THE SCORER THAT WOULD BE USED TO FIND THE BEST PARAMETERS FOR REFITTING THE ESTIMATOR AT THE END

WHERE THERE ARE CONSIDERATIONS OTHER THAN MAXIMUM SCORE IN CHOOSING A BEST ESTIMATOR REFIT CAN BE SET TO A FUNCTION WHICH RETURNS THE SELECTED BESTINDEX GIVEN THE CVRESULTS

THE REFITTED ESTIMATOR IS MADE AVAILABLE AT THE BESTESTIMATOR ATTRIBUTE AND PERMITS USINGPREDICT DIRECTLY ON THIS RANDOMIZEDSEARCHCV INSTANCE

ALSO FOR MULTIPLE METRIC EVALUATION THE ATTRIBUTES BESTINDEX BESTSCORE AND BESTPARAMS WILL ONLY BE AVAILABLE IF REFIT IS SET AND ALL OF THEM WILL BE DETERMINED WRT THIS SPECIFIC SCORER WHEN REFIT IS CALLABLE BESTSCORE IS DISABLED

SEESCORING PARAMETER TO KNOW MORE ABOUT MULTIPLE METRIC EVALUATION

CHANGED IN VERSION 020 SUPPORT FOR CALLABLE ADDED

VERBOSE INTEGER CONTROLS THE VERBOSITY THE HIGHER THE MORE MESSAGES

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE PSEUDO RANDOM NUMBER GENERATOR STATE USED FOR RANDOM UNIFORM SAMPLING FROM LISTS OF POSSIBLE VALUES INSTEAD OF SCIPYSTATS DISTRIBUTIONS IF INT RANDOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NPRANDOM ERRORSORE 'RAISE' OR NUMERIC VALUE TO ASSIGN TO THE SCORE IF AN ERROR OCCURS IN ESTIMATOR FITTING IF SET TO 'RAISE' THE ERROR IS RAISED IF A NUMERIC VALUE IS GIVEN FITFAILEDWARNING IS RAISED THIS PARAMETER DOES NOT AFFECT THE REFIT STEP WHICH WILL ALWAYS RAISE THE ERROR DEFAULT IS 'RAISE' BUT FROM VERSION 022 IT WILL CHANGE TO NPNAN

6265KLEARNMODELSELECTION MODEL SELECTION 2105

SCIKITLEARN USER GUIDE RELEASE 0213

RETURNTRAINSCORE BOOLEAN DEFAULTFALSE IF FALSE THECVRESULTS ATTRIBUTE WILL NOT INCLUDE TRAINING SCORES COMPUTING TRAINING SCORES IS USED TO GET INSIGHTS ON HOW DIFFERENT PARAMETER SETTINGS IMPACT THE OVERFITTINGUNDERFITTING TRADEOFF HOWEVER COMPUTING THE SCORES ON THE TRAINING SET CAN BE COMPUTATIONALLY EXPENSIVE AND IS NOT STRICTLY REQUIRED TO SELECT THE PARAMETERS THAT YIELD THE BEST GENERALIZATION PERFORMANCE

ATTRIBUTES

CVRESULTS DICT OF NUMPY MASKED NDARRAYS A DICT WITH KEYS AS COLUMN HEADERS AND VALUES AS COLUMNS THAT CAN BE IMPORTED INTO A PANDAS DATAFRAME

FOR INSTANCE THE BELOW GIVEN TABLE

PARAMKERNEL	PARAMGAMMA	SPLIT0TESTSCORE	RANKTESTSCORE
'RBF'	01 080	2	
'RBF'	02 090	1	
'RBF'	03 070	1	

WILL BE REPRESENTED BY A CVRESULTS DICT OF

PARAMKERNEL	MASKEDARRAYDATA	RBF	RBF	RBF
MASK	FALSE			
PARAMGAMMA	MASKEDARRAYDATA	01	02	03
MASK	FALSE			
SPLIT0TESTSCORE	080	090	070	
SPLIT1TESTSCORE	082	050	070	
MEANTESTSCORE	081	070	070	
STDTESTSCORE	001	020	000	
RANKTESTSCORE	3	1	1	
SPLIT0TRAINSCORE	080	092	070	
SPLIT1TRAINSCORE	082	055	070	
MEANTRAINSCORE	081	074	070	
STDTRAINSCORE	001	019	000	
MEANFITTIME	073	063	043	
STDFITTIME	001	002	001	
MEANSCORETIME	001	006	004	
STDSCORETIME	000	000	000	
PARAMS	KERNEL	RBF	GAMMA	01

NOTE

THE KEYPARAMS IS USED TO STORE A LIST OF PARAMETER SETTINGS DICTS FOR ALL THE PARAMETER CANDIDATES

THEMEANFITTIME STDFITTIME MEANSCORETIME AND STDSCORETIME ARE ALL IN SECONDS

FOR MULTIMETRIC EVALUATION THE SCORES FOR ALL THE SCORERS ARE AVAILABLE IN THE CVRESULTS DICT AT THE KEYS ENDING WITH THAT SCORER'S NAME SCORERNAME

INSTEAD OFSCORE SHOWN ABOVE 'SPLIT0TESTPRECISION' 'MEANTRAINPRECISION' ETC

BESTESTIMATOR ESTIMATOR OR DICT ESTIMATOR THAT WAS CHOSEN BY THE SEARCH IE ESTIMATOR WHICH GAVE HIGHEST SCORE OR SMALLEST LOSS IF SPECIFIED ON THE LEFT OUT DATA NOT AVAILABLE IF REFITFALSE

FOR MULTIMETRIC EVALUATION THIS ATTRIBUTE IS PRESENT ONLY IF REFIT IS SPECIFIED

SEEREFIT PARAMETER FOR MORE INFORMATION ON ALLOWED VALUES

2106 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

BESTSCORE FLOAT MEAN CROSSVALIDATED SCORE OF THE BESTESTIMATOR  
FOR MULTIMETRIC EVALUATION THIS IS NOT AVAILABLE IF REFIT ISFALSE SEEREFIT PARAMETER  
FOR MORE INFORMATION

BESTPARAMS DICT PARAMETER SETTING THAT GAVE THE BEST RESULTS ON THE HOLD OUT DATA  
FOR MULTIMETRIC EVALUATION THIS IS NOT AVAILABLE IF REFIT ISFALSE SEEREFIT PARAMETER  
FOR MORE INFORMATION

BESTINDEX INT THE INDEX OF THE CVRESULTS ARRAYS WHICH CORRESPONDS TO THE BEST CAN  
DIDATE PARAMETER SETTING

THE DICT AT SEARCHCVRESULTSPARAMSSEARCHBESTINDEX GIVES  
THE PARAMETER SETTING FOR THE BEST MODEL THAT GIVES THE HIGHEST MEAN SCORE SEARCH  
BESTSCORE  
FOR MULTIMETRIC EVALUATION THIS IS NOT AVAILABLE IF REFIT ISFALSE SEEREFIT PARAMETER  
FOR MORE INFORMATION

SCORER FUNCTION OR A DICT SCORER FUNCTION USED ON THE HELD OUT DATA TO CHOOSE THE BEST PARAM  
ETERS FOR THE MODEL  
FOR MULTIMETRIC EVALUATION THIS ATTRIBUTE HOLDS THE VALIDATED SCORING DICT WHICH MAPS  
THE SCORER KEY TO THE SCORER CALLABLE

NSPLITS INT THE NUMBER OF CROSSVALIDATION SPLITS FOLDSITERATIONS

REFITTIME FLOAT SECONDS USED FOR REFITTING THE BEST MODEL ON THE WHOLE DATASET  
THIS IS PRESENT ONLY IF REFIT IS NOT FALSE  
SEE ALSO

GRIDSEARCHCV DOES EXHAUSTIVE SEARCH OVER A GRID OF PARAMETERS

PARAMETERSAMPLER A GENERATOR OVER PARAMETER SETTINGS CONSTRUCTED FROM PARAMDISTRIBUTIONS

NOTES

THE PARAMETERS SELECTED ARE THOSE THAT MAXIMIZE THE SCORE OF THE HELDOUT DATA ACCORDING TO THE SCORING PARAM  
ETER

IFNJOBS WAS SET TO A VALUE HIGHER THAN ONE THE DATA IS COPIED FOR EACH PARAMETER SETTINGAND NOT NJOBS TIMES  
THIS IS DONE FOR EFFICIENCY REASONS IF INDIVIDUAL JOBS TAKE VERY LITTLE TIME BUT MAY RAISE ERRORS IF THE DATASET  
IS LARGE AND NOT ENOUGH MEMORY IS AVAILABLE A WORKAROUND IN THIS CASE IS TO SET PREDISPATCH THEN  
THE MEMORY IS COPIED ONLY PREDISPATCH MANY TIMES A REASONABLE VALUE FOR PREDISPATCH IS2

NJOBS

METHODS

DECISIONFUNCTION SELF X CALL DECISIONFUNCTION ON THE ESTIMATOR WITH THE BEST  
FOUND PARAMETERS

FITSELF X Y GROUPS RUN FIT WITH ALL SETS OF PARAMETERS

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

INVERSETRANSFORM SELF XT CALL INVERSETRANSFORM ON THE ESTIMATOR WITH THE BEST  
FOUND PARAMS

CONTINUED ON NEXT PAGE

626SKLEARNMODELSELECTION MODEL SELECTION 2107

SCIKITLEARN USER GUIDE RELEASE 0213

TABLE 6206 – CONTINUED FROM PREVIOUS PAGE

PREDICT SELF X CALL PREDICT ON THE ESTIMATOR WITH THE BEST FOUND PARAMETERS

PREDICTLOGPROBA SELF X CALL PREDICTLOGPROBA ON THE ESTIMATOR WITH THE BEST FOUND PARAMETERS

PREDICTPROBA SELF X CALL PREDICTPROBA ON THE ESTIMATOR WITH THE BEST FOUND PARAMETERS

SCORE SELF X Y RETURNS THE SCORE ON THE GIVEN DATA IF THE ESTIMATOR HAS BEEN REFIT

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF X CALL TRANSFORM ON THE ESTIMATOR WITH THE BEST FOUND PARAMETERS

INIT SELFESTIMATOR PARAMDISTRIBUTIONS NITER10 SCORINGNONE NJOBSNONE IID'WARN' REFITTRUE CV'WARN' VERBOSE0 PREDISPATCH'2NJOBS' RANDOMSTATENONE ERRORSORE'RAISEDEPRECATING' RETURNTRAINSCOREFALSE

DECISIONFUNCTION SELF X CALL DECISIONFUNCTION ON THE ESTIMATOR WITH THE BEST FOUND PARAMETERS ONLY AVAILABLE IF REFITTRUE AND THE UNDERLYING ESTIMATOR SUPPORTS DECISIONFUNCTION

XINDEXABLE LENGTH NSAMPLES MUST FULFILL THE INPUT ASSUMPTIONS OF THE UNDERLYING ESTIMATOR

FITSELFXYNONE GROUPSNONE FITPARAMS RUN FIT WITH ALL SETS OF PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING VECTOR WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUT OPTIONAL TARGET RELATIVE TO X FOR CLASSIFICATION OR REGRESSION NONE FOR UNSUPERVISED LEARNING

GROUPS ARRAYLIKE WITH SHAPE NSAMPLES OPTIONAL GROUP LABELS FOR THE SAMPLES USED WHILE SPLITTING THE DATASET INTO TRAINTEST SET ONLY USED IN CONJUNCTION WITH A "GROUP" CV INSTANCE EG GROUPKFOLD

FITPARAMS DICT OF STRING OBJECT PARAMETERS PASSED TO THE FIT METHOD OF THE ESTIMATOR

GETPARAMS SELFDEEPTREE GET PARAMETERS FOR THIS ESTIMATOR

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

INVERSETRANSFORM SELF X CALL INVERSETRANSFORM ON THE ESTIMATOR WITH THE BEST FOUND PARAMETERS ONLY AVAILABLE IF THE UNDERLYING ESTIMATOR IMPLEMENTS INVERSETRANSFORM ANDREFITTRUE

2108 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

XTINDEXABLE LENGTH NSAMPLES MUST FULFILL THE INPUT ASSUMPTIONS OF THE UNDERLYING ESTI  
MATOR

PREDICTSELF  
CALL PREDICT ON THE ESTIMATOR WITH THE BEST FOUND PARAMETERS  
ONLY AVAILABLE IF REFITTRUE AND THE UNDERLYING ESTIMATOR SUPPORTS PREDICT  
PARAMETERS

XINDEXABLE LENGTH NSAMPLES MUST FULFILL THE INPUT ASSUMPTIONS OF THE UNDERLYING ESTI  
MATOR

PREDICTLOGPROBA SELF  
CALL PREDICTLOGPROBA ON THE ESTIMATOR WITH THE BEST FOUND PARAMETERS  
ONLY AVAILABLE IF REFITTRUE AND THE UNDERLYING ESTIMATOR SUPPORTS PREDICTLOGPROBA  
PARAMETERS

XINDEXABLE LENGTH NSAMPLES MUST FULFILL THE INPUT ASSUMPTIONS OF THE UNDERLYING ESTI  
MATOR

PREDICTPROBA SELF  
CALL PREDICTPROBA ON THE ESTIMATOR WITH THE BEST FOUND PARAMETERS  
ONLY AVAILABLE IF REFITTRUE AND THE UNDERLYING ESTIMATOR SUPPORTS PREDICTPROBA  
PARAMETERS

XINDEXABLE LENGTH NSAMPLES MUST FULFILL THE INPUT ASSUMPTIONS OF THE UNDERLYING ESTI  
MATOR

SCORESELFXYNONE  
RETURNS THE SCORE ON THE GIVEN DATA IF THE ESTIMATOR HAS BEEN REFIT  
THIS USES THE SCORE DEFINED BY SCORING WHERE PROVIDED AND THE BESTESTIMATORSORE METHOD  
OTHERWISE  
PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES INPUT DATA WHERE NSAMPLES IS THE NUMBER  
OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUT OPTIONAL TARGET RELATIVE TO X  
FOR CLASSIFICATION OR REGRESSION NONE FOR UNSUPERVISED LEARNING

RETURNS  
SCORE FLOAT

SETPARAMS SELFPARAMS  
SET THE PARAMETERS OF THIS ESTIMATOR  
THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE  
PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT  
OF A NESTED OBJECT

RETURNS  
SELF

TRANSFORM SELF  
CALL TRANSFORM ON THE ESTIMATOR WITH THE BEST FOUND PARAMETERS  
ONLY AVAILABLE IF THE UNDERLYING ESTIMATOR SUPPORTS TRANSFORM ANDREFITTRUE

626SKLEARNMODELSELECTION MODEL SELECTION 2109

SCIKITLEARN USER GUIDE RELEASE 0213

PARAMETERS

XINDEXABLE LENGTH NSAMPLES MUST FULFILL THE INPUT ASSUMPTIONS OF THE UNDERLYING ESTI  
MATOR

EXAMPLES USING SKLEARNMODELSELECTIONRANDOMIZEDSEARCHCV

- COMPARING RANDOMIZED SEARCH AND GRID SEARCH FOR HYPERPARAMETER ESTIMATION

MODELSELECTIONFITGRIDPOINT X Y

    RUN FIT ON ONE SET OF PARAMETERS

SKLEARNMODELSELECTION FITGRIDPOINT

SKLEARNMODELSELECTION FITGRIDPOINT XYESTIMATOR PARAMETERS TRAIN TESTSCORER

VERBOSE ERRORSORE'RAISEDEPRECATING'

FITPARAMS

    RUN FIT ON ONE SET OF PARAMETERS

PARAMETERS

XARRAYLIKE SPARSE MATRIX OR LIST INPUT DATA

YARRAYLIKE OR NONE TARGETS FOR INPUT DATA

ESTIMATOR ESTIMATOR OBJECT A OBJECT OF THAT TYPE IS INSTANTIATED FOR EACH GRID POINT THIS IS  
ASSUMED TO IMPLEMENT THE SCIKITLEARN ESTIMATOR INTERFACE EITHER ESTIMATOR NEEDS TO PROVIDE  
ASCORE FUNCTION OR SCORING MUST BE PASSED

PARAMETERS DICT PARAMETERS TO BE SET ON ESTIMATOR FOR THIS GRID POINT

TRAIN NDARRAY DTYPE INT OR BOOL BOOLEAN MASK OR INDICES FOR TRAINING SET

TEST NDARRAY DTYPE INT OR BOOL BOOLEAN MASK OR INDICES FOR TEST SET

SCORER CALLABLE OR NONE THE SCORER CALLABLE OBJECT FUNCTION MUST HAVE ITS SIGNATURE AS  
SCORERESTIMATOR X Y

IFNONE THE ESTIMATOR'S SCORE METHOD IS USED

VERBOSE INT VERBOSITY LEVEL

FITPARAMS KWARGS ADDITIONAL PARAMETER PASSED TO THE FIT FUNCTION OF THE ESTIMATOR

ERRORSCORE 'RAISE' OR NUMERIC VALUE TO ASSIGN TO THE SCORE IF AN ERROR OCCURS IN ESTIMATOR  
FITTING IF SET TO 'RAISE' THE ERROR IS RAISED IF A NUMERIC VALUE IS GIVEN FITFAILEDWARNING  
IS RAISED THIS PARAMETER DOES NOT AFFECT THE REFIT STEP WHICH WILL ALWAYS RAISE THE ERROR  
DEFAULT IS 'RAISE' BUT FROM VERSION 022 IT WILL CHANGE TO NPNAN

RETURNS

SCORE FLOAT SCORE OF THIS PARAMETER SETTING ON GIVEN TEST SPLIT

PARAMETERS DICT THE PARAMETERS THAT HAVE BEEN EVALUATED

NSAMPLESTEST INT NUMBER OF TEST SAMPLES IN THIS SPLIT

6264 MODEL VALIDATION

2110 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213  
MODELSELECTIONCROSSVALIDATE ESTIMATOR  
XEVALUATE METRICS BY CROSSVALIDATION AND ALSO RECORD  
FITS CORE TIMES  
MODELSELECTIONCROSSVALPREDICT ESTIMATOR  
XGENERATE CROSSVALIDATED ESTIMATES FOR EACH INPUT DATA POINT  
MODELSELECTIONCROSSVALSCORE ESTIMATOR  
XEVALUATE A SCORE BY CROSSVALIDATION  
MODELSELECTIONLEARNINGCURVE ESTIMATOR X  
YLEARNING CURVE  
MODELSELECTIONPERMUTATIONTESTSCORE EVALUATE THE SIGNIFICANCE OF A CROSSVALIDATED SCORE WITH  
PERMUTATIONS  
MODELSELECTIONVALIDATIONCURVE ESTIMATOR  
VALIDATION CURVE  
SKLEARNMODELSELECTION CROSSVALIDATE  
SKLEARNMODELSELECTION CROSSVALIDATE ESTIMATOR XYNONE GROUPSNONE SCOR  
INGNONE CV'WARN' NJOBSNONE VERBOSE0  
FITPARAMSNONE PREDISPATCH'2NJOBS' RE  
TURNTRAINSCOREFALSE RETURNESTIMATORFALSE  
ERRORSCORE'RAISEDEPRECATING'  
EVALUATE METRICS BY CROSSVALIDATION AND ALSO RECORD FITSCORE TIMES  
READ MORE IN THE USER GUIDE  
PARAMETERS  
ESTIMATOR ESTIMATOR OBJECT IMPLEMENTING 'FIT' THE OBJECT TO USE TO FIT THE DATA  
XARRAYLIKE THE DATA TO FIT CAN BE FOR EXAMPLE A LIST OR AN ARRAY  
YARRAYLIKE OPTIONAL DEFAULT NONE THE TARGET VARIABLE TO TRY TO PREDICT IN THE CASE OF SUPER  
VISED LEARNING  
GROUPS ARRAYLIKE WITH SHAPE NSAMPLES OPTIONAL GROUP LABELS FOR THE SAMPLES USED WHILE  
SPLITTING THE DATASET INTO TRAI NTEST SET ONLY USED IN CONJUNCTION WITH A "GROUP" CVINSTANCE  
EGGROU PKFOLD  
SCORING STRING CALLABLE LISTTUPLE DICT OR NONE DEFAULT NONE A SINGLE STRING SEE THE SCORING  
PARAMETER DEFINING MODEL EVALUATION RULES OR A CALLABLE SEE DEFINING YOUR SCORING STRATEGY  
FROM METRIC FUNCTIONS TO EVALUATE THE PREDICTIONS ON THE TEST SET  
FOR EVALUATING MULTIPLE METRICS EITHER GIVE A LIST OF UNIQUE STRINGS OR A DICT WITH NAMES AS  
KEYS AND CALLABLES AS VALUES  
NOTE THAT WHEN USING CUSTOM SCORERS EACH SCORER SHOULD RETURN A SINGLE VALUE METRIC  
FUNCTIONS RETURNING A LISTARRAY OF VALUES CAN BE WRAPPED INTO MULTIPLE SCORERS THAT RETURN  
ONE VALUE EACH  
SEESPECIFYING MULTIPLE METRICS FOR EVALUATION FOR AN EXAMPLE  
IF NONE THE ESTIMATOR'S SCORE METHOD IS USED  
CVINT CROSSVALIDATION GENERATOR OR AN ITERABLE OPTIONAL DETERMINES THE CROSSVALIDATION SPLIT  
TING STRATEGY POSSIBLE INPUTS FOR CV ARE  
• NONE TO USE THE DEFAULT 3FOLD CROSS VALIDATION  
• INTEGER TO SPECIFY THE NUMBER OF FOLDS IN A STRATIFIEDKFOLD  
•CV SPLITTER  
626SKLEARNMODELSELECTION MODEL SELECTION 2111

SCIKITLEARN USER GUIDE RELEASE 0213

- AN ITERABLE YIELDING TRAIN TEST SPLITS AS ARRAYS OF INDICES
- FOR INTEGERNONE INPUTS IF THE ESTIMATOR IS A CLASSIFIER AND YIS EITHER BINARY OR MULTICLASS  
STRATIFIEDKFOLD IS USED IN ALL OTHER CASES KFOLD IS USED

REFER USER GUIDE FOR THE VARIOUS CROSSVALIDATION STRATEGIES THAT CAN BE USED HERE  
CHANGED IN VERSION 020 CVDEFAULT VALUE IF NONE WILL CHANGE FROM 3FOLD TO 5FOLD IN  
V022

NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF CPUS TO USE TO DO THE COMPUTA  
TIONNONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING

ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

VERBOSE INTEGER OPTIONAL THE VERBOSITY LEVEL

FITPARAMS DICT OPTIONAL PARAMETERS TO PASS TO THE FIT METHOD OF THE ESTIMATOR

PREDISPATCH INT OR STRING OPTIONAL CONTROLS THE NUMBER OF JOBS THAT GET DISPATCHED DURING  
PARALLEL EXECUTION REDUCING THIS NUMBER CAN BE USEFUL TO AVOID AN EXPLOSION OF MEMORY  
CONSUMPTION WHEN MORE JOBS GET DISPATCHED THAN CPUS CAN PROCESS THIS PARAMETER CAN  
BE

- NONE IN WHICH CASE ALL THE JOBS ARE IMMEDIATELY CREATED AND SPAWNED USE THIS FOR  
LIGHTWEIGHT AND FASTRUNNING JOBS TO AVOID DELAYS DUE TO ONDEMAND SPAWNING OF THE JOBS
- AN INT GIVING THE EXACT NUMBER OF TOTAL JOBS THAT ARE SPAWNED
- A STRING GIVING AN EXPRESSION AS A FUNCTION OF NJOBS AS IN '2NJOBS'

RETURNTRAINSCORE BOOLEAN DEFAULTFALSE WHETHER TO INCLUDE TRAIN SCORES COMPUTING TRAIN  
ING SCORES IS USED TO GET INSIGHTS ON HOW DIFFERENT PARAMETER SETTINGS IMPACT THE OVERFIT  
TINGUNDERFITTING TRADEOFF HOWEVER COMPUTING THE SCORES ON THE TRAINING SET CAN BE COM  
PUTATIONALLY EXPENSIVE AND IS NOT STRICTLY REQUIRED TO SELECT THE PARAMETERS THAT YIELD THE BEST  
GENERALIZATION PERFORMANCE

RETURNESTIMATOR BOOLEAN DEFAULT FALSE WHETHER TO RETURN THE ESTIMATORS FITTED ON EACH SPLIT  
ERRORSCORE 'RAISE' 'RAISEDEPRECATING' OR NUMERIC VALUE TO ASSIGN TO THE SCORE IF AN ERROR  
OCCURS IN ESTIMATOR FITTING IF SET TO 'RAISE' THE ERROR IS RAISED IF SET TO 'RAISEDEPRECATING'  
A FUTUREWARNING IS PRINTED BEFORE THE ERROR IS RAISED IF A NUMERIC VALUE IS GIVEN FITFAILED  
WARNING IS RAISED THIS PARAMETER DOES NOT AFFECT THE REFIT STEP WHICH WILL ALWAYS RAISE THE  
ERROR DEFAULT IS 'RAISEDEPRECATING' BUT FROM VERSION 022 IT WILL CHANGE TO NPNAN

RETURNS

SCORES DICT OF FLOAT ARRAYS OF SHAPENSPLITS ARRAY OF SCORES OF THE ESTIMATOR FOR EACH RUN OF  
THE CROSS VALIDATION

A DICT OF ARRAYS CONTAINING THE SCORETIME ARRAYS FOR EACH SCORER IS RETURNED THE POSSIBLE  
KEYS FOR THIS DICT ARE

TESTSCORE THE SCORE ARRAY FOR TEST SCORES ON EACH CV SPLIT

TRAINSCORE THE SCORE ARRAY FOR TRAIN SCORES ON EACH CV SPLIT THIS IS AVAILABLE  
ONLY IFRETURNTRAINSCORE PARAMETER IS TRUE

FITTIME THE TIME FOR FITTING THE ESTIMATOR ON THE TRAIN SET FOR EACH CV SPLIT

SCORETIME THE TIME FOR SCORING THE ESTIMATOR ON THE TEST SET FOR EACH CV

SPLIT NOTE TIME FOR SCORING ON THE TRAIN SET IS NOT INCLUDED EVEN IF

RETURNTRAINSCORE IS SET TOTRUE

2112 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

ESTIMATOR THE ESTIMATOR OBJECTS FOR EACH CV SPLIT THIS IS AVAILABLE ONLY IF  
RETURNESTIMATOR PARAMETER IS SET TO TRUE

SEE ALSO

SKLEARNMODELSELECTIONCROSSVALSCORE RUN CROSSVALIDATION FOR SINGLE METRIC EVALUATION  
SKLEARNMODELSELECTIONCROSSVALPREDICT GET PREDICTIONS FROM EACH SPLIT OF CROSS  
VALIDATION FOR DIAGNOSTIC PURPOSES  
SKLEARNMETRICSMAKESCORER MAKE A SCORER FROM A PERFORMANCE METRIC OR LOSS FUNCTION

EXAMPLES

FROM SKLEARN IMPORT DATASETS LINEARMODEL  
FROM SKLEARNMODELSELECTION IMPORT CROSSVALIDATE  
FROM SKLEARNMETRICSSCORER IMPORT MAKESCORER  
FROM SKLEARNMETRICS IMPORT CONFUSIONMATRIX  
FROM SKLEARN SVM IMPORT LINEARSVC

DIABETES DATASETSLOADDIABETES  
X DIABETESDATA150  
Y DIABETESTARGET150

LASSO LINEARMODELLASSO

SINGLE METRIC EVALUATION USING CROSSVALIDATE  
CVRESULTS CROSSVALIDATELASSO X Y CV3  
SORTEDCVRESULTSKEYS  
FITTIME SCORETIME TESTSCORE  
CVRESULTSTESTSCORE  
ARRAY033150734 008022311 003531764

MULTIPLE METRIC EVALUATION USING CROSSVALIDATE PLEASE REFER THE SCORING PARAMETER DOC FOR MORE IN  
FORMATION

SCORES CROSSVALIDATELASSO X Y CV3  
SCORINGR2 NEGMEANSQUAREDERROR  
RETURNTRAINSCORE TRUE  
PRINTSCORETESTNEGMEANSQUAREDERROR  
36355 35733 61147  
PRINTSCORESTRAINR2  
028010158 039088426 022784852

SKLEARNMODELSELECTION CROSSVALPREDICT  
SKLEARNMODELSELECTION CROSSVALPREDICT ESTIMATOR XYNONE GROUPSNONE  
CV'WARN' NJOBSNONE VERBOSE0  
FITPARAMSNONE PREDISPATCH'2NJOBS'  
METHOD'PREDICT'

GENERATE CROSSVALIDATED ESTIMATES FOR EACH INPUT DATA POINT

THE DATA IS SPLIT ACCORDING TO THE CV PARAMETER EACH SAMPLE BELONGS TO EXACTLY ONE TEST SET AND ITS PREDICTION IS  
COMPUTED WITH AN ESTIMATOR FITTED ON THE CORRESPONDING TRAINING SET

PASSING THESE PREDICTIONS INTO AN EVALUATION METRIC MAY NOT BE A VALID WAY TO MEASURE GENERALIZATION PERFOR  
MANCE RESULTS CAN DIFFER FROM CROSSVALIDATE ANDCROSSVALSCORE UNLESS ALL TESTS SETS HAVE EQUAL  
SIZE AND THE METRIC DECOMPOSES OVER SAMPLES

626SKLEARNMODELSELECTION MODEL SELECTION 2113

SCIKITLEARN USER GUIDE RELEASE 0213

READ MORE IN THE USER GUIDE

PARAMETERS

ESTIMATOR ESTIMATOR OBJECT IMPLEMENTING ‘FIT’ AND ‘PREDICT’ THE OBJECT TO USE TO FIT THE DATA  
XARRAYLIKE THE DATA TO FIT CAN BE FOR EXAMPLE A LIST OR AN ARRAY AT LEAST 2D

YARRAYLIKE OPTIONAL DEFAULT NONE THE TARGET VARIABLE TO TRY TO PREDICT IN THE CASE OF SUPER  
VISED LEARNING

GROUPS ARRAYLIKE WITH SHAPE NSAMPLES OPTIONAL GROUP LABELS FOR THE SAMPLES USED WHILE  
SPLITTING THE DATASET INTO TRAINTEST SET ONLY USED IN CONJUNCTION WITH A “GROUP” CVINSTANCE  
EGGROUPEKFOLD

CVINT CROSSVALIDATION GENERATOR OR AN ITERABLE OPTIONAL DETERMINES THE CROSSVALIDATION SPLIT  
TING STRATEGY POSSIBLE INPUTS FOR CV ARE

- NONE TO USE THE DEFAULT 3FOLD CROSS VALIDATION
- INTEGER TO SPECIFY THE NUMBER OF FOLDS IN A STRATIFIEDKFOLD
- CV SPLITTER
- AN ITERABLE YIELDING TRAIN TEST SPLITS AS ARRAYS OF INDICES

FOR INTEGERNONE INPUTS IF THE ESTIMATOR IS A CLASSIFIER AND YIS EITHER BINARY OR MULTICLASS  
STRATIFIEDKFOLD IS USED IN ALL OTHER CASES KFOLD IS USED

REFER USER GUIDE FOR THE VARIOUS CROSSVALIDATION STRATEGIES THAT CAN BE USED HERE  
CHANGED IN VERSION 020 CVDEFAULT VALUE IF NONE WILL CHANGE FROM 3FOLD TO 5FOLD IN  
V022

NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF CPUS TO USE TO DO THE COMPUTA  
TIONNONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING  
ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

VERBOSE INTEGER OPTIONAL THE VERBOSITY LEVEL

FITPARAMS DICT OPTIONAL PARAMETERS TO PASS TO THE FIT METHOD OF THE ESTIMATOR

PREDISPATCH INT OR STRING OPTIONAL CONTROLS THE NUMBER OF JOBS THAT GET DISPATCHED DURING  
PARALLEL EXECUTION REDUCING THIS NUMBER CAN BE USEFUL TO AVOID AN EXPLOSION OF MEMORY  
CONSUMPTION WHEN MORE JOBS GET DISPATCHED THAN CPUS CAN PROCESS THIS PARAMETER CAN  
BE

- NONE IN WHICH CASE ALL THE JOBS ARE IMMEDIATELY CREATED AND SPAWNED USE THIS FOR  
LIGHTWEIGHT AND FASTRUNNING JOBS TO AVOID DELAYS DUE TO ONDEMAND SPAWNING OF THE JOBS
- AN INT GIVING THE EXACT NUMBER OF TOTAL JOBS THAT ARE SPAWNED
- A STRING GIVING AN EXPRESSION AS A FUNCTION OF NJOBS AS IN ‘2NJOBS’

METHOD STRING OPTIONAL DEFAULT ‘PREDICT’ INVOKES THE PASSED METHOD NAME OF THE PASSED  
ESTIMATOR FOR METHOD‘PREDICTPROBA’ THE COLUMNS CORRESPOND TO THE CLASSES IN SORTED  
ORDER

RETURNS

PREDICTIONS NDARRAY THIS IS THE RESULT OF CALLING METHOD

SEE ALSO

CROSSVALSCORE CALCULATE SCORE FOR EACH CV SPLIT

2114 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

CROSSVALIDATE CALCULATE ONE OR MORE SCORES AND TIMINGS FOR EACH CV SPLIT

NOTES

IN THE CASE THAT ONE OR MORE CLASSES ARE ABSENT IN A TRAINING PORTION A DEFAULT SCORE NEEDS TO BE ASSIGNED TO ALL INSTANCES FOR THAT CLASS IF METHOD PRODUCES COLUMNS PER CLASS AS IN 'DECISIONFUNCTION' 'PREDICTPROBA' 'PREDICTLOGPROBA' FOR PREDICTPROBA THIS VALUE IS 0 IN ORDER TO ENSURE FINITE OUTPUT WE APPROXIMATE NEGATIVE INFINITY BY THE MINIMUM FINITE FLOAT VALUE FOR THE DTYPE IN OTHER CASES

EXAMPLES

```
FROM SKLEARN IMPORT DATASETS LINEARMODEL
FROM SKLEARNMODELSELECTION IMPORT CROSSVALPREDICT
DIABETES DATASETSLOADDIABETES
X DIABETESDATA150
Y DIABETESTARGET150
LASSO LINEARMODELLASSO
YPRED CROSSVALPREDICTLASSO X Y CV3
EXAMPLES USING SKLEARNMODELSELECTIONCROSSVALPREDICT
•PLOTING CROSSVALIDATED PREDICTIONS
SKLEARNMODELSELECTION CROSSVALSCORE
SKLEARNMODELSELECTION CROSSVALSCORE ESTIMATOR XYNONE GROUPSNONE SCOR
INGNONE CV'WARN' NJOBSNONE VERBOSE0
FITPARAMSNONE PREDISPATCH'2NJOBS'
ERRORSCORE'RAISEDEPRECATING'
EVALUATE A SCORE BY CROSSVALIDATION
READ MORE IN THE USER GUIDE
PARAMETERS
ESTIMATOR ESTIMATOR OBJECT IMPLEMENTING 'FIT' THE OBJECT TO USE TO FIT THE DATA
XARRAYLIKE THE DATA TO FIT CAN BE FOR EXAMPLE A LIST OR AN ARRAY
YARRAYLIKE OPTIONAL DEFAULT NONE THE TARGET VARIABLE TO TRY TO PREDICT IN THE CASE OF SUPER
VISED LEARNING
GROUPS ARRAYLIKE WITH SHAPE NSAMPLES OPTIONAL GROUP LABELS FOR THE SAMPLES USED WHILE
SPLITTING THE DATASET INTO TRAI NTTEST SET ONLY USED IN CONJUNCTION WITH A "GROUP" CVINSTANCE
EGGROUPKFOLD
SCORING STRING CALLABLE OR NONE OPTIONAL DEFAULT NONE A STRING SEE MODEL EVALUATION DOC
UMENTATION OR A SCORER CALLABLE OBJECT FUNCTION WITH SIGNATURE SCORERESTIMATOR
X Y WHICH SHOULD RETURN ONLY A SINGLE VALUE
SIMILAR TOCROSSVALIDATE BUT ONLY A SINGLE METRIC IS PERMITTED
IF NONE THE ESTIMATOR'S DEFAULT SCORER IF AVAILABLE IS USED
626SKLEARNMODELSELECTION MODEL SELECTION 2115
```

SCIKITLEARN USER GUIDE RELEASE 0213

CVINT CROSSVALIDATION GENERATOR OR AN ITERABLE OPTIONAL DETERMINES THE CROSSVALIDATION SPLITTING STRATEGY POSSIBLE INPUTS FOR CV ARE

- NONE TO USE THE DEFAULT 3FOLD CROSS VALIDATION
- INTEGER TO SPECIFY THE NUMBER OF FOLDS IN A STRATIFIEDKFOLD
- CV SPLITTER
- AN ITERABLE YIELDING TRAIN TEST SPLITS AS ARRAYS OF INDICES

FOR INTEGERNONE INPUTS IF THE ESTIMATOR IS A CLASSIFIER AND YIS EITHER BINARY OR MULTICLASS STRATIFIEDKFOLD IS USED IN ALL OTHER CASES KFOLD IS USED

REFER USER GUIDE FOR THE VARIOUS CROSSVALIDATION STRATEGIES THAT CAN BE USED HERE

CHANGED IN VERSION 020 CVDEFAULT VALUE IF NONE WILL CHANGE FROM 3FOLD TO 5FOLD IN V022

NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF CPUS TO USE TO DO THE COMPUTATIONNONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

VERBOSE INTEGER OPTIONAL THE VERBOSITY LEVEL

FITPARAMS DICT OPTIONAL PARAMETERS TO PASS TO THE FIT METHOD OF THE ESTIMATOR

PREDISPATCH INT OR STRING OPTIONAL CONTROLS THE NUMBER OF JOBS THAT GET DISPATCHED DURING PARALLEL EXECUTION REDUCING THIS NUMBER CAN BE USEFUL TO AVOID AN EXPLOSION OF MEMORY CONSUMPTION WHEN MORE JOBS GET DISPATCHED THAN CPUS CAN PROCESS THIS PARAMETER CAN BE

- NONE IN WHICH CASE ALL THE JOBS ARE IMMEDIATELY CREATED AND SPAWNED USE THIS FOR LIGHTWEIGHT AND FASTRUNNING JOBS TO AVOID DELAYS DUE TO ONDEMAND SPAWNING OF THE JOBS
- AN INT GIVING THE EXACT NUMBER OF TOTAL JOBS THAT ARE SPAWNED
- A STRING GIVING AN EXPRESSION AS A FUNCTION OF NJOBS AS IN '2NJOBS'

ERRORSCORE 'RAISE' 'RAISEDEPRECATING' OR NUMERIC VALUE TO ASSIGN TO THE SCORE IF AN ERROR OCCURS IN ESTIMATOR FITTING IF SET TO 'RAISE' THE ERROR IS RAISED IF SET TO 'RAISEDEPRECATING' A FUTUREWARNING IS PRINTED BEFORE THE ERROR IS RAISED IF A NUMERIC VALUE IS GIVEN FITFAILED WARNING IS RAISED THIS PARAMETER DOES NOT AFFECT THE REFIT STEP WHICH WILL ALWAYS RAISE THE ERROR DEFAULT IS 'RAISEDEPRECATING' BUT FROM VERSION 022 IT WILL CHANGE TO NPNAN

RETURNS

SCORES ARRAY OF FLOAT SHAPELENLISTCV ARRAY OF SCORES OF THE ESTIMATOR FOR EACH RUN OF THE CROSS VALIDATION

SEE ALSO

SKLEARNMODELSELECTIONCROSSVALIDATE TO RUN CROSSVALIDATION ON MULTIPLE METRICS AND ALSO TO RETURN TRAIN SCORES FIT TIMES AND SCORE TIMES

SKLEARNMODELSELECTIONCROSSVALPREDICT GET PREDICTIONS FROM EACH SPLIT OF CROSS VALIDATION FOR DIAGNOSTIC PURPOSES

SKLEARNMETRICSMAKESCORER MAKE A SCORER FROM A PERFORMANCE METRIC OR LOSS FUNCTION

2116 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES

FROM SKLEARN IMPORT DATASETS LINEARMODEL  
FROM SKLEARNMODELSELECTION IMPORT CROSSVALSCORE

DIABETES DATASETSLOADDIABETES

X DIABETESDATA150

Y DIABETESTARGET150

LASSO LINEARMODELLASSO

PRINTCROSSVALSCORELASSO X Y CV3

033150734 008022311 003531764

EXAMPLES USING SKLEARNMODELSELECTIONCROSSVALSCORE

- MODEL SELECTION WITH PROBABILISTIC PCA AND FACTOR ANALYSIS FA
- CROSSVALIDATION ON DIGITS DATASET EXERCISE
- IMPUTING MISSING VALUES WITH VARIANTS OF ITERATIVEIMPUTER
- IMPUTING MISSING VALUES BEFORE BUILDING AN ESTIMATOR
- UNDERFITTING VS OVERFITTING
- NESTED VERSUS NONNESTED CROSSVALIDATION
- SVMANOVA SVM WITH UNIVARIATE FEATURE SELECTION

SKLEARNMODELSELECTION LEARNINGCURVE

SKLEARNMODELSELECTION LEARNINGCURVE ESTIMATOR X Y GROUPSNONE

TRAINSIZESARRAY01 0325 0550775

1 CV'WARN' SCORINGNONE EX

PLOITINCREMENTALLEARNINGFALSE NJOBSNONE

PREDISPATCH'ALL' VERBOSE0 SHUFFLEFALSE RAN

DOMSTATENONE ERRORSORE'RAISEDEPRECATING'

LEARNING CURVE

DETERMINES CROSSVALIDATED TRAINING AND TEST SCORES FOR DIFFERENT TRAINING SET SIZES

A CROSSVALIDATION GENERATOR SPLITS THE WHOLE DATASET K TIMES IN TRAINING AND TEST DATA SUBSETS OF THE TRAINING SET  
WITH VARYING SIZES WILL BE USED TO TRAIN THE ESTIMATOR AND A SCORE FOR EACH TRAINING SUBSET SIZE AND THE TEST SET  
WILL BE COMPUTED AFTERWARDS THE SCORES WILL BE AVERAGED OVER ALL K RUNS FOR EACH TRAINING SUBSET SIZE  
READ MORE IN THE USER GUIDE

PARAMETERS

ESTIMATOR OBJECT TYPE THAT IMPLEMENTS THE "FIT" AND "PREDICT" METHODS AN OBJECT OF THAT TYPE  
WHICH IS CLONED FOR EACH VALIDATION

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING VECTOR WHERE NSAMPLES IS THE NUMBER  
OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NFEATURES OPTIONAL TARGET RELATIVE TO X FOR  
CLASSIFICATION OR REGRESSION NONE FOR UNSUPERVISED LEARNING

GROUPS ARRAYLIKE WITH SHAPE NSAMPLES OPTIONAL GROUP LABELS FOR THE SAMPLES USED WHILE  
SPLITTING THE DATASET INTO TRAINTEST SET ONLY USED IN CONJUNCTION WITH A "GROUP" CVINSTANCE  
EGGROUPKFOLD

626SKLEARNMODELSELECTION MODEL SELECTION 2117

SCIKITLEARN USER GUIDE RELEASE 0213

TRAINSIZES ARRAYLIKE SHAPE NTICKS DTYPE FLOAT OR INT RELATIVE OR ABSOLUTE NUMBERS OF TRAINING EXAMPLES THAT WILL BE USED TO GENERATE THE LEARNING CURVE IF THE DTYPE IS FLOAT IT IS REGARDED AS A FRACTION OF THE MAXIMUM SIZE OF THE TRAINING SET THAT IS DETERMINED BY THE SELECTED VALIDATION METHOD IE IT HAS TO BE WITHIN 0 1 OTHERWISE IT IS INTERPRETED AS ABSOLUTE SIZES OF THE TRAINING SETS NOTE THAT FOR CLASSIFICATION THE NUMBER OF SAMPLES USUALLY HAVE TO BE BIG ENOUGH TO CONTAIN AT LEAST ONE SAMPLE FROM EACH CLASS DEFAULT  
NPLINSPACE01 10 5

CVINT CROSSVALIDATION GENERATOR OR AN ITERABLE OPTIONAL DETERMINES THE CROSSVALIDATION SPLITTING STRATEGY POSSIBLE INPUTS FOR CV ARE

- NONE TO USE THE DEFAULT 3FOLD CROSS VALIDATION
- INTEGER TO SPECIFY THE NUMBER OF FOLDS IN A STRATIFIEDKFOLD
- CV SPLITTER
- AN ITERABLE YIELDING TRAIN TEST SPLITS AS ARRAYS OF INDICES

FOR INTEGERNONE INPUTS IF THE ESTIMATOR IS A CLASSIFIER AND YIS EITHER BINARY OR MULTICLASS STRATIFIEDKFOLD IS USED IN ALL OTHER CASES KFOLD IS USED

REFER USER GUIDE FOR THE VARIOUS CROSSVALIDATION STRATEGIES THAT CAN BE USED HERE

CHANGED IN VERSION 020 CVDEFAULT VALUE IF NONE WILL CHANGE FROM 3FOLD TO 5FOLD IN V022

SCORING STRING CALLABLE OR NONE OPTIONAL DEFAULT NONE A STRING SEE MODEL EVALUATION DOCUMENTATION OR A SCORER CALLABLE OBJECT FUNCTION WITH SIGNATURE SCORERESTIMATOR  
X Y

EXPLOITINCREMENTAL LEARNING BOOLEAN OPTIONAL DEFAULT FALSE IF THE ESTIMATOR SUPPORTS INCREMENTAL LEARNING THIS WILL BE USED TO SPEED UP FITTING FOR DIFFERENT TRAINING SET SIZES

NJOBS INT OR NONE OPTIONAL DEFAULTNONE NUMBER OF JOBS TO RUN IN PARALLEL NONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

PREDISPATCH INTEGER OR STRING OPTIONAL NUMBER OF PREDISPATCHED JOBS FOR PARALLEL EXECUTION DEFAULT IS ALL THE OPTION CAN REDUCE THE ALLOCATED MEMORY THE STRING CAN BE AN EXPRESSION LIKE '2NJOBS'

VERBOSE INTEGER OPTIONAL CONTROLS THE VERBOSITY THE HIGHER THE MORE MESSAGES

SHUFFLE BOOLEAN OPTIONAL WHETHER TO SHUFFLE TRAINING DATA BEFORE TAKING PREFIXES OF IT BASED ON "TRAINSIZES"

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RANDOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NRANDOM USED WHEN SHUFFLE IS TRUE

ERRORSCORE 'RAISE' 'RAISEDEPRECATING' OR NUMERIC VALUE TO ASSIGN TO THE SCORE IF AN ERROR OCCURS IN ESTIMATOR FITTING IF SET TO 'RAISE' THE ERROR IS RAISED IF SET TO 'RAISEDEPRECATING' A FUTUREWARNING IS PRINTED BEFORE THE ERROR IS RAISED IF A NUMERIC VALUE IS GIVEN FITFAILED WARNING IS RAISED THIS PARAMETER DOES NOT AFFECT THE REFIT STEP WHICH WILL ALWAYS RAISE THE ERROR DEFAULT IS 'RAISEDEPRECATING' BUT FROM VERSION 022 IT WILL CHANGE TO NPNAN

RETURNS

TRAINSIZESABS ARRAY SHAPE NUNIQUE TICKS DTYPE INT NUMBERS OF TRAINING EXAMPLES THAT HAS BEEN USED TO GENERATE THE LEARNING CURVE NOTE THAT THE NUMBER OF TICKS MIGHT BE LESS THAN NTICKS BECAUSE DUPLICATE ENTRIES WILL BE REMOVED

2118 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213  
TRAINSCORES ARRAY SHAPE NTICKS NCVFOLDS SCORES ON TRAINING SETS  
TESTSCORES ARRAY SHAPE NTICKS NCVFOLDS SCORES ON TEST SET  
NOTES  
SEEEEXAMPLESMODELSELECTIONPLOTLEARNINGCURVEPY  
EXAMPLES USING SKLEARNMODELSELECTIONLEARNINGCURVE  
•COMPARISON OF KERNEL RIDGE REGRESSION AND SVR  
•PLOTTING LEARNING CURVES  
SKLEARNMODELSELECTION PERMUTATIONTESTSCORE  
SKLEARNMODELSELECTION PERMUTATIONTESTSCORE ESTIMATOR XY GROUPSNONE  
CV'WARN' NPERMUTATIONS100  
NJOBSNONE RANDOMSTATE0 VER  
BOSE0 SCORINGNONE  
EVALUATE THE SIGNIFICANCE OF A CROSSVALIDATED SCORE WITH PERMUTATIONS  
READ MORE IN THE USER GUIDE  
PARAMETERS  
ESTIMATOR ESTIMATOR OBJECT IMPLEMENTING 'FIT' THE OBJECT TO USE TO FIT THE DATA  
XARRAYLIKE OF SHAPE AT LEAST 2D THE DATA TO FIT  
YARRAYLIKE THE TARGET VARIABLE TO TRY TO PREDICT IN THE CASE OF SUPERVISED LEARNING  
GROUPS ARRAYLIKE WITH SHAPE NSAMPLES OPTIONAL LABELS TO CONSTRAIN PERMUTATION WITHIN  
GROUPS IE YVALUES ARE PERMUTED AMONG SAMPLES WITH THE SAME GROUP IDENTIFIER WHEN NOT  
SPECIFIEDYVALUES ARE PERMUTED AMONG ALL SAMPLES  
WHEN A GROUPED CROSSVALIDATOR IS USED THE GROUP LABELS ARE ALSO PASSED ON TO THE SPLIT  
METHOD OF THE CROSSVALIDATOR THE CROSSVALIDATOR USES THEM FOR GROUPING THE SAMPLES WHILE  
SPLITTING THE DATASET INTO TRAI NTEST SET  
SCORING STRING CALLABLE OR NONE OPTIONAL DEFAULT NONE A SINGLE STRING SEE THE SCORING  
PARAMETER DEFINING MODEL EVALUATION RULES OR A CALLABLE SEE DEFINING YOUR SCORING STRATEGY  
FROM METRIC FUNCTIONS TO EVALUATE THE PREDICTIONS ON THE TEST SET  
IF NONE THE ESTIMATOR'S SCORE METHOD IS USED  
CVINT CROSSVALIDATION GENERATOR OR AN ITERABLE OPTIONAL DETERMINES THE CROSSVALIDATION SPLIT  
TING STRATEGY POSSIBLE INPUTS FOR CV ARE  
• NONE TO USE THE DEFAULT 3FOLD CROSS VALIDATION  
• INTEGER TO SPECIFY THE NUMBER OF FOLDS IN A STRATIFIEDKFOLD  
•CV SPLITTER  
• AN ITERABLE YIELDING TRAIN TEST SPLITS AS ARRAYS OF INDICES  
626SKLEARNMODELSELECTION MODEL SELECTION 2119

SCIKITLEARN USER GUIDE RELEASE 0213

FOR INTEGERNONE INPUTS IF THE ESTIMATOR IS A CLASSIFIER AND YIS EITHER BINARY OR MULTICLASS  
STRATIFIEDKFOLD IS USED IN ALL OTHER CASES KFOLD IS USED

REFER USER GUIDE FOR THE VARIOUS CROSSVALIDATION STRATEGIES THAT CAN BE USED HERE

CHANGED IN VERSION 020 CVDEFAULT VALUE IF NONE WILL CHANGE FROM 3FOLD TO 5FOLD IN  
V022

NPERMUTATIONS INTEGER OPTIONAL NUMBER OF TIMES TO PERMUTE Y

NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF CPUS TO USE TO DO THE COMPUTA  
TIONNONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING  
ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULT0 IF INT RANDOMSTATE  
IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS  
THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE  
INSTANCE USED BY NPRANDOM

VERBOSE INTEGER OPTIONAL THE VERBOSITY LEVEL

RETURNS

SCORE FLOAT THE TRUE SCORE WITHOUT PERMUTING TARGETS

PERMUTATIONSCORES ARRAY SHAPE NPERMUTATIONS THE SCORES OBTAINED FOR EACH PERMUTA  
TIONS

PVALUE FLOAT THE PVALUE WHICH APPROXIMATES THE PROBABILITY THAT THE SCORE WOULD BE OBTAINED  
BY CHANCE THIS IS CALCULATED AS

$C = \frac{1}{NPERMUTATIONS + 1}$

WHERE C IS THE NUMBER OF PERMUTATIONS WHOSE SCORE  $\geq$  THE TRUE SCORE  
THE BEST POSSIBLE PVALUE IS  $\frac{1}{NPERMUTATIONS + 1}$  THE WORST IS 10

NOTES

THIS FUNCTION IMPLEMENTS TEST 1 IN

OJALA AND GARRIGA PERMUTATION TESTS FOR STUDYING CLASSIFIER PERFORMANCE THE JOURNAL OF MACHINE  
LEARNING RESEARCH 2010 VOL 11

EXAMPLES USING SKLEARNMODELSELECTIONPERMUTATIONTESTSCORE

- TEST WITH PERMUTATIONS THE SIGNIFICANCE OF A CLASSIFICATION SCORE

SKLEARNMODELSELECTION VALIDATIONCURVE

SKLEARNMODELSELECTION VALIDATIONCURVE ESTIMATOR XYPARAMNAME PARAMRANGE

GROUPSNONE 'CV'WARN' SCORINGNONE

NJOBSNONE PREDISPATCH'ALL' VERBOSE0

ERRORSCORE'RAISEDEPRECATING'

VALIDATION CURVE

DETERMINE TRAINING AND TEST SCORES FOR VARYING PARAMETER VALUES

2120 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213  
COMPUTE SCORES FOR AN ESTIMATOR WITH DIFFERENT VALUES OF A SPECIFIED PARAMETER THIS IS SIMILAR TO GRID SEARCH WITH ONE PARAMETER HOWEVER THIS WILL ALSO COMPUTE TRAINING SCORES AND IS MERELY A UTILITY FOR PLOTTING THE RESULTS  
READ MORE IN THE USER GUIDE  
PARAMETERS  
ESTIMATOR OBJECT TYPE THAT IMPLEMENTS THE “FIT” AND “PREDICT” METHODS AN OBJECT OF THAT TYPE WHICH IS CLONED FOR EACH VALIDATION  
XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING VECTOR WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES  
YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NFEATURES OPTIONAL TARGET RELATIVE TO X FOR CLASSIFICATION OR REGRESSION NONE FOR UNSUPERVISED LEARNING  
PARAMNAME STRING NAME OF THE PARAMETER THAT WILL BE VARIED  
PARAMRANGE ARRAYLIKE SHAPE NVALUES THE VALUES OF THE PARAMETER THAT WILL BE EVALUATED  
GROUPS ARRAYLIKE WITH SHAPE NSAMPLES OPTIONAL GROUP LABELS FOR THE SAMPLES USED WHILE SPLITTING THE DATASET INTO TRAI NTEST SET ONLY USED IN CONJUNCTION WITH A “GROUP” CVINSTANCE  
EGGROU PKFOLD  
CVINT CROSSVALIDATION GENERATOR OR AN ITERABLE OPTIONAL DETERMINES THE CROSSVALIDATION SPLITTING STRATEGY POSSIBLE INPUTS FOR CV ARE  
• NONE TO USE THE DEFAULT 3FOLD CROSS VALIDATION  
• INTEGER TO SPECIFY THE NUMBER OF FOLDS IN A STRATIFIEDKFOLD  
• CV SPLITTER  
• AN ITERABLE YIELDING TRAIN TEST SPLITS AS ARRAYS OF INDICES  
FOR INTEGERNONE INPUTS IF THE ESTIMATOR IS A CLASSIFIER AND YIS EITHER BINARY OR MULTICLASS STRATIFIEDKFOLD IS USED IN ALL OTHER CASES KFOLD IS USED  
REFER USER GUIDE FOR THE VARIOUS CROSSVALIDATION STRATEGIES THAT CAN BE USED HERE  
CHANGED IN VERSION 020 CVDEFAULT VALUE IF NONE WILL CHANGE FROM 3FOLD TO 5FOLD IN V022  
SCORING STRING CALLABLE OR NONE OPTIONAL DEFAULT NONE A STRING SEE MODEL EVALUATION DOCUMENTATION OR A SCORER CALLABLE OBJECT FUNCTION WITH SIGNATURE SCORERESTIMATOR  
X Y  
NJOBS INT OR NONE OPTIONAL DEFAULTNONE NUMBER OF JOBS TO RUN IN PARALLEL NONE MEANS 1 UNLESS IN A JOBLIBPARALLELBKEND CONTEXT1MEANS USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS  
PREDISPATCH INTEGER OR STRING OPTIONAL NUMBER OF PREDISPATCHED JOBS FOR PARALLEL EXECUTION  
DEFAULT IS ALL THE OPTION CAN REDUCE THE ALLOCATED MEMORY THE STRING CAN BE AN EXPRESSION LIKE ‘2NJOBS’  
VERBOSE INTEGER OPTIONAL CONTROLS THE VERBOSITY THE HIGHER THE MORE MESSAGES  
ERRORSCORE ‘RAISE’ ‘RAISEDEPRECATING’ OR NUMERIC VALUE TO ASSIGN TO THE SCORE IF AN ERROR OCCURS IN ESTIMATOR FITTING IF SET TO ‘RAISE’ THE ERROR IS RAISED IF SET TO ‘RAISEDEPRECATING’ A FUTUREWARNING IS PRINTED BEFORE THE ERROR IS RAISED IF A NUMERIC VALUE IS GIVEN FITFAILED WARNING IS RAISED THIS PARAMETER DOES NOT AFFECT THE REFIT STEP WHICH WILL ALWAYS RAISE THE ERROR  
DEFAULT IS ‘RAISEDEPRECATING’ BUT FROM VERSION 022 IT WILL CHANGE TO NPNAN  
RETURNS  
626SKLEARNMODELSELECTION MODEL SELECTION 2121

SCIKITLEARN USER GUIDE RELEASE 0213

TRAINSCORES ARRAY SHAPE NTICKS NCVFOLDS SCORES ON TRAINING SETS

TESTSCORES ARRAY SHAPE NTICKS NCVFOLDS SCORES ON TEST SET

NOTES

SEEPLOTTING VALIDATION CURVES

EXAMPLES USING SKLEARNMODELSELECTIONVALIDATIONCURVE

- PLOTING VALIDATION CURVES

6275KLEARNMULTICLASS MULTICLASS AND MULTILABEL CLASSIFICATION

6271 MULTICLASS AND MULTILABEL CLASSIFICATION STRATEGIES

THIS MODULE IMPLEMENTS MULTICLASS LEARNING ALGORITHMS

- ONEVSTHEREST ONEVSALL
- ONEVSONE
- ERROR CORRECTING OUTPUT CODES

THE ESTIMATORS PROVIDED IN THIS MODULE ARE METAESTIMATORS THEY REQUIRE A BASE ESTIMATOR TO BE PROVIDED IN THEIR CONSTRUCTOR FOR EXAMPLE IT IS POSSIBLE TO USE THESE ESTIMATORS TO TURN A BINARY CLASSIFIER OR A REGRESSOR INTO A MULTICLASS CLASSIFIER IT IS ALSO POSSIBLE TO USE THESE ESTIMATORS WITH MULTICLASS ESTIMATORS IN THE HOPE THAT THEIR ACCURACY OR RUNTIME PERFORMANCE IMPROVES

ALL CLASSIFIERS IN SCIKITLEARN IMPLEMENT MULTICLASS CLASSIFICATION YOU ONLY NEED TO USE THIS MODULE IF YOU WANT TO EXPERIMENT WITH CUSTOM MULTICLASS STRATEGIES

THE ONEVSTHEREST METAClassIFIER ALSO IMPLEMENTS A PREDICTPROBA METHOD SO LONG AS SUCH A METHOD IS IMPLEMENTED BY THE BASE CLASSIFIER THIS METHOD RETURNS PROBABILITIES OF CLASS MEMBERSHIP IN BOTH THE SINGLE LABEL AND MULTILABEL CASE NOTE THAT IN THE MULTILABEL CASE PROBABILITIES ARE THE MARGINAL PROBABILITY THAT A GIVEN SAMPLE FALLS IN THE GIVEN CLASS AS SUCH IN THE MULTILABEL CASE THE SUM OF THESE PROBABILITIES OVER ALL POSSIBLE LABELS FOR A GIVEN SAMPLE WILL NOT SUM TO UNITY AS THEY DO IN THE SINGLE LABEL CASE

USER GUIDE SEE THE MULTICLASS AND MULTILABEL ALGORITHMS SECTION FOR FURTHER DETAILS

MULTICLASSONEVSRESTCLASSIFIER ESTIMATOR

- ONEVSTHEREST OVR MULTICLASSMULTILABEL STRATEGY

MULTICLASSONEVSONECLASSIFIER ESTIMATOR

- ONEVSONE MULTICLASS STRATEGY

MULTICLASSOUTPUTCODECLASSIFIER ESTIMATOR

- ERRORCORRECTING OUTPUTCODE MULTICLASS STRATEGY

6272SKLEARNMULTICLASS ONEVSRESTCLASSIFIER

CLASSSSKLEARNMULTICLASS ONEVSRESTCLASSIFIER ESTIMATOR NJOBSNONE

ONEVSTHEREST OVR MULTICLASSMULTILABEL STRATEGY

ALSO KNOWN AS ONEVSALL THIS STRATEGY CONSISTS IN FITTING ONE CLASSIFIER PER CLASS FOR EACH CLASSIFIER THE CLASS IS

2122 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

FITTED AGAINST ALL THE OTHER CLASSES IN ADDITION TO ITS COMPUTATIONAL EFFICIENCY ONLY NCLASSES CLASSIFIERS ARE NEEDED ONE ADVANTAGE OF THIS APPROACH IS ITS INTERPRETABILITY SINCE EACH CLASS IS REPRESENTED BY ONE AND ONE CLASSIFIER ONLY IT IS POSSIBLE TO GAIN KNOWLEDGE ABOUT THE CLASS BY INSPECTING ITS CORRESPONDING CLASSIFIER THIS IS THE MOST COMMONLY USED STRATEGY FOR MULTICLASS CLASSIFICATION AND IS A FAIR DEFAULT CHOICE THIS STRATEGY CAN ALSO BE USED FOR MULTILABEL LEARNING WHERE A CLASSIFIER IS USED TO PREDICT MULTIPLE LABELS FOR INSTANCE BY FITTING ON A 2D MATRIX IN WHICH CELL I J IS 1 IF SAMPLE I HAS LABEL J AND 0 OTHERWISE IN THE MULTILABEL LEARNING LITERATURE OVR IS ALSO KNOWN AS THE BINARY RELEVANCE METHOD READ MORE IN THE USER GUIDE

PARAMETERS

ESTIMATOR ESTIMATOR OBJECT AN ESTIMATOR OBJECT IMPLEMENTING FIT AND ONE OF DECISIONFUNCTION ORPREDICTPROBA NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF JOBS TO USE FOR THE COMPUTATION NONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

ATTRIBUTES

ESTIMATORS LIST OFNCLASSES ESTIMATORS ESTIMATORS USED FOR PREDICTIONS CLASSES ARRAY SHAPE NCLASSES CLASS LABELS LABELBINARIZER LABELBINARIZER OBJECT OBJECT USED TO TRANSFORM MULTICLASS LABELS TO BINARY LABELS AND VICEVERSA MULTILABEL BOOLEAN WHETHER THIS IS A MULTILABEL CLASSIFIER

METHODS

DECISIONFUNCTION SELF X RETURNS THE DISTANCE OF EACH SAMPLE FROM THE DECISION BOUNDARY FOR EACH CLASS FITSELF X Y FIT UNDERLYING ESTIMATORS GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR PARTIALFIT SELF X Y CLASSES PARTIALLY FIT UNDERLYING ESTIMATORS PREDICT SELF X PREDICT MULTICLASS TARGETS USING UNDERLYING ESTIMATORS PREDICTPROBA SELF X PROBABILITY ESTIMATES SCORE SELF X Y SAMPLEWEIGHT RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR INIT SELFESTIMATOR NJOBSNONE DECISIONFUNCTION SELF X RETURNS THE DISTANCE OF EACH SAMPLE FROM THE DECISION BOUNDARY FOR EACH CLASS THIS CAN ONLY BE USED WITH ESTIMATORS WHICH IMPLEMENT THE DECISIONFUNCTION METHOD PARAMETERS XARRAYLIKE SHAPE NSAMPLES NFEATURES RETURNS TARRAYLIKE SHAPE NSAMPLES NCLASSES FITSELFXY

SCIKITLEARN USER GUIDE RELEASE 0213

FIT UNDERLYING ESTIMATORS

PARAMETERS

XSPARSE ARRAYLIKE SHAPE NSAMPLES NFEATURES DATA

YSPARSE ARRAYLIKE SHAPE NSAMPLES NSAMPLES NCLASSES MULTICLASS TARGETS

AN INDICATOR MATRIX TURNS ON MULTILABEL CLASSIFICATION

RETURNS

SELF

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

MULTILABEL

WHETHER THIS IS A MULTILABEL CLASSIFIER

PARTIALFIT SELFXYCLASSESNONE

PARTIALLY FIT UNDERLYING ESTIMATORS

SHOULD BE USED WHEN MEMORY IS INEFFICIENT TO TRAIN ALL DATA CHUNKS OF DATA CAN BE PASSED IN SEVERAL

ITERATION

PARAMETERS

XSPARSE ARRAYLIKE SHAPE NSAMPLES NFEATURES DATA

YSPARSE ARRAYLIKE SHAPE NSAMPLES NSAMPLES NCLASSES MULTICLASS TARGETS

AN INDICATOR MATRIX TURNS ON MULTILABEL CLASSIFICATION

CLASSES ARRAY SHAPE NCLASSES CLASSES ACROSS ALL CALLS TO PARTIALFIT CAN BE OBTAINED

VIANPUNIQUEYALL WHERE YALL IS THE TARGET VECTOR OF THE ENTIRE DATASET THIS

ARGUMENT IS ONLY REQUIRED IN THE FIRST CALL OF PARTIALFIT AND CAN BE OMITTED IN THE SUBSEQUENT

CALLS

RETURNS

SELF

PREDICTSELF

PREDICT MULTICLASS TARGETS USING UNDERLYING ESTIMATORS

PARAMETERS

XSPARSE ARRAYLIKE SHAPE NSAMPLES NFEATURES DATA

RETURNS

YSPARSE ARRAYLIKE SHAPE NSAMPLES NSAMPLES NCLASSES PREDICTED MULTICLASS

TARGETS

PREDICTPROBA SELF

PROBABILITY ESTIMATES

THE RETURNED ESTIMATES FOR ALL CLASSES ARE ORDERED BY LABEL OF CLASSES

2124 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

NOTE THAT IN THE MULTILABEL CASE EACH SAMPLE CAN HAVE ANY NUMBER OF LABELS THIS RETURNS THE MARGINAL PROBABILITY THAT THE GIVEN SAMPLE HAS THE LABEL IN QUESTION FOR EXAMPLE IT IS ENTIRELY CONSISTENT THAT TWO LABELS BOTH HAVE A 90 PROBABILITY OF APPLYING TO A GIVEN SAMPLE  
IN THE SINGLE LABEL MULTICLASS CASE THE ROWS OF THE RETURNED MATRIX SUM TO 1

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES

RETURNS

TSPARSE ARRAYLIKE SHAPE NSAMPLES NCLASSES RETURNS THE PROBABILITY OF THE SAMPLE FOR EACH CLASS IN THE MODEL WHERE CLASSES ARE ORDERED AS THEY ARE IN SELFCLASSES

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS  
IN MULTILABEL CLASSIFICATION THIS IS THE SUBSET ACCURACY WHICH IS A HARSH METRIC SINCE YOU REQUIRE FOR EACH SAMPLE THAT EACH LABEL SET BE CORRECTLY PREDICTED

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE LABELS FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT MEAN ACCURACY OF SELFpredictX wrt Y

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

EXAMPLES USING SKLEARNMULTICLASSONEVSRESTCLASSIFIER

- MULTILABEL CLASSIFICATION
- RECEIVER OPERATING CHARACTERISTIC ROC
- PRECISIONRECALL
- CLASSIFIER CHAIN

6273SKLEARNMULTICLASS ONEVSONECLASSIFIER

CLASSSSKLEARNMULTICLASS ONEVSONECLASSIFIER ESTIMATOR NJOBSNONE  
ONEVSONE MULTICLASS STRATEGY

THIS STRATEGY CONSISTS IN FITTING ONE CLASSIFIER PER CLASS PAIR AT PREDICTION TIME THE CLASS WHICH RECEIVED THE MOST VOTES IS SELECTED SINCE IT REQUIRES TO FIT NCLASSES NCLASSES 1 2 CLASSIFIERS THIS METHOD  
IS USUALLY SLOWER THAN ONEVSTHEREST DUE TO ITS ONCLASSES2 COMPLEXITY HOWEVER THIS METHOD MAY BE  
ADVANTAGEOUS FOR ALGORITHMS SUCH AS KERNEL ALGORITHMS WHICH DON'T SCALE WELL WITH NSAMPLES THIS IS BECAUSE  
6275SKLEARNMULTICLASS MULTICLASS AND MULTILABEL CLASSIFICATION 2125

SCIKITLEARN USER GUIDE RELEASE 0213

EACH INDIVIDUAL LEARNING PROBLEM ONLY INVOLVES A SMALL SUBSET OF THE DATA WHEREAS WITH ONEVSTHEREST THE COMPLETE DATASET IS USED NCLASSES TIMES

READ MORE IN THE USER GUIDE

PARAMETERS

ESTIMATOR ESTIMATOR OBJECT AN ESTIMATOR OBJECT IMPLEMENTING FIT AND ONE OF DECISIONFUNCTION ORPREDICTPROBA

NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF JOBS TO USE FOR THE COMPUTATION

NONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

ATTRIBUTES

ESTIMATORS LIST OFNCLASSES NCLASSES 1 2 ESTIMATORS ESTIMATORS USED FOR PREDICTIONS

CLASSES NUMPY ARRAY OF SHAPE NCLASSES ARRAY CONTAINING LABELS

PAIRWISEINDICES LIST LENGTH LENESTIMATORS ORNONE INDICES OF SAMPLES USED WHEN TRAINING THE ESTIMATORS NONE WHENESTIMATOR DOES NOT HAVE PAIRWISE AT

TRIBUTE

METHODS

DECISIONFUNCTION SELF X DECISION FUNCTION FOR THE ONEVSONECLASSIFIER

FITSELF X Y FIT UNDERLYING ESTIMATORS

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PARTIALFIT SELF X Y CLASSES PARTIALLY FIT UNDERLYING ESTIMATORS

PREDICT SELF X ESTIMATE THE BEST CLASS LABEL FOR EACH SAMPLE IN X

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELFESTIMATOR NJOBSNONE

DECISIONFUNCTION SELF X

DECISION FUNCTION FOR THE ONEVSONECLASSIFIER

THE DECISION VALUES FOR THE SAMPLES ARE COMPUTED BY ADDING THE NORMALIZED SUM OF PAIRWISE CLASSIFICATION CONFIDENCE LEVELS TO THE VOTES IN ORDER TO DISAMBIGUATE BETWEEN THE DECISION VALUES WHEN THE VOTES FOR ALL THE CLASSES ARE EQUAL LEADING TO A TIE

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES

RETURNS

YARRAYLIKE SHAPE NSAMPLES NCLASSES

FITSELFXY

FIT UNDERLYING ESTIMATORS

PARAMETERS

XSPARSE ARRAYLIKE SHAPE NSAMPLES NFEATURES DATA

YARRAYLIKE SHAPE NSAMPLES MULTICLASS TARGETS

2126 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

RETURNS  
SELF

GETPARAMS SELFDEEPTREE  
GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS  
DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBOBJECTS THAT ARE ESTIMATORS

RETURNS  
PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PARTIALFIT SELFXYCLASSESNONE  
PARTIALLY FIT UNDERLYING ESTIMATORS  
SHOULD BE USED WHEN MEMORY IS INEFFICIENT TO TRAIN ALL DATA CHUNKS OF DATA CAN BE PASSED IN SEVERAL  
ITERATION WHERE THE FIRST CALL SHOULD HAVE AN ARRAY OF ALL TARGET VARIABLES

PARAMETERS  
XSPARSE ARRAYLIKE SHAPE NSAMPLES NFEATURES DATA  
YARRAYLIKE SHAPE NSAMPLES MULTICLASS TARGETS  
CLASSES ARRAY SHAPE NCLASSES CLASSES ACROSS ALL CALLS TO PARTIALFIT CAN BE OBTAINED  
VIANPUNIQUEYALL WHERE YALL IS THE TARGET VECTOR OF THE ENTIRE DATASET THIS  
ARGUMENT IS ONLY REQUIRED IN THE FIRST CALL OF PARTIALFIT AND CAN BE OMITTED IN THE SUBSEQUENT  
CALLS

RETURNS  
SELF

PREDICTSELF  
ESTIMATE THE BEST CLASS LABEL FOR EACH SAMPLE IN X  
THIS IS IMPLEMENTED AS ARGMAXDECISIONFUNCTIONX AXIS1 WHICH WILL RETURN THE LABEL  
OF THE CLASS WITH MOST VOTES BY ESTIMATORS PREDICTING THE OUTCOME OF A DECISION FOR EACH POSSIBLE CLASS PAIR

PARAMETERS  
XSPARSE ARRAYLIKE SHAPE NSAMPLES NFEATURES DATA

RETURNS  
YNUMPY ARRAY OF SHAPE NSAMPLES PREDICTED MULTICLASS TARGETS

SCORESELFXYSAMPLEWEIGHTNONE  
RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS  
IN MULTILABEL CLASSIFICATION THIS IS THE SUBSET ACCURACY WHICH IS A HARSH METRIC SINCE YOU REQUIRE FOR EACH  
SAMPLE THAT EACH LABEL SET BE CORRECTLY PREDICTED

PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES  
YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE LABELS FOR X  
SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS  
SCORE FLOAT MEAN ACCURACY OF SELF PREDICTX WRT Y

627SKLEARNMULTICLASS MULTICLASS AND MULTILABEL CLASSIFICATION 2127

SCIKITLEARN USER GUIDE RELEASE 0213

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

62745KLEARNMULTICLASS OUTPUTCODECLASSIFIER

CLASSSSKLEARNMULTICLASS OUTPUTCODECLASSIFIER ESTIMATOR CODESIZE15 RAN

DOMSTATENONE NJOBSNONE

ERRORCORRECTING OUTPUTCODE MULTICLASS STRATEGY

OUTPUTCODE BASED STRATEGIES CONSIST IN REPRESENTING EACH CLASS WITH A BINARY CODE AN ARRAY OF 0S AND 1S AT FITTING TIME ONE BINARY CLASSIFIER PER BIT IN THE CODE BOOK IS FITTED AT PREDICTION TIME THE CLASSIFIERS ARE USED TO PROJECT NEW POINTS IN THE CLASS SPACE AND THE CLASS CLOSEST TO THE POINTS IS CHOSEN THE MAIN ADVANTAGE OF THESE STRATEGIES IS THAT THE NUMBER OF CLASSIFIERS USED CAN BE CONTROLLED BY THE USER EITHER FOR COMPRESSING THE MODEL 0 CODESIZE 1 OR FOR MAKING THE MODEL MORE ROBUST TO ERRORS CODESIZE 1 SEE THE DOCUMENTATION FOR MORE DETAILS

READ MORE IN THE USER GUIDE

PARAMETERS

ESTIMATOR ESTIMATOR OBJECT AN ESTIMATOR OBJECT IMPLEMENTING FIT AND ONE OF DECISIONFUNCTION ORPREDICTPROBA

CODESIZE FLOAT PERCENTAGE OF THE NUMBER OF CLASSES TO BE USED TO CREATE THE CODE BOOK A NUMBER BETWEEN 0 AND 1 WILL REQUIRE FEWER CLASSIFIERS THAN ONEVSTHEREST A NUMBER GREATER THAN 1 WILL REQUIRE MORE CLASSIFIERS THAN ONEVSTHEREST

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULT NONE THE GENERATOR USED TO INITIALIZE THE CODEBOOK IF INT RANDOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NPRANDOM

NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF JOBS TO USE FOR THE COMPUTATION NONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

ATTRIBUTES

ESTIMATORS LIST OFINTNCLASSES CODESIZE ESTIMATORS ESTIMATORS USED FOR PRE

DICTIONS

CLASSES NUMPY ARRAY OF SHAPE NCLASSES ARRAY CONTAINING LABELS

CODEBOOK NUMPY ARRAY OF SHAPE NCLASSES CODESIZE BINARY ARRAY CONTAINING THE CODE OF EACH CLASS

REFERENCES

R2EDDAEEC08491 R2EDDAEEC08492 R2EDDAEEC08493

2128 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

METHODS

FITSELF X Y FIT UNDERLYING ESTIMATORS

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PREDICT SELF X PREDICT MULTICLASS TARGETS USING UNDERLYING ESTIMATORS

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELFESTIMATOR CODESIZE15 RANDOMSTATENONE NJOBSNONE

FITSELFXY

FIT UNDERLYING ESTIMATORS

PARAMETERS

XSPARSE ARRAYLIKE SHAPE NSAMPLES NFEATURES DATA

YNUMPY ARRAY OF SHAPE NSAMPLES MULTICLASS TARGETS

RETURNS

SELF

GETPARAMS SELFDEEPTREE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELF

PREDICT MULTICLASS TARGETS USING UNDERLYING ESTIMATORS

PARAMETERS

XSPARSE ARRAYLIKE SHAPE NSAMPLES NFEATURES DATA

RETURNS

YNUMPY ARRAY OF SHAPE NSAMPLES PREDICTED MULTICLASS TARGETS

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS

IN MULTILABEL CLASSIFICATION THIS IS THE SUBSET ACCURACY WHICH IS A HARSH METRIC SINCE YOU REQUIRE FOR EACH SAMPLE THAT EACH LABEL SET BE CORRECTLY PREDICTED

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE LABELS FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT MEAN ACCURACY OF SELF PREDICTX WRT Y

627SKLEARNMULTICLASS MULTICLASS AND MULTILABEL CLASSIFICATION 2129

SCIKITLEARN USER GUIDE RELEASE 0213

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

628SKLEARNMULTIOUTPUT MULTIOUTPUT REGRESSION AND CLASSIFICATION

THIS MODULE IMPLEMENTS MULTIOUTPUT REGRESSION AND CLASSIFICATION

THE ESTIMATORS PROVIDED IN THIS MODULE ARE METAESTIMATORS THEY REQUIRE A BASE ESTIMATOR TO BE PROVIDED IN THEIR CONSTRUCTOR THE METAESTIMATOR EXTENDS SINGLE OUTPUT ESTIMATORS TO MULTIOUTPUT ESTIMATORS

USER GUIDE SEE THE MULTICLASS AND MULTILABEL ALGORITHMS SECTION FOR FURTHER DETAILS

MULTIOUTPUTCLASSIFIERCHAIN BASEESTIMATOR A MULTILABEL MODEL THAT ARRANGES BINARY CLASSIFIERS INTO A CHAIN

MULTIOUTPUTMULTIOUTPUTREGRESSOR ESTIMATOR MULTI TARGET REGRESSION

MULTIOUTPUTMULTIOUTPUTCLASSIFIER ESTIMATOR MULTI TARGET CLASSIFICATION

MULTIOUTPUTREGRESSORCHAIN BASEESTIMATOR

A MULTILABEL MODEL THAT ARRANGES REGRESSIONS INTO A CHAIN

6281SKLEARNMULTIOUTPUT CLASSIFIERCHAIN

CLASSSKLEARNMULTIOUTPUT CLASSIFIERCHAIN BASEESTIMATOR ORDERNONE CVNONE RAN

DOMSTATENONE

A MULTILABEL MODEL THAT ARRANGES BINARY CLASSIFIERS INTO A CHAIN

EACH MODEL MAKES A PREDICTION IN THE ORDER SPECIFIED BY THE CHAIN USING ALL OF THE AVAILABLE FEATURES PROVIDED TO THE MODEL PLUS THE PREDICTIONS OF MODELS THAT ARE EARLIER IN THE CHAIN

READ MORE IN THE USER GUIDE

PARAMETERS

BASEESTIMATOR ESTIMATOR THE BASE ESTIMATOR FROM WHICH THE CLASSIFIER CHAIN IS BUILT

ORDER ARRAYLIKE SHAPENOUTPUTS OR 'RANDOM' OPTIONAL BY DEFAULT THE ORDER WILL BE DETERMINED BY THE ORDER OF COLUMNS IN THE LABEL MATRIX Y

ORDER 0 1 2 YSHAPE1 1

THE ORDER OF THE CHAIN CAN BE EXPLICITLY SET BY PROVIDING A LIST OF INTEGERS FOR EXAMPLE FOR A CHAIN OF LENGTH 5

ORDER 1 3 2 4 0

MEANS THAT THE FIRST MODEL IN THE CHAIN WILL MAKE PREDICTIONS FOR COLUMN 1 IN THE Y MATRIX

THE SECOND MODEL WILL MAKE PREDICTIONS FOR COLUMN 3 ETC

2130 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

IF ORDER IS ‘RANDOM’ A RANDOM ORDERING WILL BE USED

CVINT CROSSVALIDATION GENERATOR OR AN ITERABLE OPTIONAL DEFAULTNONE DETERMINES WHETHER TO USE CROSS VALIDATED PREDICTIONS OR TRUE LABELS FOR THE RESULTS OF PREVIOUS ESTIMATORS IN THE CHAIN IF CV IS NONE THE TRUE LABELS ARE USED WHEN FITTING OTHERWISE POSSIBLE INPUTS FOR CV ARE

- INTEGER TO SPECIFY THE NUMBER OF FOLDS IN A STRATIFIEDKFOLD
- CV SPLITTER
- AN ITERABLE YIELDING TRAIN TEST SPLITS AS ARRAYS OF INDICES

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RAN

DOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE

RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS

THE RANDOMSTATE INSTANCE USED BY NPRANDOM

THE RANDOM NUMBER GENERATOR IS USED TO GENERATE RANDOM CHAIN ORDERS

ATTRIBUTES

CLASSES LIST A LIST OF ARRAYS OF LENGTH LENESTIMATORS CONTAINING THE CLASS LABELS FOR EACH ESTIMATOR IN THE CHAIN

ESTIMATORS LIST A LIST OF CLONES OF BASEESTIMATOR

ORDER LIST THE ORDER OF LABELS IN THE CLASSIFIER CHAIN

SEE ALSO

REGRESSORCHAIN EQUIVALENT FOR REGRESSION

MULTIOUTPUTCLASSIFIER CLASSIFIES EACH OUTPUT INDEPENDENTLY RATHER THAN CHAINING

REFERENCES

JESSE READ BERNHARD PFAHRINGER GEOFF HOLMES EIBE FRANK “CLASSIFIER CHAINS FOR MULTILABEL CLASSIFICATION” 2009

METHODS

DECISIONFUNCTION SELF X EVALUATE THE DECISIONFUNCTION OF THE MODELS IN THE CHAIN

FITSELF X Y FIT THE MODEL TO DATA MATRIX X AND TARGETS Y

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PREDICT SELF X PREDICT ON THE DATA MATRIX X USING THE CLASSIFIERCHAIN

MODEL

PREDICTPROBA SELF X PREDICT PROBABILITY ESTIMATES

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELFBASEESTIMATOR ORDERNONE CVNONE RANDOMSTATENONE

DECISIONFUNCTION SELF X

EVALUATE THE DECISIONFUNCTION OF THE MODELS IN THE CHAIN

628SKLEARNMULTIOUTPUT MULTIOUTPUT REGRESSION AND CLASSIFICATION 2131

SCIKITLEARN USER GUIDE RELEASE 0213  
 PARAMETERS  
 XARRAYLIKE SHAPE NSAMPLES NFEATURES  
 RETURNS  
 YDECISION ARRAYLIKE SHAPE NSAMPLES NCLASSES RETURNS THE DECISION FUNCTION OF THE  
 SAMPLE FOR EACH MODEL IN THE CHAIN  
 FITSELFXY  
 FIT THE MODEL TO DATA MATRIX X AND TARGETS Y  
 PARAMETERS  
 XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT DATA  
 YARRAYLIKE SHAPE NSAMPLES NCLASSES THE TARGET VALUES  
 RETURNS  
 SELF OBJECT  
 GETPARAMS SELFDEEPTREE  
 GET PARAMETERS FOR THIS ESTIMATOR  
 PARAMETERS  
 DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
 SUBOBJECTS THAT ARE ESTIMATORS  
 RETURNS  
 PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES  
 PREDICTSELF  
 PREDICT ON THE DATA MATRIX X USING THE CLASSIFIERCHAIN MODEL  
 PARAMETERS  
 XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT DATA  
 RETURNS  
 YPRED ARRAYLIKE SHAPE NSAMPLES NCLASSES THE PREDICTED VALUES  
 PREDICTPROBA SELF  
 PREDICT PROBABILITY ESTIMATES  
 PARAMETERS  
 XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES  
 RETURNS  
 YPROB ARRAYLIKE SHAPE NSAMPLES NCLASSES  
 SCORESELFXYSAMPLEWEIGHTNONE  
 RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS  
 IN MULTILABEL CLASSIFICATION THIS IS THE SUBSET ACCURACY WHICH IS A HARSH METRIC SINCE YOU REQUIRE FOR EACH  
 SAMPLE THAT EACH LABEL SET BE CORRECTLY PREDICTED  
 PARAMETERS  
 XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES  
 YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE LABELS FOR X  
 SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS  
 2132 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

RETURNS

SCORE FLOAT MEAN ACCURACY OF SELF PREDICT X WRT Y

SETPARAMS SELF PARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENT PARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

EXAMPLES USING SKLEARN MULTI OUTPUT CLASSIFIER CHAIN

- CLASSIFIER CHAIN

6282 SKLEARN MULTI OUTPUT MULTI OUTPUT REGRESSOR

CLASS SKLEARN MULTI OUTPUT MULTI OUTPUT REGRESSOR ESTIMATOR NJOBS NONE

MULTI TARGET REGRESSION

THIS STRATEGY CONSISTS OF FITTING ONE REGRESSOR PER TARGET THIS IS A SIMPLE STRATEGY FOR EXTENDING REGRESSORS THAT DO NOT NATIVELY SUPPORT MULTI TARGET REGRESSION

PARAMETERS

ESTIMATOR ESTIMATOR OBJECT AN ESTIMATOR OBJECT IMPLEMENTING FIT AND PREDICT

NJOBS INT OR NONE OPTIONAL DEFAULT NONE THE NUMBER OF JOBS TO RUN IN PARALLEL FOR FIT

NONE MEANS 1 UNLESS IN A JOBLIB PARALLEL BACKEND CONTEXT 1 MEANS USING ALL

PROCESSORS SEE GLOSSARY FOR MORE DETAILS

WHEN INDIVIDUAL ESTIMATORS ARE FAST TO TRAIN OR PREDICT USING NJOBS 1 CAN RESULT IN SLOWER PERFORMANCE DUE TO THE OVERHEAD OF SPAWNING PROCESSES

ATTRIBUTES

ESTIMATORS LIST OF N OUTPUT ESTIMATORS ESTIMATORS USED FOR PREDICTIONS

METHODS

FIT SELF X Y SAMPLE WEIGHT FIT THE MODEL TO DATA

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PARTIAL FIT SELF X Y SAMPLE WEIGHT INCREMENTALLY FIT THE MODEL TO DATA

PREDICT SELF X PREDICT MULTI OUTPUT VARIABLE USING A MODEL TRAINED FOR EACH TARGET VARIABLE

SCORE SELF X Y SAMPLE WEIGHT RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELF ESTIMATOR NJOBS NONE

FIT SELF X Y SAMPLE WEIGHT NONE

FIT THE MODEL TO DATA FIT A SEPARATE MODEL FOR EACH OUTPUT VARIABLE

628 SKLEARN MULTI OUTPUT MULTI OUTPUT REGRESSION AND CLASSIFICATION 2133

SCIKITLEARN USER GUIDE RELEASE 0213

PARAMETERS

XSPARSE ARRAYLIKE SHAPE NSAMPLES NFEATURES DATA

YSPARSE ARRAYLIKE SHAPE NSAMPLES NOUTPUTS MULTIOUTPUT TARGETS AN INDICATOR MATRIX TURNS ON MULTILABEL ESTIMATION

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OR NONE SAMPLE WEIGHTS IF NONE THEN SAMPLES ARE EQUALLY WEIGHTED ONLY SUPPORTED IF THE UNDERLYING REGRESSOR SUPPORTS SAMPLE WEIGHTS

RETURNS

SELF OBJECT

GETPARAMS SELFDEEPTREE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PARTIALFIT SELFXYSAMPLEWEIGHTNONE

INCREMENTALLY FIT THE MODEL TO DATA FIT A SEPARATE MODEL FOR EACH OUTPUT VARIABLE

PARAMETERS

XSPARSE ARRAYLIKE SHAPE NSAMPLES NFEATURES DATA

YSPARSE ARRAYLIKE SHAPE NSAMPLES NOUTPUTS MULTIOUTPUT TARGETS

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OR NONE SAMPLE WEIGHTS IF NONE THEN SAMPLES ARE EQUALLY WEIGHTED ONLY SUPPORTED IF THE UNDERLYING REGRESSOR SUPPORTS SAMPLE WEIGHTS

RETURNS

SELF OBJECT

PREDICTSELF

PREDICT MULTIOUTPUT VARIABLE USING A MODEL TRAINED FOR EACH TARGET VARIABLE

PARAMETERS

XSPARSE ARRAYLIKE SHAPE NSAMPLES NFEATURES DATA

RETURNS

YSPARSE ARRAYLIKE SHAPE NSAMPLES NOUTPUTS MULTIOUTPUT TARGETS PREDICTED ACROSS MULTIPLE PREDICTORS NOTE SEPARATE MODELS ARE GENERATED FOR EACH PREDICTOR

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{U}{V}$  WHERE U IS THE RESIDUAL SUM OF SQUARES  $\sum (y_{true} - y_{pred})^2$  AND V IS THE REGRESSION SUM OF SQUARES  $\sum (y_{true} - y_{true\_mean})^2$  SUM OF SQUARES YTRUE YTRUEMEAN 2SUM BEST POSSIBLE SCORE IS 10 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 00

PARAMETERS

2134 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE VALUES FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT R2 OF SELF PREDICT X WRT Y

NOTES

R2 IS CALCULATED BY WEIGHTING ALL THE TARGETS EQUALLY USING MULTIOUTPUT UNIFORM AVERAGE

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENT PARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

EXAMPLES USING SKLEARN MULTIOUTPUT MULTIOUTPUT REGRESSOR

- COMPARING RANDOM FORESTS AND THE MULTIOUTPUT META ESTIMATOR

6283 SKLEARN MULTIOUTPUT MULTIOUTPUT CLASSIFIER

CLASS SKLEARN MULTIOUTPUT MULTIOUTPUT CLASSIFIER ESTIMATOR NJOBS NONE

MULTI TARGET CLASSIFICATION

THIS STRATEGY CONSISTS OF FITTING ONE CLASSIFIER PER TARGET THIS IS A SIMPLE STRATEGY FOR EXTENDING CLASSIFIERS THAT DO NOT NATIVELY SUPPORT MULTITARGET CLASSIFICATION

PARAMETERS

ESTIMATOR ESTIMATOR OBJECT AN ESTIMATOR OBJECT IMPLEMENTING FIT SCORE AND PREDICT PROBA

NJOBS INT OR NONE OPTIONAL DEFAULT NONE THE NUMBER OF JOBS TO USE FOR THE COMPUTATION IT DOES EACH TARGET VARIABLE IN Y IN PARALLEL NONE MEANS 1 UNLESS IN A JOBLIB PARALLEL BACKEND CONTEXT 1 MEANS USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

ATTRIBUTES

ESTIMATORS LIST OF N OUTPUT ESTIMATORS ESTIMATORS USED FOR PREDICTIONS

METHODS

FIT SELF X Y SAMPLEWEIGHT FIT THE MODEL TO DATA

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PARTIALFIT SELF X Y CLASSES SAMPLEWEIGHT INCREMENTALLY FIT THE MODEL TO DATA

CONTINUED ON NEXT PAGE

6283 SKLEARN MULTIOUTPUT MULTIOUTPUT REGRESSION AND CLASSIFICATION 2135

SCIKITLEARN USER GUIDE RELEASE 0213

TABLE 6216 – CONTINUED FROM PREVIOUS PAGE

PREDICT SELF X PREDICT MULTIOUTPUT VARIABLE USING A MODEL TRAINED FOR EACH TARGET VARIABLE

PREDICTPROBA SELF X PROBABILITY ESTIMATES

SCORE SELF X Y RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELFESTIMATOR NJOBSNONE

FITSELFXYSAMPLEWEIGHTNONE

FIT THE MODEL TO DATA FIT A SEPARATE MODEL FOR EACH OUTPUT VARIABLE

PARAMETERS

XSPARSE ARRAYLIKE SHAPE NSAMPLES NFEATURES DATA

YSPARSE ARRAYLIKE SHAPE NSAMPLES NOUTPUTS MULTIOUTPUT TARGETS AN INDICATOR MATRIX TURNS ON MULTILABEL ESTIMATION

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OR NONE SAMPLE WEIGHTS IF NONE THEN SAMPLES ARE EQUALLY WEIGHTED ONLY SUPPORTED IF THE UNDERLYING REGRESSOR SUPPORTS SAMPLE WEIGHTS

RETURNS

SELF OBJECT

GETPARAMS SELFDEEPTREE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PARTIALFIT SELFXYCLASSESNONE SAMPLEWEIGHTNONE

INCREMENTALLY FIT THE MODEL TO DATA FIT A SEPARATE MODEL FOR EACH OUTPUT VARIABLE

PARAMETERS

XSPARSE ARRAYLIKE SHAPE NSAMPLES NFEATURES DATA

YSPARSE ARRAYLIKE SHAPE NSAMPLES NOUTPUTS MULTIOUTPUT TARGETS

CLASSES LIST OF NUMPY ARRAYS SHAPE NOUTPUTS EACH ARRAY IS UNIQUE CLASSES FOR ONE OUTPUT

IN STRINT CAN BE OBTAINED BY VIA NPUNIQUEY I FOR I IN RANGEY

SHAPE1 WHERE Y IS THE TARGET MATRIX OF THE ENTIRE DATASET THIS ARGUMENT IS REQUIRED FOR THE FIRST CALL TO PARTIALFIT AND CAN BE OMITTED IN THE SUBSEQUENT CALLS NOTE THAT Y DOESN'T NEED TO CONTAIN ALL LABELS IN CLASSES

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OR NONE SAMPLE WEIGHTS IF NONE THEN SAMPLES ARE EQUALLY WEIGHTED ONLY SUPPORTED IF THE UNDERLYING REGRESSOR SUPPORTS SAMPLE WEIGHTS

RETURNS

SELF OBJECT

PREDICTSELFXY

2136 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

PREDICT MULTIOUTPUT VARIABLE USING A MODEL TRAINED FOR EACH TARGET VARIABLE

PARAMETERS

XSPARSE ARRAYLIKE SHAPE NSAMPLES NFEATURES DATA

RETURNS

YSPARSE ARRAYLIKE SHAPE NSAMPLES NOUTPUTS MULTIOUTPUT TARGETS PREDICTED ACROSS MULTIPLE PREDICTORS NOTE SEPARATE MODELS ARE GENERATED FOR EACH PREDICTOR

PREDICTPROBA SELF

PROBABILITY ESTIMATES RETURNS PREDICTION PROBABILITIES FOR EACH CLASS OF EACH OUTPUT THIS METHOD WILL RAISE A VALUEERROR IF ANY OF THE ESTIMATORS DO NOT HAVE PREDICTPROBA

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES DATA

RETURNS

PARRAY OF SHAPE NSAMPLES NCLASSES OR A LIST OF NOUTPUTS SUCH ARRAYS IF NOUTPUTS 1 THE CLASS PROBABILITIES OF THE INPUT SAMPLES THE ORDER OF THE CLASSES CORRESPONDS TO THAT IN THE ATTRIBUTE CLASSES

SCORESELFXY

RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES

YARRAYLIKE SHAPE NSAMPLES NOUTPUTS TRUE VALUES FOR X

RETURNS

SCORES FLOAT ACCURACYScore OF SELFpredictX VERSUS Y

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

6284SKLEARNMULTIOUTPUT REGRESSORCHAIN

CLASSSKLEARNMULTIOUTPUT REGRESSORCHAIN BASEESTIMATOR ORDERNONE CVNONE RAN

DOMSTATENONE

A MULTILABEL MODEL THAT ARRANGES REGRESSIONS INTO A CHAIN

EACH MODEL MAKES A PREDICTION IN THE ORDER SPECIFIED BY THE CHAIN USING ALL OF THE AVAILABLE FEATURES PROVIDED TO THE MODEL PLUS THE PREDICTIONS OF MODELS THAT ARE EARLIER IN THE CHAIN

READ MORE IN THE USER GUIDE

PARAMETERS

BASEESTIMATOR ESTIMATOR THE BASE ESTIMATOR FROM WHICH THE CLASSIFIER CHAIN IS BUILT

628SKLEARNMULTIOUTPUT MULTIOUTPUT REGRESSION AND CLASSIFICATION 2137

SCIKITLEARN USER GUIDE RELEASE 0213

ORDER ARRAYLIKE SHAPENOUTPUTS OR ‘RANDOM’ OPTIONAL BY DEFAULT THE ORDER WILL BE DETERMINED BY THE ORDER OF COLUMNS IN THE LABEL MATRIX Y

ORDER 0 1 2 YSHAPE1 1

THE ORDER OF THE CHAIN CAN BE EXPLICITLY SET BY PROVIDING A LIST OF INTEGERS FOR EXAMPLE FOR A CHAIN OF LENGTH 5

ORDER 1 3 2 4 0

MEANS THAT THE FIRST MODEL IN THE CHAIN WILL MAKE PREDICTIONS FOR COLUMN 1 IN THE Y MATRIX THE SECOND MODEL WILL MAKE PREDICTIONS FOR COLUMN 3 ETC

IF ORDER IS ‘RANDOM’ A RANDOM ORDERING WILL BE USED

CVINT CROSSVALIDATION GENERATOR OR AN ITERABLE OPTIONAL DEFAULTNONE DETERMINES WHETHER TO USE CROSS VALIDATED PREDICTIONS OR TRUE LABELS FOR THE RESULTS OF PREVIOUS ESTIMATORS IN THE CHAIN IF CV IS NONE THE TRUE LABELS ARE USED WHEN FITTING OTHERWISE POSSIBLE INPUTS FOR CV ARE

- INTEGER TO SPECIFY THE NUMBER OF FOLDS IN A STRATIFIEDKFOLD
- CV SPLITTER
- AN ITERABLE YIELDING TRAIN TEST SPLITS AS ARRAYS OF INDICES

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RANDOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NPRANDOM

THE RANDOM NUMBER GENERATOR IS USED TO GENERATE RANDOM CHAIN ORDERS

ATTRIBUTES

ESTIMATORS LIST A LIST OF CLONES OF BASEESTIMATOR

ORDER LIST THE ORDER OF LABELS IN THE CLASSIFIER CHAIN

SEE ALSO

CLASSIFIERCHAIN EQUIVALENT FOR CLASSIFICATION

MULTIOUTPUTREGRESSOR LEARNS EACH OUTPUT INDEPENDENTLY RATHER THAN CHAINING

METHODS

FITSELF X Y FIT THE MODEL TO DATA MATRIX X AND TARGETS Y

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PREDICT SELF X PREDICT ON THE DATA MATRIX X USING THE CLASSIFIERCHAIN

MODEL

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELFBASEESTIMATOR ORDERNONE CVNONE RANDOMSTATENONE

FITSELFXY

FIT THE MODEL TO DATA MATRIX X AND TARGETS Y

2138 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT DATA

YARRAYLIKE SHAPE NSAMPLES NCLASSES THE TARGET VALUES

RETURNS

SELF OBJECT

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELF

PREDICT ON THE DATA MATRIX X USING THE CLASSIFIERCHAIN MODEL

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT DATA

RETURNS

YPRED ARRAYLIKE SHAPE NSAMPLES NCLASSES THE PREDICTED VALUES

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{U}{V}$  WHERE U IS THE RESIDUAL SUM OF SQUARES YTRUE YPRED

2SUM AND V IS THE TOTAL SUM OF SQUARES YTRUE YTRUEMEAN 2SUM THE BEST POSSIBLE SCORE

IS 10 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS

PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 00

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES FOR SOME ESTIMATORS THIS MAY

BE A PRECOMPUTED KERNEL MATRIX INSTEAD SHAPE NSAMPLES NSAMPLESFITTED WHERE

NSAMPLESFITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE VALUES FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT R2 OF SELF PREDICTX WRT Y

NOTES

THE R2 SCORE USED WHEN CALLING SCORE ON A REGRESSOR WILL USE MULTIOUTPUTUNIFORMAVERAGE

FROM VERSION 023 TO KEEP CONSISTENT WITH METRICSR2SCORE THIS WILL INFLUENCE THE SCORE

METHOD OF ALL THE MULTIOUTPUT REGRESSORS EXCEPT FOR MULTIOUTPUTMULTIOUTPUTREGRESSOR

TO SPECIFY THE DEFAULT VALUE MANUALLY AND AVOID THE WARNING PLEASE EITHER CALL METRICSR2SCORE

DIRECTLY OR MAKE A CUSTOM SCORER WITH METRICSMAKESCORER THE BUILTIN SCORER R2 USES

MULTIOUTPUTUNIFORMAVERAGE

628SKLEARNMULTIOUTPUT MULTIOUTPUT REGRESSION AND CLASSIFICATION 2139

SCIKITLEARN USER GUIDE RELEASE 0213

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

629SKLEARNNAIVEBAYES NAIVE BAYES

THESKLEARNNAIVEBAYES MODULE IMPLEMENTS NAIVE BAYES ALGORITHMS THESE ARE SUPERVISED LEARNING METHODS BASED ON APPLYING BAYES' THEOREM WITH STRONG NAIVE FEATURE INDEPENDENCE ASSUMPTIONS

USER GUIDE SEE THE NAIVE BAYES SECTION FOR FURTHER DETAILS

NAIVEBAYESBERNOULLINB ALPHA BINARIZE NAIVE BAYES CLASSIFIER FOR MULTIVARIATE BERNOULLI MODELS

NAIVEBAYESGAUSSIANNB PRIORS VARSMOOTHING GAUSSIAN NAIVE BAYES GAUSSIANNB

NAIVEBAYESMULTINOMIALNB ALPHA NAIVE BAYES CLASSIFIER FOR MULTINOMIAL MODELS

NAIVEBAYESCOMPLEMENTNB ALPHA FITPRIOR THE COMPLEMENT NAIVE BAYES CLASSIFIER DESCRIBED IN REN NIE ET AL

6291SKLEARNNAIVEBAYES BERNOULLINB

CLASSSSKLEARNNAIVEBAYES BERNOULLINB ALPHA10 BINARIZE00 FITPRIORTRUE

CLASSPRIORNONE

NAIVE BAYES CLASSIFIER FOR MULTIVARIATE BERNOULLI MODELS

LIKE MULTINOMIALNB THIS CLASSIFIER IS SUITABLE FOR DISCRETE DATA THE DIFFERENCE IS THAT WHILE MULTINOMIALNB WORKS WITH OCCURRENCE COUNTS BERNOULLINB IS DESIGNED FOR BINARYBOOLEAN FEATURES

READ MORE IN THE USER GUIDE

PARAMETERS

ALPHA FLOAT OPTIONAL DEFAULT10 ADDITIVE LAPLACELIDSTONE SMOOTHING PARAMETER 0 FOR NO SMOOTHING

BINARIZE FLOAT OR NONE OPTIONAL DEFAULT00 THRESHOLD FOR BINARIZING MAPPING TO BOOLEANS OF SAMPLE FEATURES IF NONE INPUT IS PRESUMED TO ALREADY CONSIST OF BINARY VECTORS

FITPRIOR BOOLEAN OPTIONAL DEFAULTTRUE WHETHER TO LEARN CLASS PRIOR PROBABILITIES OR NOT IF FALSE A UNIFORM PRIOR WILL BE USED

CLASSPRIOR ARRAYLIKE SIZENCLASSES OPTIONAL DEFAULTNONE PRIOR PROBABILITIES OF THE CLASSES IF SPECIFIED THE PRIORS ARE NOT ADJUSTED ACCORDING TO THE DATA

ATTRIBUTES

CLASSLOGPRIOR ARRAY SHAPE NCLASSES LOG PROBABILITY OF EACH CLASS SMOOTHED

FEATURELOGPROB ARRAY SHAPE NCLASSES NFEATURES EMPIRICAL LOG PROBABILITY OF FEATURES GIVEN A CLASS PXIY

CLASSCOUNT ARRAY SHAPE NCLASSES NUMBER OF SAMPLES ENCOUNTERED FOR EACH CLASS DURING FITTING THIS VALUE IS WEIGHTED BY THE SAMPLE WEIGHT WHEN PROVIDED

2140 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

FEATURECOUNT ARRAY SHAPE NCLASSES NFEATURES NUMBER OF SAMPLES ENCOUNTERED FOR EACH CLASS FEATURE DURING FITTING THIS VALUE IS WEIGHTED BY THE SAMPLE WEIGHT WHEN PROVIDED

REFERENCES

CD MANNING P RAGHAVAN AND H SCHUETZE 2008 INTRODUCTION TO INFORMATION RETRIEVAL CAMBRIDGE UNIVERSITY PRESS PP 234265 [HTTPS://NLPSTANFORDEDUIRBOOKHTMLHTMLEDITIONTHEBERNOULLIMODEL1HTML](https://nlp.stanford.edu/IR-book/html/htmledition/the-bernoulli-model-1.html)

A MCCALLUM AND K NIGAM 1998 A COMPARISON OF EVENT MODELS FOR NAIVE BAYES TEXT CLASSIFICATION PROC AAAIICML98 WORKSHOP ON LEARNING FOR TEXT CATEGORIZATION PP 4148

V METSIS I ANDROUTSOPOULOS AND G PALIOURAS 2006 SPAM FILTERING WITH NAIVE BAYES - WHICH NAIVE BAYES 3RD CONF ON EMAIL AND ANTISPAM CEAS

EXAMPLES

```
import numpy as np
X = np.random.randint(2, size=(6, 100))
Y = np.array([2, 3, 4, 4, 5])
from sklearn.naive_bayes import BernoulliNB
clf = BernoulliNB()
clf.fit(X, Y)
BernoulliNB(alpha=10, binarize=0, class_prior=None, fit_prior=True)
print(clf.predict(X[2:3]))
```

METHODS

FITSELF X Y SAMPLEWEIGHT FIT NAIVE BAYES CLASSIFIER ACCORDING TO X Y

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PARTIALFIT SELF X Y CLASSES SAMPLEWEIGHT INCREMENTAL FIT ON A BATCH OF SAMPLES

PREDICT SELF X PERFORM CLASSIFICATION ON AN ARRAY OF TEST VECTORS X

PREDICTLOGPROBA SELF X RETURN LOGPROBABILITY ESTIMATES FOR THE TEST VECTOR X

PREDICTPROBA SELF X RETURN PROBABILITY ESTIMATES FOR THE TEST VECTOR X

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELFALPHA10 BINARIZE00 FITPRIORTRUE CLASSPRIORNONE

FITSELFXYSAMPLEWEIGHTNONE

FIT NAIVE BAYES CLASSIFIER ACCORDING TO X Y

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES TRAINING VECTORS WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

YARRAYLIKE SHAPE NSAMPLES TARGET VALUES

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES DEFAULTNONE WEIGHTS APPLIED TO INDIVIDUAL SAMPLES 1 FOR UNWEIGHTED

629SKLEARNNAIVEBAYES NAIVE BAYES 2141

SCIKITLEARN USER GUIDE RELEASE 0213

RETURNS  
SELF OBJECT

GETPARAMS SELFDEEPTREE  
GET PARAMETERS FOR THIS ESTIMATOR  
PARAMETERS  
DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBOBJECTS THAT ARE ESTIMATORS  
RETURNS  
PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PARTIALFIT SELFXYCLASSES NONE SAMPLEWEIGHT NONE  
INCREMENTAL FIT ON A BATCH OF SAMPLES  
THIS METHOD IS EXPECTED TO BE CALLED SEVERAL TIMES CONSECUTIVELY ON DIFFERENT CHUNKS OF A DATASET SO AS TO  
IMPLEMENT OUTFCORE OR ONLINE LEARNING  
THIS IS ESPECIALLY USEFUL WHEN THE WHOLE DATASET IS TOO BIG TO FIT IN MEMORY AT ONCE  
THIS METHOD HAS SOME PERFORMANCE OVERHEAD HENCE IT IS BETTER TO CALL PARTIALFIT ON CHUNKS OF DATA THAT ARE  
AS LARGE AS POSSIBLE AS LONG AS FITTING IN THE MEMORY BUDGET TO HIDE THE OVERHEAD

PARAMETERS  
X ARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING VECTORS WHERE  
NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES  
Y ARRAYLIKE SHAPE NSAMPLES TARGET VALUES  
CLASSES ARRAYLIKE SHAPE NCLASSES DEFAULT NONE LIST OF ALL THE CLASSES THAT CAN POS  
SIBLY APPEAR IN THE Y VECTOR  
MUST BE PROVIDED AT THE FIRST CALL TO PARTIALFIT CAN BE OMITTED IN SUBSEQUENT CALLS  
SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES DEFAULT NONE WEIGHTS APPLIED TO INDIVIDUAL SAMPLES 1 FOR UNWEIGHTED

RETURNS  
SELF OBJECT

PREDICTSELF  
PERFORM CLASSIFICATION ON AN ARRAY OF TEST VECTORS X  
PARAMETERS  
X ARRAYLIKE SHAPE NSAMPLES NFEATURES

RETURNS  
CARRAY SHAPE NSAMPLES PREDICTED TARGET VALUES FOR X  
PREDICTLOGPROBA SELF  
RETURN LOGPROBABILITY ESTIMATES FOR THE TEST VECTOR X  
PARAMETERS  
X ARRAYLIKE SHAPE NSAMPLES NFEATURES

RETURNS

2142 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

CARRAYLIKE SHAPE NSAMPLES NCLASSES RETURNS THE LOGPROBABILITY OF THE SAMPLES  
FOR EACH CLASS IN THE MODEL THE COLUMNS CORRESPOND TO THE CLASSES IN SORTED ORDER AS THEY  
APPEAR IN THE ATTRIBUTE CLASSES

PREDICTPROBA SELF X  
RETURN PROBABILITY ESTIMATES FOR THE TEST VECTOR X

PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES

RETURNS  
CARRAYLIKE SHAPE NSAMPLES NCLASSES RETURNS THE PROBABILITY OF THE SAMPLES FOR  
EACH CLASS IN THE MODEL THE COLUMNS CORRESPOND TO THE CLASSES IN SORTED ORDER AS THEY  
APPEAR IN THE ATTRIBUTE CLASSES

SCORESELFXY SAMPLEWEIGHT NONE  
RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS  
IN MULTILABEL CLASSIFICATION THIS IS THE SUBSET ACCURACY WHICH IS A HARSH METRIC SINCE YOU REQUIRE FOR EACH  
SAMPLE THAT EACH LABEL SET BE CORRECTLY PREDICTED

PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES  
YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE LABELS FOR X  
SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS  
SCORE FLOAT MEAN ACCURACY OF SELF PREDICT X WRT Y

SETPARAMS SELF PARAMS  
SET THE PARAMETERS OF THIS ESTIMATOR  
THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE  
PARAMETERS OF THE FORM COMPONENT PARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT  
OF A NESTED OBJECT

RETURNS  
SELF

EXAMPLES USING SKLEARN NAIVE BAYES BERNOLLINB

- HASHING FEATURE TRANSFORMATION USING TOTALLY RANDOM TREES
- CLASSIFICATION OF TEXT DOCUMENTS USING SPARSE FEATURES

6292 SKLEARN NAIVE BAYES GAUSSIANNB  
CLASS SKLEARN NAIVE BAYES GAUSSIANNB PRIORS NONE VARSMOOTHING 1E09  
GAUSSIAN NAIVE BAYES GAUSSIANNB  
CAN PERFORM ONLINE UPDATES TO MODEL PARAMETERS VIA PARTIAL FIT METHOD FOR DETAILS ON ALGORITHM USED TO  
UPDATE FEATURE MEANS AND VARIANCE ONLINE SEE STANFORD CS TECH REPORT STANCS79773 BY CHAN GOLUB AND  
LEVEQUE  
HTTP://STANFORD.EDU/PUB/CS/TECH/REPORTS/CS-TR-79-773/CS-TR-79-773.PDF

6295 SKLEARN NAIVE BAYES NAIVE BAYES 2143

SCIKITLEARN USER GUIDE RELEASE 0213  
READ MORE IN THE USER GUIDE  
PARAMETERS  
PRIORS ARRAYLIKE SHAPE NCLASSES PRIOR PROBABILITIES OF THE CLASSES IF SPECIFIED THE PRIORS  
ARE NOT ADJUSTED ACCORDING TO THE DATA  
VARSMOOTHING FLOAT OPTIONAL DEFAULT1E9 PORTION OF THE LARGEST VARIANCE OF ALL FEATURES  
THAT IS ADDED TO VARIANCES FOR CALCULATION STABILITY  
ATTRIBUTES  
CLASSPRIOR ARRAY SHAPE NCLASSES PROBABILITY OF EACH CLASS  
CLASSCOUNT ARRAY SHAPE NCLASSES NUMBER OF TRAINING SAMPLES OBSERVED IN EACH CLASS  
THETA ARRAY SHAPE NCLASSES NFEATURES MEAN OF EACH FEATURE PER CLASS  
SIGMA ARRAY SHAPE NCLASSES NFEATURES VARIANCE OF EACH FEATURE PER CLASS  
EPSILON FLOAT ABSOLUTE ADDITIVE VALUE TO VARIANCES  
EXAMPLES  
IMPORT NUMPY AS NP  
X NPARRAY1 1 2 1 3 2 1 1 2 1 3 2  
Y NPARRAY1 1 1 2 2 2  
FROM SKLEARNNAIVEBAYES IMPORT GAUSSIANNB  
CLF GAUSSIANNB  
CLFFITX Y  
GAUSSIANNBPRIORSNONE VARSMOOTHING1E09  
PRINTCLFPREDICT08 1  
1  
CLFPF GAUSSIANNB  
CLFPFPARTIALFITX Y NPUNIQUEY  
GAUSSIANNBPRIORSNONE VARSMOOTHING1E09  
PRINTCLFPFPREDICT08 1  
1  
METHODS  
FITSELF X Y SAMPLEWEIGHT FIT GAUSSIAN NAIVE BAYES ACCORDING TO X Y  
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR  
PARTIALFIT SELF X Y CLASSES SAMPLEWEIGHT INCREMENTAL FIT ON A BATCH OF SAMPLES  
PREDICT SELF X PERFORM CLASSIFICATION ON AN ARRAY OF TEST VECTORS X  
PREDICTLOGPROBA SELF X RETURN LOGPROBABILITY ESTIMATES FOR THE TEST VECTOR X  
PREDICTPROBA SELF X RETURN PROBABILITY ESTIMATES FOR THE TEST VECTOR X  
SCORE SELF X Y SAMPLEWEIGHT RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND  
LABELS  
SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR  
INIT SELFPRIORSNONE VARSMOOTHING1E09  
FITSELFXYSAMPLEWEIGHTNONE  
FIT GAUSSIAN NAIVE BAYES ACCORDING TO X Y  
PARAMETERS  
2144 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING VECTORS WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

YARRAYLIKE SHAPE NSAMPLES TARGET VALUES

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL DEFAULTNONE WEIGHTS APPLIED TO INDIVIDUAL SAMPLES 1 FOR UNWEIGHTED

NEW IN VERSION 017 GAUSSIAN NAIVE BAYES SUPPORTS FITTING WITH SAMPLEWEIGHT

RETURNS

SELF OBJECT

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PARTIALFIT SELFXYCLASSESNONE SAMPLEWEIGHTNONE

INCREMENTAL FIT ON A BATCH OF SAMPLES

THIS METHOD IS EXPECTED TO BE CALLED SEVERAL TIMES CONSECUTIVELY ON DIFFERENT CHUNKS OF A DATASET SO AS TO IMPLEMENT OUTFCORE OR ONLINE LEARNING

THIS IS ESPECIALLY USEFUL WHEN THE WHOLE DATASET IS TOO BIG TO FIT IN MEMORY AT ONCE

THIS METHOD HAS SOME PERFORMANCE AND NUMERICAL STABILITY OVERHEAD HENCE IT IS BETTER TO CALL PARTIALFIT ON CHUNKS OF DATA THAT ARE AS LARGE AS POSSIBLE AS LONG AS FITTING IN THE MEMORY BUDGET TO HIDE THE OVERHEAD

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING VECTORS WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

YARRAYLIKE SHAPE NSAMPLES TARGET VALUES

CLASSES ARRAYLIKE SHAPE NCLASSES OPTIONAL DEFAULTNONE LIST OF ALL THE CLASSES THAT CAN POSSIBLY APPEAR IN THE Y VECTOR

MUST BE PROVIDED AT THE FIRST CALL TO PARTIALFIT CAN BE OMITTED IN SUBSEQUENT CALLS

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL DEFAULTNONE WEIGHTS APPLIED TO INDIVIDUAL SAMPLES 1 FOR UNWEIGHTED

NEW IN VERSION 017

RETURNS

SELF OBJECT

PREDICTSELF

PERFORM CLASSIFICATION ON AN ARRAY OF TEST VECTORS X

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES

RETURNS

CARRAY SHAPE NSAMPLES PREDICTED TARGET VALUES FOR X

629SKLEARNNAIVEBAYES NAIVE BAYES 2145

SCIKITLEARN USER GUIDE RELEASE 0213

PREDICTLOGPROBA SELF X

RETURN LOGPROBABILITY ESTIMATES FOR THE TEST VECTOR X

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES

RETURNS

CARRAYLIKE SHAPE NSAMPLES NCLASSES RETURNS THE LOGPROBABILITY OF THE SAMPLES FOR EACH CLASS IN THE MODEL THE COLUMNS CORRESPOND TO THE CLASSES IN SORTED ORDER AS THEY APPEAR IN THE ATTRIBUTE CLASSES

PREDICTPROBA SELF X

RETURN PROBABILITY ESTIMATES FOR THE TEST VECTOR X

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES

RETURNS

CARRAYLIKE SHAPE NSAMPLES NCLASSES RETURNS THE PROBABILITY OF THE SAMPLES FOR EACH CLASS IN THE MODEL THE COLUMNS CORRESPOND TO THE CLASSES IN SORTED ORDER AS THEY APPEAR IN THE ATTRIBUTE CLASSES

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS

IN MULTILABEL CLASSIFICATION THIS IS THE SUBSET ACCURACY WHICH IS A HARSH METRIC SINCE YOU REQUIRE FOR EACH SAMPLE THAT EACH LABEL SET BE CORRECTLY PREDICTED

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE LABELS FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT MEAN ACCURACY OF SELF PREDICT X WRT Y

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

EXAMPLES USING SKLEARNNAIVEBAYESGAUSSIANNB

- COMPARISON OF CALIBRATION OF CLASSIFIERS
- PROBABILITY CALIBRATION CURVES
- PROBABILITY CALIBRATION OF CLASSIFIERS
- CLASSIFIER COMPARISON
- PLOT CLASS PROBABILITIES CALCULATED BY THE VOTINGCLASSIFIER

2146 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

- PLOTting LEARNING CURVES
- IMPORTANCE OF FEATURE SCALING

6293SKLEARNNAIVEBAYES MULTINOMIALNB

CLASSSKLEARNNAIVEBAYES MULTINOMIALNB ALPHA10 FITPRIORTRUE CLASSPRIORNONE

NAIVE BAYES CLASSIFIER FOR MULTINOMIAL MODELS

THE MULTINOMIAL NAIVE BAYES CLASSIFIER IS SUITABLE FOR CLASSIFICATION WITH DISCRETE FEATURES EG WORD COUNTS FOR TEXT CLASSIFICATION THE MULTINOMIAL DISTRIBUTION NORMALLY REQUIRES INTEGER FEATURE COUNTS HOWEVER IN PRACTICE FRACTIONAL COUNTS SUCH AS TFIDF MAY ALSO WORK

READ MORE IN THE USER GUIDE

PARAMETERS

ALPHA FLOAT OPTIONAL DEFAULT10 ADDITIVE LAPLACELIDSTONE SMOOTHING PARAMETER 0 FOR NO SMOOTHING

FITPRIOR BOOLEAN OPTIONAL DEFAULTTRUE WHETHER TO LEARN CLASS PRIOR PROBABILITIES OR NOT IF FALSE A UNIFORM PRIOR WILL BE USED

CLASSPRIOR ARRAYLIKE SIZE NCLASSES OPTIONAL DEFAULTNONE PRIOR PROBABILITIES OF THE CLASSES IF SPECIFIED THE PRIORS ARE NOT ADJUSTED ACCORDING TO THE DATA

ATTRIBUTES

CLASSLOGPRIOR ARRAY SHAPE NCLASSES SMOOTHED EMPIRICAL LOG PROBABILITY FOR EACH CLASS

INTERCEPT ARRAY SHAPE NCLASSES MIRRORS CLASSLOGPRIOR FOR INTERPRETING MULTI

NOMIALNB AS A LINEAR MODEL

FEATURELOGPROB ARRAY SHAPE NCLASSES NFEATURES EMPIRICAL LOG PROBABILITY OF FEATURES GIVEN A CLASS PXIY

COEF ARRAY SHAPE NCLASSES NFEATURES MIRRORS FEATURELOGPROB FOR INTERPRETING MULTINOMIALNB AS A LINEAR MODEL

CLASSCOUNT ARRAY SHAPE NCLASSES NUMBER OF SAMPLES ENCOUNTERED FOR EACH CLASS DURING FITTING THIS VALUE IS WEIGHTED BY THE SAMPLE WEIGHT WHEN PROVIDED

FEATURECOUNT ARRAY SHAPE NCLASSES NFEATURES NUMBER OF SAMPLES ENCOUNTERED FOR EACH CLASS FEATURE DURING FITTING THIS VALUE IS WEIGHTED BY THE SAMPLE WEIGHT WHEN PROVIDED

NOTES

FOR THE RATIONALE BEHIND THE NAMES COEF ANDINTERCEPT IE NAIVE BAYES AS A LINEAR CLASSIFIER SEE J RENNIE ET AL 2003 TACKLING THE POOR ASSUMPTIONS OF NAIVE BAYES TEXT CLASSIFIERS ICML

REFERENCES

CD MANNING P RAGHAVAN AND H SCHUETZE 2008 INTRODUCTION TO INFORMATION RETRIEVAL CAMBRIDGE UNIVER SITY PRESS PP 234265 [HTTPS://NLPSTANFORD.EDUIRBOOKHTMLHTML/EDITIONNAIVEBAYESTEXTCLASSIFICATION1 HTML](https://nlp.stanford.edu/IR-book/html/htmledition/naivebayes-text-classification-1.html)

6293SKLEARNNAIVEBAYES NAIVE BAYES 2147

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES

```
import numpy as np
X = np.random.randint(5, size=(6, 100))
Y = np.array(1, 2, 3, 4, 5, 6)
from sklearn.naive_bayes import MultinomialNB
clf = MultinomialNB()
clf.fit(X, Y)
MultinomialNB(alpha=10, class_prior=None, fit_prior=True)
print(clf.predict(X[2:3]))
```

METHODS

`fit(X, Y)` SAMPLEWEIGHT FIT NAIVE BAYES CLASSIFIER ACCORDING TO X Y

`get_params()` SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

`partial_fit(X, Y, classes)` SAMPLEWEIGHT INCREMENTAL FIT ON A BATCH OF SAMPLES

`predict(X)` PERFORM CLASSIFICATION ON AN ARRAY OF TEST VECTORS X

`predict_log_proba(X)` SELF X RETURN LOGPROBABILITY ESTIMATES FOR THE TEST VECTOR X

`predict_proba(X)` SELF X RETURN PROBABILITY ESTIMATES FOR THE TEST VECTOR X

`score(X, Y)` SAMPLEWEIGHT RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS

`set_params()` SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

`init(self, alpha=10, fit_prior=True, class_prior=None)`

`fit(X, Y, sample_weight=None)`

`fit_naive_bayes()` CLASSIFIER ACCORDING TO X Y

PARAMETERS

`X` ARRAYLIKE SPARSE MATRIX SHAPE (NSAMPLES, NFEATURES) TRAINING VECTORS WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

`Y` ARRAYLIKE SHAPE (NSAMPLES,) TARGET VALUES

`SAMPLEWEIGHT` ARRAYLIKE SHAPE (NSAMPLES,) DEFAULT NONE WEIGHTS APPLIED TO INDIVIDUAL SAMPLES 1 FOR UNWEIGHTED

RETURNS

SELF OBJECT

`get_params()` SELF DEEP TRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

`partial_fit(X, Y, classes=None, sample_weight=None)`

INCREMENTAL FIT ON A BATCH OF SAMPLES

2148 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

THIS METHOD IS EXPECTED TO BE CALLED SEVERAL TIMES CONSECUTIVELY ON DIFFERENT CHUNKS OF A DATASET SO AS TO IMPLEMENT OUTFCORE OR ONLINE LEARNING

THIS IS ESPECIALLY USEFUL WHEN THE WHOLE DATASET IS TOO BIG TO FIT IN MEMORY AT ONCE

THIS METHOD HAS SOME PERFORMANCE OVERHEAD HENCE IT IS BETTER TO CALL PARTIALFIT ON CHUNKS OF DATA THAT ARE AS LARGE AS POSSIBLE AS LONG AS FITTING IN THE MEMORY BUDGET TO HIDE THE OVERHEAD

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES TRAINING VECTORS WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

YARRAYLIKE SHAPE NSAMPLES TARGET VALUES

CLASSES ARRAYLIKE SHAPE NCLASSES DEFAULTNONE LIST OF ALL THE CLASSES THAT CAN POSSIBLY APPEAR IN THE Y VECTOR

MUST BE PROVIDED AT THE FIRST CALL TO PARTIALFIT CAN BE OMITTED IN SUBSEQUENT CALLS

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES DEFAULTNONE WEIGHTS APPLIED TO INDIVIDUAL SAMPLES 1 FOR UNWEIGHTED

RETURNS

SELF OBJECT

PREDICTSELF

PERFORM CLASSIFICATION ON AN ARRAY OF TEST VECTORS X

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES

RETURNS

CARRAY SHAPE NSAMPLES PREDICTED TARGET VALUES FOR X

PREDICTLOGPROBA SELF

RETURN LOGPROBABILITY ESTIMATES FOR THE TEST VECTOR X

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES

RETURNS

CARRAYLIKE SHAPE NSAMPLES NCLASSES RETURNS THE LOGPROBABILITY OF THE SAMPLES FOR EACH CLASS IN THE MODEL THE COLUMNS CORRESPOND TO THE CLASSES IN SORTED ORDER AS THEY APPEAR IN THE ATTRIBUTE CLASSES

PREDICTPROBA SELF

RETURN PROBABILITY ESTIMATES FOR THE TEST VECTOR X

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES

RETURNS

CARRAYLIKE SHAPE NSAMPLES NCLASSES RETURNS THE PROBABILITY OF THE SAMPLES FOR EACH CLASS IN THE MODEL THE COLUMNS CORRESPOND TO THE CLASSES IN SORTED ORDER AS THEY APPEAR IN THE ATTRIBUTE CLASSES

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS

629SKLEARNNAIVEBAYES NAIVE BAYES 2149

SCIKITLEARN USER GUIDE RELEASE 0213

IN MULTILABEL CLASSIFICATION THIS IS THE SUBSET ACCURACY WHICH IS A HARSH METRIC SINCE YOU REQUIRE FOR EACH SAMPLE THAT EACH LABEL SET BE CORRECTLY PREDICTED

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE LABELS FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT MEAN ACCURACY OF SELF PREDICTX WRT Y

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

EXAMPLES USING SKLEARNNAIVEBAYESMULTINOMIALNB

- OUTOFCORE CLASSIFICATION OF TEXT DOCUMENTS
- CLASSIFICATION OF TEXT DOCUMENTS USING SPARSE FEATURES

6294SKLEARNNAIVEBAYES COMPLEMENTNB

CLASSSSKLEARNNAIVEBAYES COMPLEMENTNB ALPHA10 FITPRIORTRUE CLASSPRIORNONE

NORMFALSE

THE COMPLEMENT NAIVE BAYES CLASSIFIER DESCRIBED IN RENNIE ET AL 2003

THE COMPLEMENT NAIVE BAYES CLASSIFIER WAS DESIGNED TO CORRECT THE "SEVERE ASSUMPTIONS" MADE BY THE STANDARD MULTINOMIAL NAIVE BAYES CLASSIFIER IT IS PARTICULARLY SUITED FOR IMBALANCED DATA SETS

READ MORE IN THE USER GUIDE

PARAMETERS

ALPHA FLOAT OPTIONAL DEFAULT10 ADDITIVE LAPLACELIDSTONE SMOOTHING PARAMETER 0 FOR NO SMOOTHING

FITPRIOR BOOLEAN OPTIONAL DEFAULTTRUE ONLY USED IN EDGE CASE WITH A SINGLE CLASS IN THE TRAINING SET

CLASSPRIOR ARRAYLIKE SIZE NCLASSES OPTIONAL DEFAULTNONE PRIOR PROBABILITIES OF THE CLASSES NOT USED

NORM BOOLEAN OPTIONAL DEFAULTFALSE WHETHER OR NOT A SECOND NORMALIZATION OF THE WEIGHTS IS PERFORMED THE DEFAULT BEHAVIOR MIRRORS THE IMPLEMENTATIONS FOUND IN MAHOUT AND WEKA WHICH DO NOT FOLLOW THE FULL ALGORITHM DESCRIBED IN TABLE 9 OF THE PAPER

ATTRIBUTES

CLASSLOGPRIOR ARRAY SHAPE NCLASSES SMOOTHED EMPIRICAL LOG PROBABILITY FOR EACH CLASS ONLY USED IN EDGE CASE WITH A SINGLE CLASS IN THE TRAINING SET

2150 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

FEATURELOGPROB ARRAY SHAPE NCLASSES NFEATURES EMPIRICAL WEIGHTS FOR CLASS COMPLEMENTS

CLASSCOUNT ARRAY SHAPE NCLASSES NUMBER OF SAMPLES ENCOUNTERED FOR EACH CLASS DURING FITTING THIS VALUE IS WEIGHTED BY THE SAMPLE WEIGHT WHEN PROVIDED

FEATURECOUNT ARRAY SHAPE NCLASSES NFEATURES NUMBER OF SAMPLES ENCOUNTERED FOR EACH CLASS FEATURE DURING FITTING THIS VALUE IS WEIGHTED BY THE SAMPLE WEIGHT WHEN PROVIDED

FEATUREALL ARRAY SHAPE NFEATURES NUMBER OF SAMPLES ENCOUNTERED FOR EACH FEATURE DURING FITTING THIS VALUE IS WEIGHTED BY THE SAMPLE WEIGHT WHEN PROVIDED

REFERENCES

RENNIE J D SHIH L TEEVAN J KARGER D R 2003 TACKLING THE POOR ASSUMPTIONS OF NAIVE BAYES TEXT CLASSIFIERS IN ICML V OL 3 PP 616623 [HTTPS://PEOPLE.CSAIL.MIT.EDU/JRENNIEPAPERS/ICML03NB.PDF](https://people.csail.mit.edu/jrennie/papers/icml03nb.pdf)

EXAMPLES

```
import numpy as np
X = np.random.randint(5, size=(6, 100))
Y = np.array(1, 2, 3, 4, 5, 6)
from sklearn.naive_bayes import ComplementNB
clf = ComplementNB()
clf.fit(X, Y)
complementnb.alpha = 10
class_prior = None
fit_prior = True
norm = False
print(clf.predict(X[2:3]))
```

METHODS

fit(X, Y, sample\_weight) FIT NAIVE BAYES CLASSIFIER ACCORDING TO X, Y

get\_params() SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

partial\_fit(X, Y, classes, sample\_weight) INCREMENTAL FIT ON A BATCH OF SAMPLES

predict(X) SELF X PERFORM CLASSIFICATION ON AN ARRAY OF TEST VECTORS X

predict\_log\_proba(X) SELF X RETURN LOGPROBABILITY ESTIMATES FOR THE TEST VECTOR X

predict\_proba(X) SELF X RETURN PROBABILITY ESTIMATES FOR THE TEST VECTOR X

score(X, Y, sample\_weight) RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS

set\_params(\*\*kwargs) SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

init(alpha=10, fit\_prior=True, class\_prior=None, norm=False)

fit(X, Y, sample\_weight=None)

FIT NAIVE BAYES CLASSIFIER ACCORDING TO X, Y

PARAMETERS

X: ARRAY-LIKE SPARSE MATRIX SHAPE (NSAMPLES, NFEATURES) TRAINING VECTORS WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

Y: ARRAY-LIKE SHAPE (NSAMPLES,) TARGET VALUES

629SKLEARNNAIVEBAYES NAIVE BAYES 2151

SCIKITLEARN USER GUIDE RELEASE 0213

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES DEFAULTNONE WEIGHTS APPLIED TO INDIVIDUAL SAMPLES 1 FOR UNWEIGHTED

RETURNS

SELF OBJECT

GETPARAMS SELFDEEPTREE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PARTIALFIT SELFXYCLASSESNONE SAMPLEWEIGHTNONE

INCREMENTAL FIT ON A BATCH OF SAMPLES

THIS METHOD IS EXPECTED TO BE CALLED SEVERAL TIMES CONSECUTIVELY ON DIFFERENT CHUNKS OF A DATASET SO AS TO IMPLEMENT OUTOF CORE OR ONLINE LEARNING

THIS IS ESPECIALLY USEFUL WHEN THE WHOLE DATASET IS TOO BIG TO FIT IN MEMORY AT ONCE

THIS METHOD HAS SOME PERFORMANCE OVERHEAD HENCE IT IS BETTER TO CALL PARTIALFIT ON CHUNKS OF DATA THAT ARE AS LARGE AS POSSIBLE AS LONG AS FITTING IN THE MEMORY BUDGET TO HIDE THE OVERHEAD

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES TRAINING VECTORS WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

YARRAYLIKE SHAPE NSAMPLES TARGET VALUES

CLASSES ARRAYLIKE SHAPE NCLASSES DEFAULTNONE LIST OF ALL THE CLASSES THAT CAN POSSIBLY APPEAR IN THE Y VECTOR

MUST BE PROVIDED AT THE FIRST CALL TO PARTIALFIT CAN BE OMITTED IN SUBSEQUENT CALLS

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES DEFAULTNONE WEIGHTS APPLIED TO INDIVIDUAL SAMPLES 1 FOR UNWEIGHTED

RETURNS

SELF OBJECT

PREDICTSELF

PERFORM CLASSIFICATION ON AN ARRAY OF TEST VECTORS X

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES

RETURNS

CARRAY SHAPE NSAMPLES PREDICTED TARGET VALUES FOR X

PREDICTLOGPROB SELF

RETURN LOGPROBABILITY ESTIMATES FOR THE TEST VECTOR X

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES

RETURNS

2152 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

CARRAYLIKE SHAPE NSAMPLES NCLASSES RETURNS THE LOGPROBABILITY OF THE SAMPLES FOR EACH CLASS IN THE MODEL THE COLUMNS CORRESPOND TO THE CLASSES IN SORTED ORDER AS THEY APPEAR IN THE ATTRIBUTE CLASSES

PREDICTPROBA SELF X

RETURN PROBABILITY ESTIMATES FOR THE TEST VECTOR X

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES

RETURNS

CARRAYLIKE SHAPE NSAMPLES NCLASSES RETURNS THE PROBABILITY OF THE SAMPLES FOR EACH CLASS IN THE MODEL THE COLUMNS CORRESPOND TO THE CLASSES IN SORTED ORDER AS THEY APPEAR IN THE ATTRIBUTE CLASSES

SCORESELFXY SAMPLEWEIGHT NONE

RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS

IN MULTILABEL CLASSIFICATION THIS IS THE SUBSET ACCURACY WHICH IS A HARSH METRIC SINCE YOU REQUIRE FOR EACH SAMPLE THAT EACH LABEL SET BE CORRECTLY PREDICTED

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE LABELS FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT MEAN ACCURACY OF SELF PREDICT X WRT Y

SETPARAMS SELF PARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENT PARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

EXAMPLES USING SKLEARN NNAIVE BAYES COMPLEMENT NB

- CLASSIFICATION OF TEXT DOCUMENTS USING SPARSE FEATURES

630 SKLEARN NEIGHBORS NEAREST NEIGHBORS

THE SKLEARN NEIGHBORS MODULE IMPLEMENTS THE K NEAREST NEIGHBORS ALGORITHM

USER GUIDE SEE THE NEAREST NEIGHBORS SECTION FOR FURTHER DETAILS

NEIGHBORS BALL TREE BALL TREE FOR FAST GENERALIZED N POINT PROBLEMS

NEIGHBORS DISTANCE METRIC DISTANCE METRIC CLASS

NEIGHBORS KD TREE KD TREE FOR FAST GENERALIZED N POINT PROBLEMS

CONTINUED ON NEXT PAGE

630 SKLEARN NEIGHBORS NEAREST NEIGHBORS 2153

SCIKITLEARN USER GUIDE RELEASE 0213

TABLE 6223 – CONTINUED FROM PREVIOUS PAGE

NEIGHBORSKERNELDENSITY BANDWIDTH    KERNEL DENSITY ESTIMATION

NEIGHBORSKNEIGHBORSCLASSIFIER    CLASSIFIER IMPLEMENTING THE KNEAREST NEIGHBORS VOTE

NEIGHBORSKNEIGHBORSREGRESSOR NNEIGHBORS

    REGRESSION BASED ON KNEAREST NEIGHBORS

NEIGHBORSLOCALOUTLIERFACTOR NNEIGHBORS

    UNSUPERVISED OUTLIER DETECTION USING LOCAL OUTLIER FACTOR

LOF

NEIGHBORSRADIUSNEIGHBORSCLASSIFIER    CLASSIFIER IMPLEMENTING A VOTE AMONG NEIGHBORS WITHIN A GIVEN RADIUS

NEIGHBORSRADIUSNEIGHBORSREGRESSOR RADIUS

    REGRESSION BASED ON NEIGHBORS WITHIN A FIXED RADIUS

NEIGHBORSNEARESTCENTROID METRIC    NEAREST CENTROID CLASSIFIER

NEIGHBORSNEARESTNEIGHBORS NNEIGHBORS

    UNSUPERVISED LEARNER FOR IMPLEMENTING NEIGHBOR SEARCHES

NEIGHBORSNEIGHBORHOODCOMPONENTSANALYSIS    NEIGHBORHOOD COMPONENTS ANALYSIS

6301SKLEARNNEIGHBORS BALLTREE

CLASSSSKLEARNNEIGHBORS BALLTREE

BALLTREE FOR FAST GENERALIZED NPOINT PROBLEMS

BALLTREEX LEAFSIZE40 METRIC'MINKOWSKI' KWARGS

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES NSAMPLES IS THE NUMBER OF POINTS IN THE DATA SET AND NFEATURES IS THE DIMENSION OF THE PARAMETER SPACE NOTE IF X IS A CCONTIGUOUS ARRAY OF DOUBLES THEN DATA WILL NOT BE COPIED OTHERWISE AN INTERNAL COPY WILL BE MADE

LEAFSIZE POSITIVE INTEGER DEFAULT 40 NUMBER OF POINTS AT WHICH TO SWITCH TO BRUTEFORCE

CHANGING LEAFSIZE WILL NOT AFFECT THE RESULTS OF A QUERY BUT CAN SIGNIFICANTLY IMPACT THE SPEED OF A QUERY AND THE MEMORY REQUIRED TO STORE THE CONSTRUCTED TREE THE AMOUNT OF MEMORY NEEDED TO STORE THE TREE SCALES AS APPROXIMATELY NSAMPLES LEAFSIZE FOR A SPEC IFIEDLEAFSIZE A LEAF NODE IS GUARANTEED TO SATISFY LEAFSIZE NPOINTS

2LEAFSIZE EXCEPT IN THE CASE THAT NSAMPLES LEAFSIZE

METRIC STRING OR DISTANCEMETRIC OBJECT THE DISTANCE METRIC TO USE FOR THE TREE DE FAULT'MINKOWSKI' WITH P2 THAT IS A EUCLIDEAN METRIC SEE THE DOCUMENTATION OF THE DISTANCEMETRIC CLASS FOR A LIST OF AVAILABLE METRICS BALLTREEVALIDMETRICS GIVES A LIST OF THE METRICS WHICH ARE VALID FOR BALLTREE

ADDITIONAL KEYWORDS ARE PASSED TO THE DISTANCE METRIC CLASS

ATTRIBUTES

DATA MEMORY VIEW THE TRAINING DATA

EXAMPLES

QUERY FOR KNEAREST NEIGHBORS

    IMPORT NUMPY AS NP

RNG NPRANDOMRANDOMSTATE0

X RNGRANDOMSAMPLE10 3 10 POINTS IN 3 DIMENSIONS

TREE BALLTREEX LEAFSIZE2

DIST IND TREEQUERYX1 K3

2154 CHAPTER 6 API REFERENCE



```
SCIKITLEARN USER GUIDE RELEASE 0213
PRINTIND INDICES OF 3 CLOSEST NEIGHBORS
0 3 1
PRINTDIST DISTANCES TO 3 CLOSEST NEIGHBORS
0 019662693 029473397
PICKLE AND UNPICKLE A TREE NOTE THAT THE STATE OF THE TREE IS SAVED IN THE PICKLE OPERATION THE TREE NEEDS NOT BE
REBUILT UPON UNPICKLING
IMPORT NUMPY AS NP
IMPORT PICKLE
RNG NPRANDOMRANDOMSTATE0
X RNGRANDOMSAMPLE10 3 10 POINTS IN 3 DIMENSIONS
TREE BALLTREEX LEAFSIZE2
S PICKLEDUMPSTREE
TRECOPY PICKLELOADSS
DIST IND TRECOPYQUERYX1 K3
PRINTIND INDICES OF 3 CLOSEST NEIGHBORS
0 3 1
PRINTDIST DISTANCES TO 3 CLOSEST NEIGHBORS
0 019662693 029473397
QUERY FOR NEIGHBORS WITHIN A GIVEN RADIUS
IMPORT NUMPY AS NP
RNG NPRANDOMRANDOMSTATE0
X RNGRANDOMSAMPLE10 3 10 POINTS IN 3 DIMENSIONS
TREE BALLTREEX LEAFSIZE2
PRINTTREEQUERYRADIUSX1 R03 COUNTONLY TRUE
3
IND TREEQUERYRADIUSX1 R03
PRINTIND INDICES OF NEIGHBORS WITHIN DISTANCE 03
3 0 1
COMPUTE A GAUSSIAN KERNEL DENSITY ESTIMATE
IMPORT NUMPY AS NP
RNG NPRANDOMRANDOMSTATE42
X RNGRANDOMSAMPLE100 3
TREE BALLTREEX
TREEKERNELDENSITYX3 H01 KERNELGAUSSIAN
ARRAY 694114649 783281226 72071716
COMPUTE A TWOPOINT AUTOCORRELATION FUNCTION
IMPORT NUMPY AS NP
RNG NPRANDOMRANDOMSTATE0
X RNGRANDOMSAMPLE30 3
R NPLinspace0 1 5
TREE BALLTREEX
TREETWOPOINTCORRELATIONX R
ARRAY 30 62 278 580 820
METHODS
630SKLEARNNEIGHBORS NEAREST NEIGHBORS 2155
```

SCIKITLEARN USER GUIDE RELEASE 0213

KERNELDENSITY SELF X H KERNEL ATOL COMPUTE THE KERNEL DENSITY ESTIMATE AT POINTS X WITH THE GIVEN KERNEL USING THE DISTANCE METRIC SPECIFIED AT TREE CREATION

QUERY X K RETURNDISTANCE DUALTREE QUERY THE TREE FOR THE K NEAREST NEIGHBORS

QUERYRADIUS QUERYRADIUSSELF X R COUNTONLY FALSE

TWOPOINTCORRELATION COMPUTE THE TWOPOINT CORRELATION FUNCTION

GETARRAYS

GETNCALLS

GETTREESTATS

RESETNCALLS

INIT SELFARGS KWARGS

INITIALIZE SELF SEE HELPTYPESELF FOR ACCURATE SIGNATURE

KERNELDENSITY SELFHXHKERNEL'GAUSSIAN' ATOL0 RTOL1E8 BREADTHFIRSTTRUE RETURNLOGFALSE

COMPUTE THE KERNEL DENSITY ESTIMATE AT POINTS X WITH THE GIVEN KERNEL USING THE DISTANCE METRIC SPECIFIED AT TREE CREATION

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES AN ARRAY OF POINTS TO QUERY LAST DIMENSION SHOULD MATCH DIMENSION OF TRAINING DATA

HFLOAT THE BANDWIDTH OF THE KERNEL

KERNEL STRING SPECIFY THE KERNEL TO USE OPTIONS ARE 'GAUSSIAN' 'TOPHAT' 'EPANECHNIKOV' 'EXPONENTIAL' 'LINEAR' 'COSINE' DEFAULT IS KERNEL 'GAUSSIAN'

ATOL RTOL FLOAT DEFAULT 0 SPECIFY THE DESIRED RELATIVE AND ABSOLUTE TOLERANCE OF THE RESULT IF THE TRUE RESULT IS KTRUE THEN THE RETURNED RESULT KRET SATISFIES ABSKTRUE

KRET ATOL RTOL KRET THE DEFAULT IS ZERO IE MACHINE PRECISION FOR BOTH

BREADTHFIRST BOOLEAN DEFAULT FALSE IF TRUE USE A BREADTHFIRST SEARCH IF FALSE DEFAULT USE A DEPTHFIRST SEARCH BREADTHFIRST IS GENERALLY FASTER FOR COMPACT KERNELS ANDOR HIGH TOLERANCES

RETURNLOG BOOLEAN DEFAULT FALSE RETURN THE LOGARITHM OF THE RESULT THIS CAN BE MORE ACCURATE THAN RETURNING THE RESULT ITSELF FOR NARROW KERNELS

RETURNS

DENSITY NDARRAY THE ARRAY OF LOGDENSITY EVALUATIONS SHAPE XSHAPE1

QUERYXK1RETURNDISTANCETRUE DUALTREEFALSE BREADTHFIRSTFALSE

QUERY THE TREE FOR THE K NEAREST NEIGHBORS

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES AN ARRAY OF POINTS TO QUERY

KINTEGER DEFAULT 1 THE NUMBER OF NEAREST NEIGHBORS TO RETURN

RETURNDISTANCE BOOLEAN DEFAULT TRUE IF TRUE RETURN A TUPLE D I OF DISTANCES AND INDICES IF FALSE RETURN ARRAY I

2156 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

DUALTREE BOOLEAN DEFAULT FALSE IF TRUE USE THE DUAL TREE FORMALISM FOR THE QUERY A TREE IS BUILT FOR THE QUERY POINTS AND THE PAIR OF TREES IS USED TO EFFICIENTLY SEARCH THIS SPACE THIS CAN LEAD TO BETTER PERFORMANCE AS THE NUMBER OF POINTS GROWS LARGE

BREADTHFIRST BOOLEAN DEFAULT FALSE IF TRUE THEN QUERY THE NODES IN A BREADTHFIRST MANNER OTHERWISE QUERY THE NODES IN A DEPTHFIRST MANNER

SORTRESULTS BOOLEAN DEFAULT TRUE IF TRUE THEN DISTANCES AND INDICES OF EACH POINT ARE SORTED ON RETURN SO THAT THE FIRST COLUMN CONTAINS THE CLOSEST POINTS OTHERWISE NEIGHBORS ARE RETURNED IN AN ARBITRARY ORDER

RETURNS

IIF RETURNDISTANCE FALSE

DI IF RETURNDISTANCE TRUE

DARRAY OF DOUBLES SHAPE XSHAPE1 K EACH ENTRY GIVES THE LIST OF DISTANCES TO THE NEIGHBORS OF THE CORRESPONDING POINT

IARRAY OF INTEGERS SHAPE XSHAPE1 K EACH ENTRY GIVES THE LIST OF INDICES OF NEIGHBORS OF THE CORRESPONDING POINT

QUERYRADIUS

QUERYRADIUSSELF X R COUNTONLY FALSE

QUERY THE TREE FOR NEIGHBORS WITHIN A RADIUS R

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES AN ARRAY OF POINTS TO QUERY

RDISTANCE WITHIN WHICH NEIGHBORS ARE RETURNED R CAN BE A SINGLE VALUE OR AN ARRAY OF VALUES OF SHAPE XSHAPE1 IF DIFFERENT RADII ARE DESIRED FOR EACH POINT

RETURNDISTANCE BOOLEAN DEFAULT FALSE IF TRUE RETURN DISTANCES TO NEIGHBORS OF EACH POINT IF FALSE RETURN ONLY NEIGHBORS NOTE THAT UNLIKE THE QUERY METHOD SETTING RETURNDISTANCETRUE HERE ADDS TO THE COMPUTATION TIME NOT ALL DISTANCES NEED TO BE CALCULATED EXPLICITLY FOR RETURNDISTANCEFALSE RESULTS ARE NOT SORTED BY DEFAULT SEE SORTRESULTS KEYWORD

COUNTONLY BOOLEAN DEFAULT FALSE IF TRUE RETURN ONLY THE COUNT OF POINTS WITHIN DISTANCE R IF FALSE RETURN THE INDICES OF ALL POINTS WITHIN DISTANCE R IF RETURNDISTANCETRUE SETTING COUNTONLYTRUE WILL RESULT IN AN ERROR

SORTRESULTS BOOLEAN DEFAULT FALSE IF TRUE THE DISTANCES AND INDICES WILL BE SORTED BEFORE BEING RETURNED IF FALSE THE RESULTS WILL NOT BE SORTED IF RETURNDISTANCE FALSE SETTING SORTRESULTS TRUE WILL RESULT IN AN ERROR

RETURNS

COUNT IF COUNTONLY TRUE

IND IF COUNTONLY FALSE AND RETURNDISTANCE FALSE

IND DIST IF COUNTONLY FALSE AND RETURNDISTANCE TRUE

COUNT ARRAY OF INTEGERS SHAPE XSHAPE1 EACH ENTRY GIVES THE NUMBER OF NEIGHBORS WITHIN A DISTANCE R OF THE CORRESPONDING POINT

IND ARRAY OF OBJECTS SHAPE XSHAPE1 EACH ELEMENT IS A NUMPY INTEGER ARRAY LISTING THE INDICES OF NEIGHBORS OF THE CORRESPONDING POINT NOTE THAT UNLIKE THE RESULTS OF A K NEIGHBORS QUERY THE RETURNED NEIGHBORS ARE NOT SORTED BY DISTANCE BY DEFAULT

630SKLEARNNEIGHBORS NEAREST NEIGHBORS 2157

SCIKITLEARN USER GUIDE RELEASE 0213

DIST ARRAY OF OBJECTS SHAPE XSHAPE1 EACH ELEMENT IS A NUMPY DOUBLE ARRAY LISTING THE DISTANCES CORRESPONDING TO INDICES IN I

TWOPOINTCORRELATION

COMPUTE THE TWOPOINT CORRELATION FUNCTION

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES AN ARRAY OF POINTS TO QUERY LAST DIMENSION SHOULD MATCH DIMENSION OF TRAINING DATA

RARRAYLIKE A ONEDIMENSIONAL ARRAY OF DISTANCES

DUALTREE BOOLEAN DEFAULT FALSE IF TRUE USE A DUALTREE ALGORITHM OTHERWISE USE A SINGLETREE ALGORITHM DUAL TREE ALGORITHMS CAN HAVE BETTER SCALING FOR LARGE N

RETURNS

COUNTS NDARRAY COUNTSI CONTAINS THE NUMBER OF PAIRS OF POINTS WITH DISTANCE LESS THAN OR EQUAL TO RI

6302SKLEARNNEIGHBORS DISTANCEMETRIC

CLASSSSKLEARNNEIGHBORS DISTANCEMETRIC

DISTANCEMETRIC CLASS

THIS CLASS PROVIDES A UNIFORM INTERFACE TO FAST DISTANCE METRIC FUNCTIONS THE VARIOUS METRICS CAN BE ACCESSED VIA THEGETMETRIC CLASS METHOD AND THE METRIC STRING IDENTIFIER SEE BELOW FOR EXAMPLE TO USE THE EUCLIDEAN DISTANCE

DIST DISTANCEMETRICGETMETRICEUCLIDEAN

X 0 1 2

3 4 5

DISTPAIRWISEX

ARRAY 0 519615242

519615242 0

AVAILABLE METRICS

THE FOLLOWING LISTS THE STRING METRIC IDENTIFIERS AND THE ASSOCIATED DISTANCE METRIC CLASSES

METRICS INTENDED FOR REALVALUED VECTOR SPACES

2158 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

IDENTIFIER CLASS NAME ARGS DISTANCE FUNCTION

“EUCLIDEAN” EUCLIDEANDISTANCE

- SQRTSUMX

Y2

“MANHATTAN” MANHATTANDISTANCE

- SUMX Y

“CHEBYSHEV” CHEBYSHEVDISTANCE

- MAXX Y

“MINKOWSKI” MINKOWSKIDISTANCE P SUMX

YP1P

“WMINKOWSKI” WMINKOWSKIDISTANCE P W SUMWX

YP1P

“SEUCLIDEAN” SEUCLIDEANDISTANCE V SQRTSUMX

Y2 V

“MAHALANOBIS” MAHALANOBISDISTANCE V OR VI SQRTX Y

V1 X Y

METRICS INTENDED FOR TWODIMENSIONAL VECTOR SPACES NOTE THAT THE HAVERSINE DISTANCE METRIC REQUIRES DATA IN THE FORM OF LATITUDE LONGITUDE AND BOTH INPUTS AND OUTPUTS ARE IN UNITS OF RADIAN

IDENTIFIER CLASS NAME DISTANCE FUNCTION

“HAVER

SINE”HAVERSINEDIS

TANCE2 ARCSINSQRTSIN205 DX

COSX1COSX2SIN205 DY

METRICS INTENDED FOR INTEGERVERVALUED VECTOR SPACES THOUGH INTENDED FOR INTEGERVERVALUED VECTORS THESE ARE ALSO VALID METRICS IN THE CASE OF REALVALUED VECTORS

IDENTIFIER CLASS NAME DISTANCE FUNCTION

“HAMMING” HAMMINGDISTANCE NUNEQUALX Y NTOT

“CANBERRA” CANBERRADISTANCE SUMX Y X Y

“BRAYCURTIS” BRAYCURTISDISTANCE SUMX Y SUMX SUMY

METRICS INTENDED FOR BOOLEANVALUED VECTOR SPACES ANY NONZERO ENTRY IS EVALUATED TO “TRUE” IN THE LISTINGS

BELOW THE FOLLOWING ABBREVIATIONS ARE USED

- N NUMBER OF DIMENSIONS
- NTT NUMBER OF DIMS IN WHICH BOTH VALUES ARE TRUE
- NTF NUMBER OF DIMS IN WHICH THE FIRST VALUE IS TRUE SECOND IS FALSE
- NFT NUMBER OF DIMS IN WHICH THE FIRST VALUE IS FALSE SECOND IS TRUE
- NFF NUMBER OF DIMS IN WHICH BOTH VALUES ARE FALSE
- NNEQ NUMBER OF NONEQUAL DIMENSIONS NNEQ NTF NFT
- NNZ NUMBER OF NONZERO DIMENSIONS NNZ NTF NFT NTT

630SKLEARNNEIGHBORS NEAREST NEIGHBORS 2159

SCIKITLEARN USER GUIDE RELEASE 0213

IDENTIFIER CLASS NAME DISTANCE FUNCTION

“JACCARD” JACCARDDISTANCE NNEQ NNZ

“MATCHING” MATCHINGDISTANCE NNEQ N

“DICE” DICEDISTANCE NNEQ NTT NNZ

“KULSINSKI” KULSINSKIDISTANCE NNEQ N NTT NNEQ N

“ROGERSTANIMOTO” ROGERSTANIMOTODISTANCE 2 NNEQ N NNEQ

“RUSSELLRAO” RUSSELLRAODISTANCE NNZ N

“SOKALMICHENER” SOKALMICHENERDISTANCE 2 NNEQ N NNEQ

“SOKALSNEATH” SOKALSNEATHDISTANCE NNEQ NNEQ 05 NTT

USERDEFINED DISTANCE

IDENTIFIER CLASS NAME ARGS

“PYFUNC” PYFUNCDISTANCE FUNC

HEREFUNC IS A FUNCTION WHICH TAKES TWO ONEDIMENSIONAL NUMPY ARRAYS AND RETURNS A DISTANCE NOTE THAT IN ORDER TO BE USED WITHIN THE BALLTREE THE DISTANCE MUST BE A TRUE METRIC IE IT MUST SATISFY THE FOLLOWING PROPERTIES

1 NONNEGATIVITY DX Y 0

2 IDENTITY DX Y 0 IF AND ONLY IF X Y

3 SYMMETRY DX Y DY X

4 TRIANGLE INEQUALITY DX Y DY Z DX Z

BECAUSE OF THE PYTHON OBJECT OVERHEAD INVOLVED IN CALLING THE PYTHON FUNCTION THIS WILL BE FAIRLY SLOW BUT IT WILL HAVE THE SAME SCALING AS OTHER DISTANCES

METHODS

DISTTORDIST CONVERT THE TRUE DISTANCE TO THE REDUCED DISTANCE

GETMETRIC GET THE GIVEN DISTANCE METRIC FROM THE STRING IDENTIFIER

PAIRWISE COMPUTE THE PAIRWISE DISTANCES BETWEEN X AND Y

RDISTTODIST CONVERT THE REDUCED DISTANCE TO THE TRUE DISTANCE

INIT SELFARGS KWARGS

INITIALIZE SELF SEE HELPTYPESELF FOR ACCURATE SIGNATURE

DISTTORDIST

CONVERT THE TRUE DISTANCE TO THE REDUCED DISTANCE

THE REDUCED DISTANCE DEFINED FOR SOME METRICS IS A COMPUTATIONALLY MORE EFFICIENT MEASURE WHICH PRESERVES THE RANK OF THE TRUE DISTANCE FOR EXAMPLE IN THE EUCLIDEAN DISTANCE METRIC THE REDUCED DISTANCE IS THE SQUAREDEUCLIDEAN DISTANCE

GETMETRIC

GET THE GIVEN DISTANCE METRIC FROM THE STRING IDENTIFIER

SEE THE DOCSTRING OF DISTANCEMETRIC FOR A LIST OF AVAILABLE METRICS

PARAMETERS

METRIC STRING OR CLASS NAME THE DISTANCE METRIC TO USE

2160 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

KWARGS ADDITIONAL ARGUMENTS WILL BE PASSED TO THE REQUESTED METRIC

PAIRWISE

COMPUTE THE PAIRWISE DISTANCES BETWEEN X AND Y

THIS IS A CONVENIENCE ROUTINE FOR THE SAKE OF TESTING FOR MANY METRICS THE UTILITIES IN SCIPYSPATIALDISTANCECDIST AND SCIPYSPATIALDISTANCEPDIST WILL BE FASTER

PARAMETERS

XARRAYLIKE ARRAY OF SHAPE NX D REPRESENTING NX POINTS IN D DIMENSIONS

YARRAYLIKE OPTIONAL ARRAY OF SHAPE NY D REPRESENTING NY POINTS IN D DIMENSIONS

IF NOT SPECIFIED THEN YX

RETURNS

---

DIST NDARRAY THE SHAPE NX NY ARRAY OF PAIRWISE DISTANCES BETWEEN POINTS IN X AND Y

RDISTTODIST

CONVERT THE REDUCED DISTANCE TO THE TRUE DISTANCE

THE REDUCED DISTANCE DEFINED FOR SOME METRICS IS A COMPUTATIONALLY MORE EFFICIENT MEASURE WHICH PRE SERVES THE RANK OF THE TRUE DISTANCE FOR EXAMPLE IN THE EUCLIDEAN DISTANCE METRIC THE REDUCED DISTANCE IS THE SQUAREDEUCLIDEAN DISTANCE

6303SKLEARNNEIGHBORS KDTree

CLASSSKLEARNNEIGHBORS KDTree

KDTree FOR FAST GENERALIZED NPOINT PROBLEMS

KDTreeX LEAFSize40 METRIC'MINKOWSKI' KWARGS

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES NSAMPLES IS THE NUMBER OF POINTS IN THE DATA SET AND NFEATURES IS THE DIMENSION OF THE PARAMETER SPACE NOTE IF X IS A CCONTIGUOUS ARRAY OF DOUBLES THEN DATA WILL NOT BE COPIED OTHERWISE AN INTERNAL COPY WILL BE MADE

LEAFSize POSITIVE INTEGER DEFAULT 40 NUMBER OF POINTS AT WHICH TO SWITCH TO BRUTEFORCE

CHANGING LEAFSize WILL NOT AFFECT THE RESULTS OF A QUERY BUT CAN SIGNIFICANTLY IMPACT THE SPEED OF A QUERY AND THE MEMORY REQUIRED TO STORE THE CONSTRUCTED TREE THE AMOUNT OF MEMORY NEEDED TO STORE THE TREE SCALES AS APPROXIMATELY NSAMPLES LEAFSize FOR A SPEC IFIEDLEAFSize A LEAF NODE IS GUARANTEED TO SATISFY LEAFSize NPOINTS

2LEAFSize EXCEPT IN THE CASE THAT NSAMPLES LEAFSize

METRIC STRING OR DISTANCEMETRIC OBJECT THE DISTANCE METRIC TO USE FOR THE TREE DE FAULT'MINKOWSKI' WITH P2 THAT IS A EUCLIDEAN METRIC SEE THE DOCUMENTATION OF THE DISTANCEMETRIC CLASS FOR A LIST OF AVAILABLE METRICS KDTreeVALIDMETRICS GIVES A LIST OF THE METRICS WHICH ARE VALID FOR KDTree

ADDITIONAL KEYWORDS ARE PASSED TO THE DISTANCE METRIC CLASS

ATTRIBUTES

DATA MEMORY VIEW THE TRAINING DATA

6303SKLEARNNEIGHBORS NEAREST NEIGHBORS 2161

```
SCIKITLEARN USER GUIDE RELEASE 0213
EXAMPLES
QUERY FOR KNEAREST NEIGHBORS
  IMPORT NUMPY AS NP
RNG  NPRANDOMRANDOMSTATE0
X  RNGRANDOMSAMPLE10 3  10 POINTS IN 3 DIMENSIONS
TREE  KDTreeX LEAFSize2
DIST IND  TREEQUERYX1 K3
PRINTIND  INDICES OF 3 CLOSEST NEIGHBORS
0 3 1
PRINTDIST  DISTANCES TO 3 CLOSEST NEIGHBORS
  0 019662693 029473397
PICKLE AND UNPICKLE A TREE NOTE THAT THE STATE OF THE TREE IS SAVED IN THE PICKLE OPERATION THE TREE NEEDS NOT BE
REBUILT UPON UNPICKLING
  IMPORT NUMPY AS NP
  IMPORT PICKLE
RNG  NPRANDOMRANDOMSTATE0
X  RNGRANDOMSAMPLE10 3  10 POINTS IN 3 DIMENSIONS
TREE  KDTreeX LEAFSize2
S  PICKLEDUMPSTREE
TREECOPY  PICKLELOADSS
DIST IND  TREECOPYQUERYX1 K3
PRINTIND  INDICES OF 3 CLOSEST NEIGHBORS
0 3 1
PRINTDIST  DISTANCES TO 3 CLOSEST NEIGHBORS
  0 019662693 029473397
QUERY FOR NEIGHBORS WITHIN A GIVEN RADIUS
  IMPORT NUMPY AS NP
RNG  NPRANDOMRANDOMSTATE0
X  RNGRANDOMSAMPLE10 3  10 POINTS IN 3 DIMENSIONS
TREE  KDTreeX LEAFSize2
PRINTTREEQUERYRADIUSX1 R03 COUNTONLY TRUE
3
IND  TREEQUERYRADIUSX1 R03
PRINTIND  INDICES OF NEIGHBORS WITHIN DISTANCE 03
3 0 1
COMPUTE A GAUSSIAN KERNEL DENSITY ESTIMATE
  IMPORT NUMPY AS NP
RNG  NPRANDOMRANDOMSTATE42
X  RNGRANDOMSAMPLE100 3
TREE  KDTreeX
TREEKERNELDENSITYX3 H01 KERNELGAUSSIAN
ARRAY 694114649 783281226 72071716
COMPUTE A TWOPOINT AUTOCORRELATION FUNCTION
  IMPORT NUMPY AS NP
RNG  NPRANDOMRANDOMSTATE0
X  RNGRANDOMSAMPLE30 3
R  NPLinspace0 1 5
TREE  KDTreeX
2162 CHAPTER 6 API REFERENCE
```



SCIKITLEARN USER GUIDE RELEASE 0213

TREETWOPOINTCORRELATIONX R

ARRAY 30 62 278 580 820

METHODS

KERNELDENSITY SELF X H KERNEL ATOL COMPUTE THE KERNEL DENSITY ESTIMATE AT POINTS X WITH THE GIVEN KERNEL USING THE DISTANCE METRIC SPECIFIED AT

TREE CREATION

QUERY X K RETURNDISTANCE DUALTREE QUERY THE TREE FOR THE K NEAREST NEIGHBORS

QUERYRADIUS QUERYRADIUSSELF X R COUNTONLY FALSE

TWOPOINTCORRELATION COMPUTE THE TWOPOINT CORRELATION FUNCTION

GETARRAYS

GETNCALLS

GETTREESTATS

RESETNCALLS

INIT SELFARGS KWARGS

INITIALIZE SELF SEE HELPTYPESELF FOR ACCURATE SIGNATURE

KERNELDENSITY SELFHXKERNEL'GAUSSIAN' ATOL0 RTOL1E8 BREADTHFIRSTTRUE RE

TURNLOGFALSE

COMPUTE THE KERNEL DENSITY ESTIMATE AT POINTS X WITH THE GIVEN KERNEL USING THE DISTANCE METRIC SPECIFIED AT TREE CREATION

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES AN ARRAY OF POINTS TO QUERY LAST DIMENSION

SHOULD MATCH DIMENSION OF TRAINING DATA

HFLOAT THE BANDWIDTH OF THE KERNEL

KERNEL STRING SPECIFY THE KERNEL TO USE OPTIONS ARE 'GAUSSIAN' 'TOPHAT' 'EPANECHNIKOV' 'EXPONENTIAL' 'LINEAR' 'COSINE' DEFAULT IS KERNEL 'GAUSSIAN'

ATOL RTOL FLOAT DEFAULT 0 SPECIFY THE DESIRED RELATIVE AND ABSOLUTE TOLERANCE OF THE RESULT IF THE TRUE RESULT IS KTRUE THEN THE RETURNED RESULT KRET SATISFIES ABSKTRUE

KRET ATOL RTOL KRET THE DEFAULT IS ZERO IE MACHINE PRECISION FOR

BOTH

BREADTHFIRST BOOLEAN DEFAULT FALSE IF TRUE USE A BREADTHFIRST SEARCH IF FALSE DEFAULT USE A DEPTHFIRST SEARCH BREADTHFIRST IS GENERALLY FASTER FOR COMPACT KERNELS ANDOR HIGH TOLERANCES

RETURNLOG BOOLEAN DEFAULT FALSE RETURN THE LOGARITHM OF THE RESULT THIS CAN BE MORE ACCURATE THAN RETURNING THE RESULT ITSELF FOR NARROW KERNELS

RETURNS

DENSITY NDARRAY THE ARRAY OF LOGDENSITY EVALUATIONS SHAPE XSHAPE1

QUERYXK1RETURNDISTANCETRUE DUALTREEFALSE BREADTHFIRSTFALSE

QUERY THE TREE FOR THE K NEAREST NEIGHBORS

PARAMETERS

630SKLEARNNEIGHBORS NEAREST NEIGHBORS 2163

SCIKITLEARN USER GUIDE RELEASE 0213

XARRAYLIKE SHAPE NSAMPLES NFEATURES AN ARRAY OF POINTS TO QUERY

KINTEGER DEFAULT 1 THE NUMBER OF NEAREST NEIGHBORS TO RETURN

RETURNDISTANCE BOOLEAN DEFAULT TRUE IF TRUE RETURN A TUPLE D I OF DISTANCES AND INDICES IF FALSE RETURN ARRAY I

DUALTREE BOOLEAN DEFAULT FALSE IF TRUE USE THE DUAL TREE FORMALISM FOR THE QUERY A TREE IS BUILT FOR THE QUERY POINTS AND THE PAIR OF TREES IS USED TO EFFICIENTLY SEARCH THIS SPACE THIS CAN LEAD TO BETTER PERFORMANCE AS THE NUMBER OF POINTS GROWS LARGE

BREADTHFIRST BOOLEAN DEFAULT FALSE IF TRUE THEN QUERY THE NODES IN A BREADTHFIRST MANNER OTHERWISE QUERY THE NODES IN A DEPTHFIRST MANNER

SORTRESULTS BOOLEAN DEFAULT TRUE IF TRUE THEN DISTANCES AND INDICES OF EACH POINT ARE SORTED ON RETURN SO THAT THE FIRST COLUMN CONTAINS THE CLOSEST POINTS OTHERWISE NEIGHBORS ARE RETURNED IN AN ARBITRARY ORDER

RETURNS

IIF RETURNDISTANCE FALSE

DI IF RETURNDISTANCE TRUE

DARRAY OF DOUBLES SHAPE XSHAPE1 K EACH ENTRY GIVES THE LIST OF DISTANCES TO THE NEIGHBORS OF THE CORRESPONDING POINT

IARRAY OF INTEGERS SHAPE XSHAPE1 K EACH ENTRY GIVES THE LIST OF INDICES OF NEIGHBORS OF THE CORRESPONDING POINT

QUERYRADIUS

QUERYRADIUSSELF X R COUNTONLY FALSE

QUERY THE TREE FOR NEIGHBORS WITHIN A RADIUS R

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES AN ARRAY OF POINTS TO QUERY

RDISTANCE WITHIN WHICH NEIGHBORS ARE RETURNED R CAN BE A SINGLE VALUE OR AN ARRAY OF VALUES OF SHAPE XSHAPE1 IF DIFFERENT RADII ARE DESIRED FOR EACH POINT

RETURNDISTANCE BOOLEAN DEFAULT FALSE IF TRUE RETURN DISTANCES TO NEIGHBORS OF EACH POINT IF FALSE RETURN ONLY NEIGHBORS NOTE THAT UNLIKE THE QUERY METHOD SETTING RETURNDISTANCETRUE HERE ADDS TO THE COMPUTATION TIME NOT ALL DISTANCES NEED TO BE CALCULATED EXPLICITLY FOR RETURNDISTANCEFALSE RESULTS ARE NOT SORTED BY DEFAULT SEE SORTRESULTS KEYWORD

COUNTONLY BOOLEAN DEFAULT FALSE IF TRUE RETURN ONLY THE COUNT OF POINTS WITHIN DISTANCE R IF FALSE RETURN THE INDICES OF ALL POINTS WITHIN DISTANCE R IF RETURNDISTANCETRUE SETTING COUNTONLYTRUE WILL RESULT IN AN ERROR

SORTRESULTS BOOLEAN DEFAULT FALSE IF TRUE THE DISTANCES AND INDICES WILL BE SORTED BEFORE BEING RETURNED IF FALSE THE RESULTS WILL NOT BE SORTED IF RETURNDISTANCE FALSE SETTING SORTRESULTS TRUE WILL RESULT IN AN ERROR

RETURNS

COUNT IF COUNTONLY TRUE

IND IF COUNTONLY FALSE AND RETURNDISTANCE FALSE

IND DIST IF COUNTONLY FALSE AND RETURNDISTANCE TRUE

2164 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

COUNT ARRAY OF INTEGERS SHAPE XSHAPE1 EACH ENTRY GIVES THE NUMBER OF NEIGHBORS WITHIN A DISTANCE R OF THE CORRESPONDING POINT  
IND ARRAY OF OBJECTS SHAPE XSHAPE1 EACH ELEMENT IS A NUMPY INTEGER ARRAY LISTING THE INDICES OF NEIGHBORS OF THE CORRESPONDING POINT NOTE THAT UNLIKE THE RESULTS OF A K NEIGHBORS QUERY THE RETURNED NEIGHBORS ARE NOT SORTED BY DISTANCE BY DEFAULT  
DIST ARRAY OF OBJECTS SHAPE XSHAPE1 EACH ELEMENT IS A NUMPY DOUBLE ARRAY LISTING THE DISTANCES CORRESPONDING TO INDICES IN I

TWOPOINTCORRELATION

COMPUTE THE TWOPOINT CORRELATION FUNCTION

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES AN ARRAY OF POINTS TO QUERY LAST DIMENSION SHOULD MATCH DIMENSION OF TRAINING DATA  
RARRAYLIKE A ONEDIMENSIONAL ARRAY OF DISTANCES  
DUALTREE BOOLEAN DEFAULT FALSE IF TRUE USE A DUALTREE ALGORITHM OTHERWISE USE A SINGLETREE ALGORITHM DUAL TREE ALGORITHMS CAN HAVE BETTER SCALING FOR LARGE N

RETURNS

COUNTS NDARRAY COUNTSI CONTAINS THE NUMBER OF PAIRS OF POINTS WITH DISTANCE LESS THAN OR EQUAL TO RI

6304SKLEARNNEIGHBORS KERNELDENSITY

CLASSSSKLEARNNEIGHBORS KERNELDENSITY BANDWIDTH10 ALGORITHM'AUTO' KERNEL'GAUSSIAN' METRIC'EUCLIDEAN' ATOL0 RTOL0 BREADTHFIRSTTRUE

LEAFSIZE40 METRICPARAMSNONE

KERNEL DENSITY ESTIMATION

READ MORE IN THE USER GUIDE

PARAMETERS

BANDWIDTH FLOAT THE BANDWIDTH OF THE KERNEL

ALGORITHM STRING THE TREE ALGORITHM TO USE VALID OPTIONS ARE 'KDTREE''BALLTREE''AUTO' DEFAULT IS 'AUTO'

KERNEL STRING THE KERNEL TO USE VALID KERNELS ARE 'GAUS

SIAN''TOPHAT''EPANECHNIKOV''EXPONENTIAL''LINEAR''COSINE' DEFAULT IS 'GAUSSIAN'

METRIC STRING THE DISTANCE METRIC TO USE NOTE THAT NOT ALL METRICS ARE VALID WITH ALL ALGORITHMS

REFER TO THE DOCUMENTATION OF BALLTREE ANDKDTREE FOR A DESCRIPTION OF AVAILABLE ALGO

RITHMS NOTE THAT THE NORMALIZATION OF THE DENSITY OUTPUT IS CORRECT ONLY FOR THE EUCLIDEAN

DISTANCE METRIC DEFAULT IS 'EUCLIDEAN'

ATOL FLOAT THE DESIRED ABSOLUTE TOLERANCE OF THE RESULT A LARGER TOLERANCE WILL GENERALLY LEAD TO FASTER EXECUTION DEFAULT IS 0

RTOL FLOAT THE DESIRED RELATIVE TOLERANCE OF THE RESULT A LARGER TOLERANCE WILL GENERALLY LEAD TO FASTER EXECUTION DEFAULT IS 1E8

BREADTHFIRST BOOLEAN IF TRUE DEFAULT USE A BREADTHFIRST APPROACH TO THE PROBLEM OTHERWISE USE A DEPTHFIRST APPROACH

LEAFSIZE INT SPECIFY THE LEAF SIZE OF THE UNDERLYING TREE SEE BALLTREE ORKDTREE FOR

DETAILS DEFAULT IS 40

630SKLEARNNEIGHBORS NEAREST NEIGHBORS 2165

SCIKITLEARN USER GUIDE RELEASE 0213

METRICPARAMS DICT ADDITIONAL PARAMETERS TO BE PASSED TO THE TREE FOR USE WITH THE METRIC  
FOR MORE INFORMATION SEE THE DOCUMENTATION OF BALLTREE ORKDTREE

METHODS

FITSELF X Y SAMPLEWEIGHT FIT THE KERNEL DENSITY MODEL ON THE DATA

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

SAMPLE SELF NSAMPLES RANDOMSTATE GENERATE RANDOM SAMPLES FROM THE MODEL

SCORE SELF X Y COMPUTE THE TOTAL LOG PROBABILITY DENSITY UNDER THE  
MODEL

SCORESAMPLES SELF X EVALUATE THE DENSITY MODEL ON THE DATA

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELF BANDWIDTH10 ALGORITHM'AUTO' KERNEL'GAUSSIAN' METRIC'EUCLIDEAN' ATOL0  
RTOL0 BREADTHFIRSTTRUE LEAFSIZE40 METRICPARAMSNONE

FITSELFXYNONE SAMPLEWEIGHTNONE

FIT THE KERNEL DENSITY MODEL ON THE DATA

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES LIST OF NFEATURESDIMENSIONAL DATA POINTS  
EACH ROW CORRESPONDS TO A SINGLE DATA POINT

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL LIST OF SAMPLE WEIGHTS ATTACHED  
TO THE DATA X

GETPARAMS SELFDEEPTTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

SAMPLESELFNSAMPLES1 RANDOMSTATENONE

GENERATE RANDOM SAMPLES FROM THE MODEL

CURRENTLY THIS IS IMPLEMENTED ONLY FOR GAUSSIAN AND TOPHAT KERNELS

PARAMETERS

NSAMPLES INT OPTIONAL NUMBER OF SAMPLES TO GENERATE DEFAULTS TO 1

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE DEFAULT TO NONE IF INT RANDOMSTATE IS  
THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS  
THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE  
INSTANCE USED BY NPRANDOM

RETURNS

XARRAYLIKE SHAPE NSAMPLES NFEATURES LIST OF SAMPLES

SCORESELFXYNONE

COMPUTE THE TOTAL LOG PROBABILITY DENSITY UNDER THE MODEL

PARAMETERS

2166 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

XARRAYLIKE SHAPE NSAMPLES NFEATURES LIST OF NFEATURESDIMENSIONAL DATA POINTS  
EACH ROW CORRESPONDS TO A SINGLE DATA POINT

RETURNS

LOGPROB FLOAT TOTAL LOGLIKELIHOOD OF THE DATA IN X THIS IS NORMALIZED TO BE A PROBABILITY  
DENSITY SO THE VALUE WILL BE LOW FOR HIGHDIMENSIONAL DATA

SCORESAMPLES SELF  
EVALUATE THE DENSITY MODEL ON THE DATA

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES AN ARRAY OF POINTS TO QUERY LAST DIMENSION  
SHOULD MATCH DIMENSION OF TRAINING DATA NFEATURES

RETURNS

DENSITY NDARRAY SHAPE NSAMPLES THE ARRAY OF LOGDENSITY EVALUATIONS THESE ARE NOR  
MALIZED TO BE PROBABILITY DENSITIES SO VALUES WILL BE LOW FOR HIGHDIMENSIONAL DATA

SETPARAMS SELF  
SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE  
PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT  
OF A NESTED OBJECT

RETURNS

SELF

EXAMPLES USING SKLEARNNEIGHBORSKERNELDENSITY

- KERNEL DENSITY ESTIMATION
- KERNEL DENSITY ESTIMATE OF SPECIES DISTRIBUTIONS
- SIMPLE 1D KERNEL DENSITY ESTIMATION

6305SKLEARNNEIGHBORS KNEIGHBORSCLASSIFIER

CLASSSSKLEARNNEIGHBORS KNEIGHBORSCLASSIFIER NNEIGHBORS5 WEIGHTS'UNIFORM' AL

GORITHM'AUTO' LEAFSIZE30 P2MET

RIC'MINKOWSKI' METRICPARAMSNONE

NJOBSNONE KWARGS

CLASSIFIER IMPLEMENTING THE KNEAREST NEIGHBORS VOTE

READ MORE IN THE USER GUIDE

PARAMETERS

NNEIGHBORS INT OPTIONAL DEFAULT 5 NUMBER OF NEIGHBORS TO USE BY DEFAULT FOR

KNEIGHBORS QUERIES

WEIGHTS STR OR CALLABLE OPTIONAL DEFAULT 'UNIFORM' WEIGHT FUNCTION USED IN PREDICTION POS  
SIBLE VALUES

- 'UNIFORM' UNIFORM WEIGHTS ALL POINTS IN EACH NEIGHBORHOOD ARE WEIGHTED EQUALLY
- 'DISTANCE' WEIGHT POINTS BY THE INVERSE OF THEIR DISTANCE IN THIS CASE CLOSER NEIGHBORS  
OF A QUERY POINT WILL HAVE A GREATER INFLUENCE THAN NEIGHBORS WHICH ARE FURTHER AWAY

630SKLEARNNEIGHBORS NEAREST NEIGHBORS 2167

SCIKITLEARN USER GUIDE RELEASE 0213

- CALLABLE A USERDEFINED FUNCTION WHICH ACCEPTS AN ARRAY OF DISTANCES AND RETURNS AN ARRAY OF THE SAME SHAPE CONTAINING THE WEIGHTS
- ALGORITHM 'AUTO' 'BALLTREE' 'KDTREE' 'BRUTE' OPTIONAL ALGORITHM USED TO COMPUTE THE NEAREST NEIGHBORS
- 'BALLTREE' WILL USE BALLTREE
- 'KDTREE' WILL USE KDTREE
- 'BRUTE' WILL USE A BRUTEFORCE SEARCH
- 'AUTO' WILL ATTEMPT TO DECIDE THE MOST APPROPRIATE ALGORITHM BASED ON THE VALUES PASSED

TOFIT METHOD

NOTE FITTING ON SPARSE INPUT WILL OVERRIDE THE SETTING OF THIS PARAMETER USING BRUTE FORCE  
LEAFSIZE INT OPTIONAL DEFAULT 30 LEAF SIZE PASSED TO BALLTREE OR KDTREE THIS CAN AFFECT  
THE SPEED OF THE CONSTRUCTION AND QUERY AS WELL AS THE MEMORY REQUIRED TO STORE THE TREE  
THE OPTIMAL VALUE DEPENDS ON THE NATURE OF THE PROBLEM

PINTEGER OPTIONAL DEFAULT 2 POWER PARAMETER FOR THE MINKOWSKI METRIC WHEN P 1 THIS  
IS EQUIVALENT TO USING MANHATTANDISTANCE L1 AND EUCLIDEANDISTANCE L2 FOR P 2 FOR

ARBITRARY P MINKOWSKIDISTANCE LP IS USED

METRIC STRING OR CALLABLE DEFAULT 'MINKOWSKI' THE DISTANCE METRIC TO USE FOR THE TREE THE  
DEFAULT METRIC IS MINKOWSKI AND WITH P2 IS EQUIVALENT TO THE STANDARD EUCLIDEAN METRIC  
SEE THE DOCUMENTATION OF THE DISTANCEMETRIC CLASS FOR A LIST OF AVAILABLE METRICS  
METRICPARAMS DICT OPTIONAL DEFAULT NONE ADDITIONAL KEYWORD ARGUMENTS FOR THE METRIC  
FUNCTION

NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF PARALLEL JOBS TO RUN FOR NEIGHBORS  
SEARCHNONE MEANS 1 UNLESS IN A JOBLIBPARALLELBKEND CONTEXT1MEANS

USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS DOESN'T AFFECT FIT METHOD

SEE ALSO

RADIUSNEIGHBORSCLASSIFIER

KNEIGHBORSREGRESSOR

RADIUSNEIGHBORSREGRESSOR

NEARESTNEIGHBORS

NOTES

SEE NEAREST NEIGHBORS IN THE ONLINE DOCUMENTATION FOR A DISCUSSION OF THE CHOICE OF ALGORITHM AND

LEAFSIZE

WARNING REGARDING THE NEAREST NEIGHBORS ALGORITHMS IF IT IS FOUND THAT TWO NEIGHBORS NEIGHBOR K1 AND  
K HAVE IDENTICAL DISTANCES BUT DIFFERENT LABELS THE RESULTS WILL DEPEND ON THE ORDERING OF THE TRAINING DATA

[HTTPS://ENWIKIPEDIA.ORG/WIKI/KNEARESTNEIGHBORALGORITHM](https://en.wikipedia.org/wiki/Nearest_neighbor_algorithm)

2168 CHAPTER 6 API REFERENCE

```
SCIKITLEARN USER GUIDE RELEASE 0213
EXAMPLES
X 0 1 2 3
Y 0 0 1 1
FROM SKLEARNNEIGHBORS IMPORT KNEIGHBORSCLASSIFIER
NEIGH KNEIGHBORSCLASSIFIERNNEIGHBORS3
NEIGHFITX Y
KNEIGHBORSCLASSIFIER
PRINTNEIGHPREDICT11
0
PRINTNEIGHPREDICTPROBA09
066666667 033333333
METHODS
FITSELF X Y FIT THE MODEL USING X AS TRAINING DATA AND Y AS TARGET
VALUES
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR
KNEIGHBORS SELF X NNEIGHBORS FINDS THE KNEIGHBORS OF A POINT
KNEIGHBORSGRAPH SELF X NNEIGHBORS MODE COMPUTES THE WEIGHTED GRAPH OF KNEIGHBORS FOR
POINTS IN X
PREDICT SELF X PREDICT THE CLASS LABELS FOR THE PROVIDED DATA
PREDICTPROBA SELF X RETURN PROBABILITY ESTIMATES FOR THE TEST DATA X
SCORE SELF X Y SAMPLEWEIGHT RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND
LABELS
SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR
INIT SELFNNEIGHBORS5 WEIGHTS'UNIFORM' ALGORITHM'AUTO' LEAFSIZE30 P2MET
RIC'MINKOWSKI' METRICPARAMSNONE NJOBSNONE KWARGS
FITSELFXY
FIT THE MODEL USING X AS TRAINING DATA AND Y AS TARGET VALUES
PARAMETERS
XARRAYLIKE SPARSE MATRIX BALLTREE KDTree TRAINING DATA IF ARRAY OR MATRIX SHAPE
NSAMPLES NFEATURES OR NSAMPLES NSAMPLES IF METRIC'PRECOMPUTED'
YARRAYLIKE SPARSE MATRIX TARGET VALUES OF SHAPE NSAMPLES OR NSAMPLES
NOUTPUTS
GETPARAMS SELFDEEPTREE
GET PARAMETERS FOR THIS ESTIMATOR
PARAMETERS
DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED
SUBOBJECTS THAT ARE ESTIMATORS
RETURNS
PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES
KNEIGHBORS SELFXNONE NNEIGHBORSNONE RETURNDISTANCETRUE
FINDS THE KNEIGHBORS OF A POINT RETURNS INDICES OF AND DISTANCES TO THE NEIGHBORS OF EACH POINT
PARAMETERS
630SKLEARNNEIGHBORS NEAREST NEIGHBORS 2169
```

SCIKITLEARN USER GUIDE RELEASE 0213

XARRAYLIKE SHAPE NQUERY NFEATURES OR NQUERY NINDEXED IF METRIC 'PRECOMPUTED' THE QUERY POINT OR POINTS IF NOT PROVIDED NEIGHBORS OF EACH INDEXED POINT ARE RETURNED IN THIS CASE THE QUERY POINT IS NOT CONSIDERED ITS OWN NEIGHBOR NNEIGHBORS INT NUMBER OF NEIGHBORS TO GET DEFAULT IS THE VALUE PASSED TO THE CONSTRUCTOR RETURNEDISTANCE BOOLEAN OPTIONAL DEFAULTS TO TRUE IF FALSE DISTANCES WILL NOT BE RETURNED

RETURNS  
DIST ARRAY ARRAY REPRESENTING THE LENGTHS TO POINTS ONLY PRESENT IF RETURNEDISTANCETRUE  
IND ARRAY INDICES OF THE NEAREST POINTS IN THE POPULATION MATRIX

EXAMPLES  
IN THE FOLLOWING EXAMPLE WE CONSTRUCT A NEIGHBORSCLASSIFIER CLASS FROM AN ARRAY REPRESENTING OUR DATA SET AND ASK WHO'S THE CLOSEST POINT TO 111  
SAMPLES 0 0 0 0 5 0 1 1 5  
FROM SKLEARNNEIGHBORS IMPORT NEARESTNEIGHBORS  
NEIGH NEARESTNEIGHBORSNNEIGHBORS1  
NEIGHFITSAMPLES  
NEARESTNEIGHBORSALGORITHM AUTO LEAF SIZE30  
PRINTNEIGHKNEIGHBORS1 1 1  
ARRAY05 ARRAY2

AS YOU CAN SEE IT RETURNS 05 AND 2 WHICH MEANS THAT THE ELEMENT IS AT DISTANCE 05 AND IS THE THIRD ELEMENT OF SAMPLES INDEXES START AT 0 YOU CAN ALSO QUERY FOR MULTIPLE POINTS  
X 0 1 0 1 0 1  
NEIGHKNEIGHBORSX RETURNEDISTANCE FALSE  
ARRAY1  
2

KNEIGHBORSGRAPH SELF XNONE NNEIGHBORSNONE MODE'CONNECTIVITY'  
COMPUTES THE WEIGHTED GRAPH OF KNEIGHBORS FOR POINTS IN X  
PARAMETERS

XARRAYLIKE SHAPE NQUERY NFEATURES OR NQUERY NINDEXED IF METRIC 'PRECOMPUTED' THE QUERY POINT OR POINTS IF NOT PROVIDED NEIGHBORS OF EACH INDEXED POINT ARE RETURNED IN THIS CASE THE QUERY POINT IS NOT CONSIDERED ITS OWN NEIGHBOR NNEIGHBORS INT NUMBER OF NEIGHBORS FOR EACH SAMPLE DEFAULT IS VALUE PASSED TO THE CONSTRUCTOR  
MODE 'CONNECTIVITY' 'DISTANCE' OPTIONAL TYPE OF RETURNED MATRIX 'CONNECTIVITY' WILL RETURN THE CONNECTIVITY MATRIX WITH ONES AND ZEROS IN 'DISTANCE' THE EDGES ARE EUCLIDEAN DISTANCE BETWEEN POINTS

RETURNS  
ASPARSE MATRIX IN CSR FORMAT SHAPE NSAMPLES NSAMPLESFIT NSAMPLESFIT IS THE NUMBER OF SAMPLES IN THE FITTED DATA A I J IS ASSIGNED THE WEIGHT OF EDGE THAT CONNECTS I TO J  
SEE ALSO  
2170 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213  
 NEARESTNEIGHBORSRADIUSNEIGHBORSGRAPH  
 EXAMPLES  
 X 0 3 1  
 FROM SKLEARNNEIGHBORS IMPORT NEARESTNEIGHBORS  
 NEIGH NEARESTNEIGHBORSNNEIGHBORS2  
 NEIGHFITX  
 NEARESTNEIGHBORSALGORITHMAUTO LEAFSIZE30  
 A NEIGHKNEIGHBORSGRAPHX  
 ATOARRAY  
 ARRAY1 0 1  
 0 1 1  
 1 0 1  
 PREDICTSELF  
 PREDICT THE CLASS LABELS FOR THE PROVIDED DATA  
 PARAMETERS  
 XARRAYLIKE SHAPE NQUERY NFEATURES OR NQUERY NINDEXED IF METRIC 'PRECOM  
 PUTED' TEST SAMPLES  
 RETURNS  
 YARRAY OF SHAPE NSAMPLES OR NSAMPLES NOUTPUTS CLASS LABELS FOR EACH DATA SAMPLE  
 PREDICTPROBA SELF  
 RETURN PROBABILITY ESTIMATES FOR THE TEST DATA X  
 PARAMETERS  
 XARRAYLIKE SHAPE NQUERY NFEATURES OR NQUERY NINDEXED IF METRIC 'PRECOM  
 PUTED' TEST SAMPLES  
 RETURNS  
 PARRAY OF SHAPE NSAMPLES NCLASSES OR A LIST OF NOUTPUTS OF SUCH ARRAYS IF NOUTPUTS  
 1 THE CLASS PROBABILITIES OF THE INPUT SAMPLES CLASSES ARE ORDERED BY LEXICOGRAPHIC  
 ORDER  
 SCORESELFXYSAMPLEWEIGHTNONE  
 RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS  
 IN MULTILABEL CLASSIFICATION THIS IS THE SUBSET ACCURACY WHICH IS A HARSH METRIC SINCE YOU REQUIRE FOR EACH  
 SAMPLE THAT EACH LABEL SET BE CORRECTLY PREDICTED  
 PARAMETERS  
 XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES  
 YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE LABELS FOR X  
 SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS  
 RETURNS  
 SCORE FLOAT MEAN ACCURACY OF SELF PREDICTX WRT Y  
 SETPARAMS SELFPARAMS  
 SET THE PARAMETERS OF THIS ESTIMATOR  
 630SKLEARNNEIGHBORS NEAREST NEIGHBORS 2171

SCIKITLEARN USER GUIDE RELEASE 0213

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS  
SELF

EXAMPLES USING SKLEARNNEIGHBORSKNEIGHBORSCLASSIFIER

- CLASSIFIER COMPARISON
- PLOT THE DECISION BOUNDARIES OF A VOTINGCLASSIFIER
- DIGITS CLASSIFICATION EXERCISE
- NEAREST NEIGHBORS CLASSIFICATION
- COMPARING NEAREST NEIGHBORS WITH AND WITHOUT NEIGHBORHOOD COMPONENTS ANALYSIS
- DIMENSIONALITY REDUCTION WITH NEIGHBORHOOD COMPONENTS ANALYSIS
- CLASSIFICATION OF TEXT DOCUMENTS USING SPARSE FEATURES

6306SKLEARNNEIGHBORS KNEIGHBORSREGRESSOR

CLASSSSKLEARNNEIGHBORS KNEIGHBORSREGRESSOR NNEIGHBORS5 WEIGHTS'UNIFORM' ALGO

RITHM'AUTO' LEAFSIZE30 P2MET

RIC'MINKOWSKI' METRICPARAMSNONE

NJOBSNONE KWARGS

REGRESSION BASED ON KNEAREST NEIGHBORS

THE TARGET IS PREDICTED BY LOCAL INTERPOLATION OF THE TARGETS ASSOCIATED OF THE NEAREST NEIGHBORS IN THE TRAINING SET

READ MORE IN THE USER GUIDE

PARAMETERS

NNEIGHBORS INT OPTIONAL DEFAULT 5 NUMBER OF NEIGHBORS TO USE BY DEFAULT FOR

KNEIGHBORS QUERIES

WEIGHTS STR OR CALLABLE WEIGHT FUNCTION USED IN PREDICTION POSSIBLE VALUES

- 'UNIFORM' UNIFORM WEIGHTS ALL POINTS IN EACH NEIGHBORHOOD ARE WEIGHTED EQUALLY
- 'DISTANCE' WEIGHT POINTS BY THE INVERSE OF THEIR DISTANCE IN THIS CASE CLOSER NEIGHBORS OF A QUERY POINT WILL HAVE A GREATER INFLUENCE THAN NEIGHBORS WHICH ARE FURTHER AWAY
- CALLABLE A USERDEFINED FUNCTION WHICH ACCEPTS AN ARRAY OF DISTANCES AND RETURNS AN ARRAY OF THE SAME SHAPE CONTAINING THE WEIGHTS

UNIFORM WEIGHTS ARE USED BY DEFAULT

ALGORITHM 'AUTO' 'BALLTREE' 'KDTree' 'BRUTE' OPTIONAL ALGORITHM USED TO COMPUTE THE

NEAREST NEIGHBORS

- 'BALLTREE' WILL USE BALLTREE
- 'KDTree' WILL USE KDTree
- 'BRUTE' WILL USE A BRUTEFORCE SEARCH
- 'AUTO' WILL ATTEMPT TO DECIDE THE MOST APPROPRIATE ALGORITHM BASED ON THE VALUES PASSED

TOFIT METHOD

2172 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

NOTE FITTING ON SPARSE INPUT WILL OVERRIDE THE SETTING OF THIS PARAMETER USING BRUTE FORCE  
LEAFSIZE INT OPTIONAL DEFAULT 30 LEAF SIZE PASSED TO BALLTREE OR KDTree THIS CAN AFFECT  
THE SPEED OF THE CONSTRUCTION AND QUERY AS WELL AS THE MEMORY REQUIRED TO STORE THE TREE  
THE OPTIMAL VALUE DEPENDS ON THE NATURE OF THE PROBLEM

P INTEGER OPTIONAL DEFAULT 2 POWER PARAMETER FOR THE MINKOWSKI METRIC WHEN P 1 THIS  
IS EQUIVALENT TO USING MANHATTANDISTANCE L1 AND EUCLIDEANDISTANCE L2 FOR P 2 FOR  
ARBITRARY P MINKOWSKI DISTANCE LP IS USED

METRIC STRING OR CALLABLE DEFAULT 'MINKOWSKI' THE DISTANCE METRIC TO USE FOR THE TREE THE  
DEFAULT METRIC IS MINKOWSKI AND WITH P2 IS EQUIVALENT TO THE STANDARD EUCLIDEAN METRIC  
SEE THE DOCUMENTATION OF THE DISTANCEMETRIC CLASS FOR A LIST OF AVAILABLE METRICS  
METRICPARAMS DICT OPTIONAL DEFAULT NONE ADDITIONAL KEYWORD ARGUMENTS FOR THE METRIC  
FUNCTION

NJOBS INT OR NONE OPTIONAL DEFAULT NONE THE NUMBER OF PARALLEL JOBS TO RUN FOR NEIGHBORS  
SEARCH NONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT 1 MEANS  
USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS DOESN'T AFFECT FIT METHOD

SEE ALSO

NEARESTNEIGHBORS  
RADIUSNEIGHBORSREGRESSOR  
KNEIGHBORSCLASSIFIER  
RADIUSNEIGHBORSCLASSIFIER

NOTES

SEE NEAREST NEIGHBORS IN THE ONLINE DOCUMENTATION FOR A DISCUSSION OF THE CHOICE OF ALGORITHM AND  
LEAFSIZE

WARNING REGARDING THE NEAREST NEIGHBORS ALGORITHMS IF IT IS FOUND THAT TWO NEIGHBORS NEIGHBOR K1 AND  
K HAVE IDENTICAL DISTANCES BUT DIFFERENT LABELS THE RESULTS WILL DEPEND ON THE ORDERING OF THE TRAINING DATA  
[HTTPSENWIKIPEDIAORGWIKINEARESTNEIGHBORALGORITHM](http://en.wikipedia.org/wiki/Nearest_neighbor_algorithm)

EXAMPLES

```
X 0 1 2 3
Y 0 0 1 1

from sklearn.neighbors import KNeighborsRegressor
neigh = KNeighborsRegressor(n_neighbors=2)
neigh.fit(X, Y)
KNeighborsRegressor()
print(neigh.predict(1.5))
05
630SKLEARNNEIGHBORS NEAREST NEIGHBORS 2173
```

METHODS

FITSELF X Y FIT THE MODEL USING X AS TRAINING DATA AND Y AS TARGET VALUES

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

KNEIGHBORS SELF X NNEIGHBORS FINDS THE KNEIGHBORS OF A POINT

KNEIGHBORSGRAPH SELF X NNEIGHBORS MODE COMPUTES THE WEIGHTED GRAPH OF KNEIGHBORS FOR POINTS IN X

PREDICT SELF X PREDICT THE TARGET FOR THE PROVIDED DATA

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELF NNEIGHBORS5 WEIGHTS'UNIFORM' ALGORITHM'AUTO' LEAFSIZE30 P2MET

RIC'MINKOWSKI' METRICPARAMSNONE NJOBSNONE KWARGS

FITSELFXY

FIT THE MODEL USING X AS TRAINING DATA AND Y AS TARGET VALUES

PARAMETERS

XARRAYLIKE SPARSE MATRIX BALLTREE KDTree TRAINING DATA IF ARRAY OR MATRIX SHAPE

NSAMPLES NFEATURES OR NSAMPLES NSAMPLES IF METRIC'PRECOMPUTED'

YARRAYLIKE SPARSE MATRIX

TARGET VALUES ARRAY OF FLOAT VALUES SHAPE NSAMPLES OR NSAMPLES NOUTPUTS

GETPARAMS SELFDEEPTREE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

KNEIGHBORS SELF X NONE NNEIGHBORS NONE RETURN DISTANCE TRUE

FINDS THE KNEIGHBORS OF A POINT RETURNS INDICES OF AND DISTANCES TO THE NEIGHBORS OF EACH POINT

PARAMETERS

XARRAYLIKE SHAPE NQUERY NFEATURES OR NQUERY NINDEXED IF METRIC 'PRECOM

PUTED' THE QUERY POINT OR POINTS IF NOT PROVIDED NEIGHBORS OF EACH INDEXED POINT ARE

RETURNED IN THIS CASE THE QUERY POINT IS NOT CONSIDERED ITS OWN NEIGHBOR

NNEIGHBORS INT NUMBER OF NEIGHBORS TO GET DEFAULT IS THE VALUE PASSED TO THE CONSTRUCTOR

RETURN DISTANCE BOOLEAN OPTIONAL DEFAULTS TO TRUE IF FALSE DISTANCES WILL NOT BE RETURNED

RETURNS

DIST ARRAY ARRAY REPRESENTING THE LENGTHS TO POINTS ONLY PRESENT IF RETURN DISTANCE TRUE

IND ARRAY INDICES OF THE NEAREST POINTS IN THE POPULATION MATRIX

```
SCIKITLEARN USER GUIDE RELEASE 0213
EXAMPLES
IN THE FOLLOWING EXAMPLE WE CONSTRUCT A NEIGHBORSCLASSIFIER CLASS FROM AN ARRAY REPRESENTING OUR DATA SET
AND ASK WHO'S THE CLOSEST POINT TO 111
SAMPLES 0 0 0 0 5 0 1 1 5
FROM SKLEARNNEIGHBORS IMPORT NEARESTNEIGHBORS
NEIGH NEARESTNEIGHBORSNNEIGHBORS1
NEIGHFITSAMPLES
NEARESTNEIGHBORSALGORITHMAUTO LEAFSIZE30
PRINTNEIGHKNEIGHBORS1 1 1
ARRAY05 ARRAY2
AS YOU CAN SEE IT RETURNS 05 AND 2 WHICH MEANS THAT THE ELEMENT IS AT DISTANCE 05 AND IS THE THIRD
ELEMENT OF SAMPLES INDEXES START AT 0 YOU CAN ALSO QUERY FOR MULTIPLE POINTS
X 0 1 0 1 0 1
NEIGHKNEIGHBORSX RETURNDISTANCE FALSE
ARRAY1
2
KNEIGHBORSGRAPH SELFXXNONE NNEIGHBORSNONE MODE'CONNECTIVITY'
COMPUTES THE WEIGHTED GRAPH OF KNEIGHBORS FOR POINTS IN X
PARAMETERS
XARRAYLIKE SHAPE NQUERY NFEATURES OR NQUERY NINDEXED IF METRIC 'PRECOM
PUTED' THE QUERY POINT OR POINTS IF NOT PROVIDED NEIGHBORS OF EACH INDEXED POINT ARE
RETURNED IN THIS CASE THE QUERY POINT IS NOT CONSIDERED ITS OWN NEIGHBOR
NNEIGHBORS INT NUMBER OF NEIGHBORS FOR EACH SAMPLE DEFAULT IS VALUE PASSED TO THE
CONSTRUCTOR
MODE 'CONNECTIVITY' 'DISTANCE' OPTIONAL TYPE OF RETURNED MATRIX 'CONNECTIVITY' WILL
RETURN THE CONNECTIVITY MATRIX WITH ONES AND ZEROS IN 'DISTANCE' THE EDGES ARE EUCLIDEAN
DISTANCE BETWEEN POINTS
RETURNS
ASPARSE MATRIX IN CSR FORMAT SHAPE NSAMPLES NSAMPLESFIT NSAMPLESFIT IS THE
NUMBER OF SAMPLES IN THE FITTED DATA AI J IS ASSIGNED THE WEIGHT OF EDGE THAT CONNECTS I
TO J
SEE ALSO
NEARESTNEIGHBORSRADIUSNEIGHBORSGRAPH
EXAMPLES
X 0 3 1
FROM SKLEARNNEIGHBORS IMPORT NEARESTNEIGHBORS
NEIGH NEARESTNEIGHBORSNNEIGHBORS2
NEIGHFITX
NEARESTNEIGHBORSALGORITHMAUTO LEAFSIZE30
A NEIGHKNEIGHBORSGRAPHX
ATOARRAY
ARRAY1 0 1
630SKLEARNNEIGHBORS NEAREST NEIGHBORS 2175
```

SCIKITLEARN USER GUIDE RELEASE 0213

0 1 1  
1 0 1

PREDICTSELF  
PREDICT THE TARGET FOR THE PROVIDED DATA

PARAMETERS  
XARRAYLIKE SHAPE NQUERY NFEATURES OR NQUERY NINDEXED IF METRIC 'PRECOMPUTED' TEST SAMPLES

RETURNS  
YARRAY OF INT SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TARGET VALUES

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{U}{V}$  WHERE U IS THE RESIDUAL SUM OF SQUARES  $Y_{TRUE} - Y_{PRED}$   $2SUM$  AND V IS THE TOTAL SUM OF SQUARES  $Y_{TRUE} - Y_{TRUEMEAN}$   $2SUM$  THE BEST POSSIBLE SCORE IS 10 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 00

PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES FOR SOME ESTIMATORS THIS MAY BE A PRECOMPUTED KERNEL MATRIX INSTEAD SHAPE NSAMPLES NSAMPLESFITTED WHERE NSAMPLESFITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE VALUES FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS  
SCORE FLOAT R2 OF SELF PREDICTX WRT Y

NOTES  
THE R2 SCORE USED WHEN CALLING SCORE ON A REGRESSOR WILL USE MULTIOUTPUTUNIFORMAVERAGE FROM VERSION 023 TO KEEP CONSISTENT WITH METRICSR2SCORE THIS WILL INFLUENCE THE SCORE METHOD OF ALL THE MULTIOUTPUT REGRESSORS EXCEPT FOR MULTIOUTPUTMULTIOUTPUTREGRESSOR TO SPECIFY THE DEFAULT VALUE MANUALLY AND AVOID THE WARNING PLEASE EITHER CALL METRICSR2SCORE DIRECTLY OR MAKE A CUSTOM SCORER WITH METRICSMAKESCORER THE BUILTIN SCORER R2 USES MULTIOUTPUTUNIFORMAVERAGE

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS  
SELF

EXAMPLES USING SKLEARNNEIGHBORSKNEIGHBORSREGRESSOR

- FACE COMPLETION WITH A MULTIOUTPUT ESTIMATORS

2176 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

- IMPUTING MISSING VALUES WITH VARIANTS OF ITERATIVEIMPUTER
- NEAREST NEIGHBORS REGRESSION

6307SKLEARNNEIGHBORS LOCALOUTLIERFACTOR

CLASSSKLEARNNEIGHBORS LOCALOUTLIERFACTOR NNEIGHBORS20 ALGORITHM'AUTO'

LEAFSIZE30 METRIC'MINKOWSKI' P2MET

RICPARAMSNONE CONTAMINATION'LEGACY'

NOVELTYFALSE NJOBSNONE

UNSUPERVISED OUTLIER DETECTION USING LOCAL OUTLIER FACTOR LOF

THE ANOMALY SCORE OF EACH SAMPLE IS CALLED LOCAL OUTLIER FACTOR IT MEASURES THE LOCAL DEVIATION OF DENSITY OF A GIVEN SAMPLE WITH RESPECT TO ITS NEIGHBORS IT IS LOCAL IN THAT THE ANOMALY SCORE DEPENDS ON HOW ISOLATED THE OBJECT IS WITH RESPECT TO THE SURROUNDING NEIGHBORHOOD MORE PRECISELY LOCALITY IS GIVEN BY KNEAREST NEIGHBORS WHOSE DISTANCE IS USED TO ESTIMATE THE LOCAL DENSITY BY COMPARING THE LOCAL DENSITY OF A SAMPLE TO THE LOCAL DENSITIES OF ITS NEIGHBORS ONE CAN IDENTIFY SAMPLES THAT HAVE A SUBSTANTIALLY LOWER DENSITY THAN THEIR NEIGHBORS THESE ARE CONSIDERED OUTLIERS

PARAMETERS

NNEIGHBORS INT OPTIONAL DEFAULT20 NUMBER OF NEIGHBORS TO USE BY DEFAULT FOR KNEIGHBORS QUERIES IF NNEIGHBORS IS LARGER THAN THE NUMBER OF SAMPLES PROVIDED ALL SAMPLES WILL BE USED

ALGORITHM 'AUTO' 'BALLTREE' 'KDTree' 'BRUTE' OPTIONAL ALGORITHM USED TO COMPUTE THE NEAREST NEIGHBORS

- 'BALLTREE' WILL USE BALLTREE
- 'KDTree' WILL USE KDTree
- 'BRUTE' WILL USE A BRUTEFORCE SEARCH
- 'AUTO' WILL ATTEMPT TO DECIDE THE MOST APPROPRIATE ALGORITHM BASED ON THE VALUES PASSED

TOFIT METHOD

NOTE FITTING ON SPARSE INPUT WILL OVERRIDE THE SETTING OF THIS PARAMETER USING BRUTE FORCE

LEAFSIZE INT OPTIONAL DEFAULT30 LEAF SIZE PASSED TO BALLTREE ORKDTree THIS CAN AFFECT THE SPEED OF THE CONSTRUCTION AND QUERY AS WELL AS THE MEMORY REQUIRED TO STORE THE TREE THE OPTIMAL VALUE DEPENDS ON THE NATURE OF THE PROBLEM

METRIC STRING OR CALLABLE DEFAULT 'MINKOWSKI' METRIC USED FOR THE DISTANCE COMPUTATION ANY METRIC FROM SCIKITLEARN OR SCIPYSPATIALDISTANCE CAN BE USED

IF 'PRECOMPUTED' THE TRAINING INPUT X IS EXPECTED TO BE A DISTANCE MATRIX

IF METRIC IS A CALLABLE FUNCTION IT IS CALLED ON EACH PAIR OF INSTANCES ROWS AND THE RESULTING VALUE RECORDED THE CALLABLE SHOULD TAKE TWO ARRAYS AS INPUT AND RETURN ONE VALUE INDICATING THE DISTANCE BETWEEN THEM THIS WORKS FOR SCIPY'S METRICS BUT IS LESS EFFICIENT THAN PASSING

THE METRIC NAME AS A STRING

VALID VALUES FOR METRIC ARE

- FROM SCIKITLEARN 'CITYBLOCK' 'COSINE' 'EUCLIDEAN' 'L1' 'L2' 'MANHATTAN'
- FROM SCIPYSPATIALDISTANCE 'BRAYCURTIS' 'CANBERRA' 'CHEBYSHEV' 'CORRELATION' 'DICE' 'HAMMING' 'JACCARD' 'KULSINSKI' 'MAHALANOBIS' 'MINKOWSKI' 'ROGERSTANIMOTO' 'RUSSELLRAO' 'SEUCLIDEAN' 'SOKALMICHENER' 'SOKALSNEATH' 'SQUEUCLIDEAN' 'YULE'

630SKLEARNNEIGHBORS NEAREST NEIGHBORS 2177

SCIKITLEARN USER GUIDE RELEASE 0213  
SEE THE DOCUMENTATION FOR SCIPYSPATIALDISTANCE FOR DETAILS ON THESE METRICS HTTPSDOCS  
SCIPYORGDOCS SCIPYREFERENCESPATIALDISTANCEHTML  
PINTEGER OPTIONAL DEFAULT2 PARAMETER FOR THE MINKOWSKI METRIC FROM SKLEARN  
METRICSPAIRWISEPAIRWISEDISTANCES WHEN P 1 THIS IS EQUIVALENT TO  
USING MANHATTANDISTANCE L1 AND EUCLIDEANDISTANCE L2 FOR P 2 FOR ARBITRARY P  
MINKOWSKIDISTANCE LP IS USED  
METRICPARAMS DICT OPTIONAL DEFAULTNONE ADDITIONAL KEYWORD ARGUMENTS FOR THE METRIC  
FUNCTION  
CONTAMINATION FLOAT IN 0 05 OPTIONAL DEFAULT01 THE AMOUNT OF CONTAMINATION OF THE  
DATA SET IE THE PROPORTION OF OUTLIERS IN THE DATA SET WHEN FITTING THIS IS USED TO DEFINE THE  
THRESHOLD ON THE DECISION FUNCTION IF “AUTO” THE DECISION FUNCTION THRESHOLD IS DETERMINED  
AS IN THE ORIGINAL PAPER  
CHANGED IN VERSION 020 THE DEFAULT VALUE OF CONTAMINATION WILL CHANGE FROM 01 IN  
020 TOAUTO IN 022  
NOVELTY BOOLEAN DEFAULT FALSE BY DEFAULT LOCALOUTLIERFACTOR IS ONLY MEANT TO BE USED FOR  
OUTLIER DETECTION NOVELTYFALSE SET NOVELTY TO TRUE IF YOU WANT TO USE LOCALOUTLIERFACTOR  
FOR NOVELTY DETECTION IN THIS CASE BE AWARE THAT THAT YOU SHOULD ONLY USE PREDICT DECI  
SIONFUNCTION AND SCORESAMPLES ON NEW UNSEEN DATA AND NOT ON THE TRAINING SET  
NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF PARALLEL JOBS TO RUN FOR NEIGH  
BORS SEARCH NONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1  
MEANS USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS AFFECTS ONLY KNEIGHBORS AND  
KNEIGHBORSGRAPH METHODS  
ATTRIBUTES  
NEGATIVEOUTLIERFACTOR NUMPY ARRAY SHAPE NSAMPLES THE OPPOSITE LOF OF THE TRAIN  
ING SAMPLES THE HIGHER THE MORE NORMAL INLIERS TEND TO HAVE A LOF SCORE CLOSE TO 1  
NEGATIVEOUTLIERFACTOR CLOSE TO 1 WHILE OUTLIERS TEND TO HAVE A LARGER LOF  
SCORE  
THE LOCAL OUTLIER FACTOR LOF OF A SAMPLE CAPTURES ITS SUPPOSED ‘DEGREE OF ABNORMALITY’  
IT IS THE AVERAGE OF THE RATIO OF THE LOCAL REACHABILITY DENSITY OF A SAMPLE AND THOSE OF ITS  
KNEAREST NEIGHBORS  
NNEIGHBORS INTEGER THE ACTUAL NUMBER OF NEIGHBORS USED FOR KNEIGHBORS QUERIES  
OFFSET FLOAT OFFSET USED TO OBTAIN BINARY LABELS FROM THE RAW SCORES OBSERVATIONS HAVING A  
NEGATIVEOUTLIERFACTOR SMALLER THAN OFFSET ARE DETECTED AS ABNORMAL THE OFFSET IS SET TO  
15 INLIERS SCORE AROUND 1 EXCEPT WHEN A CONTAMINATION PARAMETER DIFFERENT THAN “AUTO”  
IS PROVIDED IN THAT CASE THE OFFSET IS DEFINED IN SUCH A WAY WE OBTAIN THE EXPECTED NUMBER  
OF OUTLIERS IN TRAINING  
REFERENCES  
RCA479BB498411  
METHODS  
FITSELF X Y FIT THE MODEL USING X AS TRAINING DATA  
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR  
CONTINUED ON NEXT PAGE  
2178 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

TABLE 6230 – CONTINUED FROM PREVIOUS PAGE

KNEIGHBORS SELF X NNEIGHBORS FINDS THE KNEIGHBORS OF A POINT

KNEIGHBORSGRAPH SELF X NNEIGHBORS MODE COMPUTES THE WEIGHTED GRAPH OF KNEIGHBORS FOR POINTS IN X

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELFNNEIGHBORS20 ALGORITHM‘AUTO’ LEAFSIZE30 METRIC‘MINKOWSKI’ P2MET

RICPARAMSNONE CONTAMINATION‘LEGACY’ NOVELTYFALSE NJOBSNONE

DECISIONFUNCTION

SHIFTED OPPOSITE OF THE LOCAL OUTLIER FACTOR OF X

BIGGER IS BETTER IE LARGE VALUES CORRESPOND TO INLIERS

THE SHIFT OFFSET ALLOWS A ZERO THRESHOLD FOR BEING AN OUTLIER ONLY AVAILABLE FOR NOVELTY DETECTION WHEN

NOVELTY IS SET TO TRUE THE ARGUMENT X IS SUPPOSED TO CONTAIN NEW DATA IF X CONTAINS A POINT FROM

TRAINING IT CONSIDERS THE LATER IN ITS OWN NEIGHBORHOOD ALSO THE SAMPLES IN X ARE NOT CONSIDERED IN THE

NEIGHBORHOOD OF ANY POINT

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES THE QUERY SAMPLE OR SAMPLES TO COMPUTE THE

LOCAL OUTLIER FACTOR WRT THE TRAINING SAMPLES

RETURNS

SHIFTEDOPPOSITELOFSCORES ARRAY SHAPE NSAMPLES THE SHIFTED OPPOSITE OF THE LOCAL

OUTLIER FACTOR OF EACH INPUT SAMPLES THE LOWER THE MORE ABNORMAL NEGATIVE SCORES

REPRESENT OUTLIERS POSITIVE SCORES REPRESENT INLIERS

FITSELFXYNONE

FIT THE MODEL USING X AS TRAINING DATA

PARAMETERS

XARRAYLIKE SPARSE MATRIX BALLTREE KDTree TRAINING DATA IF ARRAY OR MATRIX SHAPE

NSAMPLES NFEATURES OR NSAMPLES NSAMPLES IF METRIC‘PRECOMPUTED’

YIGNORED NOT USED PRESENT FOR API CONSISTENCY BY CONVENTION

RETURNS

SELF OBJECT

FITPREDICT

“FITS THE MODEL TO THE TRAINING SET X AND RETURNS THE LABELS

LABEL IS 1 FOR AN INLIER AND 1 FOR AN OUTLIER ACCORDING TO THE LOF SCORE AND THE CONTAMINATION PARAMETER

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES DEFAULTNONE THE QUERY SAMPLE OR SAMPLES

TO COMPUTE THE LOCAL OUTLIER FACTOR WRT TO THE TRAINING SAMPLES

YIGNORED NOT USED PRESENT FOR API CONSISTENCY BY CONVENTION

RETURNS

ISINLIER ARRAY SHAPE NSAMPLES RETURNS 1 FOR ANOMALIESOUTLIERS AND 1 FOR INLIERS

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

630SKLEARNNEIGHBORS NEAREST NEIGHBORS 2179

SCIKITLEARN USER GUIDE RELEASE 0213

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

KNEIGHBORS SELFXTNONE NNEIGHBORSNONE RETURNDISTANCETRUE

FINDS THE KNEIGHBORS OF A POINT RETURNS INDICES OF AND DISTANCES TO THE NEIGHBORS OF EACH POINT

PARAMETERS

XARRAYLIKE SHAPE NQUERY NFEATURES OR NQUERY NINDEXED IF METRIC 'PRECOMPUTED' THE QUERY POINT OR POINTS IF NOT PROVIDED NEIGHBORS OF EACH INDEXED POINT ARE RETURNED IN THIS CASE THE QUERY POINT IS NOT CONSIDERED ITS OWN NEIGHBOR

NNEIGHBORS INT NUMBER OF NEIGHBORS TO GET DEFAULT IS THE VALUE PASSED TO THE CONSTRUCTOR

RETURNDISTANCE BOOLEAN OPTIONAL DEFAULTS TO TRUE IF FALSE DISTANCES WILL NOT BE RETURNED

RETURNS

DIST ARRAY ARRAY REPRESENTING THE LENGTHS TO POINTS ONLY PRESENT IF RETURNDISTANCETRUE

IND ARRAY INDICES OF THE NEAREST POINTS IN THE POPULATION MATRIX

EXAMPLES

IN THE FOLLOWING EXAMPLE WE CONSTRUCT A NEIGHBORSCLASSIFIER CLASS FROM AN ARRAY REPRESENTING OUR DATA SET AND ASK WHO'S THE CLOSEST POINT TO 111

SAMPLES 0 0 0 0 5 0 1 1 5

FROM SKLEARNNEIGHBORS IMPORT NEARESTNEIGHBORS

NEIGH NEARESTNEIGHBORSNNEIGHBORS1

NEIGHFITSAMPLES

NEARESTNEIGHBORSALGORITHM AUTO LEAF SIZE30

PRINTNEIGHKNEIGHBORS1 1 1

ARRAY05 ARRAY2

AS YOU CAN SEE IT RETURNS 05 AND 2 WHICH MEANS THAT THE ELEMENT IS AT DISTANCE 05 AND IS THE THIRD ELEMENT OF SAMPLES INDEXES START AT 0 YOU CAN ALSO QUERY FOR MULTIPLE POINTS

X 0 1 0 1 0 1

NEIGHKNEIGHBORSX RETURNDISTANCE FALSE

ARRAY1

2

KNEIGHBORSGRAPH SELFXTNONE NNEIGHBORSNONE MODE'CONNECTIVITY'

COMPUTES THE WEIGHTED GRAPH OF KNEIGHBORS FOR POINTS IN X

PARAMETERS

XARRAYLIKE SHAPE NQUERY NFEATURES OR NQUERY NINDEXED IF METRIC 'PRECOMPUTED' THE QUERY POINT OR POINTS IF NOT PROVIDED NEIGHBORS OF EACH INDEXED POINT ARE RETURNED IN THIS CASE THE QUERY POINT IS NOT CONSIDERED ITS OWN NEIGHBOR

NNEIGHBORS INT NUMBER OF NEIGHBORS FOR EACH SAMPLE DEFAULT IS VALUE PASSED TO THE CONSTRUCTOR

2180 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

MODE 'CONNECTIVITY' 'DISTANCE' OPTIONAL TYPE OF RETURNED MATRIX 'CONNECTIVITY' WILL RETURN THE CONNECTIVITY MATRIX WITH ONES AND ZEROS IN 'DISTANCE' THE EDGES ARE EUCLIDEAN DISTANCE BETWEEN POINTS

RETURNS

SPARSE MATRIX IN CSR FORMAT SHAPE (NSAMPLES, NSAMPLES) FIT IS THE NUMBER OF SAMPLES IN THE FITTED DATA A[I, J] IS ASSIGNED THE WEIGHT OF EDGE THAT CONNECTS I TO J

SEE ALSO

NEARESTNEIGHBORSRADIUSNEIGHBORSGRAPH

EXAMPLES

```
X = [[0, 3, 1],
      [1, 0, 1]]
from sklearn.neighbors import NearestNeighbors
neigh = NearestNeighbors(n_neighbors=2)
neigh.fit(X)
nearest_neighbors_algorithm = NearestNeighbors(n_neighbors=2,
                                              algorithm='auto',
                                              leaf_size=30)
A = neigh.kneighbors_graph(X)
A.toarray()
array([[0, 1],
       [1, 0],
       [1, 0]])
```

PREDICT

PREDICT THE LABELS 1 INLIER 1 OUTLIER OF X ACCORDING TO LOF

THIS METHOD ALLOWS TO GENERALIZE PREDICTION TO NEW OBSERVATIONS NOT IN THE TRAINING SET ONLY AVAILABLE FOR NOVELTY DETECTION WHEN NOVELTY IS SET TO TRUE

PARAMETERS

X: ARRAY-LIKE SHAPE (NSAMPLES, NFEATURES) THE QUERY SAMPLE OR SAMPLES TO COMPUTE THE LOCAL OUTLIER FACTOR WRT TO THE TRAINING SAMPLES

RETURNS

IS\_INLIER: ARRAY SHAPE (NSAMPLES) RETURNS 1 FOR ANOMALIES/OUTLIERS AND 0 FOR INLIERS

SCORES: SAMPLES OPPOSITE OF THE LOCAL OUTLIER FACTOR OF X

IT IS THE OPPOSITE AS BIGGER IS BETTER IE LARGE VALUES CORRESPOND TO INLIERS

ONLY AVAILABLE FOR NOVELTY DETECTION WHEN NOVELTY IS SET TO TRUE THE ARGUMENT X IS SUPPOSED TO CONTAIN NEW DATA IF X CONTAINS A POINT FROM TRAINING IT CONSIDERS THE LOCAL IN ITS OWN NEIGHBORHOOD ALSO THE SAMPLES IN X ARE NOT CONSIDERED IN THE NEIGHBORHOOD OF ANY POINT THE SCORES ON TRAINING DATA IS AVAILABLE BY CONSIDERING THE NEGATIVE OUTLIER FACTOR ATTRIBUTE

PARAMETERS

X: ARRAY-LIKE SHAPE (NSAMPLES, NFEATURES) THE QUERY SAMPLE OR SAMPLES TO COMPUTE THE LOCAL OUTLIER FACTOR WRT THE TRAINING SAMPLES

RETURNS

OPPOSITE\_OF\_SCORES: ARRAY SHAPE (NSAMPLES) THE OPPOSITE OF THE LOCAL OUTLIER FACTOR OF EACH INPUT SAMPLES THE LOWER THE MORE ABNORMAL

630SKLEARNNEIGHBORS NEAREST NEIGHBORS 2181

SCIKITLEARN USER GUIDE RELEASE 0213

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

EXAMPLES USING SKLEARNNEIGHBORSLOCALOUTLIERFACTOR

- COMPARING ANOMALY DETECTION ALGORITHMS FOR OUTLIER DETECTION ON TOY DATASETS
- OUTLIER DETECTION WITH LOCAL OUTLIER FACTOR LOF
- NOVELTY DETECTION WITH LOCAL OUTLIER FACTOR LOF

6308SKLEARNNEIGHBORS RADIUSNEIGHBORSCLASSIFIER

CLASSSSKLEARNNEIGHBORS RADIUSNEIGHBORSCLASSIFIER RADIUS10 WEIGHTS'UNIFORM' ALGO

RITHM'AUTO' LEAFSIZE30 P2MET

RIC'MINKOWSKI' OUTLIERLABELNONE

METRICPARAMSNONE NJOBSNONE

KWARGS

CLASSIFIER IMPLEMENTING A VOTE AMONG NEIGHBORS WITHIN A GIVEN RADIUS

READ MORE IN THE USER GUIDE

PARAMETERS

RADIUS FLOAT OPTIONAL DEFAULT 10 RANGE OF PARAMETER SPACE TO USE BY DEFAULT FOR

RADIUSNEIGHBORS QUERIES

WEIGHTS STR OR CALLABLE WEIGHT FUNCTION USED IN PREDICTION POSSIBLE VALUES

- 'UNIFORM' UNIFORM WEIGHTS ALL POINTS IN EACH NEIGHBORHOOD ARE WEIGHTED EQUALLY
- 'DISTANCE' WEIGHT POINTS BY THE INVERSE OF THEIR DISTANCE IN THIS CASE CLOSER NEIGHBORS OF A QUERY POINT WILL HAVE A GREATER INFLUENCE THAN NEIGHBORS WHICH ARE FURTHER AWAY
- CALLABLE A USERDEFINED FUNCTION WHICH ACCEPTS AN ARRAY OF DISTANCES AND RETURNS AN ARRAY OF THE SAME SHAPE CONTAINING THE WEIGHTS

UNIFORM WEIGHTS ARE USED BY DEFAULT

ALGORITHM 'AUTO' 'BALLTREE' 'KDTREE' 'BRUTE' OPTIONAL ALGORITHM USED TO COMPUTE THE NEAREST NEIGHBORS

- 'BALLTREE' WILL USE BALLTREE
- 'KDTREE' WILL USE KDTREE
- 'BRUTE' WILL USE A BRUTEFORCE SEARCH
- 'AUTO' WILL ATTEMPT TO DECIDE THE MOST APPROPRIATE ALGORITHM BASED ON THE VALUES PASSED

TOFIT METHOD

NOTE FITTING ON SPARSE INPUT WILL OVERRIDE THE SETTING OF THIS PARAMETER USING BRUTE FORCE

2182 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

LEAFSIZE INT OPTIONAL DEFAULT 30 LEAF SIZE PASSED TO BALLTREE OR KDTree THIS CAN AFFECT THE SPEED OF THE CONSTRUCTION AND QUERY AS WELL AS THE MEMORY REQUIRED TO STORE THE TREE THE OPTIMAL VALUE DEPENDS ON THE NATURE OF THE PROBLEM

PINTEGRER OPTIONAL DEFAULT 2 POWER PARAMETER FOR THE MINKOWSKI METRIC WHEN P 1 THIS IS EQUIVALENT TO USING MANHATTANDISTANCE L1 AND EUCLIDEANDISTANCE L2 FOR P 2 FOR ARBITRARY P MINKOWSKIDISTANCE LP IS USED

METRIC STRING OR CALLABLE DEFAULT 'MINKOWSKI' THE DISTANCE METRIC TO USE FOR THE TREE THE DEFAULT METRIC IS MINKOWSKI AND WITH P2 IS EQUIVALENT TO THE STANDARD EUCLIDEAN METRIC SEE THE DOCUMENTATION OF THE DISTANCEMETRIC CLASS FOR A LIST OF AVAILABLE METRICS

OUTLIERLABEL INT OPTIONAL DEFAULT NONE LABEL WHICH IS GIVEN FOR OUTLIER SAMPLES SAMPLES WITH NO NEIGHBORS ON GIVEN RADIUS IF SET TO NONE VALUEERROR IS RAISED WHEN OUTLIER IS DETECTED

METRICPARAMS DICT OPTIONAL DEFAULT NONE ADDITIONAL KEYWORD ARGUMENTS FOR THE METRIC FUNCTION

NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF PARALLEL JOBS TO RUN FOR NEIGHBORS SEARCHNONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

SEE ALSO

KNEIGHBORSCLASSIFIER

RADIUSNEIGHBORSREGRESSOR

KNEIGHBORSREGRESSOR

NEARESTNEIGHBORS

NOTES

SEE NEAREST NEIGHBORS IN THE ONLINE DOCUMENTATION FOR A DISCUSSION OF THE CHOICE OF ALGORITHM AND LEAFSIZE

[HTTPSENWIKIPEDIAORGWIKIKNEARESTNEIGHBORALGORITHM](https://en.wikipedia.org/wiki/K-nearest_neighbors_algorithm)

EXAMPLES

```
X 0 1 2 3
Y 0 0 1 1
```

```
FROM SKLEARNNEIGHBORS IMPORT RADIUSNEIGHBORSCLASSIFIER
NEIGH RADIUSNEIGHBORSCLASSIFIERRADIUS10
NEIGHFITX Y
RADIUSNEIGHBORSCLASSIFIER
PRINTNEIGHPREDICT15
0
METHODS
630SKLEARNNEIGHBORS NEAREST NEIGHBORS 2183
```

SCIKITLEARN USER GUIDE RELEASE 0213

FITSELF X Y FIT THE MODEL USING X AS TRAINING DATA AND Y AS TARGET VALUES

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PREDICT SELF X PREDICT THE CLASS LABELS FOR THE PROVIDED DATA

RADIUSNEIGHBORS SELF X RADIUS FINDS THE NEIGHBORS WITHIN A GIVEN RADIUS OF A POINT OR POINTS

RADIUSNEIGHBORSGRAPH SELF X RADIUS

MODECOMPUTES THE WEIGHTED GRAPH OF NEIGHBORS FOR POINTS

IN X

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELF RADIUS10 WEIGHTS'UNIFORM' ALGORITHM'AUTO' LEAFSIZE30 P2METRIC'MINKOWSKI' OUTLIERLABELNONE METRICPARAMSNONE NJOBSNONE KWARGS

FITSELFXY

FIT THE MODEL USING X AS TRAINING DATA AND Y AS TARGET VALUES

PARAMETERS

XARRAYLIKE SPARSE MATRIX BALLTREE KDTree TRAINING DATA IF ARRAY OR MATRIX SHAPE

NSAMPLES NFEATURES OR NSAMPLES NSAMPLES IF METRIC'PRECOMPUTED'

YARRAYLIKE SPARSE MATRIX TARGET VALUES OF SHAPE NSAMPLES OR NSAMPLES

NOUTPUTS

GETPARAMS SELFDEEPTREE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELFXY

PREDICT THE CLASS LABELS FOR THE PROVIDED DATA

PARAMETERS

XARRAYLIKE SHAPE NQUERY NFEATURES OR NQUERY NINDEXED IF METRIC 'PRECOMPUTED' TEST SAMPLES

RETURNS

YARRAY OF SHAPE NSAMPLES OR NSAMPLES NOUTPUTS CLASS LABELS FOR EACH DATA SAMPLE

RADIUSNEIGHBORS SELF X NONE RADIUS NONE RETURN DISTANCE TRUE

FINDS THE NEIGHBORS WITHIN A GIVEN RADIUS OF A POINT OR POINTS

RETURN THE INDICES AND DISTANCES OF EACH POINT FROM THE DATASET LYING IN A BALL WITH SIZE RADIUS AROUND

THE POINTS OF THE QUERY ARRAY POINTS LYING ON THE BOUNDARY ARE INCLUDED IN THE RESULTS

THE RESULT POINTS ARE NOT NECESSARILY SORTED BY DISTANCE TO THEIR QUERY POINT

PARAMETERS

2184 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

XARRAYLIKE NSAMPLES NFEATURES OPTIONAL THE QUERY POINT OR POINTS IF NOT PROVIDED  
NEIGHBORS OF EACH INDEXED POINT ARE RETURNED IN THIS CASE THE QUERY POINT IS NOT CONSIDERED  
ITS OWN NEIGHBOR  
RADIUS FLOAT LIMITING DISTANCE OF NEIGHBORS TO RETURN DEFAULT IS THE VALUE PASSED TO THE  
CONSTRUCTOR  
RETURNDISTANCE BOOLEAN OPTIONAL DEFAULTS TO TRUE IF FALSE DISTANCES WILL NOT BE RE  
TURNED  
RETURNS  
DIST ARRAY SHAPE NSAMPLES OF ARRAYS ARRAY REPRESENTING THE DISTANCES TO EACH POINT  
ONLY PRESENT IF RETURNDISTANCETRUE THE DISTANCE VALUES ARE COMPUTED ACCORDING TO THE  
METRIC CONSTRUCTOR PARAMETER  
IND ARRAY SHAPE NSAMPLES OF ARRAYS AN ARRAY OF ARRAYS OF INDICES OF THE APPROXIMATE  
NEAREST POINTS FROM THE POPULATION MATRIX THAT LIE WITHIN A BALL OF SIZE RADIUS AROUND THE  
QUERY POINTS  
NOTES  
BECAUSE THE NUMBER OF NEIGHBORS OF EACH POINT IS NOT NECESSARILY EQUAL THE RESULTS FOR MULTIPLE QUERY  
POINTS CANNOT BE FIT IN A STANDARD DATA ARRAY FOR EFFICIENCY RADIUSNEIGHBORS RETURNS ARRAYS OF  
OBJECTS WHERE EACH OBJECT IS A 1D ARRAY OF INDICES OR DISTANCES  
EXAMPLES  
IN THE FOLLOWING EXAMPLE WE CONSTRUCT A NEIGHBORSCLASSIFIER CLASS FROM AN ARRAY REPRESENTING OUR DATA SET  
AND ASK WHO'S THE CLOSEST POINT TO 1 1 1  
IMPORT NUMPY AS NP  
SAMPLES 0 0 0 0 5 0 1 1 5  
FROM SKLEARNNEIGHBORS IMPORT NEARESTNEIGHBORS  
NEIGH NEARESTNEIGHBORSRADIUS16  
NEIGHFITSAMPLES  
NEARESTNEIGHBORSALGORITHMAUTO LEAFSIZE30  
RNG NEIGHRADIUSNEIGHBORS1 1 1  
PRINTNPASARRAYRNG00  
15 05  
PRINTNPASARRAYRNG10  
1 2  
THE FIRST ARRAY RETURNED CONTAINS THE DISTANCES TO ALL POINTS WHICH ARE CLOSER THAN 16 WHILE THE SECOND  
ARRAY RETURNED CONTAINS THEIR INDICES IN GENERAL MULTIPLE POINTS CAN BE QUERIED AT THE SAME TIME  
RADIUSNEIGHBORSGRAPH SELFXTNONE RADIUSNONE MODE'CONNECTIVITY'  
COMPUTES THE WEIGHTED GRAPH OF NEIGHBORS FOR POINTS IN X  
NEIGHBORHOODS ARE RESTRICTED THE POINTS AT A DISTANCE LOWER THAN RADIUS  
PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES OPTIONAL THE QUERY POINT OR POINTS IF NOT  
PROVIDED NEIGHBORS OF EACH INDEXED POINT ARE RETURNED IN THIS CASE THE QUERY POINT IS NOT  
CONSIDERED ITS OWN NEIGHBOR  
RADIUS FLOAT RADIUS OF NEIGHBORHOODS DEFAULT IS THE VALUE PASSED TO THE CONSTRUCTOR  
630SKLEARNNEIGHBORS NEAREST NEIGHBORS 2185

SCIKITLEARN USER GUIDE RELEASE 0213

MODE ‘CONNECTIVITY’ ‘DISTANCE’ OPTIONAL TYPE OF RETURNED MATRIX ‘CONNECTIVITY’ WILL RETURN THE CONNECTIVITY MATRIX WITH ONES AND ZEROS IN ‘DISTANCE’ THE EDGES ARE EUCLIDEAN DISTANCE BETWEEN POINTS

RETURNS

SPARSE MATRIX IN CSR FORMAT SHAPE (NSAMPLES, NSAMPLES) A[i, j] IS ASSIGNED THE WEIGHT OF EDGE THAT CONNECTS i TO j

SEE ALSO

KNEIGHBORSGRAPH

EXAMPLES

```
X = [[0, 3, 1],
      [1, 0, 0],
      [0, 1, 0]]
from sklearn.neighbors import NearestNeighbors
neigh = NearestNeighbors(radius=15)
neigh.fit(X)
nearest_neighbors_algorithm = NearestNeighbors(radius=15)
A = neigh.radius_neighbors_graph(X)
A.toarray()
array([[1, 0, 1],
       [0, 1, 0],
       [1, 0, 1]])
```

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS

IN MULTILABEL CLASSIFICATION THIS IS THE SUBSET ACCURACY WHICH IS A HARSH METRIC SINCE YOU REQUIRE FOR EACH SAMPLE THAT EACH LABEL SET BE CORRECTLY PREDICTED

PARAMETERS

X: ARRAYLIKE SHAPE (NSAMPLES, NFEATURES) TEST SAMPLES

Y: ARRAYLIKE SHAPE (NSAMPLES, OR NSAMPLES) OUTPUTS TRUE LABELS FOR X

SAMPLEWEIGHT: ARRAYLIKE SHAPE (NSAMPLES) OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE: FLOAT MEAN ACCURACY OF SELF-PREDICT X WRT Y

SETPARAMS: SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENT.PARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

2186 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

63095KLEARNNEIGHBORS RADIUSNEIGHBORSREGRESSOR

CLASSSSKLEARNNEIGHBORS RADIUSNEIGHBORSREGRESSOR RADIUS10 WEIGHTS'UNIFORM' ALGO

RITHM'AUTO' LEAFSIZE30 P2MET

RIC'MINKOWSKI' METRICPARAMSNONE

NJOBSNONE KWARGS

REGRESSION BASED ON NEIGHBORS WITHIN A FIXED RADIUS

THE TARGET IS PREDICTED BY LOCAL INTERPOLATION OF THE TARGETS ASSOCIATED OF THE NEAREST NEIGHBORS IN THE TRAINING SET

READ MORE IN THE USER GUIDE

PARAMETERS

RADIUS FLOAT OPTIONAL DEFAULT 10 RANGE OF PARAMETER SPACE TO USE BY DEFAULT FOR

RADIUSNEIGHBORS QUERIES

WEIGHTS STR OR CALLABLE WEIGHT FUNCTION USED IN PREDICTION POSSIBLE VALUES

- 'UNIFORM' UNIFORM WEIGHTS ALL POINTS IN EACH NEIGHBORHOOD ARE WEIGHTED EQUALLY
- 'DISTANCE' WEIGHT POINTS BY THE INVERSE OF THEIR DISTANCE IN THIS CASE CLOSER NEIGHBORS OF A QUERY POINT WILL HAVE A GREATER INFLUENCE THAN NEIGHBORS WHICH ARE FURTHER AWAY
- CALLABLE A USERDEFINED FUNCTION WHICH ACCEPTS AN ARRAY OF DISTANCES AND RETURNS AN ARRAY OF THE SAME SHAPE CONTAINING THE WEIGHTS

UNIFORM WEIGHTS ARE USED BY DEFAULT

ALGORITHM 'AUTO' 'BALLTREE' 'KDTree' 'BRUTE' OPTIONAL ALGORITHM USED TO COMPUTE THE NEAREST NEIGHBORS

- 'BALLTREE' WILL USE BALLTREE
- 'KDTree' WILL USE KDTree
- 'BRUTE' WILL USE A BRUTEFORCE SEARCH
- 'AUTO' WILL ATTEMPT TO DECIDE THE MOST APPROPRIATE ALGORITHM BASED ON THE VALUES PASSED

TOFIT METHOD

NOTE FITTING ON SPARSE INPUT WILL OVERRIDE THE SETTING OF THIS PARAMETER USING BRUTE FORCE

LEAFSIZE INT OPTIONAL DEFAULT 30 LEAF SIZE PASSED TO BALLTREE OR KDTree THIS CAN AFFECT THE SPEED OF THE CONSTRUCTION AND QUERY AS WELL AS THE MEMORY REQUIRED TO STORE THE TREE

THE OPTIMAL VALUE DEPENDS ON THE NATURE OF THE PROBLEM

PINTEGER OPTIONAL DEFAULT 2 POWER PARAMETER FOR THE MINKOWSKI METRIC WHEN P 1 THIS IS EQUIVALENT TO USING MANHATTANDISTANCE L1 AND EUCLIDEANDISTANCE L2 FOR P 2 FOR ARBITRARY P MINKOWSKIDISTANCE LP IS USED

METRIC STRING OR CALLABLE DEFAULT 'MINKOWSKI' THE DISTANCE METRIC TO USE FOR THE TREE THE DEFAULT METRIC IS MINKOWSKI AND WITH P2 IS EQUIVALENT TO THE STANDARD EUCLIDEAN METRIC

SEE THE DOCUMENTATION OF THE DISTANCEMETRIC CLASS FOR A LIST OF AVAILABLE METRICS

METRICPARAMS DICT OPTIONAL DEFAULT NONE ADDITIONAL KEYWORD ARGUMENTS FOR THE METRIC FUNCTION

NJOBS INT OR NONE OPTIONAL DEFAULTNONE

THE NUMBER OF PARALLEL JOBS TO RUN FOR NEIGHBORS SEARCH NONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT

1MEANS USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

SEE ALSO

6305KLEARNNEIGHBORS NEAREST NEIGHBORS 2187

SCIKITLEARN USER GUIDE RELEASE 0213

NEARESTNEIGHBORS

KNEIGHBORSREGRESSOR

KNEIGHBORSCLASSIFIER

RADIUSNEIGHBORSCLASSIFIER

NOTES

SEE NEAREST NEIGHBORS IN THE ONLINE DOCUMENTATION FOR A DISCUSSION OF THE CHOICE OF ALGORITHM AND LEAFSIZE

[HTTPSENWIKIPEDIAORGWIKIKNEARESTNEIGHBORALGORITHM](http://en.wikipedia.org/wiki/K-nearest_neighbors_algorithm)

EXAMPLES

```
X 0 1 2 3
Y 0 0 1 1
```

```
FROM SKLEARNNEIGHBORS IMPORT RADIUSNEIGHBORSREGRESSOR
NEIGH RADIUSNEIGHBORSREGRESSORRADIUS10
NEIGHFITX Y
RADIUSNEIGHBORSREGRESSOR
PRINTNEIGHPREDICT15
05
```

METHODS

FITSELF X Y FIT THE MODEL USING X AS TRAINING DATA AND Y AS TARGET VALUES

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PREDICT SELF X PREDICT THE TARGET FOR THE PROVIDED DATA

RADIUSNEIGHBORS SELF X RADIUS FINDS THE NEIGHBORS WITHIN A GIVEN RADIUS OF A POINT OR POINTS

RADIUSNEIGHBORSGRAPH SELF X RADIUS

MODECOMPUTES THE WEIGHTED GRAPH OF NEIGHBORS FOR POINTS

IN X

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELF RADIUS10 WEIGHTS'UNIFORM' ALGORITHM'AUTO' LEAFSIZE30 P2METRIC'MINKOWSKI' METRICPARAMSNONE NJOBSNONE KWARGS

FITSELFXY

FIT THE MODEL USING X AS TRAINING DATA AND Y AS TARGET VALUES

PARAMETERS

XARRAYLIKE SPARSE MATRIX BALLTREE KDTree TRAINING DATA IF ARRAY OR MATRIX SHAPE

NSAMPLES NFEATURES OR NSAMPLES NSAMPLES IF METRIC'PRECOMPUTED'

YARRAYLIKE SPARSE MATRIX

TARGET VALUES ARRAY OF FLOAT VALUES SHAPE NSAMPLES OR NSAMPLES NOUTPUTS

2188 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

GETPARAMS SELFDEEPTURE  
GET PARAMETERS FOR THIS ESTIMATOR  
PARAMETERS  
DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBOBJECTS THAT ARE ESTIMATORS  
RETURNS  
PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES  
PREDICTSELF  
PREDICT THE TARGET FOR THE PROVIDED DATA  
PARAMETERS  
XARRAYLIKE SHAPE NQUERY NFEATURES OR NQUERY NINDEXED IF METRIC 'PRECOM  
PUTED' TEST SAMPLES  
RETURNS  
YARRAY OF FLOAT SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TARGET VALUES  
RADIUSNEIGHBORS SELFXTNONE RADIUSNONE RETURNDISTANCETURE  
FINDS THE NEIGHBORS WITHIN A GIVEN RADIUS OF A POINT OR POINTS  
RETURN THE INDICES AND DISTANCES OF EACH POINT FROM THE DATASET LYING IN A BALL WITH SIZE RADIUS AROUND  
THE POINTS OF THE QUERY ARRAY POINTS LYING ON THE BOUNDARY ARE INCLUDED IN THE RESULTS  
THE RESULT POINTS ARE NOTNECESSARILY SORTED BY DISTANCE TO THEIR QUERY POINT  
PARAMETERS  
XARRAYLIKE NSAMPLES NFEATURES OPTIONAL THE QUERY POINT OR POINTS IF NOT PROVIDED  
NEIGHBORS OF EACH INDEXED POINT ARE RETURNED IN THIS CASE THE QUERY POINT IS NOT CONSIDERED  
ITS OWN NEIGHBOR  
RADIUS FLOAT LIMITING DISTANCE OF NEIGHBORS TO RETURN DEFAULT IS THE VALUE PASSED TO THE  
CONSTRUCTOR  
RETURNDISTANCE BOOLEAN OPTIONAL DEFAULTS TO TRUE IF FALSE DISTANCES WILL NOT BE RE  
TURNED  
RETURNS  
DIST ARRAY SHAPE NSAMPLES OF ARRAYS ARRAY REPRESENTING THE DISTANCES TO EACH POINT  
ONLY PRESENT IF RETURNDISTANCETURE THE DISTANCE VALUES ARE COMPUTED ACCORDING TO THE  
METRIC CONSTRUCTOR PARAMETER  
IND ARRAY SHAPE NSAMPLES OF ARRAYS AN ARRAY OF ARRAYS OF INDICES OF THE APPROXIMATE  
NEAREST POINTS FROM THE POPULATION MATRIX THAT LIE WITHIN A BALL OF SIZE RADIUS AROUND THE  
QUERY POINTS  
NOTES  
BECAUSE THE NUMBER OF NEIGHBORS OF EACH POINT IS NOT NECESSARILY EQUAL THE RESULTS FOR MULTIPLE QUERY  
POINTS CANNOT BE FIT IN A STANDARD DATA ARRAY FOR EFFICIENCY RADIUSNEIGHBORS RETURNS ARRAYS OF  
OBJECTS WHERE EACH OBJECT IS A 1D ARRAY OF INDICES OR DISTANCES  
6305KLEARNNEIGHBORS NEAREST NEIGHBORS 2189

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES  
IN THE FOLLOWING EXAMPLE WE CONSTRUCT A NEIGHBORSCLASSIFIER CLASS FROM AN ARRAY REPRESENTING OUR DATA SET  
AND ASK WHO'S THE CLOSEST POINT TO 1 1 1

```
import numpy as np
samples = 0 0 0 0 5 0 1 1 5
from sklearn.neighbors import NearestNeighbors
neigh = NearestNeighbors(radius=16)
neigh.fit(samples)
nearest_neighbors_algorithm = LeafSize(30)
rng = NeighRadiusNeighbors(1, 1)
print(np.asarray(rng[0][15, 0:5]))
print(np.asarray(rng[10][1, 2]))
```

THE FIRST ARRAY RETURNED CONTAINS THE DISTANCES TO ALL POINTS WHICH ARE CLOSER THAN 16 WHILE THE SECOND  
ARRAY RETURNED CONTAINS THEIR INDICES IN GENERAL MULTIPLE POINTS CAN BE QUERIED AT THE SAME TIME  
RADIUSNEIGHBORSGRAPH SELFKNONE RADIUSNONE MODE'CONNECTIVITY'  
COMPUTES THE WEIGHTED GRAPH OF NEIGHBORS FOR POINTS IN X  
NEIGHBORHOODS ARE RESTRICTED THE POINTS AT A DISTANCE LOWER THAN RADIUS  
PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES OPTIONAL THE QUERY POINT OR POINTS IF NOT  
PROVIDED NEIGHBORS OF EACH INDEXED POINT ARE RETURNED IN THIS CASE THE QUERY POINT IS NOT  
CONSIDERED ITS OWN NEIGHBOR  
RADIUS FLOAT RADIUS OF NEIGHBORHOODS DEFAULT IS THE VALUE PASSED TO THE CONSTRUCTOR  
MODE 'CONNECTIVITY' 'DISTANCE' OPTIONAL TYPE OF RETURNED MATRIX 'CONNECTIVITY' WILL  
RETURN THE CONNECTIVITY MATRIX WITH ONES AND ZEROS IN 'DISTANCE' THE EDGES ARE EUCLIDEAN  
DISTANCE BETWEEN POINTS

RETURNS  
ASPARSE MATRIX IN CSR FORMAT SHAPE NSAMPLES NSAMPLES AI J IS ASSIGNED THE  
WEIGHT OF EDGE THAT CONNECTS I TO J

SEE ALSO  
KNEIGHBORSGRAPH  
EXAMPLES

```
x = 0 3 1
from sklearn.neighbors import NearestNeighbors
neigh = NearestNeighbors(radius=15)
neigh.fit(x)
nearest_neighbors_algorithm = LeafSize(30)
a = NeighRadiusNeighborsGraph(x)
a.toarray()
array([0 1
       1 0
       1 0])
```

SCIKITLEARN USER GUIDE RELEASE 0213

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{U}{V}$  WHERE U IS THE RESIDUAL SUM OF SQUARES  $\sum (y_{true} - y_{pred})^2$  AND V IS THE TOTAL SUM OF SQUARES  $\sum (y_{true} - y_{true\_mean})^2$  THE BEST POSSIBLE SCORE IS 1.0 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 0.0

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES FOR SOME ESTIMATORS THIS MAY BE A PRECOMPUTED KERNEL MATRIX INSTEAD SHAPE NSAMPLES NSAMPLESFITTED WHERE NSAMPLESFITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE VALUES FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT R2 OF SELF PREDICT X WRT Y

NOTES

THE R2 SCORE USED WHEN CALLING SCORE ON A REGRESSOR WILL USE MULTIOUTPUTUNIFORMAVERAGE FROM VERSION 0.23 TO KEEP CONSISTENT WITH METRICS R2 SCORE THIS WILL INFLUENCE THE SCORE METHOD OF ALL THE MULTIOUTPUT REGRESSORS EXCEPT FOR MULTIOUTPUTMULTIOUTPUTREGRESSOR TO SPECIFY THE DEFAULT VALUE MANUALLY AND AVOID THE WARNING PLEASE EITHER CALL METRICS R2 SCORE DIRECTLY OR MAKE A CUSTOM SCORER WITH METRICS MAKE SCORER THE BUILT IN SCORER R2 USES MULTIOUTPUTUNIFORMAVERAGE

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENT PARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

63010 SKLEARN NEIGHBORS NEAREST CENTROID

CLASS SKLEARN NEIGHBORS NEAREST CENTROID METRIC 'EUCLIDEAN' SHRINK THRESHOLD NONE

NEAREST CENTROID CLASSIFIER

EACH CLASS IS REPRESENTED BY ITS CENTROID WITH TEST SAMPLES CLASSIFIED TO THE CLASS WITH THE NEAREST CENTROID

READ MORE IN THE USER GUIDE

PARAMETERS

METRIC STRING OR CALLABLE THE METRIC TO USE WHEN CALCULATING DISTANCE BETWEEN INSTANCES IN A FEATURE ARRAY IF METRIC IS A STRING OR CALLABLE IT MUST BE ONE OF THE OPTIONS ALLOWED BY METRICS PAIRWISE DISTANCES FOR ITS METRIC PARAMETER THE CENTROIDS FOR THE SAMPLES CORRESPONDING TO EACH CLASS IS THE POINT FROM WHICH THE SUM OF THE DISTANCES ACCORDING TO THE METRIC OF ALL SAMPLES THAT BELONG TO THAT PARTICULAR CLASS ARE MINIMIZED IF THE "MANHATTAN" METRIC IS PROVIDED THIS CENTROID IS THE MEDIAN AND FOR ALL OTHER METRICS THE CENTROID IS NOW SET TO BE THE MEAN

630 SKLEARN NEIGHBORS NEAREST NEIGHBORS 2191

SCIKITLEARN USER GUIDE RELEASE 0213

SHRINKTHRESHOLD FLOAT OPTIONAL DEFAULT NONE THRESHOLD FOR SHRINKING CENTROIDS TO RE  
MOVE FEATURES

ATTRIBUTES

CENTROIDS ARRAYLIKE SHAPE NCLASSES NFEATURES CENTROID OF EACH CLASS

SEE ALSO

SKLEARNNEIGHBORSKNEIGHBORSCLASSIFIER NEAREST NEIGHBORS CLASSIFIER

NOTES

WHEN USED FOR TEXT CLASSIFICATION WITH TFIDF VECTORS THIS CLASSIFIER IS ALSO KNOWN AS THE ROCCHIO CLASSIFIER

REFERENCES

TIBSHIRANI R HASTIE T NARASIMHAN B CHU G 2002 DIAGNOSIS OF MULTIPLE CANCER TYPES BY SHRUNKEN  
CENTROIDS OF GENE EXPRESSION PROCEEDINGS OF THE NATIONAL ACADEMY OF SCIENCES OF THE UNITED STATES OF AMERICA  
9910 65676572 THE NATIONAL ACADEMY OF SCIENCES

EXAMPLES

```
FROM SKLEARNNEIGHBORSNEARESTCENTROID IMPORT NEARESTCENTROID
IMPORT NUMPY AS NP
X NPARRAY1 1 2 1 3 2 1 1 2 1 3 2
Y NPARRAY1 1 1 2 2 2
CLF NEARESTCENTROID
CLFFITX Y
NEARESTCENTROIDMETRICEUCLIDEAN SHRINKTHRESHOLDNONE
PRINTCLFPREDICT08 1
1
```

METHODS

FITSELF X Y FIT THE NEARESTCENTROID MODEL ACCORDING TO THE GIVEN  
TRAINING DATA

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PREDICT SELF X PERFORM CLASSIFICATION ON AN ARRAY OF TEST VECTORS X

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND  
LABELS

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELFMETRIC'EUCLIDEAN' SHRINKTHRESHOLDNONE

FITSELFXY

FIT THE NEARESTCENTROID MODEL ACCORDING TO THE GIVEN TRAINING DATA

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES TRAINING VECTOR WHERE  
NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES NOTE THAT  
2192 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

CENTROID SHRINKING CANNOT BE USED WITH SPARSE MATRICES

YARRAY SHAPE NSAMPLES TARGET VALUES INTEGERS

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELF

PERFORM CLASSIFICATION ON AN ARRAY OF TEST VECTORS X

THE PREDICTED CLASS C FOR EACH SAMPLE IN X IS RETURNED

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES

RETURNS

CARRAY SHAPE NSAMPLES

NOTES

IF THE METRIC CONSTRUCTOR PARAMETER IS “PRECOMPUTED” X IS ASSUMED TO BE THE DISTANCE MATRIX BETWEEN THE

DATA TO BE PREDICTED AND SELFCENTROIDS

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS

IN MULTILABEL CLASSIFICATION THIS IS THE SUBSET ACCURACY WHICH IS A HARSH METRIC SINCE YOU REQUIRE FOR EACH

SAMPLE THAT EACH LABEL SET BE CORRECTLY PREDICTED

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE LABELS FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT MEAN ACCURACY OF SELF PREDICTX WRT Y

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE

PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT

OF A NESTED OBJECT

RETURNS

SELF

630SKLEARNNEIGHBORS NEAREST NEIGHBORS 2193

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES USING SKLEARNNEIGHBORSNEARESTCENTROID

- NEAREST CENTROID CLASSIFICATION
- CLASSIFICATION OF TEXT DOCUMENTS USING SPARSE FEATURES

63011SKLEARNNEIGHBORS NEARESTNEIGHBORS

CLASSSSKLEARNNEIGHBORS NEARESTNEIGHBORS NNEIGHBORS5 RADIUS10 ALGORITHM'AUTO'

LEAFSIZE30 METRIC'MINKOWSKI' P2MET

RICPARAMSNONE NJOBSNONE KWARGS

UNSUPERVISED LEARNER FOR IMPLEMENTING NEIGHBOR SEARCHES

READ MORE IN THE USER GUIDE

PARAMETERS

NNEIGHBORS INT OPTIONAL DEFAULT 5 NUMBER OF NEIGHBORS TO USE BY DEFAULT FOR KNEIGHBORS QUERIES

RADIUS FLOAT OPTIONAL DEFAULT 10 RANGE OF PARAMETER SPACE TO USE BY DEFAULT FOR RADIUSNEIGHBORS QUERIES

ALGORITHM 'AUTO' 'BALLTREE' 'KDTREE' 'BRUTE' OPTIONAL ALGORITHM USED TO COMPUTE THE NEAREST NEIGHBORS

- 'BALLTREE' WILL USE BALLTREE
- 'KDTREE' WILL USE KDTREE
- 'BRUTE' WILL USE A BRUTEFORCE SEARCH
- 'AUTO' WILL ATTEMPT TO DECIDE THE MOST APPROPRIATE ALGORITHM BASED ON THE VALUES PASSED

TOFIT METHOD

NOTE FITTING ON SPARSE INPUT WILL OVERRIDE THE SETTING OF THIS PARAMETER USING BRUTE FORCE

LEAFSIZE INT OPTIONAL DEFAULT 30 LEAF SIZE PASSED TO BALLTREE OR KDTREE THIS CAN AFFECT THE SPEED OF THE CONSTRUCTION AND QUERY AS WELL AS THE MEMORY REQUIRED TO STORE THE TREE THE OPTIMAL VALUE DEPENDS ON THE NATURE OF THE PROBLEM

METRIC STRING OR CALLABLE DEFAULT 'MINKOWSKI' METRIC TO USE FOR DISTANCE COMPUTATION ANY METRIC FROM SCIKITLEARN OR SCIPYSPATIALDISTANCE CAN BE USED

IF METRIC IS A CALLABLE FUNCTION IT IS CALLED ON EACH PAIR OF INSTANCES ROWS AND THE RESULTING VALUE RECORDED THE CALLABLE SHOULD TAKE TWO ARRAYS AS INPUT AND RETURN ONE VALUE INDICATING THE DISTANCE BETWEEN THEM THIS WORKS FOR SCIPY'S METRICS BUT IS LESS EFFICIENT THAN PASSING THE METRIC NAME AS A STRING

DISTANCE MATRICES ARE NOT SUPPORTED

VALID VALUES FOR METRIC ARE

- FROM SCIKITLEARN 'CITYBLOCK' 'COSINE' 'EUCLIDEAN' 'L1' 'L2' 'MANHATTAN'
- FROM SCIPYSPATIALDISTANCE 'BRAYCURTIS' 'CANBERRA' 'CHEBYSHEV' 'CORRELATION' 'DICE' 'HAMMING' 'JACCARD' 'KULSINSKI' 'MAHALANOBIS' 'MINKOWSKI' 'ROGERSTANIMOTO' 'RUSSELLRAO' 'SEUCLIDEAN' 'SOKALMICHENER' 'SOKALSNEATH' 'SQUEUCLIDEAN' 'YULE'

SEE THE DOCUMENTATION FOR SCIPYSPATIALDISTANCE FOR DETAILS ON THESE METRICS

2194 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

PINTEGER OPTIONAL DEFAULT 2 PARAMETER FOR THE MINKOWSKI METRIC FROM  
SKLEARNMETRICSPAIRWISEPAIRWISEDISTANCES WHEN P 1 THIS IS EQUIVALENT TO US  
ING MANHATTANDISTANCE L1 AND EUCLIDEANDISTANCE L2 FOR P 2 FOR ARBITRARY P  
MINKOWSKIDISTANCE LP IS USED

METRICPARAMS DICT OPTIONAL DEFAULT NONE ADDITIONAL KEYWORD ARGUMENTS FOR THE METRIC  
FUNCTION

NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF PARALLEL JOBS TO RUN FOR NEIGHBORS  
SEARCHNONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS  
USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

SEE ALSO

KNEIGHBORSCLASSIFIER  
RADIUSNEIGHBORSCLASSIFIER  
KNEIGHBORSREGRESSOR  
RADIUSNEIGHBORSREGRESSOR  
BALLTREE

NOTES

SEE NEAREST NEIGHBORS IN THE ONLINE DOCUMENTATION FOR A DISCUSSION OF THE CHOICE OF ALGORITHM AND  
LEAFSIZE

[HTTPSENWIKIPEDIAORGWIKIKNEARESTNEIGHBORALGORITHM](https://en.wikipedia.org/wiki/Nearest_neighbor_algorithm)

EXAMPLES

```
import numpy as np
from sklearn.neighbors import NearestNeighbors
samples = 0 0 2 1 0 0 0 0 1
neigh = NearestNeighbors(2, 0.4)
neigh.fit(samples)
neigh.kneighbors(0 0 13 2, return_distance=False)
```

ARRAY2 0  
NBRS NEIGHRADIUSNEIGHBORS0 0 13 0.4, return\_distance=False  
NPASARRAYNBRS00  
ARRAY2  
METHODS  
fit(self, X, Y) fit the model using X as training data  
continued on next page  
630sklearn.neighbors.NearestNeighbors 2195

SCIKITLEARN USER GUIDE RELEASE 0213

TABLE 6234 – CONTINUED FROM PREVIOUS PAGE

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

KNEIGHBORS SELF X NNEIGHBORS FINDS THE KNEIGHBORS OF A POINT

KNEIGHBORSGRAPH SELF X NNEIGHBORS MODE COMPUTES THE WEIGHTED GRAPH OF KNEIGHBORS FOR POINTS IN X

RADIUSNEIGHBORS SELF X RADIUS FINDS THE NEIGHBORS WITHIN A GIVEN RADIUS OF A POINT OR POINTS

RADIUSNEIGHBORSGRAPH SELF X RADIUS

MODECOMPUTES THE WEIGHTED GRAPH OF NEIGHBORS FOR POINTS

IN X

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELFNNNEIGHBORS5 RADIUS10 ALGORITHM‘AUTO’ LEAFSIZE30 METRIC‘MINKOWSKI’

P2METRICPARAMSNONE NJOBSNONE KWARGS

FITSELFXYNONE

FIT THE MODEL USING X AS TRAINING DATA

PARAMETERS

XARRAYLIKE SPARSE MATRIX BALLTREE KDTREE TRAINING DATA IF ARRAY OR MATRIX SHAPE

NSAMPLES NFEATURES OR NSAMPLES NSAMPLES IF METRIC‘PRECOMPUTED’

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

KNEIGHBORS SELFXYNONE NNEIGHBORSNONE RETURNDISTANCETURE

FINDS THE KNEIGHBORS OF A POINT RETURNS INDICES OF AND DISTANCES TO THE NEIGHBORS OF EACH POINT

PARAMETERS

XARRAYLIKE SHAPE NQUERY NFEATURES OR NQUERY NINDEXED IF METRIC ‘PRECOM

PUTED’ THE QUERY POINT OR POINTS IF NOT PROVIDED NEIGHBORS OF EACH INDEXED POINT ARE

RETURNED IN THIS CASE THE QUERY POINT IS NOT CONSIDERED ITS OWN NEIGHBOR

NNEIGHBORS INT NUMBER OF NEIGHBORS TO GET DEFAULT IS THE VALUE PASSED TO THE CONSTRUCT

TOR

RETURNDISTANCE BOOLEAN OPTIONAL DEFAULTS TO TRUE IF FALSE DISTANCES WILL NOT BE RETURNED

RETURNS

DIST ARRAY ARRAY REPRESENTING THE LENGTHS TO POINTS ONLY PRESENT IF RETURNDISTANCETURE

IND ARRAY INDICES OF THE NEAREST POINTS IN THE POPULATION MATRIX

EXAMPLES

IN THE FOLLOWING EXAMPLE WE CONSTRUCT A NEIGHBORSCLASSIFIER CLASS FROM AN ARRAY REPRESENTING OUR DATA SET

AND ASK WHO’S THE CLOSEST POINT TO 111

2196 CHAPTER 6 API REFERENCE

```
SCIKITLEARN USER GUIDE RELEASE 0213
SAMPLES 0 0 0 0 5 0 1 1 5
FROM SKLEARNNEIGHBORS IMPORT NEARESTNEIGHBORS
NEIGH NEARESTNEIGHBORSNNEIGHBORS1
NEIGHFITSAMPLES
NEARESTNEIGHBORSALGORITHM AUTO LEAF SIZE30
PRINTNEIGHKNEIGHBORS1 1 1
ARRAY05 ARRAY2
AS YOU CAN SEE IT RETURNS 05 AND 2 WHICH MEANS THAT THE ELEMENT IS AT DISTANCE 05 AND IS THE THIRD
ELEMENT OF SAMPLES INDEXES START AT 0 YOU CAN ALSO QUERY FOR MULTIPLE POINTS
X 0 1 0 1 0 1
NEIGHKNEIGHBORSX RETURN DISTANCE FALSE
ARRAY1
2
KNEIGHBORSGRAPH SELF X NONE NNEIGHBORS NONE MODE 'CONNECTIVITY'
COMPUTES THE WEIGHTED GRAPH OF KNEIGHBORS FOR POINTS IN X
PARAMETERS
X ARRAY LIKE SHAPE NQUERY NFEATURES OR NQUERY NINDEXED IF METRIC 'PRECOM
PUTED' THE QUERY POINT OR POINTS IF NOT PROVIDED NEIGHBORS OF EACH INDEXED POINT ARE
RETURNED IN THIS CASE THE QUERY POINT IS NOT CONSIDERED ITS OWN NEIGHBOR
NNEIGHBORS INT NUMBER OF NEIGHBORS FOR EACH SAMPLE DEFAULT IS VALUE PASSED TO THE
CONSTRUCTOR
MODE 'CONNECTIVITY' 'DISTANCE' OPTIONAL TYPE OF RETURNED MATRIX 'CONNECTIVITY' WILL
RETURN THE CONNECTIVITY MATRIX WITH ONES AND ZEROS IN 'DISTANCE' THE EDGES ARE EUCLIDEAN
DISTANCE BETWEEN POINTS
RETURNS
A SPARSE MATRIX IN CSR FORMAT SHAPE NSAMPLES NSAMPLESFIT NSAMPLESFIT IS THE
NUMBER OF SAMPLES IN THE FITTED DATA A I J IS ASSIGNED THE WEIGHT OF EDGE THAT CONNECTS I
TO J
SEE ALSO
NEARESTNEIGHBORSRADIUSNEIGHBORSGRAPH
EXAMPLES
X 0 3 1
FROM SKLEARNNEIGHBORS IMPORT NEARESTNEIGHBORS
NEIGH NEARESTNEIGHBORSNNEIGHBORS2
NEIGH FITX
NEARESTNEIGHBORSALGORITHM AUTO LEAF SIZE30
A NEIGHKNEIGHBORSGRAPHX
A TO ARRAY
ARRAY1 0 1
0 1 1
1 0 1
RADIUSNEIGHBORS SELF X NONE RADIUS NONE RETURN DISTANCE TRUE
FINDS THE NEIGHBORS WITHIN A GIVEN RADIUS OF A POINT OR POINTS
630 SKLEARNNEIGHBORS NEAREST NEIGHBORS 2197
```

SCIKITLEARN USER GUIDE RELEASE 0213

RETURN THE INDICES AND DISTANCES OF EACH POINT FROM THE DATASET LYING IN A BALL WITH SIZE RADIUS AROUND THE POINTS OF THE QUERY ARRAY POINTS LYING ON THE BOUNDARY ARE INCLUDED IN THE RESULTS  
THE RESULT POINTS ARE NOTNECESSARILY SORTED BY DISTANCE TO THEIR QUERY POINT

PARAMETERS

XARRAYLIKE NSAMPLES NFEATURES OPTIONAL THE QUERY POINT OR POINTS IF NOT PROVIDED  
NEIGHBORS OF EACH INDEXED POINT ARE RETURNED IN THIS CASE THE QUERY POINT IS NOT CONSIDERED ITS OWN NEIGHBOR

RADIUS FLOAT LIMITING DISTANCE OF NEIGHBORS TO RETURN DEFAULT IS THE VALUE PASSED TO THE CONSTRUCTOR

RETURNDISTANCE BOOLEAN OPTIONAL DEFAULTS TO TRUE IF FALSE DISTANCES WILL NOT BE RETURNED

RETURNS

DIST ARRAY SHAPE NSAMPLES OF ARRAYS ARRAY REPRESENTING THE DISTANCES TO EACH POINT  
ONLY PRESENT IF RETURNDISTANCETRUE THE DISTANCE VALUES ARE COMPUTED ACCORDING TO THE METRIC CONSTRUCTOR PARAMETER

IND ARRAY SHAPE NSAMPLES OF ARRAYS AN ARRAY OF ARRAYS OF INDICES OF THE APPROXIMATE NEAREST POINTS FROM THE POPULATION MATRIX THAT LIE WITHIN A BALL OF SIZE RADIUS AROUND THE QUERY POINTS

NOTES

BECAUSE THE NUMBER OF NEIGHBORS OF EACH POINT IS NOT NECESSARILY EQUAL THE RESULTS FOR MULTIPLE QUERY POINTS CANNOT BE FIT IN A STANDARD DATA ARRAY FOR EFFICIENCY RADIUSNEIGHBORS RETURNS ARRAYS OF OBJECTS WHERE EACH OBJECT IS A 1D ARRAY OF INDICES OR DISTANCES

EXAMPLES

IN THE FOLLOWING EXAMPLE WE CONSTRUCT A NEIGHBORSCLASSIFIER CLASS FROM AN ARRAY REPRESENTING OUR DATA SET AND ASK WHO'S THE CLOSEST POINT TO 1 1 1

```
import numpy as np
samples = 0 0 0 0 5 0 1 1 5
from sklearn.neighbors import NearestNeighbors
neigh = NearestNeighbors(radius=16)
neigh.fit(samples)
nearest_neighbors_algorithm = auto(leafsize=30)
rng = NeighRadiusNeighbors(1, 1, 1)
print(np.asarray(rng[0]))
15 05
print(np.asarray(rng[10]))
1 2
```

THE FIRST ARRAY RETURNED CONTAINS THE DISTANCES TO ALL POINTS WHICH ARE CLOSER THAN 16 WHILE THE SECOND ARRAY RETURNED CONTAINS THEIR INDICES IN GENERAL MULTIPLE POINTS CAN BE QUERIED AT THE SAME TIME

RADIUSNEIGHBORSGRAPH SELFXTNONE RADIUSNONE MODE'CONNECTIVITY'  
COMPUTES THE WEIGHTED GRAPH OF NEIGHBORS FOR POINTS IN X  
NEIGHBORHOODS ARE RESTRICTED THE POINTS AT A DISTANCE LOWER THAN RADIUS

PARAMETERS

SCIKITLEARN USER GUIDE RELEASE 0213

XARRAYLIKE SHAPE NSAMPLES NFEATURES OPTIONAL THE QUERY POINT OR POINTS IF NOT PROVIDED NEIGHBORS OF EACH INDEXED POINT ARE RETURNED IN THIS CASE THE QUERY POINT IS NOT CONSIDERED ITS OWN NEIGHBOR

RADIUS FLOAT RADIUS OF NEIGHBORHOODS DEFAULT IS THE VALUE PASSED TO THE CONSTRUCTOR

MODE 'CONNECTIVITY' 'DISTANCE' OPTIONAL TYPE OF RETURNED MATRIX 'CONNECTIVITY' WILL RETURN THE CONNECTIVITY MATRIX WITH ONES AND ZEROS IN 'DISTANCE' THE EDGES ARE EUCLIDEAN DISTANCE BETWEEN POINTS

RETURNS

ASPARSE MATRIX IN CSR FORMAT SHAPE NSAMPLES NSAMPLES AI J IS ASSIGNED THE WEIGHT OF EDGE THAT CONNECTS I TO J

SEE ALSO

KNEIGHBORSGRAPH

EXAMPLES

```
X 0 3 1
FROM SKLEARNNEIGHBORS IMPORT NEARESTNEIGHBORS
NEIGH NEARESTNEIGHBORSRADIUS15
NEIGHFITX
NEARESTNEIGHBORSALGORITHMAUTO LEAFSIZE30
A NEIGHRADIUSNEIGHBORSGRAPHX
ATOARRAY
ARRAY1 0 1
0 1 0
1 0 1
SETPARAMS SELFPARAMS
SET THE PARAMETERS OF THIS ESTIMATOR
THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT
RETURNS
SELF
63012SKLEARNNEIGHBORS NEIGHBORHOODCOMPONENTSANALYSIS
CLASSSKLEARNNEIGHBORS NEIGHBORHOODCOMPONENTSANALYSIS NCOMPONENTSNONE
INIT'AUTO' WARMSTARTFALSE
MAXITER50 TOL1E05
CALLBACKNONE VERBOSE0
RANDOMSTATENONE
NEIGHBORHOOD COMPONENTS ANALYSIS
NEIGHBORHOOD COMPONENT ANALYSIS NCA IS A MACHINE LEARNING ALGORITHM FOR METRIC LEARNING IT LEARNS A LINEAR TRANSFORMATION IN A SUPERVISED FASHION TO IMPROVE THE CLASSIFICATION ACCURACY OF A STOCHASTIC NEAREST NEIGHBORS RULE IN THE TRANSFORMED SPACE
READ MORE IN THE USER GUIDE
630SKLEARNNEIGHBORS NEAREST NEIGHBORS 2199
```

SCIKITLEARN USER GUIDE RELEASE 0213

PARAMETERS

NCOMPONENTS INT OPTIONAL DEFAULTNONE PREFERRED DIMENSIONALITY OF THE PROJECTED SPACE  
IF NONE IT WILL BE SET TO NFEATURES

INIT STRING OR NUMPY ARRAY OPTIONAL DEFAULT'AUTO' INITIALIZATION OF THE LINEAR TRANSFORMATION  
POSSIBLE OPTIONS ARE 'AUTO' 'PCA' 'LDA' 'IDENTITY' 'RANDOM' AND A NUMPY ARRAY OF SHAPE

NFEATURESA NFEATURESB

'AUTO' DEPENDING ON NCOMPONENTS THE MOST REASONABLE INITIALIZATION WILL BE CHOSEN

IFNCOMPONENTS NCLASSES WE USE 'LDA' AS IT USES LABELS INFORMATION IF

NOT BUTNCOMPONENTS MINNFEATURES NSAMPLES WE USE 'PCA' AS IT

PROJECTS DATA IN MEANINGFUL DIRECTIONS THOSE OF HIGHER VARIANCE OTHERWISE WE JUST USE  
'IDENTITY'

'PCA'NCOMPONENTS PRINCIPAL COMPONENTS OF THE INPUTS PASSED TO FIT WILL BE USED TO  
INITIALIZE THE TRANSFORMATION SEE DECOMPOSITIONPCA

'LDA'MINNCOMPONENTS NCLASSES MOST DISCRIMINATIVE COMPONENTS OF THE

INPUTS PASSED TO FIT WILL BE USED TO INITIALIZE THE TRANSFORMATION IF

NCOMPONENTS NCLASSES THE REST OF THE COMPONENTS WILL BE ZERO SEE

DISCRIMINANTANALYSISLINEARDISCRIMINANTANALYSIS

'IDENTITY' IFNCOMPONENTS IS STRICTLY SMALLER THAN THE DIMENSIONALITY OF THE INPUTS

PASSED TOFIT THE IDENTITY MATRIX WILL BE TRUNCATED TO THE FIRST NCOMPONENTS ROWS

'RANDOM' THE INITIAL TRANSFORMATION WILL BE A RANDOM ARRAY OF SHAPE NCOMPONENTS

NFEATURES EACH VALUE IS SAMPLED FROM THE STANDARD NORMAL DISTRIBUTION

NUMPY ARRAY NFEATURESB MUST MATCH THE DIMENSIONALITY OF THE INPUTS PASSED TO FIT

AND NFEATURESA MUST BE LESS THAN OR EQUAL TO THAT IF NCOMPONENTS IS NOT NONE

NFEATURESA MUST MATCH IT

WARMSTART BOOL OPTIONAL DEFAULTFALSE IF TRUE AND FIT HAS BEEN CALLED BEFORE THE SOLU  
TION OF THE PREVIOUS CALL TO FIT IS USED AS THE INITIAL LINEAR TRANSFORMATION NCOMPONENTS

ANDINIT WILL BE IGNORED

MAXITER INT OPTIONAL DEFAULT50 MAXIMUM NUMBER OF ITERATIONS IN THE OPTIMIZATION

TOLFLOAT OPTIONAL DEFAULT1E5 CONVERGENCE TOLERANCE FOR THE OPTIMIZATION

CALLBACK CALLABLE OPTIONAL DEFAULTNONE IF NOT NONE THIS FUNCTION IS CALLED AFTER EVERY  
ITERATION OF THE OPTIMIZER TAKING AS ARGUMENTS THE CURRENT SOLUTION FLATTENED TRANSFORMATION  
MATRIX AND THE NUMBER OF ITERATIONS THIS MIGHT BE USEFUL IN CASE ONE WANTS TO EXAMINE OR  
STORE THE TRANSFORMATION FOUND AFTER EACH ITERATION

VERBOSE INT OPTIONAL DEFAULT0 IF 0 NO PROGRESS MESSAGES WILL BE PRINTED IF 1 PROGRESS  
MESSAGES WILL BE PRINTED TO STDOUT IF 1 PROGRESS MESSAGES WILL BE PRINTED AND THE DISP

PARAMETER OF SCIPYOPTIMIZE MINIMIZE WILL BE SET TO VERBOSE 2

RANDOMSTATE INT OR NUMPYRANDOMSTATE OR NONE OPTIONAL DEFAULTNONE A PSEUDO RAN

DOM NUMBER GENERATOR OBJECT OR A SEED FOR IT IF INT IF INITRANDOM RANDOMSTATE

IS USED TO INITIALIZE THE RANDOM TRANSFORMATION IF INITPCA RANDOMSTATE IS

PASSED AS AN ARGUMENT TO PCA WHEN INITIALIZING THE TRANSFORMATION

ATTRIBUTES

COMPONENTS ARRAY SHAPE NCOMPONENTS NFEATURES THE LINEAR TRANSFORMATION LEARNED DUR  
ING FITTING

NITER INT COUNTS THE NUMBER OF ITERATIONS PERFORMED BY THE OPTIMIZER

2200 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

REFERENCES

RF9B6BAEE82291 RF9B6BAEE82292

EXAMPLES

FROM SKLEARNNEIGHBORSNCA IMPORT NEIGHBORHOODCOMPONENTSANALYSIS

FROM SKLEARNNEIGHBORS IMPORT KNEIGHBORSCLASSIFIER

FROM SKLEARNDATASETS IMPORT LOADIRIS

FROM SKLEARNMODELSELECTION IMPORT TRAINTESTSPLIT

X Y LOADIRISRETURNXY TRUE

XTRAIN XTEST YTRAIN YTEST TRAINTESTSPLITX Y

STRATIFY TESTSIZE07 RANDOMSTATE42

NCA NEIGHBORHOODCOMPONENTSANALYSISRANDOMSTATE42

NCAFITXTRAIN YTRAIN

NEIGHBORHOODCOMPONENTSANALYSIS

KNN KNEIGHBORSCLASSIFIERNNEIGHBORS3

KNNFITXTRAIN YTRAIN

KNEIGHBORSCLASSIFIER

PRINTKNNSCOREXTEST YTEST

0933333

KNNFITNCATTRANSFORMXTRAIN YTRAIN

KNEIGHBORSCLASSIFIER

PRINTKNNSCORENCATTRANSFORMXTEST YTEST

0961904

METHODS

FITSELF X Y FIT THE MODEL ACCORDING TO THE GIVEN TRAINING DATA

FITTRANSFORM SELF X Y FIT TO DATA THEN TRANSFORM IT

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF X APPLIES THE LEARNED TRANSFORMATION TO THE GIVEN DATA

INIT SELFNCOMPONENTSNONE INIT'AUTO' WARMSTARTFALSE MAXITER50 TOL1E05 CALL

BACKNONE VERBOSE0 RANDOMSTATENONE

FITSELFXY

FIT THE MODEL ACCORDING TO THE GIVEN TRAINING DATA

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES THE TRAINING SAMPLES

YARRAYLIKE SHAPE NSAMPLES THE CORRESPONDING TRAINING LABELS

RETURNS

SELF OBJECT RETURNS A TRAINED NEIGHBORHOODCOMPONENTSANALYSIS MODEL

FITTRANSFORM SELFXYNONE FITPARAMS

FIT TO DATA THEN TRANSFORM IT

FITS TRANSFORMER TO X AND Y WITH OPTIONAL PARAMETERS FITPARAMS AND RETURNS A TRANSFORMED VERSION OF X

PARAMETERS

630SKLEARNNEIGHBORS NEAREST NEIGHBORS 2201

SCIKITLEARN USER GUIDE RELEASE 0213

XNUMPY ARRAY OF SHAPE NSAMPLES NFEATURES TRAINING SET

YNUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES

RETURNS

XNEW NUMPY ARRAY OF SHAPE NSAMPLES NFEATURESNEW TRANSFORMED ARRAY

GETPARAMS SELFDEEPTTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE

PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT

OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELF

APPLIES THE LEARNED TRANSFORMATION TO THE GIVEN DATA

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES DATA SAMPLES

RETURNS

XEMBEDDED ARRAY SHAPE NSAMPLES NCOMPONENTS THE DATA SAMPLES TRANSFORMED

RAISES

NOTFITTEDERROR IFFIT HAS NOT BEEN CALLED BEFORE

EXAMPLES USING SKLEARNNEIGHBORSNEIGHBORHOODCOMPONENTSANALYSIS

- MANIFOLD LEARNING ON HANDWRITTEN DIGITS LOCALLY LINEAR EMBEDDING ISOMAP
- COMPARING NEAREST NEIGHBORS WITH AND WITHOUT NEIGHBORHOOD COMPONENTS ANALYSIS
- DIMENSIONALITY REDUCTION WITH NEIGHBORHOOD COMPONENTS ANALYSIS
- NEIGHBORHOOD COMPONENTS ANALYSIS ILLUSTRATION

NEIGHBORSKNEIGHBORSGRAPH X NNEIGHBORS

COMPUTES THE WEIGHTED GRAPH OF KNEIGHBORS FOR POINTS

IN X

NEIGHBORSRADIUSNEIGHBORSGRAPH X RADIUS COMPUTES THE WEIGHTED GRAPH OF NEIGHBORS FOR POINTS IN

X

2202 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

63013SKLEARNNEIGHBORS KNEIGHBORSGRAPH

SKLEARNNEIGHBORS KNEIGHBORSGRAPH X NNEIGHBORS MODE'CONNECTIVITY' METRIC'MINKOWSKI' P2METRICPARAMSNONE INCLUDESELFFALSE NJOBSNONE COMPUTES THE WEIGHTED GRAPH OF KNEIGHBORS FOR POINTS IN X

READ MORE IN THE USER GUIDE

PARAMETERS

XARRAYLIKE OR BALLTREE SHAPE NSAMPLES NFEATURES SAMPLE DATA IN THE FORM OF A NUMPY ARRAY OR A PRECOMPUTED BALLTREE

NNEIGHBORS INT NUMBER OF NEIGHBORS FOR EACH SAMPLE

MODE 'CONNECTIVITY' 'DISTANCE' OPTIONAL TYPE OF RETURNED MATRIX 'CONNECTIVITY' WILL RETURN THE CONNECTIVITY MATRIX WITH ONES AND ZEROS AND 'DISTANCE' WILL RETURN THE DISTANCES BETWEEN NEIGHBORS ACCORDING TO THE GIVEN METRIC

METRIC STRING DEFAULT 'MINKOWSKI' THE DISTANCE METRIC USED TO CALCULATE THE KNEIGHBORS FOR EACH SAMPLE POINT THE DISTANCEMETRIC CLASS GIVES A LIST OF AVAILABLE METRICS THE DEFAULT DISTANCE IS 'EUCLIDEAN' 'MINKOWSKI' METRIC WITH THE P PARAM EQUAL TO 2

PINT DEFAULT 2 POWER PARAMETER FOR THE MINKOWSKI METRIC WHEN P 1 THIS IS EQUIVALENT TO USING MANHATTANDISTANCE L1 AND EUCLIDEANDISTANCE L2 FOR P 2 FOR ARBITRARY P MINKOWSKIDISTANCE LP IS USED

METRICPARAMS DICT OPTIONAL ADDITIONAL KEYWORD ARGUMENTS FOR THE METRIC FUNCTION

INCLUDESELF BOOL DEFAULTFALSE WHETHER OR NOT TO MARK EACH SAMPLE AS THE FIRST NEAREST NEIGHBOR TO ITSELF IF NONE THEN TRUE IS USED FOR MODE'CONNECTIVITY' AND FALSE FOR MODE'DISTANCE' AS THIS WILL PRESERVE BACKWARDS COMPATIBILITY

NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF PARALLEL JOBS TO RUN FOR NEIGHBORS SEARCHNONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

RETURNS

ASPARSE MATRIX IN CSR FORMAT SHAPE NSAMPLES NSAMPLES AI J IS ASSIGNED THE WEIGHT OF EDGE THAT CONNECTS I TO J

SEE ALSO

RADIUSNEIGHBORSGRAPH

EXAMPLES

```
X = 0 3 1
FROM SKLEARNNEIGHBORS IMPORT KNEIGHBORSGRAPH
A = KNEIGHBORSGRAPH(X, 2, MODE='CONNECTIVITY', INCLUDE_SELF=True)
A.toarray()
array([[0, 1, 1],
       [1, 0, 1]])
```

630SKLEARNNEIGHBORS NEAREST NEIGHBORS 2203

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES USING SKLEARNNEIGHBORSKNEIGHBORSGRAPH

- AGGLOMERATIVE CLUSTERING WITH AND WITHOUT STRUCTURE
- HIERARCHICAL CLUSTERING STRUCTURED VS UNSTRUCTURED WARD
- COMPARING DIFFERENT CLUSTERING ALGORITHMS ON TOY DATASETS

63014SKLEARNNEIGHBORS RADIUSNEIGHBORSGRAPH

SKLEARNNEIGHBORS RADIUSNEIGHBORSGRAPH X RADIUS MODE'CONNECTIVITY' METRIC'MINKOWSKI' P2METRICPARAMSNONE

INCLUDESELFFALSE NJOBSNONE

COMPUTES THE WEIGHTED GRAPH OF NEIGHBORS FOR POINTS IN X

NEIGHBORHOODS ARE RESTRICTED THE POINTS AT A DISTANCE LOWER THAN RADIUS

READ MORE IN THE USER GUIDE

PARAMETERS

XARRAYLIKE OR BALLTREE SHAPE NSAMPLES NFEATURES SAMPLE DATA IN THE FORM OF A NUMPY ARRAY OR A PRECOMPUTED BALLTREE

RADIUS FLOAT RADIUS OF NEIGHBORHOODS

MODE 'CONNECTIVITY' 'DISTANCE' OPTIONAL TYPE OF RETURNED MATRIX 'CONNECTIVITY' WILL RETURN THE CONNECTIVITY MATRIX WITH ONES AND ZEROS AND 'DISTANCE' WILL RETURN THE DISTANCES BETWEEN NEIGHBORS ACCORDING TO THE GIVEN METRIC

METRIC STRING DEFAULT 'MINKOWSKI' THE DISTANCE METRIC USED TO CALCULATE THE NEIGHBORS WITHIN A GIVEN RADIUS FOR EACH SAMPLE POINT THE DISTANCEMETRIC CLASS GIVES A LIST OF AVAILABLE METRICS THE DEFAULT DISTANCE IS 'EUCLIDEAN' 'MINKOWSKI' METRIC WITH THE PARAM EQUAL TO 2

PINT DEFAULT 2 POWER PARAMETER FOR THE MINKOWSKI METRIC WHEN P 1 THIS IS EQUIVALENT TO USING MANHATTANDISTANCE L1 AND EUCLIDEANDISTANCE L2 FOR P 2 FOR ARBITRARY P

MINKOWSKIDISTANCE LP IS USED

METRICPARAMS DICT OPTIONAL ADDITIONAL KEYWORD ARGUMENTS FOR THE METRIC FUNCTION

INCLUDESELF BOOL DEFAULTFALSE WHETHER OR NOT TO MARK EACH SAMPLE AS THE FIRST NEAREST NEIGHBOR TO ITSELF IF NONE THEN TRUE IS USED FOR MODE'CONNECTIVITY' AND FALSE FOR MODE'DISTANCE' AS THIS WILL PRESERVE BACKWARDS COMPATIBILITY

NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF PARALLEL JOBS TO RUN FOR NEIGHBORS

SEARCHNONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

RETURNS

ASPARSE MATRIX IN CSR FORMAT SHAPE NSAMPLES NSAMPLES AI J IS ASSIGNED THE WEIGHT OF EDGE THAT CONNECTS I TO J

SEE ALSO

KNEIGHBORSGRAPH

2204 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES

X 0 3 1

FROM SKLEARNNEIGHBORS IMPORT RADIUSNEIGHBORSGRAPH

A RADIUSNEIGHBORSGRAPHX 15 MODECONNECTIVITY

INCLUDESELF TRUE

ATOARRAY

ARRAY1 0 1

0 1 0

1 0 1

631SKLEARNNEURALNETWORK NEURAL NETWORK MODELS

THESKLEARNNEURALNETWORK MODULE INCLUDES MODELS BASED ON NEURAL NETWORKS

USER GUIDE SEE THE NEURAL NETWORK MODELS SUPERVISED ANDNEURAL NETWORK MODELS UNSUPERVISED SECTIONS FOR

FURTHER DETAILS

NEURALNETWORKBERNOULLIRBM NCOMPONENTS

BERNOULLI RESTRICTED BOLTZMANN MACHINE RBM

NEURALNETWORKMLPCLASSIFIER MULTILAYER PERCEPTRON CLASSIFIER

NEURALNETWORKMLPREGRESSOR MULTILAYER PERCEPTRON REGRESSOR

6311SKLEARNNEURALNETWORK BERNOULLIRBM

CLASSSKLEARNNEURALNETWORK BERNOULLIRBM NCOMPONENTS256 LEARNINGRATE01

BATCHSIZE10 NITER10 VERBOSE0 RAN

DOMSTATENONE

BERNOULLI RESTRICTED BOLTZMANN MACHINE RBM

A RESTRICTED BOLTZMANN MACHINE WITH BINARY VISIBLE UNITS AND BINARY HIDDEN UNITS PARAMETERS ARE ESTIMATED

USING STOCHASTIC MAXIMUM LIKELIHOOD SML ALSO KNOWN AS PERSISTENT CONTRASTIVE DIVERGENCE PCD 2

THE TIME COMPLEXITY OF THIS IMPLEMENTATION IS OD2ASSUMING D NFEATURES NCOMPONENTS

READ MORE IN THE USER GUIDE

PARAMETERS

NCOMPONENTS INT OPTIONAL NUMBER OF BINARY HIDDEN UNITS

LEARNINGRATE FLOAT OPTIONAL THE LEARNING RATE FOR WEIGHT UPDATES IT IS HIGHLY RECOMMENDED

TO TUNE THIS HYPERPARAMETER REASONABLE VALUES ARE IN THE 100 3 RANGE

BATCHSIZE INT OPTIONAL NUMBER OF EXAMPLES PER MINIBATCH

NITER INT OPTIONAL NUMBER OF ITERATIONSSWEEPS OVER THE TRAINING DATASET TO PERFORM DURING

TRAINING

VERBOSE INT OPTIONAL THE VERBOSITY LEVEL THE DEFAULT ZERO MEANS SILENT MODE

RANDOMSTATE INTEGER OR RANDOMSTATE OPTIONAL A RANDOM NUMBER GENERATOR INSTANCE TO DE

FINE THE STATE OF THE RANDOM PERMUTATIONS GENERATOR IF AN INTEGER IS GIVEN IT FIXES THE SEED

DEFAULTS TO THE GLOBAL NUMPY RANDOM NUMBER GENERATOR

ATTRIBUTES

631SKLEARNNEURALNETWORK NEURAL NETWORK MODELS 2205

SCIKITLEARN USER GUIDE RELEASE 0213

INTERCEPTHIDDEN ARRAYLIKE SHAPE NCOMPONENTS BIASES OF THE HIDDEN UNITS

INTERCEPTVISIBLE ARRAYLIKE SHAPE NFEATURES BIASES OF THE VISIBLE UNITS

COMPONENTS ARRAYLIKE SHAPE NCOMPONENTS NFEATURES WEIGHT MATRIX WHERE NFEATURES

IN THE NUMBER OF VISIBLE UNITS AND NCOMPONENTS IS THE NUMBER OF HIDDEN UNITS

REFERENCES

1 HINTON G E OSINDERO S AND TEH Y A FAST LEARNING ALGORITHM FOR DEEP BELIEF NETS NEURAL COMPU  
TATION 18 PP 15271554 HTTPSWWWCSTORONTOEDUHINTONABSPSFASTNCPDF

2 TIELEMAN T TRAINING RESTRICTED BOLTZMANN MACHINES USING APPROXIMATIONS TO THE LIKELIHOOD GRADI  
ENT INTERNATIONAL CONFERENCE ON MACHINE LEARNING ICML 2008

EXAMPLES

```
import numpy as np
from sklearn.neural_network import BernoulliRBM
X = np.array([0, 0, 0, 1, 1, 1, 0, 1, 1, 1])
model = BernoulliRBM(n_components=2)
model.fit(X)
BernoulliRBM(batch_size=10, learning_rate=0.1, n_components=2, n_iter=10,
              random_state=None, verbose=0)

METHODS
fit(self, X, Y) fit the model to the data X
fit_transform(self, X, Y) fit to data then transform it
get_params(self, deep=True) get parameters for this estimator
gibbs(self, V) perform one Gibbs sampling step
partial_fit(self, X, Y) fit the model to the data X which should contain a par  
tial segment of the data
score_samples(self, X) compute the pseudolikelihood of X
set_params(self, **params) set the parameters of this estimator
transform(self, X) compute the hidden layer activation probabilities
ph1vX
init(self, n_components=256, learning_rate=0.1, batch_size=10, n_iter=10, verbose=0, ran  
dom_state=None)
fit(self, X, Y=None)
fit(self, X, Y=None) fit the model to the data X
parameters
X array-like sparse matrix shape (n_samples, n_features) training data
returns
self BernoulliRBM the fitted model
fit_transform(self, X, Y=None) fit params  
fit to data then transform it
```

2206 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

FITS TRANSFORMER TO X AND Y WITH OPTIONAL PARAMETERS FITPARAMS AND RETURNS A TRANSFORMED VERSION OF X

PARAMETERS

XNUMPY ARRAY OF SHAPE NSAMPLES NFEATURES TRAINING SET

YNUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES

RETURNS

XNEW NUMPY ARRAY OF SHAPE NSAMPLES NFEATURESNEW TRANSFORMED ARRAY

GETPARAMS SELFDEEPTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

GIBBSSELFV

PERFORM ONE GIBBS SAMPLING STEP

PARAMETERS

VARRAYLIKE SHAPE NSAMPLES NFEATURES VALUES OF THE VISIBLE LAYER TO START FROM

RETURNS

VNEW ARRAYLIKE SHAPE NSAMPLES NFEATURES VALUES OF THE VISIBLE LAYER AFTER ONE GIBBS

STEP

PARTIALFIT SELFXYNONE

FIT THE MODEL TO THE DATA X WHICH SHOULD CONTAIN A PARTIAL SEGMENT OF THE DATA

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TRAINING DATA

RETURNS

SELF BERNOULLIRBM THE FITTED MODEL

SCORESAMPLES SELFXY

COMPUTE THE PSEUDOLIKELIHOOD OF X

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES VALUES OF THE VISIBLE LAYER

MUST BE ALLBOOLEAN NOT CHECKED

RETURNS

PSEUDOLIKELIHOOD ARRAYLIKE SHAPE NSAMPLES VALUE OF THE PSEUDOLIKELIHOOD PROXY

FOR LIKELIHOOD

NOTES

THIS METHOD IS NOT DETERMINISTIC IT COMPUTES A QUANTITY CALLED THE FREE ENERGY ON X THEN ON A RANDOMLY

CORRUPTED VERSION OF X AND RETURNS THE LOG OF THE LOGISTIC FUNCTION OF THE DIFFERENCE

631SKLEARNNEURALNETWORK NEURAL NETWORK MODELS 2207

SCIKITLEARN USER GUIDE RELEASE 0213

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELF

COMPUTE THE HIDDEN LAYER ACTIVATION PROBABILITIES PH1VX

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE DATA TO BE TRANSFORMED

RETURNS

HARRAY SHAPE NSAMPLES NCOMPONENTS LATENT REPRESENTATIONS OF THE DATA

EXAMPLES USING SKLEARNNEURALNETWORKBERNOULLIRBM

- RESTRICTED BOLTZMANN MACHINE FEATURES FOR DIGIT CLASSIFICATION

6312SKLEARNNEURALNETWORK MLPCLASSIFIER

CLASSSSKLEARNNEURALNETWORK MLPCLASSIFIER HIDDENLAYERSIZES100 ACTIVA

TION‘RELU’ SOLVER‘ADAM’ ALPHA00001

BATCHSIZE‘AUTO’ LEARNINGRATE‘CONSTANT’

LEARNINGRATEINIT0001 POWERT05

MAXITER200 SHUFFLETRUE RAN

DOMSTATENONE TOL00001 VER

BOSEFALSE WARMSTARTFALSE MOMEN

TUM09 NESTEROVSMOMENTUMTRUE

EARLYSTOPPINGFALSE VALIDATIONFRACTION01

BETA109 BETA20999 EPSILON1E08

NITERNOCHANGE10

MULTILAYER PERCEPTRON CLASSIFIER

THIS MODEL OPTIMIZES THE LOGLOSS FUNCTION USING LBFGS OR STOCHASTIC GRADIENT DESCENT

NEW IN VERSION 018

PARAMETERS

HIDDENLAYERSIZES TUPLE LENGTH NLAYERS 2 DEFAULT 100 THE ITH ELEMENT REPRESENTS THE NUMBER OF NEURONS IN THE ITH HIDDEN LAYER

ACTIVATION ‘IDENTITY’ ‘LOGISTIC’ ‘TANH’ ‘RELU’ DEFAULT ‘RELU’ ACTIVATION FUNCTION FOR THE HID

DEN LAYER

- ‘IDENTITY’ NOOP ACTIVATION USEFUL TO IMPLEMENT LINEAR BOTTLENECK RETURNS FX X
- ‘LOGISTIC’ THE LOGISTIC SIGMOID FUNCTION RETURNS FX 1 1 EXPX
- ‘TANH’ THE HYPERBOLIC TAN FUNCTION RETURNS FX TANHX
- ‘RELU’ THE RECTIFIED LINEAR UNIT FUNCTION RETURNS FX MAX0 X

SOLVER ‘LBFGS’ ‘SGD’ ‘ADAM’ DEFAULT ‘ADAM’ THE SOLVER FOR WEIGHT OPTIMIZATION

2208 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

- ‘LBFGS’ IS AN OPTIMIZER IN THE FAMILY OF QUASINEWTON METHODS
- ‘SGD’ REFERS TO STOCHASTIC GRADIENT DESCENT
- ‘ADAM’ REFERS TO A STOCHASTIC GRADIENTBASED OPTIMIZER PROPOSED BY KINGMA DIEDERIK AND JIMMY BA

NOTE THE DEFAULT SOLVER ‘ADAM’ WORKS PRETTY WELL ON RELATIVELY LARGE DATASETS WITH THOU SANDS OF TRAINING SAMPLES OR MORE IN TERMS OF BOTH TRAINING TIME AND VALIDATION SCORE FOR SMALL DATASETS HOWEVER ‘LBFGS’ CAN CONVERGE FASTER AND PERFORM BETTER

ALPHA FLOAT OPTIONAL DEFAULT 00001 L2 PENALTY REGULARIZATION TERM PARAMETER

BATCHSIZE INT OPTIONAL DEFAULT ‘AUTO’ SIZE OF MINIBATCHES FOR STOCHASTIC OPTIMIZERS IF THE SOLVER IS ‘LBFGS’ THE CLASSIFIER WILL NOT USE MINIBATCH WHEN SET TO “AUTO”

BATCHSIZEMIN200 NSAMPLES

LEARNINGRATE ‘CONSTANT’ ‘INVSCALING’ ‘ADAPTIVE’ DEFAULT ‘CONSTANT’ LEARNING RATE SCHEDULE FOR WEIGHT UPDATES

- ‘CONSTANT’ IS A CONSTANT LEARNING RATE GIVEN BY ‘LEARNINGRATEINIT’
- ‘INVSCALING’ GRADUALLY DECREASES THE LEARNING RATE AT EACH TIME STEP ‘T’ USING AN INVERSE SCALING EXPONENT OF ‘POWERT’ EFFECTIVELEARNINGRATE LEARNINGRATEINIT POWT POWERT
- ‘ADAPTIVE’ KEEPS THE LEARNING RATE CONSTANT TO ‘LEARNINGRATEINIT’ AS LONG AS TRAINING LOSS KEEPS DECREASING EACH TIME TWO CONSECUTIVE EPOCHS FAIL TO DECREASE TRAINING LOSS BY AT LEAST TOL OR FAIL TO INCREASE VALIDATION SCORE BY AT LEAST TOL IF ‘EARLYSTOPPING’ IS ON THE CURRENT LEARNING RATE IS DIVIDED BY 5

ONLY USED WHEN SOLVERSGD

LEARNINGRATEINIT DOUBLE OPTIONAL DEFAULT 0001 THE INITIAL LEARNING RATE USED IT CONTROLS THE STEPSIZE IN UPDATING THE WEIGHTS ONLY USED WHEN SOLVER‘SGD’ OR ‘ADAM’

POWERT DOUBLE OPTIONAL DEFAULT 05 THE EXPONENT FOR INVERSE SCALING LEARNING RATE IT IS USED IN UPDATING EFFECTIVE LEARNING RATE WHEN THE LEARNINGRATE IS SET TO ‘INVSCALING’ ONLY USED WHEN SOLVER‘SGD’

MAXITER INT OPTIONAL DEFAULT 200 MAXIMUM NUMBER OF ITERATIONS THE SOLVER ITERATES UNTIL CONVERGENCE DETERMINED BY ‘TOL’ OR THIS NUMBER OF ITERATIONS FOR STOCHASTIC SOLVERS ‘SGD’ ‘ADAM’ NOTE THAT THIS DETERMINES THE NUMBER OF EPOCHS HOW MANY TIMES EACH DATA POINT WILL BE USED NOT THE NUMBER OF GRADIENT STEPS

SHUFFLE BOOL OPTIONAL DEFAULT TRUE WHETHER TO SHUFFLE SAMPLES IN EACH ITERATION ONLY USED WHEN SOLVER‘SGD’ OR ‘ADAM’

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULT NONE IF INT RANDOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NPRANDOM

TOLFLOAT OPTIONAL DEFAULT 1E4 TOLERANCE FOR THE OPTIMIZATION WHEN THE LOSS OR SCORE IS NOT IMPROVING BY AT LEAST TOL FORNITERNOCHANGE CONSECUTIVE ITERATIONS UNLESS LEARNINGRATE IS SET TO ‘ADAPTIVE’ CONVERGENCE IS CONSIDERED TO BE REACHED AND TRAINING STOPS

VERBOSE BOOL OPTIONAL DEFAULT FALSE WHETHER TO PRINT PROGRESS MESSAGES TO STDOUT

WARMSTART BOOL OPTIONAL DEFAULT FALSE WHEN SET TO TRUE REUSE THE SOLUTION OF THE PREVIOUS CALL TO FIT AS INITIALIZATION OTHERWISE JUST ERASE THE PREVIOUS SOLUTION SEE THE GLOSSARY

631SKLEARNNEURALNETWORK NEURAL NETWORK MODELS 2209

SCIKITLEARN USER GUIDE RELEASE 0213

MOMENTUM FLOAT DEFAULT 0.9 MOMENTUM FOR GRADIENT DESCENT UPDATE SHOULD BE BETWEEN 0 AND 1 ONLY USED WHEN SOLVER='SGD'

NESTEROVMOMENTUM BOOLEAN DEFAULT TRUE WHETHER TO USE NESTEROV'S MOMENTUM ONLY USED WHEN SOLVER='SGD' AND MOMENTUM 0

EARLYSTOPPING BOOL DEFAULT FALSE WHETHER TO USE EARLY STOPPING TO TERMINATE TRAINING WHEN VALIDATION SCORE IS NOT IMPROVING IF SET TO TRUE IT WILL AUTOMATICALLY SET ASIDE 10 OF TRAINING DATA AS VALIDATION AND TERMINATE TRAINING WHEN VALIDATION SCORE IS NOT IMPROVING BY AT LEAST TOL FOR NITERNOCHANGE CONSECUTIVE EPOCHS THE SPLIT IS STRATIFIED EXCEPT IN A MULTILABEL SETTING ONLY EFFECTIVE WHEN SOLVER='SGD' OR 'ADAM'

VALIDATIONFRACTION FLOAT OPTIONAL DEFAULT 0.1 THE PROPORTION OF TRAINING DATA TO SET ASIDE AS VALIDATION SET FOR EARLY STOPPING MUST BE BETWEEN 0 AND 1 ONLY USED IF EARLYSTOPPING IS TRUE

BETA1 FLOAT OPTIONAL DEFAULT 0.9 EXPONENTIAL DECAY RATE FOR ESTIMATES OF FIRST MOMENT VECTOR IN ADAM SHOULD BE IN 0 1 ONLY USED WHEN SOLVER='ADAM'

BETA2 FLOAT OPTIONAL DEFAULT 0.999 EXPONENTIAL DECAY RATE FOR ESTIMATES OF SECOND MOMENT VECTOR IN ADAM SHOULD BE IN 0 1 ONLY USED WHEN SOLVER='ADAM'

EPSILON FLOAT OPTIONAL DEFAULT 1e-8 VALUE FOR NUMERICAL STABILITY IN ADAM ONLY USED WHEN SOLVER='ADAM'

NITERNOCHANGE INT OPTIONAL DEFAULT 10 MAXIMUM NUMBER OF EPOCHS TO NOT MEET TOL IMPROVEMENT ONLY EFFECTIVE WHEN SOLVER='SGD' OR 'ADAM'

NEW IN VERSION 0.20

ATTRIBUTES

CLASSES ARRAY OR LIST OF ARRAY OF SHAPE NCLASSES CLASS LABELS FOR EACH OUTPUT

LOSS FLOAT THE CURRENT LOSS COMPUTED WITH THE LOSS FUNCTION

COEFS LIST LENGTH NLAYERS - 1 THE ITH ELEMENT IN THE LIST REPRESENTS THE WEIGHT MATRIX CORRESPONDING TO LAYER I

INTERCEPTS LIST LENGTH NLAYERS - 1 THE ITH ELEMENT IN THE LIST REPRESENTS THE BIAS VECTOR CORRESPONDING TO LAYER I - 1

NITER INT THE NUMBER OF ITERATIONS THE SOLVER HAS RAN

NLAYERS INT NUMBER OF LAYERS

NOUTPUTS INT NUMBER OF OUTPUTS

OUTACTIVATION STRING NAME OF THE OUTPUT ACTIVATION FUNCTION

NOTES

MLPCLASSIFIER TRAINS ITERATIVELY SINCE AT EACH TIME STEP THE PARTIAL DERIVATIVES OF THE LOSS FUNCTION WITH RESPECT TO THE MODEL PARAMETERS ARE COMPUTED TO UPDATE THE PARAMETERS

IT CAN ALSO HAVE A REGULARIZATION TERM ADDED TO THE LOSS FUNCTION THAT SHRINKS MODEL PARAMETERS TO PREVENT OVER FITTING

THIS IMPLEMENTATION WORKS WITH DATA REPRESENTED AS DENSE NUMPY ARRAYS OR SPARSE SCIPY ARRAYS OF FLOATING POINT VALUES

2210 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

REFERENCES

HINTON GEOFFREY E “CONNECTIONIST LEARNING PROCEDURES” ARTIFICIAL INTELLIGENCE 401 1989 185234  
GLOROT XAVIER AND YOSHUA BENGIO “UNDERSTANDING THE DIFFICULTY OF TRAINING DEEP FEEDFORWARD NEURAL NET  
WORKS” INTERNATIONAL CONFERENCE ON ARTIFICIAL INTELLIGENCE AND STATISTICS 2010  
HE KAIMING ET AL “DELIVING DEEP INTO RECTIFIERS SURPASSING HUMANLEVEL PERFORMANCE ON IMAGENET CLASSI  
FICATION” ARXIV PREPRINT ARXIV150201852 2015  
KINGMA DIEDERIK AND JIMMY BA “ADAM A METHOD FOR STOCHASTIC OPTIMIZATION” ARXIV PREPRINT  
ARXIV14126980 2014

METHODS

FITSELF X Y FIT THE MODEL TO DATA MATRIX X AND TARGETS Y  
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR  
PREDICT SELF X PREDICT USING THE MULTILAYER PERCEPTRON CLASSIFIER  
PREDICTLOGPROBA SELF X RETURN THE LOG OF PROBABILITY ESTIMATES  
PREDICTPROBA SELF X PROBABILITY ESTIMATES  
SCORE SELF X Y SAMPLEWEIGHT RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND  
LABELS  
SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR  
INIT SELF HIDDENLAYERSIZES100 ACTIVATION'RELU' SOLVER'ADAM' AL  
PHA00001 BATCHSIZE'AUTO' LEARNINGRATE'CONSTANT' LEARNINGRATEINIT0001  
POWERT05 MAXITER200 SHUFFLETRUE RANDOMSTATENONE TOL00001 VER  
BOSEFALSE WARMSTARTFALSE MOMENTUM09 NESTEROVSMOMENTUMTRUE  
EARLYSTOPPINGFALSE VALIDATIONFRACTION01 BETA109 BETA20999 EPSILON1E08  
NITERNOCHANGE10  
FITSELFXY  
FIT THE MODEL TO DATA MATRIX X AND TARGETS Y  
PARAMETERS  
XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT DATA  
YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS THE TARGET VALUES CLASS LABELS  
IN CLASSIFICATION REAL NUMBERS IN REGRESSION  
RETURNS  
SELF RETURNS A TRAINED MLP MODEL  
GETPARAMS SELFDEEPTURE  
GET PARAMETERS FOR THIS ESTIMATOR  
PARAMETERS  
DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBOBJECTS THAT ARE ESTIMATORS  
RETURNS  
PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES  
PARTIALFIT  
UPDATE THE MODEL WITH A SINGLE ITERATION OVER THE GIVEN DATA  
631SKLEARNNEURALNETWORK NEURAL NETWORK MODELS 2211

SCIKITLEARN USER GUIDE RELEASE 0213

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT DATA

YARRAYLIKE SHAPE NSAMPLES THE TARGET VALUES

CLASSES ARRAY SHAPE NCLASSES DEFAULT NONE CLASSES ACROSS ALL CALLS TO PARTIALFIT CAN BE OBTAINED VIA NPUNIQUEYALL WHERE YALL IS THE TARGET VECTOR OF THE ENTIRE DATASET

THIS ARGUMENT IS REQUIRED FOR THE FIRST CALL TO PARTIALFIT AND CAN BE OMITTED IN THE SUBSEQUENT CALLS NOTE THAT Y DOESN'T NEED TO CONTAIN ALL LABELS IN CLASSES

RETURNS

SELF RETURNS A TRAINED MLP MODEL

PREDICTSELF

PREDICT USING THE MULTILAYER PERCEPTRON CLASSIFIER

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT DATA

RETURNS

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NCLASSES THE PREDICTED CLASSES

PREDICTLOGPROBA SELF

RETURN THE LOG OF PROBABILITY ESTIMATES

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES THE INPUT DATA

RETURNS

LOGYPROB ARRAYLIKE SHAPE NSAMPLES NCLASSES THE PREDICTED LOGPROBABILITY OF THE SAMPLE FOR EACH CLASS IN THE MODEL WHERE CLASSES ARE ORDERED AS THEY ARE IN SELFCLASSES EQUIVALENT TO LOGPREDICTPROBAX

PREDICTPROBA SELF

PROBABILITY ESTIMATES

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT DATA

RETURNS

YPROB ARRAYLIKE SHAPE NSAMPLES NCLASSES THE PREDICTED PROBABILITY OF THE SAMPLE FOR EACH CLASS IN THE MODEL WHERE CLASSES ARE ORDERED AS THEY ARE IN SELFCLASSES

SCORESELF

SCORESELFXY SAMPLEWEIGHT NONE RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS IN MULTILABEL CLASSIFICATION THIS IS THE SUBSET ACCURACY WHICH IS A HARSH METRIC SINCE YOU REQUIRE FOR EACH SAMPLE THAT EACH LABEL SET BE CORRECTLY PREDICTED

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE LABELS FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT MEAN ACCURACY OF SELF PREDICTX WRT Y

SCIKITLEARN USER GUIDE RELEASE 0213

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS  
SELF

EXAMPLES USING SKLEARNNEURALNETWORKMLPCLASSIFIER

- CLASSIFIER COMPARISON
- VISUALIZATION OF MLP WEIGHTS ON MNIST
- VARYING REGULARIZATION IN MULTILAYER PERCEPTRON
- COMPARE STOCHASTIC LEARNING STRATEGIES FOR MLPCLASSIFIER

6313SKLEARNNEURALNETWORK MLPREGRESSOR

CLASSSSKLEARNNEURALNETWORK MLPREGRESSOR HIDDENLAYERSIZES100 ACTIVATION'RELU' SOLVER'ADAM' ALPHA00001 BATCHSIZE'AUTO' LEARNINGRATE'CONSTANT' LEARNINGRATEINIT0001 POWERT05 MAXITER200 SHUFFLETRUE RANDOMSTATENONE TOL00001 VERBOSEFALSE WARMSTARTFALSE MOMENTUM09 NESTEROVSMOMENTUMTRUE EARLYSTOPPINGFALSE VALIDATIONFRACTION01 BETA109 BETA20999 EPSILON1E08 NITERNOCHANGE10

MULTILAYER PERCEPTRON REGRESSOR

THIS MODEL OPTIMIZES THE SQUAREDLOSS USING LBFGS OR STOCHASTIC GRADIENT DESCENT

NEW IN VERSION 018

PARAMETERS

HIDDENLAYERSIZES TUPLE LENGTH NLAYERS 2 DEFAULT 100 THE ITH ELEMENT REPRESENTS THE NUMBER OF NEURONS IN THE ITH HIDDEN LAYER

ACTIVATION 'IDENTITY' 'LOGISTIC' 'TANH' 'RELU' DEFAULT 'RELU' ACTIVATION FUNCTION FOR THE HIDDEN LAYER

- 'IDENTITY' NOOP ACTIVATION USEFUL TO IMPLEMENT LINEAR BOTTLENECK RETURNS  $f(x) = x$
- 'LOGISTIC' THE LOGISTIC SIGMOID FUNCTION RETURNS  $f(x) = \frac{1}{1 + \exp(-x)}$
- 'TANH' THE HYPERBOLIC TAN FUNCTION RETURNS  $f(x) = \tanh(x)$
- 'RELU' THE RECTIFIED LINEAR UNIT FUNCTION RETURNS  $f(x) = \max(0, x)$

SOLVER 'LBFGS' 'SGD' 'ADAM' DEFAULT 'ADAM' THE SOLVER FOR WEIGHT OPTIMIZATION

- 'LBFGS' IS AN OPTIMIZER IN THE FAMILY OF QUASINEWTON METHODS
- 'SGD' REFERS TO STOCHASTIC GRADIENT DESCENT

6313SKLEARNNEURALNETWORK NEURAL NETWORK MODELS 2213

SCIKITLEARN USER GUIDE RELEASE 0213

- ‘ADAM’ REFERS TO A STOCHASTIC GRADIENTBASED OPTIMIZER PROPOSED BY KINGMA DIEDERIK AND JIMMY BA

NOTE THE DEFAULT SOLVER ‘ADAM’ WORKS PRETTY WELL ON RELATIVELY LARGE DATASETS WITH THOU SANDS OF TRAINING SAMPLES OR MORE IN TERMS OF BOTH TRAINING TIME AND VALIDATION SCORE FOR SMALL DATASETS HOWEVER ‘LBFGS’ CAN CONVERGE FASTER AND PERFORM BETTER

ALPHA FLOAT OPTIONAL DEFAULT 00001 L2 PENALTY REGULARIZATION TERM PARAMETER

BATCHSIZE INT OPTIONAL DEFAULT ‘AUTO’ SIZE OF MINIBATCHES FOR STOCHASTIC OPTIMIZERS IF THE SOLVER IS ‘LBFGS’ THE CLASSIFIER WILL NOT USE MINIBATCH WHEN SET TO “AUTO”

BATCHSIZEMIN200 NSAMPLES

LEARNINGRATE ‘CONSTANT’ ‘INVSCALING’ ‘ADAPTIVE’ DEFAULT ‘CONSTANT’ LEARNING RATE SCHEDULE FOR WEIGHT UPDATES

- ‘CONSTANT’ IS A CONSTANT LEARNING RATE GIVEN BY ‘LEARNINGRATEINIT’
- ‘INVSCALING’ GRADUALLY DECREASES THE LEARNING RATE LEARNINGRATE AT EACH TIME STEP ‘T’ USING AN INVERSE SCALING EXPONENT OF ‘POWER’ EFFECTIVELEARNINGRATE LEARN

INGRATEINIT POWT POWER

- ‘ADAPTIVE’ KEEPS THE LEARNING RATE CONSTANT TO ‘LEARNINGRATEINIT’ AS LONG AS TRAINING LOSS KEEPS DECREASING EACH TIME TWO CONSECUTIVE EPOCHS FAIL TO DECREASE TRAINING LOSS BY AT LEAST TOL OR FAIL TO INCREASE VALIDATION SCORE BY AT LEAST TOL IF ‘EARLYSTOPPING’ IS ON THE CURRENT LEARNING RATE IS DIVIDED BY 5

ONLY USED WHEN SOLVER‘SGD’

LEARNINGRATEINIT DOUBLE OPTIONAL DEFAULT 0001 THE INITIAL LEARNING RATE USED IT CONTROLS THE STEPSIZE IN UPDATING THE WEIGHTS ONLY USED WHEN SOLVER‘SGD’ OR ‘ADAM’

POWER DOUBLE OPTIONAL DEFAULT 05 THE EXPONENT FOR INVERSE SCALING LEARNING RATE IT IS USED IN UPDATING EFFECTIVE LEARNING RATE WHEN THE LEARNINGRATE IS SET TO ‘INVSCALING’ ONLY USED WHEN SOLVER‘SGD’

MAXITER INT OPTIONAL DEFAULT 200 MAXIMUM NUMBER OF ITERATIONS THE SOLVER ITERATES UNTIL CONVERGENCE DETERMINED BY ‘TOL’ OR THIS NUMBER OF ITERATIONS FOR STOCHASTIC SOLVERS ‘SGD’ ‘ADAM’ NOTE THAT THIS DETERMINES THE NUMBER OF EPOCHS HOW MANY TIMES EACH DATA POINT WILL BE USED NOT THE NUMBER OF GRADIENT STEPS

SHUFFLE BOOL OPTIONAL DEFAULT TRUE WHETHER TO SHUFFLE SAMPLES IN EACH ITERATION ONLY USED WHEN SOLVER‘SGD’ OR ‘ADAM’

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULT NONE IF INT RANDOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NRANDOM

TOLFLOAT OPTIONAL DEFAULT 1E4 TOLERANCE FOR THE OPTIMIZATION WHEN THE LOSS OR SCORE IS NOT IMPROVING BY AT LEAST TOL FORNITERNOCHANGE CONSECUTIVE ITERATIONS UNLESS LEARNINGRATE IS SET TO ‘ADAPTIVE’ CONVERGENCE IS CONSIDERED TO BE REACHED AND TRAINING STOPS

VERBOSE BOOL OPTIONAL DEFAULT FALSE WHETHER TO PRINT PROGRESS MESSAGES TO STDOUT

WARMSTART BOOL OPTIONAL DEFAULT FALSE WHEN SET TO TRUE REUSE THE SOLUTION OF THE PREVIOUS CALL TO FIT AS INITIALIZATION OTHERWISE JUST ERASE THE PREVIOUS SOLUTION SEE THE GLOSSARY

MOMENTUM FLOAT DEFAULT 09 MOMENTUM FOR GRADIENT DESCENT UPDATE SHOULD BE BETWEEN 0 AND 1 ONLY USED WHEN SOLVER‘SGD’

2214 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

NESTEROVMOMENTUM BOOLEAN DEFAULT TRUE WHETHER TO USE NESTEROV’S MOMENTUM ONLY  
USED WHEN SOLVER’S ‘SGD’ AND MOMENTUM 0

EARLYSTOPPING BOOL DEFAULT FALSE WHETHER TO USE EARLY STOPPING TO TERMINATE TRAINING WHEN  
VALIDATION SCORE IS NOT IMPROVING IF SET TO TRUE IT WILL AUTOMATICALLY SET ASIDE 10 OF TRAINING  
DATA AS VALIDATION AND TERMINATE TRAINING WHEN VALIDATION SCORE IS NOT IMPROVING BY AT LEAST  
TOL FORNITERNOCCHANGE CONSECUTIVE EPOCHS ONLY EFFECTIVE WHEN SOLVER’S ‘SGD’ OR  
‘ADAM’

VALIDATIONFRACTION FLOAT OPTIONAL DEFAULT 01 THE PROPORTION OF TRAINING DATA TO SET ASIDE AS  
VALIDATION SET FOR EARLY STOPPING MUST BE BETWEEN 0 AND 1 ONLY USED IF EARLYSTOPPING IS  
TRUE

BETA1 FLOAT OPTIONAL DEFAULT 09 EXPONENTIAL DECAY RATE FOR ESTIMATES OF FIRST MOMENT VECTOR  
IN ADAM SHOULD BE IN 0 1 ONLY USED WHEN SOLVER’S ‘ADAM’

BETA2 FLOAT OPTIONAL DEFAULT 0999 EXPONENTIAL DECAY RATE FOR ESTIMATES OF SECOND MOMENT  
VECTOR IN ADAM SHOULD BE IN 0 1 ONLY USED WHEN SOLVER’S ‘ADAM’

EPSILON FLOAT OPTIONAL DEFAULT 1E8 VALUE FOR NUMERICAL STABILITY IN ADAM ONLY USED WHEN  
SOLVER’S ‘ADAM’

NITERNOCCHANGE INT OPTIONAL DEFAULT 10 MAXIMUM NUMBER OF EPOCHS TO NOT MEET TOL  
IMPROVEMENT ONLY EFFECTIVE WHEN SOLVER’S ‘SGD’ OR ‘ADAM’

NEW IN VERSION 020

ATTRIBUTES

LOSS FLOAT THE CURRENT LOSS COMPUTED WITH THE LOSS FUNCTION

COEFS LIST LENGTH NLAYERS 1 THE ITH ELEMENT IN THE LIST REPRESENTS THE WEIGHT MATRIX CORRE  
SPONDING TO LAYER I

INTERCEPTS LIST LENGTH NLAYERS 1 THE ITH ELEMENT IN THE LIST REPRESENTS THE BIAS VECTOR COR  
RESPONDING TO LAYER I 1

NITER INT THE NUMBER OF ITERATIONS THE SOLVER HAS RAN

NLAYERS INT NUMBER OF LAYERS

NOUTPUTS INT NUMBER OF OUTPUTS

OUTACTIVATION STRING NAME OF THE OUTPUT ACTIVATION FUNCTION

NOTES

MLPREGRESSOR TRAINS ITERATIVELY SINCE AT EACH TIME STEP THE PARTIAL DERIVATIVES OF THE LOSS FUNCTION WITH RESPECT  
TO THE MODEL PARAMETERS ARE COMPUTED TO UPDATE THE PARAMETERS

IT CAN ALSO HAVE A REGULARIZATION TERM ADDED TO THE LOSS FUNCTION THAT SHRINKS MODEL PARAMETERS TO PREVENT OVER  
FITTING

THIS IMPLEMENTATION WORKS WITH DATA REPRESENTED AS DENSE AND SPARSE NUMPY ARRAYS OF FLOATING POINT VALUES

REFERENCES

HINTON GEOFFREY E “CONNECTIONIST LEARNING PROCEDURES” ARTIFICIAL INTELLIGENCE 401 1989 185234

GLOROT XAVIER AND YOSHUA BENGIO “UNDERSTANDING THE DIFFICULTY OF TRAINING DEEP FEEDFORWARD NEURAL NET  
WORKS” INTERNATIONAL CONFERENCE ON ARTIFICIAL INTELLIGENCE AND STATISTICS 2010

631SKLEARNNEURALNETWORK NEURAL NETWORK MODELS 2215

SCIKITLEARN USER GUIDE RELEASE 0213

HE KAIMING ET AL “DELIVING DEEP INTO RECTIFIERS SURPASSING HUMANLEVEL PERFORMANCE ON IMAGENET CLASSIFICATION” ARXIV PREPRINT ARXIV150201852 2015

KINGMA DIEDERIK AND JIMMY BA “ADAM A METHOD FOR STOCHASTIC OPTIMIZATION” ARXIV PREPRINT ARXIV14126980 2014

METHODS

FITSELF X Y FIT THE MODEL TO DATA MATRIX X AND TARGETS Y

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PREDICT SELF X PREDICT USING THE MULTILAYER PERCEPTRON MODEL

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELF HIDDENLAYERSIZES100 ACTIVATION'RELU' SOLVER'ADAM' AL

PHA00001 BATCHSIZE'AUTO' LEARNINGRATE'CONSTANT' LEARNINGRATEINIT0001

POWERT05 MAXITER200 SHUFFLETRUE RANDOMSTATENONE TOL00001 VER

BOSEFALSE WARMSTARTFALSE MOMENTUM09 NESTEROVSMOMENTUMTRUE

EARLYSTOPPINGFALSE VALIDATIONFRACTION01 BETA109 BETA20999 EPSILON1E08

NITERNOCHANGE10

FITSELFXY

FIT THE MODEL TO DATA MATRIX X AND TARGETS Y

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT DATA

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS THE TARGET VALUES CLASS LABELS

IN CLASSIFICATION REAL NUMBERS IN REGRESSION

RETURNS

SELF RETURNS A TRAINED MLP MODEL

GETPARAMS SELFDEEPTTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PARTIALFIT

UPDATE THE MODEL WITH A SINGLE ITERATION OVER THE GIVEN DATA

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT DATA

YARRAYLIKE SHAPE NSAMPLES THE TARGET VALUES

RETURNS

SELF RETURNS A TRAINED MLP MODEL

2216 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

PREDICTSELF

PREDICT USING THE MULTILAYER PERCEPTRON MODEL

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT DATA

RETURNS

YARRAYLIKE SHAPE NSAMPLES NOUTPUTS THE PREDICTED VALUES

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{U}{V}$  WHERE U IS THE RESIDUAL SUM OF SQUARES  $Y_{TRUE} - Y_{PRED}$   $2 \sum$  AND V IS THE TOTAL SUM OF SQUARES  $Y_{TRUE} - Y_{TRUEMEAN}$   $2 \sum$  THE BEST POSSIBLE SCORE IS 10 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 00

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES FOR SOME ESTIMATORS THIS MAY BE A PRECOMPUTED KERNEL MATRIX INSTEAD SHAPE NSAMPLES NSAMPLESFITTED WHERE NSAMPLESFITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE VALUES FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT R2 OF SELF PREDICTX WRT Y

NOTES

THE R2 SCORE USED WHEN CALLING SCORE ON A REGRESSOR WILL USE MULTIOUTPUTUNIFORMAVERAGE FROM VERSION 023 TO KEEP CONSISTENT WITH METRICSR2SCORE THIS WILL INFLUENCE THE SCORE METHOD OF ALL THE MULTIOUTPUT REGRESSORS EXCEPT FOR MULTIOUTPUTMULTIOUTPUTREGRESSOR TO SPECIFY THE DEFAULT VALUE MANUALLY AND AVOID THE WARNING PLEASE EITHER CALL METRICSR2SCORE DIRECTLY OR MAKE A CUSTOM SCORER WITH METRICSMAKESCORER THE BUILTIN SCORER R2 USES MULTIOUTPUTUNIFORMAVERAGE

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

EXAMPLES USING SKLEARNNEURALNETWORKMLPREGRESSOR

- PARTIAL DEPENDENCE PLOTS

631SKLEARNNEURALNETWORK NEURAL NETWORK MODELS 2217

SCIKITLEARN USER GUIDE RELEASE 0213

632SKLEARNPIPELINE PIPELINE

THE SKLEARNPIPELINE MODULE IMPLEMENTS UTILITIES TO BUILD A COMPOSITE ESTIMATOR AS A CHAIN OF TRANSFORMS AND ESTIMATORS

PIPELINEFEATUREUNION TRANSFORMERLIST CONCATENATES RESULTS OF MULTIPLE TRANSFORMER OBJECTS

PIPELINEPIPELINE STEPS MEMORY VERBOSE PIPELINE OF TRANSFORMS WITH A FINAL ESTIMATOR

6321SKLEARNPIPELINE FEATUREUNION

CLASS SKLEARNPIPELINE FEATUREUNION TRANSFORMERLIST NJOBS NONE TRANSFORMERWEIGHTS NONE

VERBOSE FALSE

CONCATENATES RESULTS OF MULTIPLE TRANSFORMER OBJECTS

THIS ESTIMATOR APPLIES A LIST OF TRANSFORMER OBJECTS IN PARALLEL TO THE INPUT DATA THEN CONCATENATES THE RESULTS

THIS IS USEFUL TO COMBINE SEVERAL FEATURE EXTRACTION MECHANISMS INTO A SINGLE TRANSFORMER

PARAMETERS OF THE TRANSFORMERS MAY BE SET USING ITS NAME AND THE PARAMETER NAME SEPARATED BY A " " A

TRANSFORMER MAY BE REPLACED ENTIRELY BY SETTING THE PARAMETER WITH ITS NAME TO ANOTHER TRANSFORMER OR REMOVED

BY SETTING TO 'DROP' OR NONE

READ MORE IN THE USER GUIDE

PARAMETERS

TRANSFORMERLIST LIST OF STRING TRANSFORMER TUPLES LIST OF TRANSFORMER OBJECTS TO BE APPLIED

TO THE DATA THE FIRST HALF OF EACH TUPLE IS THE NAME OF THE TRANSFORMER

NJOBS INT OR NONE OPTIONAL DEFAULT NONE NUMBER OF JOBS TO RUN IN PARALLEL NONE MEANS 1

UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT 1 MEANS USING ALL PROCESSORS SEE

GLOSSARY FOR MORE DETAILS

TRANSFORMERWEIGHTS DICT OPTIONAL MULTIPLICATIVE WEIGHTS FOR FEATURES PER TRANSFORMER KEYS

ARE TRANSFORMER NAMES VALUES THE WEIGHTS

VERBOSE BOOLEAN OPTIONAL DEFAULT FALSE IF TRUE THE TIME ELAPSED WHILE FITTING EACH TRANS

FORMER WILL BE PRINTED AS IT IS COMPLETED

SEE ALSO

SKLEARNPIPELINEMAKEUNION CONVENIENCE FUNCTION FOR SIMPLIFIED FEATURE UNION CONSTRUCTION

EXAMPLES

```
FROM SKLEARNPIPELINE IMPORT FEATUREUNION
FROM SKLEARNDECOMPOSITION IMPORT PCA TRUNCATEDSVD
UNION FEATUREUNIONPCA PCANCOMPONENTS1
SVD TRUNCATEDSVDNCOMPONENTS2
X 0 1 3 2 2 5
UNIONFITTRANSFORMX
ARRAY 15 30 08
15 57 04
```

2218 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

METHODS

FITSELF X Y FIT ALL TRANSFORMERS USING X

FITTRANSFORM SELF X Y FIT ALL TRANSFORMERS TRANSFORM THE DATA AND CONCATENATE RESULTS

GETFEATURENAMES SELF GET FEATURE NAMES FROM ALL TRANSFORMERS

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

SETPARAMS SELF KWARGS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF X TRANSFORM X SEPARATELY BY EACH TRANSFORMER CONCATENATE RESULTS

INIT SELFTRANSFORMERLIST NJOBSNONE TRANSFORMERWEIGHTSNONE VERBOSEFALSE

FITSELFXYNONE

FIT ALL TRANSFORMERS USING X

PARAMETERS

XITERABLE OR ARRAYLIKE DEPENDING ON TRANSFORMERS INPUT DATA USED TO FIT TRANSFORMERS

YARRAYLIKE SHAPE NSAMPLES OPTIONAL TARGETS FOR SUPERVISED LEARNING

RETURNS

SELF FEATUREUNION THIS ESTIMATOR

FITTRANSFORM SELFXYNONE FITPARAMS

FIT ALL TRANSFORMERS TRANSFORM THE DATA AND CONCATENATE RESULTS

PARAMETERS

XITERABLE OR ARRAYLIKE DEPENDING ON TRANSFORMERS INPUT DATA TO BE TRANSFORMED

YARRAYLIKE SHAPE NSAMPLES OPTIONAL TARGETS FOR SUPERVISED LEARNING

RETURNS

XT ARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES SUMNCOMPONENTS HSTACK OF RESULTS OF TRANSFORMERS SUMNCOMPONENTS IS THE SUM OF NCOMPONENTS OUTPUT DIMENSION OVER TRANSFORMERS

GETFEATURENAMES SELF

GET FEATURE NAMES FROM ALL TRANSFORMERS

RETURNS

FEATURENAMES LIST OF STRINGS NAMES OF THE FEATURES PRODUCED BY TRANSFORM

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

SETPARAMS SELFKWARGS

SET THE PARAMETERS OF THIS ESTIMATOR

632SKLEARNPIPELINE PIPELINE 2219

SCIKITLEARN USER GUIDE RELEASE 0213

VALID PARAMETER KEYS CAN BE LISTED WITH GETPARAMS

RETURNS

SELF

TRANSFORM SELF

TRANSFORM X SEPARATELY BY EACH TRANSFORMER CONCATENATE RESULTS

PARAMETERS

X ITERABLE OR ARRAYLIKE DEPENDING ON TRANSFORMERS INPUT DATA TO BE TRANSFORMED

RETURNS

XT ARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES SUMNCOMPONENTS HSTACK OF RESULTS

OF TRANSFORMERS SUMNCOMPONENTS IS THE SUM OF NCOMPONENTS OUTPUT DIMENSION OVER TRANSFORMERS

EXAMPLES USING SKLEARNPIPELINEFEATUREUNION

- CONCATENATING MULTIPLE FEATURE EXTRACTION METHODS

6322SKLEARNPIPELINE PIPELINE

CLASSSSKLEARNPIPELINE PIPELINE STEPS MEMORYNONE VERBOSEFALSE

PIPELINE OF TRANSFORMS WITH A FINAL ESTIMATOR

SEQUENTIALLY APPLY A LIST OF TRANSFORMS AND A FINAL ESTIMATOR INTERMEDIATE STEPS OF THE PIPELINE MUST BE ‘TRANSFORMS’ THAT IS THEY MUST IMPLEMENT FIT AND TRANSFORM METHODS THE FINAL ESTIMATOR ONLY NEEDS TO IMPLEMENT FIT

THE TRANSFORMERS IN THE PIPELINE CAN BE CACHED USING MEMORY ARGUMENT

THE PURPOSE OF THE PIPELINE IS TO ASSEMBLE SEVERAL STEPS THAT CAN BE CROSSVALIDATED TOGETHER WHILE SETTING DIFFERENT PARAMETERS FOR THIS IT ENABLES SETTING PARAMETERS OF THE VARIOUS STEPS USING THEIR NAMES AND THE PARAMETER NAME SEPARATED BY A ‘ ’ AS IN THE EXAMPLE BELOW A STEP’S ESTIMATOR MAY BE REPLACED ENTIRELY BY SETTING THE PARAMETER WITH ITS NAME TO ANOTHER ESTIMATOR OR A TRANSFORMER REMOVED BY SETTING IT TO ‘PASSTHROUGH’ OR NONE

READ MORE IN THE USER GUIDE

PARAMETERS

STEPS LIST LIST OF NAME TRANSFORM TUPLES IMPLEMENTING FITTRANSFORM THAT ARE CHAINED IN THE ORDER IN WHICH THEY ARE CHAINED WITH THE LAST OBJECT AN ESTIMATOR

MEMORY NONE STR OR OBJECT WITH THE JOBLIBMEMORY INTERFACE OPTIONAL USED TO CACHE THE FITTED TRANSFORMERS OF THE PIPELINE BY DEFAULT NO CACHING IS PERFORMED IF A STRING IS GIVEN IT IS THE PATH TO THE CACHING DIRECTORY ENABLING CACHING TRIGGERS A CLONE OF THE TRANSFORMERS BEFORE FITTING THEREFORE THE TRANSFORMER INSTANCE GIVEN TO THE PIPELINE CANNOT BE INSPECTED DIRECTLY USE THE ATTRIBUTE NAMEDSTEPS OR STEPS TO INSPECT ESTIMATORS WITHIN THE PIPELINE CACHING THE TRANSFORMERS IS ADVANTAGEOUS WHEN FITTING IS TIME CONSUMING

VERBOSE BOOLEAN OPTIONAL IF TRUE THE TIME ELAPSED WHILE FITTING EACH STEP WILL BE PRINTED AS IT IS COMPLETED

ATTRIBUTES

NAMEDSTEPS BUNCH OBJECT A DICTIONARY WITH ATTRIBUTE ACCESS READONLY ATTRIBUTE TO ACCESS ANY STEP PARAMETER BY USER GIVEN NAME KEYS ARE STEP NAMES AND VALUES ARE STEPS PARAMETERS

SEE ALSO

2220 CHAPTER 6 API REFERENCE

```
SCIKITLEARN USER GUIDE RELEASE 0213
SKLEARNPIPELINEMAKEPIPELINE CONVENIENCE FUNCTION FOR SIMPLIFIED PIPELINE CONSTRUCTION
EXAMPLES
FROM SKLEARN IMPORT SVM
FROM SKLEARNDATASETS IMPORT SAMPLESGENERATOR
FROM SKLEARNFEATURESELECTION IMPORT SELECTKBEST
FROM SKLEARNFEATURESELECTION IMPORT FREGRESSION
FROM SKLEARNPIPELINE IMPORT PIPELINE
GENERATE SOME DATA TO PLAY WITH
X Y SAMPLESGENERATORMAKECLASSIFICATION
NINFORMATIVE5 NREDUNDANT0 RANDOMSTATE42
ANOVA SVMC
ANOVAFILTER SELECTKBESTFREGRESSION K5
CLF SVM SVCKERNELLINEAR
ANOVASVM PIPELINEANOVA ANOVAFILTER SVC CLF
YOU CAN SET THE PARAMETERS USING THE NAMES ISSUED
FOR INSTANCE FIT USING A K OF 10 IN THE SELECTKBEST
AND A PARAMETER C OF THE SVM
ANOVASVMSETPARAMSANOVAK10 SVCC1FITX Y

PIPELINEMEMORYNONE
STEPSANOVA SELECTKBEST
SVC SVC VERBOSEFALSE
PREDICTION ANOVASVMPREDICTX
ANOVASVMSCOREX Y
083
GETTING THE SELECTED FEATURES CHOSEN BY ANOVAFILTER
ANOVASVMANOVAGETSUPPORT

ARRAYFALSE FALSE TRUE TRUE FALSE FALSE TRUE TRUE FALSE
TRUE FALSE TRUE TRUE FALSE TRUE FALSE TRUE TRUE
FALSE FALSE
ANOTHER WAY TO GET SELECTED FEATURES CHOSEN BY ANOVAFILTER
ANOVASVMNAMEDSTEPSANOVAGETSUPPORT

ARRAYFALSE FALSE TRUE TRUE FALSE FALSE TRUE TRUE FALSE
TRUE FALSE TRUE TRUE FALSE TRUE FALSE TRUE TRUE
FALSE FALSE
INDEXING CAN ALSO BE USED TO EXTRACT A SUBPIPELINE
SUBPIPELINE ANOVASVM1
SUBPIPELINE
PIPELINEMEMORYNONE STEPSANOVA VERBOSEFALSE
COEF ANOVASVM1COEF
ANOVASVMSVC ISANOVASVM1
TRUE
COEFSHAPE
1 10
SUBPIPELINEINVERSETRANSFORMCOEFSHAPE
1 20
METHODS
632SKLEARNPIPELINE PIPELINE 2221
```

SCIKITLEARN USER GUIDE RELEASE 0213

DECISIONFUNCTION SELF X APPLY TRANSFORMS AND DECISIONFUNCTION OF THE FINAL ESTIMATOR

FITSELF X Y FIT THE MODEL

FITPREDICT SELF X Y APPLIES FITPREDICT OF LAST STEP IN PIPELINE AFTER TRANSFORMS

FITTRANSFORM SELF X Y FIT THE MODEL AND TRANSFORM WITH THE FINAL ESTIMATOR

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PREDICT SELF X PREDICTPARAMS APPLY TRANSFORMS TO THE DATA AND PREDICT WITH THE FINAL ESTIMATOR

PREDICTLOGPROBA SELF X APPLY TRANSFORMS AND PREDICTLOGPROBA OF THE FINAL ESTIMATOR

PREDICTPROBA SELF X APPLY TRANSFORMS AND PREDICTPROBA OF THE FINAL ESTIMATOR

SCORE SELF X Y SAMPLEWEIGHT APPLY TRANSFORMS AND SCORE WITH THE FINAL ESTIMATOR

SETPARAMS SELF KWARGS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELFSTEPS MEMORYNONE VERBOSEFALSE

DECISIONFUNCTION SELF X

APPLY TRANSFORMS AND DECISIONFUNCTION OF THE FINAL ESTIMATOR

PARAMETERS

XITERABLE DATA TO PREDICT ON MUST FULFILL INPUT REQUIREMENTS OF FIRST STEP OF THE PIPELINE

RETURNS

YSCORE ARRAYLIKE SHAPE NSAMPLES NCLASSES

FITSELFXYNONE FITPARAMS

FIT THE MODEL

FIT ALL THE TRANSFORMS ONE AFTER THE OTHER AND TRANSFORM THE DATA THEN FIT THE TRANSFORMED DATA USING THE FINAL ESTIMATOR

PARAMETERS

XITERABLE TRAINING DATA MUST FULFILL INPUT REQUIREMENTS OF FIRST STEP OF THE PIPELINE

YITERABLE DEFAULTNONE TRAINING TARGETS MUST FULFILL LABEL REQUIREMENTS FOR ALL STEPS OF THE PIPELINE

FITPARAMS DICT OF STRING OBJECT PARAMETERS PASSED TO THE FIT METHOD OF EACH STEP

WHERE EACH PARAMETER NAME IS PREFIXED SUCH THAT PARAMETER PREFIX HAS KEYS

RETURNS

SELF PIPELINE THIS ESTIMATOR

FITPREDICT SELFXYNONE FITPARAMS

APPLIES FITPREDICT OF LAST STEP IN PIPELINE AFTER TRANSFORMS

APPLIES FITTRANSFORMS OF A PIPELINE TO THE DATA FOLLOWED BY THE FITPREDICT METHOD OF THE FINAL ESTIMATOR IN THE PIPELINE VALID ONLY IF THE FINAL ESTIMATOR IMPLEMENTS FITPREDICT

PARAMETERS

XITERABLE TRAINING DATA MUST FULFILL INPUT REQUIREMENTS OF FIRST STEP OF THE PIPELINE

YITERABLE DEFAULTNONE TRAINING TARGETS MUST FULFILL LABEL REQUIREMENTS FOR ALL STEPS OF THE PIPELINE

2222 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

FITPARAMS DICT OF STRING OBJECT PARAMETERS PASSED TO THE FIT METHOD OF EACH STEP  
WHERE EACH PARAMETER NAME IS PREFIXED SUCH THAT PARAMETER PFOR STEPSHAS KEYS  
RETURNS  
YPRED ARRAYLIKE

FITTRANSFORM SELFXYNONE FITPARAMS

FIT THE MODEL AND TRANSFORM WITH THE FINAL ESTIMATOR

FITS ALL THE TRANSFORMS ONE AFTER THE OTHER AND TRANSFORMS THE DATA THEN USES FITTRANSFORM ON TRANSFORMED  
DATA WITH THE FINAL ESTIMATOR

PARAMETERS

XITERABLE TRAINING DATA MUST FULFILL INPUT REQUIREMENTS OF FIRST STEP OF THE PIPELINE

YITERABLE DEFAULTNONE TRAINING TARGETS MUST FULFILL LABEL REQUIREMENTS FOR ALL STEPS OF  
THE PIPELINE

FITPARAMS DICT OF STRING OBJECT PARAMETERS PASSED TO THE FIT METHOD OF EACH STEP  
WHERE EACH PARAMETER NAME IS PREFIXED SUCH THAT PARAMETER PFOR STEPSHAS KEYS  
RETURNS

XTARRAYLIKE SHAPE NSAMPLES NTRANSFORMEDFEATURES TRANSFORMED SAMPLES

GETPARAMS SELFDEEPTTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

INVERSETRANSFORM

APPLY INVERSE TRANSFORMATIONS IN REVERSE ORDER

ALL ESTIMATORS IN THE PIPELINE MUST SUPPORT INVERSETRANSFORM

PARAMETERS

XTARRAYLIKE SHAPE NSAMPLES NTRANSFORMEDFEATURES DATA SAMPLES WHERE  
NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES MUST  
FULFILL INPUT REQUIREMENTS OF LAST STEP OF PIPELINE'S INVERSETRANSFORM METHOD

RETURNS

XTARRAYLIKE SHAPE NSAMPLES NFEATURES

PREDICTSELFXPREDICTPARAMS

APPLY TRANSFORMS TO THE DATA AND PREDICT WITH THE FINAL ESTIMATOR

PARAMETERS

XITERABLE DATA TO PREDICT ON MUST FULFILL INPUT REQUIREMENTS OF FIRST STEP OF THE PIPELINE

PREDICTPARAMS DICT OF STRING OBJECT PARAMETERS TO THE PREDICT CALLED AT THE END OF  
ALL TRANSFORMATIONS IN THE PIPELINE NOTE THAT WHILE THIS MAY BE USED TO RETURN UNCERTAINTIES  
FROM SOME MODELS WITH RETURNSTD OR RETURNCOV UNCERTAINTIES THAT ARE GENERATED BY THE  
TRANSFORMATIONS IN THE PIPELINE ARE NOT PROPAGATED TO THE FINAL ESTIMATOR

RETURNS

632SKLEARNPIPELINE PIPELINE 2223

SCIKITLEARN USER GUIDE RELEASE 0213

YPRED ARRAYLIKE

PREDICTLOGPROBA SELF

APPLY TRANSFORMS AND PREDICTLOGPROBA OF THE FINAL ESTIMATOR

PARAMETERS

XITERABLE DATA TO PREDICT ON MUST FULFILL INPUT REQUIREMENTS OF FIRST STEP OF THE PIPELINE

RETURNS

YSCORE ARRAYLIKE SHAPE NSAMPLES NCLASSES

PREDICTPROBA SELF

APPLY TRANSFORMS AND PREDICTPROBA OF THE FINAL ESTIMATOR

PARAMETERS

XITERABLE DATA TO PREDICT ON MUST FULFILL INPUT REQUIREMENTS OF FIRST STEP OF THE PIPELINE

RETURNS

YPROBA ARRAYLIKE SHAPE NSAMPLES NCLASSES

SCORESELFXYNONE SAMPLEWEIGHTNONE

APPLY TRANSFORMS AND SCORE WITH THE FINAL ESTIMATOR

PARAMETERS

XITERABLE DATA TO PREDICT ON MUST FULFILL INPUT REQUIREMENTS OF FIRST STEP OF THE PIPELINE

YITERABLE DEFAULTNONE TARGETS USED FOR SCORING MUST FULFILL LABEL REQUIREMENTS FOR ALL STEPS OF THE PIPELINE

SAMPLEWEIGHT ARRAYLIKE DEFAULTNONE IF NOT NONE THIS ARGUMENT IS PASSED AS

SAMPLEWEIGHT KEYWORD ARGUMENT TO THE SCORE METHOD OF THE FINAL ESTIMATOR

RETURNS

SCORE FLOAT

SETPARAMS SELFKWARGS

SET THE PARAMETERS OF THIS ESTIMATOR

VALID PARAMETER KEYS CAN BE LISTED WITH GETPARAMS

RETURNS

SELF

TRANSFORM

APPLY TRANSFORMS AND TRANSFORM WITH THE FINAL ESTIMATOR

THIS ALSO WORKS WHERE FINAL ESTIMATOR IS NONE ALL PRIOR TRANSFORMATIONS ARE APPLIED

PARAMETERS

XITERABLE DATA TO TRANSFORM MUST FULFILL INPUT REQUIREMENTS OF FIRST STEP OF THE PIPELINE

RETURNS

XTARRAYLIKE SHAPE NSAMPLES NTRANSFORMEDFEATURES

2224 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES USING SKLEARNPIPELINEPIPELINE

- EXPLICIT FEATURE MAP APPROXIMATION FOR RBF KERNELS
- FEATURE AGGLOMERATION VS UNIVARIATE SELECTION
- CONCATENATING MULTIPLE FEATURE EXTRACTION METHODS
- PIPELINING CHAINING A PCA AND A LOGISTIC REGRESSION
- COLUMN TRANSFORMER WITH MIXED TYPES
- SELECTING DIMENSIONALITY REDUCTION WITH PIPELINE AND GRIDSEARCHCV
- COLUMN TRANSFORMER WITH HETEROGENEOUS DATA SOURCES
- UNDERFITTING VS OVERFITTING
- BALANCE MODEL COMPLEXITY AND CROSSVALIDATED SCORE
- SAMPLE PIPELINE FOR TEXT FEATURE EXTRACTION AND EVALUATION
- COMPARING NEAREST NEIGHBORS WITH AND WITHOUT NEIGHBORHOOD COMPONENTS ANALYSIS
- RESTRICTED BOLTZMANN MACHINE FEATURES FOR DIGIT CLASSIFICATION
- SVMANOVA SVM WITH UNIVARIATE FEATURE SELECTION
- CLASSIFICATION OF TEXT DOCUMENTS USING SPARSE FEATURES

PIPELINEMAKEPIPELINE STEPS KWARGS CONSTRUCT A PIPELINE FROM THE GIVEN ESTIMATORS

PIPELINEMAKEUNION TRANSFORMERS KWARGS CONSTRUCT A FEATUREUNION FROM THE GIVEN TRANSFORMERS

6323SKLEARNPIPELINE MAKEPIPELINE

SKLEARNPIPELINE MAKEPIPELINE STEPS KWARGS

CONSTRUCT A PIPELINE FROM THE GIVEN ESTIMATORS

THIS IS A SHORTHAND FOR THE PIPELINE CONSTRUCTOR IT DOES NOT REQUIRE AND DOES NOT PERMIT NAMING THE ESTIMATORS INSTEAD THEIR NAMES WILL BE SET TO THE LOWERCASE OF THEIR TYPES AUTOMATICALLY

PARAMETERS

STEPS LIST OF ESTIMATORS

MEMORY NONE STR OR OBJECT WITH THE JOBLIBMEMORY INTERFACE OPTIONAL USED TO CACHE THE FIT

TED TRANSFORMERS OF THE PIPELINE BY DEFAULT NO CACHING IS PERFORMED IF A STRING IS GIVEN

IT IS THE PATH TO THE CACHING DIRECTORY ENABLING CACHING TRIGGERS A CLONE OF THE TRANSFORM

ERS BEFORE FITTING THEREFORE THE TRANSFORMER INSTANCE GIVEN TO THE PIPELINE CANNOT BE IN

SPECTED DIRECTLY USE THE ATTRIBUTE NAMEDSTEPS ORSTEPS TO INSPECT ESTIMATORS WITHIN

THE PIPELINE CACHING THE TRANSFORMERS IS ADVANTAGEOUS WHEN FITTING IS TIME CONSUMING

VERBOSE BOOLEAN OPTIONAL IF TRUE THE TIME ELAPSED WHILE FITTING EACH STEP WILL BE PRINTED AS IT IS COMPLETED

RETURNS

PIPELINE

SEE ALSO

SKLEARNPIPELINEPIPELINE CLASS FOR CREATING A PIPELINE OF TRANSFORMS WITH A FINAL ESTIMATOR

6325SKLEARNPIPELINE PIPELINE 2225

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES

FROM SKLEARNNAIVEBAYES IMPORT GAUSSIANNB  
FROM SKLEARNPREPROCESSING IMPORT STANDARDSCALER  
MAKEPIPELINESTANDARDSCALER GAUSSIANNBPRIORS NONE

PIPELINEMEMORYNONE  
STEPSSTANDARDSCALER  
STANDARDSCALERCOPYTRUE WITHMEANTRUE WITHSTDTRUE  
GAUSSIANNB  
GAUSSIANNBPRIORSNONE VARSMOOTHING1E09  
VERBOSEFALSE

EXAMPLES USING SKLEARNPIPELINEMAKEPIPELINE

- FEATURE TRANSFORMATIONS WITH ENSEMBLES OF TREES
- PIPELINE ANOVA SVM
- IMPUTING MISSING VALUES WITH VARIANTS OF ITERATIVEIMPUTER
- IMPUTING MISSING VALUES BEFORE BUILDING AN ESTIMATOR
- POLYNOMIAL INTERPOLATION
- ROBUST LINEAR ESTIMATOR FITTING
- DIMENSIONALITY REDUCTION WITH NEIGHBORHOOD COMPONENTS ANALYSIS
- USING FUNCTIONTRANSFORMER TO SELECT COLUMNS
- IMPORTANCE OF FEATURE SCALING
- FEATURE DISCRETIZATION
- CLUSTERING TEXT DOCUMENTS USING KMEANS

6324SKLEARNPIPELINE MAKEUNION

SKLEARNPIPELINE MAKEUNION TRANSFORMERS KWARGS

CONSTRUCT A FEATUREUNION FROM THE GIVEN TRANSFORMERS

THIS IS A SHORTHAND FOR THE FEATUREUNION CONSTRUCTOR IT DOES NOT REQUIRE AND DOES NOT PERMIT NAMING THE TRANSFORMERS INSTEAD THEY WILL BE GIVEN NAMES AUTOMATICALLY BASED ON THEIR TYPES IT ALSO DOES NOT ALLOW WEIGHTING

PARAMETERS

TRANSFORMERS LIST OF ESTIMATORS

NJOBS INT OR NONE OPTIONAL DEFAULTNONE NUMBER OF JOBS TO RUN IN PARALLEL NONE MEANS 1

UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL PROCESSORS SEE

GLOSSARY FOR MORE DETAILS

VERBOSE BOOLEAN OPTIONALDEFAULTFALSE IF TRUE THE TIME ELAPSED WHILE FITTING EACH TRANS FORMER WILL BE PRINTED AS IT IS COMPLETED

RETURNS

FFEATUREUNION

2226 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

SEE ALSO

SKLEARNPIPELINEFEATUREUNION CLASS FOR CONCATENATING THE RESULTS OF MULTIPLE TRANSFORMER OBJECTS

EXAMPLES

FROM SKLEARNDECOMPOSITION IMPORT PCA TRUNCATEDSVD

FROM SKLEARNPIPELINE IMPORT MAKEUNION

MAKEUNIONPCA TRUNCATEDSVD

FEATUREUNIONNJOBSNONE

TRANSFORMERLISTPCA

PCACOPYTRUE ITERATEDPOWERAUTO

NCOMPONENTSNONE RANDOMSTATENONE

SVDSOLVERAUTO TOL00 WHITENFALSE

TRUNCATEDSVD

TRUNCATEDSVDALGORITHMRANDOMIZED

NCOMPONENTS2 NITER5

RANDOMSTATENONE TOL00

TRANSFORMERWEIGHTSNONE VERBOSEFALSE

EXAMPLES USING SKLEARNPIPELINEMAKEUNION

- IMPUTING MISSING VALUES BEFORE BUILDING AN ESTIMATOR

633SKLEARNINSPECTION INSPECTION

THESKLEARNINSPECTION MODULE INCLUDES TOOLS FOR MODEL INSPECTION

INSPECTIONPARTIALDEPENDENCE ESTIMATOR X

PARTIAL DEPENDENCE OF FEATURES

INSPECTIONPLOTPARTIALDEPENDENCE

PARTIAL DEPENDENCE PLOTS

6331SKLEARNINSPECTION PARTIALDEPENDENCE

SKLEARNINSPECTION PARTIALDEPENDENCE ESTIMATOR XFEATURES RESPONSEMETHOD'AUTO'

PERCENTILES005 095 GRIDRESOLUTION100

METHOD'AUTO'

PARTIAL DEPENDENCE OF FEATURES

PARTIAL DEPENDENCE OF A FEATURE OR A SET OF FEATURES CORRESPONDS TO THE AVERAGE RESPONSE OF AN ESTIMATOR FOR EACH POSSIBLE VALUE OF THE FEATURE

READ MORE IN THE USER GUIDE

PARAMETERS

ESTIMATOR BASEESTIMATOR A FITTED ESTIMATOR OBJECT IMPLEMENTING PREDICT PREDICTPROBA OR DECISIONFUNCTION MULTIOUTPUTMULTICLASS CLASSIFIERS ARE NOT SUPPORTED

XARRAYLIKE SHAPE NSAMPLES NFEATURES XIS USED BOTH TO GENERATE A GRID OF VALUES FOR THE FEATURES AND TO COMPUTE THE AVERAGED PREDICTIONS WHEN METHOD IS 'BRUTE'

633SKLEARNINSPECTION INSPECTION 2227

SCIKITLEARN USER GUIDE RELEASE 0213

FEATURES LIST OR ARRAYLIKE OF INT THE TARGET FEATURES FOR WHICH THE PARTIAL DEPENDENCY SHOULD BE COMPUTED

RESPONSEMETHOD 'AUTO' 'PREDICTPROBA' OR 'DECISIONFUNCTION' OPTIONAL DEFAULT'AUTO' SPECIFIES WHETHER TO USE PREDICTPROBA ORDECISIONFUNCTION AS THE TARGET RESPONSE FOR REGRESSORS THIS PARAMETER IS IGNORED AND THE RESPONSE IS ALWAYS THE OUTPUT OF PREDICT BY DEFAULT PREDICTPROBA IS TRIED FIRST AND WE REVERT TO DECISIONFUNCTION IF IT DOESN'T EXIST IF METHOD IS 'RECURSION' THE RESPONSE IS ALWAYS THE OUTPUT OF DECISIONFUNCTION PERCENTILES TUPLE OF FLOAT OPTIONAL DEFAULT005 095 THE LOWER AND UPPER PERCENTILE USED TO CREATE THE EXTREME VALUES FOR THE GRID MUST BE IN 0 1

GRIDRESOLUTION INT OPTIONAL DEFAULT100 THE NUMBER OF EQUALLY SPACED POINTS ON THE GRID FOR EACH TARGET FEATURE

METHOD STR OPTIONAL DEFAULT'AUTO' THE METHOD USED TO CALCULATE THE AVERAGED PREDICTIONS

- 'RECURSION' IS ONLY SUPPORTED FOR OBJECTS INHERITING FROM BASEGRADIENTBOOSTING BUT IS MORE EFFICIENT IN TERMS OF SPEED WITH THIS METHOD XIS ONLY USED TO BUILD THE GRID AND THE PARTIAL DEPENDENCES ARE COMPUTED USING THE TRAINING DATA THIS METHOD DOES NOT ACCOUNT FOR THE INIT PREDICOR OF THE BOOSTING PROCESS WHICH MAY LEAD TO INCORRECT VALUES SEE WARNING BELOW WITH THIS METHOD THE TARGET RESPONSE OF A CLASSIFIER IS ALWAYS THE DECISION FUNCTION NOT THE PREDICTED PROBABILITIES

- 'BRUTE' IS SUPPORTED FOR ANY ESTIMATOR BUT IS MORE COMPUTATIONALLY INTENSIVE
- IF 'AUTO' THEN 'RECURSION' WILL BE USED FOR BASEGRADIENTBOOSTING ESTIMATORS WITH INITNONE AND 'BRUTE' FOR ALL OTHER

RETURNS  
AVERAGEDPREDICTIONS NDARRAY SHAPE NOUTPUTS LENVALUES0 LENVALUES1 THE PRE DITIONS FOR ALL THE POINTS IN THE GRID AVERAGED OVER ALL SAMPLES IN X OR OVER THE TRAINING DATA IFMETHOD IS 'RECURSION' NOUTPUTS CORRESPONDS TO THE NUMBER OF CLASSES IN A MULTI CLASS SETTING OR TO THE NUMBER OF TASKS FOR MULTIOUTPUT REGRESSION FOR CLASSICAL REGRESSION AND BINARY CLASSIFICATION NOUTPUTS1 NVALUESFEATUREJ CORRESPONDS TO THE SIZEVALUESJ

VALUES SEQ OF 1D NDARRAYS THE VALUES WITH WHICH THE GRID HAS BEEN CREATED THE GENERATED GRID IS A CARTESIAN PRODUCT OF THE ARRAYS IN VALUES LENVALUES LENFEATURES THE SIZE OF EACH ARRAY VALUESJ IS EITHERGRIDRESOLUTION OR THE NUMBER OF UNIQUE VALUES IN X J WHICHEVER IS SMALLER

WARNING THE 'RECURSION' METHOD ONLY WORKS FOR GRADIENT BOOSTING ESTIMATORS AND UNLIKE THE 'BRUTE' METHOD IT DOES NOT ACCOUNT FOR THE INIT PREDICTOR OF THE BOOSTING PROCESS IN PRACTICE THIS WILL PRODUCE THE SAME VALUES AS 'BRUTE' UP TO A CONSTANT OFFSET IN THE TARGET RESPONSE PROVIDED THAT INIT IS A CONSANT ESTIMATOR WHICH IS THE DEFAULT HOWEVER AS SOON AS INIT IS NOT A CONSTANT ESTIMATOR THE PARTIAL DEPENDENCE VALUES ARE INCORRECT FOR 'RECURSION'

SEE ALSO  
SKLEARNINSPECTIONPLOTPARTIALDEPENDENCE PLOT PARTIAL DEPENDENCE  
2228 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES

X 0 0 2 1 0 0

Y 0 1

FROM SKLEARNENSEMBLE IMPORT GRADIENTBOOSTINGCLASSIFIER

GB GRADIENTBOOSTINGCLASSIFIERRANDOMSTATE0FITX Y

PARTIALDEPENDENCEGB FEATURES0 XX PERCENTILES0 1

GRIDRESOLUTION2

ARRAY452 452 ARRAY 0 1

EXAMPLES USING SKLEARNINSPECTIONPARTIALDEPENDENCE

•PARTIAL DEPENDENCE PLOTS

6332SKLEARNINSPECTION PLOTPARTIALDEPENDENCE

SKLEARNINSPECTION PLOTPARTIALDEPENDENCE ESTIMATOR XFEATURES FEATURENAMESNONE

TARGETNONE RESPONSEMETHOD'AUTO'

NCOLS3 GRIDRESOLUTION100 PER

CENTILES005 095 METHOD'AUTO'

NJOBSNONE VERBOSE0 FIGNONE

LINEKWNONE CONTOURKWNONE

PARTIAL DEPENDENCE PLOTS

THELENFEATURES PLOTS ARE ARRANGED IN A GRID WITH NCOLS COLUMNS TWOWAY PARTIAL DEPENDENCE PLOTS

ARE PLOTTED AS CONTOUR PLOTS

READ MORE IN THE USER GUIDE

PARAMETERS

ESTIMATOR BASEESTIMATOR A FITTED ESTIMATOR OBJECT IMPLEMENTING PREDICT PREDICTPROBA OR

DECISIONFUNCTION MULTIOUTPUTMULTICLASS CLASSIFIERS ARE NOT SUPPORTED

XARRAYLIKE SHAPE NSAMPLES NFEATURES THE DATA TO USE TO BUILD THE GRID OF VALUES ON

WHICH THE DEPENDENCE WILL BE EVALUATED THIS IS USUALLY THE TRAINING DATA

FEATURES LIST OF INT STR PAIR OF INT PAIR OF STR THE TARGET FEATURES FOR WHICH TO CREATE THE

PDPS IF FEATURES<sub>i</sub> IS AN INT OR A STRING A ONEWAY PDP IS CREATED IF FEATURES<sub>i</sub> IS A TUPLE A

TWOWAY PDP IS CREATED EACH TUPLE MUST BE OF SIZE 2 IF ANY ENTRY IS A STRING THEN IT MUST

BE INFEATURENAMES

FEATURENAMES SEQ OF STR SHAPE NFEATURES OPTIONAL NAME OF EACH FEATURE FEA

TURENAMES<sub>i</sub> HOLDS THE NAME OF THE FEATURE WITH INDEX <sub>i</sub> BY DEFAULT THE NAME OF THE FEATURE

CORRESPONDS TO THEIR NUMERICAL INDEX

TARGET INT OPTIONAL DEFAULTNONE

• IN A MULTICLASS SETTING SPECIFIES THE CLASS FOR WHICH THE PDPS SHOULD BE COMPUTED NOTE

THAT FOR BINARY CLASSIFICATION THE POSITIVE CLASS INDEX 1 IS ALWAYS USED

• IN A MULTIOUTPUT SETTING SPECIFIES THE TASK FOR WHICH THE PDPS SHOULD BE COMPUTED

IGNORED IN BINARY CLASSIFICATION OR CLASSICAL REGRESSION SETTINGS

RESPONSEMETHOD 'AUTO' 'PREDICTPROBA' OR 'DECISIONFUNCTION' OPTIONAL DEFAULT'AUTO'

SPECIFIES WHETHER TO USE PREDICTPROBA ORDECISIONFUNCTION AS THE TARGET RESPONSE FOR

REGRESSORS THIS PARAMETER IS IGNORED AND THE RESPONSE IS ALWAYS THE OUTPUT OF PREDICT BY

6333SKLEARNINSPECTION INSPECTION 2229

SCIKITLEARN USER GUIDE RELEASE 0213

DEFAULT PREDICTPROBA IS TRIED FIRST AND WE REVERT TO DECISIONFUNCTION IF IT DOESN'T EXIST IF METHOD IS 'RECURSION' THE RESPONSE IS ALWAYS THE OUTPUT OF DECISIONFUNCTION  
NCOLS INT OPTIONAL DEFAULT3 THE MAXIMUM NUMBER OF COLUMNS IN THE GRID PLOT  
GRIDRESOLUTION INT OPTIONAL DEFAULT100 THE NUMBER OF EQUALLY SPACED POINTS ON THE AXES OF THE PLOTS FOR EACH TARGET FEATURE

PERCENTILES TUPLE OF FLOAT OPTIONAL DEFAULT005 095 THE LOWER AND UPPER PERCENTILE USED TO CREATE THE EXTREME VALUES FOR THE PDP AXES MUST BE IN 0 1

METHOD STR OPTIONAL DEFAULT'AUTO' THE METHOD TO USE TO CALCULATE THE PARTIAL DEPENDENCE PREDICTIONS

- 'RECURSION' IS ONLY SUPPORTED FOR OBJECTS INHERITING FROM BASEGRADIENTBOOSTING BUT IS MORE EFFICIENT IN TERMS OF SPEED WITH THIS METHOD XIS OPTIONAL AND IS ONLY USED TO BUILD THE GRID AND THE PARTIAL DEPENDENCES ARE COMPUTED USING THE TRAINING DATA THIS METHOD DOES NOT ACCOUNT FOR THE INIT PREDICOR OF THE BOOSTING PROCESS WHICH MAY LEAD TO INCORRECT VALUES SEE WARNING BELOW WITH THIS METHOD THE TARGET RESPONSE OF A CLASSIFIER IS ALWAYS THE DECISION FUNCTION NOT THE PREDICTED PROBABILITIES

- 'BRUTE' IS SUPPORTED FOR ANY ESTIMATOR BUT IS MORE COMPUTATIONALLY INTENSIVE
- IF 'AUTO' THEN 'RECURSION' WILL BE USED FOR BASEGRADIENTBOOSTING ESTIMATORS WITH INITNONE AND 'BRUTE' FOR ALL OTHER

UNLIKE THE 'BRUTE' METHOD 'RECURSION' DOES NOT ACCOUNT FOR THE INIT PREDICTOR OF THE BOOSTING PROCESS IN PRACTICE THIS STILL PRODUCES THE SAME PLOTS UP TO A CONSTANT OFFSET IN THE TARGET RESPONSE

NJOBS INT OPTIONAL DEFAULTNONE THE NUMBER OF CPUS TO USE TO COMPUTE THE PARTIAL DEPENDENCESNONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

VERBOSE INT OPTIONAL DEFAULT0 VERBOSE OUTPUT DURING PD COMPUTATIONS

FIGMATPLOTLIB FIGURE OBJECT OPTIONAL DEFAULTNONE A FIGURE OBJECT ONTO WHICH THE PLOTS WILL BE DRAWN AFTER THE FIGURE HAS BEEN CLEARED BY DEFAULT A NEW ONE IS CREATED

LINEKW DICT OPTIONAL DICT WITH KEYWORDS PASSED TO THE MATPLOTLIBPYPLOTPLOT CALL FOR ONEWAY PARTIAL DEPENDENCE PLOTS

CONTOURKW DICT OPTIONAL DICT WITH KEYWORDS PASSED TO THE MATPLOTLIBPYPLOTPLOT CALL FOR TWOWAY PARTIAL DEPENDENCE PLOTS

WARNING THE 'RECURSION' METHOD ONLY WORKS FOR GRADIENT BOOSTING ESTIMATORS AND UNLIKE THE 'BRUTE' METHOD IT DOES NOT ACCOUNT FOR THE INIT PREDICTOR OF THE BOOSTING PROCESS IN PRACTICE THIS WILL PRODUCE THE SAME VALUES AS 'BRUTE' UP TO A CONSTANT OFFSET IN THE TARGET RESPONSE PROVIDED THAT INIT IS A CONSANT ESTIMATOR WHICH IS THE DEFAULT HOWEVER AS SOON AS INIT IS NOT A CONSTANT ESTIMATOR THE PARTIAL DEPENDENCE VALUES ARE INCORRECT FOR 'RECURSION'

SEE ALSO

SKLEARNINSPECTIONPARTIALDEPENDENCE RETURN RAW PARTIAL DEPENDENCE VALUES

2230 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES

```
FROM SKLEARNDATASETS IMPORT MAKEFRIEDMAN1
FROM SKLEARNENSEMBLE IMPORT GRADIENTBOOSTINGREGRESSOR
X Y MAKEFRIEDMAN1
CLF GRADIENTBOOSTINGREGRESSORNESTIMATORS10FITX Y
PLOTPARTIALDEPENDENCECLF X 0 0 1
EXAMPLES USING SKLEARNINSPECTIONPLOTPARTIALDEPENDENCE
•PARTIAL DEPENDENCE PLOTS
634SKLEARNPREPROCESSING PREPROCESSING AND NORMALIZATION
THESKLEARNPREPROCESSING MODULE INCLUDES SCALING CENTERING NORMALIZATION BINARIZATION AND IMPUTATION
METHODS
USER GUIDE SEE THE PREPROCESSING DATA SECTION FOR FURTHER DETAILS
PREPROCESSINGBINARIZER THRESHOLD COPY BINARIZE DATA SET FEATURE VALUES TO 0 OR 1 ACCORDING TO A
THRESHOLD
PREPROCESSINGFUNCTIONTRANSFORMER FUNC
    CONSTRUCTS A TRANSFORMER FROM AN ARBITRARY CALLABLE
PREPROCESSINGKBINSDISCRETIZER NBINS
    BIN CONTINUOUS DATA INTO INTERVALS
PREPROCESSINGKERNELCENTERER CENTER A KERNEL MATRIX
PREPROCESSINGLABELBINARIZER NEGLABEL
    BINARIZE LABELS IN A ONEVSALL FASHION
PREPROCESSINGLABELENCODER ENCODE LABELS WITH VALUE BETWEEN 0 AND NCLASSES1
PREPROCESSINGMULTILABELBINARIZER CLASSES
    TRANSFORM BETWEEN ITERABLE OF ITERABLES AND A MULTILABEL
FORMAT
PREPROCESSINGMAXABSSCALER COPY SCALE EACH FEATURE BY ITS MAXIMUM ABSOLUTE VALUE
PREPROCESSINGMINMAXSCALER FEATURERANGE
COPYTRANSFORMS FEATURES BY SCALING EACH FEATURE TO A GIVEN
RANGE
PREPROCESSINGNORMALIZER NORM COPY NORMALIZE SAMPLES INDIVIDUALLY TO UNIT NORM
PREPROCESSINGONEHOTENCODER NVALUES ENCODE CATEGORICAL INTEGER FEATURES AS A ONEHOT NUMERIC
ARRAY
PREPROCESSINGORDINALENCODER CATEGORIES
DTYPEENCODE CATEGORICAL FEATURES AS AN INTEGER ARRAY
PREPROCESSINGPOLYNOMIALFEATURES DEGREE
    GENERATE POLYNOMIAL AND INTERACTION FEATURES
PREPROCESSINGPOWERTRANSFORMER METHOD
    APPLY A POWER TRANSFORM FEATUREWISE TO MAKE DATA MORE
GAUSSIANLIKE
PREPROCESSINGQUANTILETRANSFORMER TRANSFORM FEATURES USING QUANTILES INFORMATION
PREPROCESSINGROBUSTSCALER WITHCENTERING
    SCALE FEATURES USING STATISTICS THAT ARE ROBUST TO OUTLIERS
PREPROCESSINGSTANDARDSCALER COPY STANDARDIZE FEATURES BY REMOVING THE MEAN AND SCALING TO
UNIT VARIANCE
634SKLEARNPREPROCESSING PREPROCESSING AND NORMALIZATION 2231
```

SCIKITLEARN USER GUIDE RELEASE 0213

6341SKLEARNPREPROCESSING BINARIZER

CLASSSSKLEARNPREPROCESSING BINARIZER THRESHOLD00 COPYTRUE

BINARIZE DATA SET FEATURE VALUES TO 0 OR 1 ACCORDING TO A THRESHOLD

VALUES GREATER THAN THE THRESHOLD MAP TO 1 WHILE VALUES LESS THAN OR EQUAL TO THE THRESHOLD MAP TO 0 WITH THE

DEFAULT THRESHOLD OF 0 ONLY POSITIVE VALUES MAP TO 1

BINARIZATION IS A COMMON OPERATION ON TEXT COUNT DATA WHERE THE ANALYST CAN DECIDE TO ONLY CONSIDER THE PRESENCE

OR ABSENCE OF A FEATURE RATHER THAN A QUANTIFIED NUMBER OF OCCURRENCES FOR INSTANCE

IT CAN ALSO BE USED AS A PREPROCESSING STEP FOR ESTIMATORS THAT CONSIDER BOOLEAN RANDOM VARIABLES EG MODELLED

USING THE BERNOULLI DISTRIBUTION IN A BAYESIAN SETTING

READ MORE IN THE USER GUIDE

PARAMETERS

THRESHOLD FLOAT OPTIONAL 00 BY DEFAULT FEATURE VALUES BELOW OR EQUAL TO THIS ARE REPLACED BY

0 ABOVE IT BY 1 THRESHOLD MAY NOT BE LESS THAN 0 FOR OPERATIONS ON SPARSE MATRICES

COPY BOOLEAN OPTIONAL DEFAULT TRUE SET TO FALSE TO PERFORM INPLACE BINARIZATION AND AVOID A

COPY IF THE INPUT IS ALREADY A NUMPY ARRAY OR A SCIPYSPARSE CSR MATRIX

SEE ALSO

BINARIZE EQUIVALENT FUNCTION WITHOUT THE ESTIMATOR API

NOTES

IF THE INPUT IS A SPARSE MATRIX ONLY THE NONZERO VALUES ARE SUBJECT TO UPDATE BY THE BINARIZER CLASS

THIS ESTIMATOR IS STATELESS BESIDES CONSTRUCTOR PARAMETERS THE FIT METHOD DOES NOTHING BUT IS USEFUL WHEN USED

IN A PIPELINE

EXAMPLES

FROM SKLEARNPREPROCESSING IMPORT BINARIZER

X 1 1 2

2 0 0

0 1 1

TRANSFORMER BINARIZERFITX FIT DOES NOTHING

TRANSFORMER

BINARIZERCOPYTRUE THRESHOLD00

TRANSFORMERTRANSFORMX

ARRAY1 0 1

1 0 0

0 1 0

METHODS

FITSELF X Y DO NOTHING AND RETURN THE ESTIMATOR UNCHANGED

FITTRANSFORM SELF X Y FIT TO DATA THEN TRANSFORM IT

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

CONTINUED ON NEXT PAGE

2232 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

TABLE 6247 – CONTINUED FROM PREVIOUS PAGE

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF X COPY BINARIZE EACH ELEMENT OF X

INIT SELFTHRESHOLD00 COPYTRUE

FITSELFXYNONE

DO NOTHING AND RETURN THE ESTIMATOR UNCHANGED

THIS METHOD IS JUST THERE TO IMPLEMENT THE USUAL API AND HENCE WORK IN PIPELINES

PARAMETERS

XARRAYLIKE

FITTRANSFORM SELFXYNONE FITPARAMS

FIT TO DATA THEN TRANSFORM IT

FITS TRANSFORMER TO X AND Y WITH OPTIONAL PARAMETERS FITPARAMS AND RETURNS A TRANSFORMED VERSION OF X

PARAMETERS

XNUMPY ARRAY OF SHAPE NSAMPLES NFEATURES TRAINING SET

YNUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES

RETURNS

XNEW NUMPY ARRAY OF SHAPE NSAMPLES NFEATURESNEW TRANSFORMED ARRAY

GETPARAMS SELFDEEPTTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE

PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT

OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELFXCOPYNONE

BINARIZE EACH ELEMENT OF X

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE DATA TO BINARIZE ELEMENT

BY ELEMENT SCIPYSPARSE MATRICES SHOULD BE IN CSR FORMAT TO AVOID AN UNNECESSARY COPY

COPY BOOL COPY THE INPUT X OR NOT

634SKLEARNPREPROCESSING PREPROCESSING AND NORMALIZATION 2233

SCIKITLEARN USER GUIDE RELEASE 0213

63425KLEARNPREPROCESSING FUNCTIONTRANSFORMER

CLASSSSKLEARNPREPROCESSING FUNCTIONTRANSFORMER FUNCNONE INVERSEFUNCNONE VAL

IDATENONE ACCEPTSPARSEFALSE

PASSY'DEPRECATED'

CHECKINVERSETRUE KWARGSNONE

INVKWARGSNONE

CONSTRUCTS A TRANSFORMER FROM AN ARBITRARY CALLABLE

A FUNCTIONTRANSFORMER FORWARDS ITS X AND OPTIONALLY Y ARGUMENTS TO A USERDEFINED FUNCTION OR FUNCTION OBJECT AND RETURNS THE RESULT OF THIS FUNCTION THIS IS USEFUL FOR STATELESS TRANSFORMATIONS SUCH AS TAKING THE LOG OF FREQUENCIES DOING CUSTOM SCALING ETC

NOTE IF A LAMBDA IS USED AS THE FUNCTION THEN THE RESULTING TRANSFORMER WILL NOT BE PICKLEABLE

NEW IN VERSION 017

READ MORE IN THE USER GUIDE

PARAMETERS

FUNC CALLABLE OPTIONAL DEFAULTNONE THE CALLABLE TO USE FOR THE TRANSFORMATION THIS WILL BE PASSED THE SAME ARGUMENTS AS TRANSFORM WITH ARGS AND KWARGS FORWARDED IF FUNC IS NONE THEN FUNC WILL BE THE IDENTITY FUNCTION

INVERSEFUNC CALLABLE OPTIONAL DEFAULTNONE THE CALLABLE TO USE FOR THE INVERSE TRANSFORMATION THIS WILL BE PASSED THE SAME ARGUMENTS AS INVERSE TRANSFORM WITH ARGS AND KWARGS FORWARDED IF INVERSEFUNC IS NONE THEN INVERSEFUNC WILL BE THE IDENTITY FUNCTION

VALIDATE BOOL OPTIONAL DEFAULTTRUE INDICATE THAT THE INPUT X ARRAY SHOULD BE CHECKED BEFORE CALLINGFUNC THE POSSIBILITIES ARE

- IF FALSE THERE IS NO INPUT VALIDATION
- IF TRUE THEN X WILL BE CONVERTED TO A 2DIMENSIONAL NUMPY ARRAY OR SPARSE MATRIX IF THE CONVERSION IS NOT POSSIBLE AN EXCEPTION IS RAISED

DEPRECATED SINCE VERSION 020 VALIDATETRUE AS DEFAULT WILL BE REPLACED BY

VALIDATEFALSE IN 022

ACCEPTSPARSE BOOLEAN OPTIONAL INDICATE THAT FUNC ACCEPTS A SPARSE MATRIX AS INPUT IF VALIDATE IS FALSE THIS HAS NO EFFECT OTHERWISE IF ACCEPTSPARSE IS FALSE SPARSE MATRIX INPUTS WILL CAUSE AN EXCEPTION TO BE RAISED

PASSY BOOL OPTIONAL DEFAULTFALSE INDICATE THAT TRANSFORM SHOULD FORWARD THE Y ARGUMENT TO THE INNER CALLABLE

DEPRECATED SINCE VERSION 019

CHECKINVERSE BOOL DEFAULTTRUE WHETHER TO CHECK THAT OR FUNC FOLLOWED BY INVERSEFUNC LEADS TO THE ORIGINAL INPUTS IT CAN BE USED FOR A SANITY CHECK RAISING A WARNING WHEN THE CONDITION IS NOT FULFILLED

NEW IN VERSION 020

KWARGS DICT OPTIONAL DICTIONARY OF ADDITIONAL KEYWORD ARGUMENTS TO PASS TO FUNC

INVKWARGS DICT OPTIONAL DICTIONARY OF ADDITIONAL KEYWORD ARGUMENTS TO PASS TO IN VERSEFUNC

2234 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

METHODS

FITSELF X Y FIT TRANSFORMER BY CHECKING X

FITTRANSFORM SELF X Y FIT TO DATA THEN TRANSFORM IT

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

INVERSETRANSFORM SELF X TRANSFORM X USING THE INVERSE FUNCTION

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF X TRANSFORM X USING THE FORWARD FUNCTION

INIT SELF FUNCNONE INVERSEFUNCNONE VALIDATENONE ACCEPTSPARSEFALSE

PASSY'DEPRECATED' CHECKINVERSETRUE KWARGSNONE INVKWARGSNONE

FITSELFXYNONE

FIT TRANSFORMER BY CHECKING X

IFVALIDATE ISTRUE XWILL BE CHECKED

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES INPUT ARRAY

RETURNS

SELF

FITTRANSFORM SELFXYNONE FITPARAMS

FIT TO DATA THEN TRANSFORM IT

FITS TRANSFORMER TO X AND Y WITH OPTIONAL PARAMETERS FITPARAMS AND RETURNS A TRANSFORMED VERSION OF X

PARAMETERS

XNUMPY ARRAY OF SHAPE NSAMPLES NFEATURES TRAINING SET

YNUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES

RETURNS

XNEW NUMPY ARRAY OF SHAPE NSAMPLES NFEATURESNEW TRANSFORMED ARRAY

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

INVERSETRANSFORM SELFXY

TRANSFORM X USING THE INVERSE FUNCTION

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES INPUT ARRAY

RETURNS

XOUT ARRAYLIKE SHAPE NSAMPLES NFEATURES TRANSFORMED INPUT

634SKLEARNPREPROCESSING PREPROCESSING AND NORMALIZATION 2235

SCIKITLEARN USER GUIDE RELEASE 0213

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELF

TRANSFORM X USING THE FORWARD FUNCTION

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES INPUT ARRAY

RETURNS

XOUT ARRAYLIKE SHAPE NSAMPLES NFEATURES TRANSFORMED INPUT

EXAMPLES USING SKLEARNPREPROCESSINGFUNCTIONTRANSFORMER

- USING FUNCTIONTRANSFORMER TO SELECT COLUMNS

6343SKLEARNPREPROCESSING KBINSDISCRETIZER

CLASSSSKLEARNPREPROCESSING KBINSDISCRETIZER NBINS5 ENCODE’ONEHOT’ STRAT

EGY’QUANTILE’

BIN CONTINUOUS DATA INTO INTERVALS

READ MORE IN THE USER GUIDE

PARAMETERS

NBINS INT OR ARRAYLIKE SHAPE NFEATURES DEFAULT5 THE NUMBER OF BINS TO PRODUCE

RAISES VALUEERROR IF NBINS 2

ENCODE ’ONEHOT’ ’ONEHOTDENSE’ ’ORDINAL’ DEFAULT’ONEHOT’ METHOD USED TO ENCODE THE TRANSFORMED RESULT

ONEHOT ENCODE THE TRANSFORMED RESULT WITH ONEHOT ENCODING AND RETURN A SPARSE MATRIX

IGNORED FEATURES ARE ALWAYS STACKED TO THE RIGHT

ONEHOTDENSE ENCODE THE TRANSFORMED RESULT WITH ONEHOT ENCODING AND RETURN A DENSE AR

RAY IGNORED FEATURES ARE ALWAYS STACKED TO THE RIGHT

ORDINAL RETURN THE BIN IDENTIFIER ENCODED AS AN INTEGER VALUE

STRATEGY ’UNIFORM’ ’QUANTILE’ ’KMEANS’ DEFAULT’QUANTILE’ STRATEGY USED TO DEFINE THE WIDTHS OF THE BINS

UNIFORM ALL BINS IN EACH FEATURE HAVE IDENTICAL WIDTHS

QUANTILE ALL BINS IN EACH FEATURE HAVE THE SAME NUMBER OF POINTS

KMEANS VALUES IN EACH BIN HAVE THE SAME NEAREST CENTER OF A 1D KMEANS CLUSTER

ATTRIBUTES

NBINS INT ARRAY SHAPE NFEATURES NUMBER OF BINS PER FEATURE BINS WHOSE WIDTH ARE TOO SMALL IE 1E8 ARE REMOVED WITH A WARNING

2236 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

BINEDGES ARRAY OF ARRAYS SHAPE NFEATURES THE EDGES OF EACH BIN CONTAIN ARRAYS OF VARYING SHAPES NBINS IGNORED FEATURES WILL HAVE EMPTY ARRAYS

SEE ALSO

SKLEARNPREPROCESSINGBINARIZER CLASS USED TO BIN VALUES AS 0OR1BASED ON A PARAMETER THRESHOLD

NOTES

IN BIN EDGES FOR FEATURE I THE FIRST AND LAST VALUES ARE USED ONLY FOR INVERSETRANSFORM DURING TRANSFORM BIN EDGES ARE EXTENDED TO

NPCONCATENATENPINF BINEDGESI11 NPINF

YOU CAN COMBINE KBINSDISCRETIZER WITHSKLEARNCOMPOSECOLUMNTRANSFORMER IF YOU ONLY WANT TO PREPROCESS PART OF THE FEATURES

KBINSDISCRETIZER MIGHT PRODUCE CONSTANT FEATURES EG WHEN ENCODE ONEHOT AND CERTAIN BINS DO NOT CONTAIN ANY DATA THESE FEATURES CAN BE REMOVED WITH FEATURE SELECTION ALGORITHMS EG SKLEARN FEATURESELECTIONVARIANCETHRESHOLD

EXAMPLES

X 2 1 4 1

1 2 3 05

0 3 2 05

1 4 1 2

EST KBINSDISCRETIZERNBINS3 ENCODEORDINAL STRATEGYUNIFORM

ESTFITX

KBINSDISCRETIZER

XT ESTTRANSFORMX

XT

ARRAY 0 0 0 0

1 1 1 0

2 2 2 1

2 2 2 2

SOMETIMES IT MAY BE USEFUL TO CONVERT THE DATA BACK INTO THE ORIGINAL FEATURE SPACE THE INVERSETRANSFORM FUNCTION CONVERTS THE BINNED DATA INTO THE ORIGINAL FEATURE SPACE EACH VALUE WILL BE EQUAL TO THE MEAN OF THE TWO BIN EDGES

ESTBINEDGES0

ARRAY2 1 0 1

ESTINVERSETRANSFORMXT

ARRAY15 15 35 05

05 25 25 05

05 35 15 05

05 35 15 15

METHODS

634SKLEARNPREPROCESSING PREPROCESSING AND NORMALIZATION 2237

SCIKITLEARN USER GUIDE RELEASE 0213

FITSELF X Y FITS THE ESTIMATOR

FITTRANSFORM SELF X Y FIT TO DATA THEN TRANSFORM IT

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

INVERSETRANSFORM SELF XT TRANSFORMS DISCRETIZED DATA BACK TO ORIGINAL FEATURE SPACE

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF X DISCRETIZES THE DATA

INIT SELFNBINS5 ENCODE'ONEHOT' STRATEGY'QUANTILE'

FITSELFXYNONE

FITS THE ESTIMATOR

PARAMETERS

XNUMERIC ARRAYLIKE SHAPE NSAMPLES NFEATURES DATA TO BE DISCRETIZED

YIGNORED

RETURNS

SELF

FITTRANSFORM SELFXYNONE FITPARAMS

FIT TO DATA THEN TRANSFORM IT

FITS TRANSFORMER TO X AND Y WITH OPTIONAL PARAMETERS FITPARAMS AND RETURNS A TRANSFORMED VERSION OF X

PARAMETERS

XNUMPY ARRAY OF SHAPE NSAMPLES NFEATURES TRAINING SET

YNUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES

RETURNS

XNEW NUMPY ARRAY OF SHAPE NSAMPLES NFEATURESNEW TRANSFORMED ARRAY

GETPARAMS SELFDEEPTTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

INVERSETRANSFORM SELFXT

TRANSFORMS DISCRETIZED DATA BACK TO ORIGINAL FEATURE SPACE

NOTE THAT THIS FUNCTION DOES NOT REGENERATE THE ORIGINAL DATA DUE TO DISCRETIZATION ROUNDING

PARAMETERS

XTNUMERIC ARRAYLIKE SHAPE NSAMPLE NFEATURES TRANSFORMED DATA IN THE BINNED SPACE

RETURNS

XINV NUMERIC ARRAYLIKE DATA IN THE ORIGINAL FEATURE SPACE

2238 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELF

DISCRETIZES THE DATA

PARAMETERS

XNUMERIC ARRAYLIKE SHAPE NSAMPLES NFEATURES DATA TO BE DISCRETIZED

RETURNS

XNUMERIC ARRAYLIKE OR SPARSE MATRIX DATA IN THE BINNED SPACE

EXAMPLES USING SKLEARNPREPROCESSINGKBINSDISCRETIZER

- USING KBINSDISCRETIZER TO DISCRETIZE CONTINUOUS FEATURES
- DEMONSTRATING THE DIFFERENT STRATEGIES OF KBINSDISCRETIZER
- FEATURE DISCRETIZATION

6344SKLEARNPREPROCESSING KERNELCENTERER

CLASSSKLEARNPREPROCESSING KERNELCENTERER

CENTER A KERNEL MATRIX

LET  $K_X Z$  BE A KERNEL DEFINED BY  $\phi(X)T \phi(Z)$  WHERE  $\phi$  IS A FUNCTION MAPPING  $X$  TO A HILBERT SPACE

KERNELCENTERER CENTERS IE NORMALIZE TO HAVE ZERO MEAN THE DATA WITHOUT EXPLICITLY COMPUTING  $\phi(X)$  IT IS EQUIVALENT TO CENTERING  $\phi(X)$  WITH SKLEARNPREPROCESSINGSTANDARDSCALERWITHSTDFALSE

READ MORE IN THE USER GUIDE

EXAMPLES

```
FROM SKLEARNPREPROCESSING IMPORT KERNELCENTERER
FROM SKLEARNMETRICSPAIRWISE IMPORT PAIRWISEKERNELS

X = [[1, 2, 2],
     [2, 1, 3],
     [4, 1, 2]]

K = PAIRWISEKERNELS(X, METRIC=LINEAR)

K
ARRAY 9 2 2
[[2 14 13]
 [2 13 21]]

TRANSFORMER = KERNELCENTERER(FIT)
TRANSFORMER
KERNELCENTERER
TRANSFORMER.TRANSFORM(K)
ARRAY 5 0 5
```

634SKLEARNPREPROCESSING PREPROCESSING AND NORMALIZATION 2239

SCIKITLEARN USER GUIDE RELEASE 0213

0 14 14

5 14 19

METHODS

FITSELF K Y FIT KERNELCENTERER

FITTRANSFORM SELF X Y FIT TO DATA THEN TRANSFORM IT

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF K COPY CENTER KERNEL MATRIX

INIT SELF

FITSELFKYNONE

FIT KERNELCENTERER

PARAMETERS

KNUMPY ARRAY OF SHAPE NSAMPLES NSAMPLES KERNEL MATRIX

RETURNS

SELF RETURNS AN INSTANCE OF SELF

FITTRANSFORM SELFXYNONE FITPARAMS

FIT TO DATA THEN TRANSFORM IT

FITS TRANSFORMER TO X AND Y WITH OPTIONAL PARAMETERS FITPARAMS AND RETURNS A TRANSFORMED VERSION OF X

PARAMETERS

XNUMPY ARRAY OF SHAPE NSAMPLES NFEATURES TRAINING SET

YNUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES

RETURNS

XNEW NUMPY ARRAY OF SHAPE NSAMPLES NFEATURESNEW TRANSFORMED ARRAY

GETPARAMS SELFDEEPTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE

PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT

OF A NESTED OBJECT

RETURNS

SELF

2240 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

TRANSFORM SELF COPY TRUE

CENTER KERNEL MATRIX

PARAMETERS

KNUMPY ARRAY OF SHAPE NSAMPLES1 NSAMPLES2 KERNEL MATRIX

COPY BOOLEAN OPTIONAL DEFAULT TRUE SET TO FALSE TO PERFORM INPLACE COMPUTATION

RETURNS

KNEW NUMPY ARRAY OF SHAPE NSAMPLES1 NSAMPLES2

6345SKLEARNPREPROCESSING LABELBINARIZER

CLASSSKLEARNPREPROCESSING LABELBINARIZER NEGLABEL0 POSLABEL1

SPARSEOUTPUTFALSE

BINARIZE LABELS IN A ONEVSALL FASHION

SEVERAL REGRESSION AND BINARY CLASSIFICATION ALGORITHMS ARE AVAILABLE IN SCIKITLEARN A SIMPLE WAY TO EXTEND THESE ALGORITHMS TO THE MULTICLASS CLASSIFICATION CASE IS TO USE THE SO CALLED ONEVSALL SCHEME

AT LEARNING TIME THIS SIMPLY CONSISTS IN LEARNING ONE REGRESSOR OR BINARY CLASSIFIER PER CLASS IN DOING SO ONE NEEDS TO CONVERT MULTICLASS LABELS TO BINARY LABELS BELONG OR DOES NOT BELONG TO THE CLASS LABELBINARIZER MAKES THIS PROCESS EASY WITH THE TRANSFORM METHOD

AT PREDICTION TIME ONE ASSIGNS THE CLASS FOR WHICH THE CORRESPONDING MODEL GAVE THE GREATEST CONFIDENCE LABELBINARIZER MAKES THIS EASY WITH THE INVERSE TRANSFORM METHOD

READ MORE IN THE USER GUIDE

PARAMETERS

NEGLABEL INT DEFAULT 0 VALUE WITH WHICH NEGATIVE LABELS MUST BE ENCODED

POSLABEL INT DEFAULT 1 VALUE WITH WHICH POSITIVE LABELS MUST BE ENCODED

SPARSEOUTPUT BOOLEAN DEFAULT FALSE TRUE IF THE RETURNED ARRAY FROM TRANSFORM IS DESIRED TO BE IN SPARSE CSR FORMAT

ATTRIBUTES

CLASSES ARRAY OF SHAPE NCLASS HOLDS THE LABEL FOR EACH CLASS

YTYPE STR REPRESENTS THE TYPE OF THE TARGET DATA AS EVALUATED BY

UTILSMULTICLASSTYPEOFTARGET POSSIBLE TYPE ARE 'CONTINUOUS' 'CONTINUOUSMULTIOUTPUT' 'BINARY' 'MULTICLASS' 'MULTICLASSMULTIOUTPUT' 'MULTILABELINDICATOR' AND 'UNKNOWN'

SPARSEINPUT BOOLEAN TRUE IF THE INPUT DATA TO TRANSFORM IS GIVEN AS A SPARSE MATRIX FALSE OTHERWISE

SEE ALSO

LABELBINARIZE FUNCTION TO PERFORM THE TRANSFORM OPERATION OF LABELBINARIZER WITH FIXED CLASSES

SKLEARNPREPROCESSINGONEHOTENCODER ENCODE CATEGORICAL FEATURES USING A ONEHOT AKA ONE OF K SCHEME

6345SKLEARNPREPROCESSING PREPROCESSING AND NORMALIZATION 2241

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES

FROM SKLEARN IMPORT PREPROCESSING

LB PREPROCESSINGLABELBINARIZER

LBFIT1 2 6 4 2

LABELBINARIZERNEGLABEL0 POSLABEL1 SPARSEOUTPUTFALSE

LBCLASSES

ARRAY1 2 4 6

LBTRANSFORM1 6

ARRAY1 0 0 0

0 0 0 1

BINARY TARGETS TRANSFORM TO A COLUMN VECTOR

LB PREPROCESSINGLABELBINARIZER

LBFITTRANSFORMYES NO NO YES

ARRAY1

0

0

1

PASSING A 2D MATRIX FOR MULTILABEL CLASSIFICATION

IMPORT NUMPY AS NP

LBFITNPARRAY0 1 1 1 0 0

LABELBINARIZERNEGLABEL0 POSLABEL1 SPARSEOUTPUTFALSE

LBCLASSES

ARRAY0 1 2

LBTRANSFORM0 1 2 1

ARRAY1 0 0

0 1 0

0 0 1

0 1 0

METHODS

FITSELF Y FIT LABEL BINARIZER

FITTRANSFORM SELF Y FIT LABEL BINARIZER AND TRANSFORM MULTICLASS LABELS TO BINARY LABELS

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

INVERSETRANSFORM SELF Y THRESHOLD TRANSFORM BINARY LABELS BACK TO MULTICLASS LABELS

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF Y TRANSFORM MULTICLASS LABELS TO BINARY LABELS

INIT SELFNEGLABEL0 POSLABEL1 SPARSEOUTPUTFALSE

FITSELFY

FIT LABEL BINARIZER

PARAMETERS

YARRAY OF SHAPE NSAMPLES OR NSAMPLES NCLASSES TARGET VALUES THE 2D MATRIX SHOULD ONLY CONTAIN 0 AND 1 REPRESENTS MULTILABEL CLASSIFICATION

RETURNS

2242 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

SELF RETURNS AN INSTANCE OF SELF

FITTRANSFORM SELF

FIT LABEL BINARIZER AND TRANSFORM MULTICLASS LABELS TO BINARY LABELS

THE OUTPUT OF TRANSFORM IS SOMETIMES REFERRED TO AS THE 10FK CODING SCHEME

PARAMETERS

YARRAY OR SPARSE MATRIX OF SHAPE NSAMPLES OR NSAMPLES NCLASSES TARGET VALUES

THE 2D MATRIX SHOULD ONLY CONTAIN 0 AND 1 REPRESENTS MULTILABEL CLASSIFICATION SPARSE

MATRIX CAN BE CSR CSC COO DOK OR LIL

RETURNS

YARRAY OR CSR MATRIX OF SHAPE NSAMPLES NCLASSES SHAPE WILL BE NSAMPLES 1 FOR

BINARY PROBLEMS

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

INVERSETRANSFORM SELFTHRESHOLDNONE

TRANSFORM BINARY LABELS BACK TO MULTICLASS LABELS

PARAMETERS

YNUMPY ARRAY OR SPARSE MATRIX WITH SHAPE NSAMPLES NCLASSES TARGET VALUES ALL

SPARSE MATRICES ARE CONVERTED TO CSR BEFORE INVERSE TRANSFORMATION

THRESHOLD FLOAT OR NONE THRESHOLD USED IN THE BINARY AND MULTILABEL CASES

USE 0 WHEN YCONTAINS THE OUTPUT OF DECISIONFUNCTION CLASSIFIER USE 05 WHEN YCON

TAINS THE OUTPUT OF PREDICTPROBA

IF NONE THE THRESHOLD IS ASSUMED TO BE HALF WAY BETWEEN NEGLABEL AND POSLABEL

RETURNS

YNUMPY ARRAY OR CSR MATRIX OF SHAPE NSAMPLES TARGET VALUES

NOTES

IN THE CASE WHEN THE BINARY LABELS ARE FRACTIONAL PROBABILISTIC INVERSETRANSFORM CHOOSES THE CLASS WITH

THE GREATEST VALUE TYPICALLY THIS ALLOWS TO USE THE OUTPUT OF A LINEAR MODEL’S DECISIONFUNCTION METHOD

DIRECTLY AS THE INPUT OF INVERSETRANSFORM

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE

PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT

OF A NESTED OBJECT

RETURNS

634SKLEARNPREPROCESSING PREPROCESSING AND NORMALIZATION 2243

SCIKITLEARN USER GUIDE RELEASE 0213

SELF

TRANSFORM SELFY

TRANSFORM MULTICLASS LABELS TO BINARY LABELS

THE OUTPUT OF TRANSFORM IS SOMETIMES REFERRED TO BY SOME AUTHORS AS THE 10FK CODING SCHEME

PARAMETERS

YARRAY OR SPARSE MATRIX OF SHAPE NSAMPLES OR NSAMPLES NCLASSES TARGET VALUES

THE 2D MATRIX SHOULD ONLY CONTAIN 0 AND 1 REPRESENTS MULTILABEL CLASSIFICATION SPARSE

MATRIX CAN BE CSR CSC COO DOK OR LIL

RETURNS

YNUMPY ARRAY OR CSR MATRIX OF SHAPE NSAMPLES NCLASSES SHAPE WILL BE NSAMPLES

1 FOR BINARY PROBLEMS

6346SKLEARNPREPROCESSING LABELENCODER

CLASSSSKLEARNPREPROCESSING LABELENCODER

ENCODE LABELS WITH VALUE BETWEEN 0 AND NCLASSES1

READ MORE IN THE USER GUIDE

ATTRIBUTES

CLASSES ARRAY OF SHAPE NCLASS HOLDS THE LABEL FOR EACH CLASS

SEE ALSO

SKLEARNPREPROCESSINGORDINALENCODER ENCODE CATEGORICAL FEATURES USING A ONEHOT OR ORDINAL

ENCODING SCHEME

EXAMPLES

LABELENCODER CAN BE USED TO NORMALIZE LABELS

FROM SKLEARN IMPORT PREPROCESSING

LE PREPROCESSINGLABELENCODER

LEFIT1 2 2 6

LABELENCODER

LECLASSES

ARRAY1 2 6

LETRANSFORM1 1 2 6

ARRAY0 0 1 2

LEINVERSETRANSFORM0 0 1 2

ARRAY1 1 2 6

IT CAN ALSO BE USED TO TRANSFORM NONNUMERICAL LABELS AS LONG AS THEY ARE HASHABLE AND COMPARABLE TO NUMERICAL

LABELS

LE PREPROCESSINGLABELENCODER

LEFITPARIS PARIS TOKYO AMSTERDAM

LABELENCODER

LISTLECLASSES

AMSTERDAM PARIS TOKYO

LETRANSFORMTOKYO TOKYO PARIS

ARRAY2 2 1

2244 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213  
LISTLEINVERSETRANSFORM2 2 1  
TOKYO TOKYO PARIS  
METHODS  
FITSELF Y FIT LABEL ENCODER  
FITTRANSFORM SELF Y FIT LABEL ENCODER AND RETURN ENCODED LABELS  
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR  
INVERSETRANSFORM SELF Y TRANSFORM LABELS BACK TO ORIGINAL ENCODING  
SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR  
TRANSFORM SELF Y TRANSFORM LABELS TO NORMALIZED ENCODING  
INIT SELFARGS KWARGS  
INITIALIZE SELF SEE HELPTYPESELF FOR ACCURATE SIGNATURE  
FITSELFY  
FIT LABEL ENCODER  
PARAMETERS  
YARRAYLIKE OF SHAPE NSAMPLES TARGET VALUES  
RETURNS  
SELF RETURNS AN INSTANCE OF SELF  
FITTRANSFORM SELFY  
FIT LABEL ENCODER AND RETURN ENCODED LABELS  
PARAMETERS  
YARRAYLIKE OF SHAPE NSAMPLES TARGET VALUES  
RETURNS  
YARRAYLIKE OF SHAPE NSAMPLES  
GETPARAMS SELFDEEPTURE  
GET PARAMETERS FOR THIS ESTIMATOR  
PARAMETERS  
DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBOBJECTS THAT ARE ESTIMATORS  
RETURNS  
PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES  
INVERSETRANSFORM SELFY  
TRANSFORM LABELS BACK TO ORIGINAL ENCODING  
PARAMETERS  
YNUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES  
RETURNS  
YNUMPY ARRAY OF SHAPE NSAMPLES  
634SKLEARNPREPROCESSING PREPROCESSING AND NORMALIZATION 2245

SCIKITLEARN USER GUIDE RELEASE 0213

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELFY

TRANSFORM LABELS TO NORMALIZED ENCODING

PARAMETERS

YARRAYLIKE OF SHAPE NSAMPLES TARGET VALUES

RETURNS

YARRAYLIKE OF SHAPE NSAMPLES

6347SKLEARNPREPROCESSING MULTILABELBINARIZER

CLASSSSKLEARNPREPROCESSING MULTILABELBINARIZER CLASSESNONE SPARSEOUTPUTFALSE

TRANSFORM BETWEEN ITERABLE OF ITERABLES AND A MULTILABEL FORMAT

ALTHOUGH A LIST OF SETS OR TUPLES IS A VERY INTUITIVE FORMAT FOR MULTILABEL DATA IT IS UNWIELDY TO PROCESS THIS

TRANSFORMER CONVERTS BETWEEN THIS INTUITIVE FORMAT AND THE SUPPORTED MULTILABEL FORMAT A SAMPLES X CLASSES

BINARY MATRIX INDICATING THE PRESENCE OF A CLASS LABEL

PARAMETERS

CLASSES ARRAYLIKE OF SHAPE NCLASSES OPTIONAL INDICATES AN ORDERING FOR THE CLASS LABELS ALL

ENTRIES SHOULD BE UNIQUE CANNOT CONTAIN DUPLICATE CLASSES

SPARSEOUTPUT BOOLEAN DEFAULT FALSE SET TO TRUE IF OUTPUT BINARY ARRAY IS DESIRED IN CSR

SPARSE FORMAT

ATTRIBUTES

CLASSES ARRAY OF LABELS A COPY OF THE CLASSES PARAMETER WHERE PROVIDED OR OTHERWISE THE SORTED SET OF CLASSES FOUND WHEN FITTING

SEE ALSO

SKLEARNPREPROCESSINGONEHOTENCODER ENCODE CATEGORICAL FEATURES USING A ONEHOT AKA ONEOFK

SCHEME

EXAMPLES

FROM SKLEARNPREPROCESSING IMPORT MULTILABELBINARIZER

MLB MULTILABELBINARIZER

MLBFITTRANSFORM1 2 3

ARRAY1 1 0

0 0 1

MLBCLASSES

ARRAY1 2 3

2246 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213  
MLBFITTRANSFORMSCIFI THRILLER COMEDY  
ARRAY0 1 1  
1 0 0  
LISTMLBCLASSES  
COMEDY SCIFI THRILLER  
METHODS  
FITSELF Y FIT THE LABEL SETS BINARIZER STORING CLASSES  
FITTRANSFORM SELF Y FIT THE LABEL SETS BINARIZER AND TRANSFORM THE GIVEN LABEL SETS  
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR  
INVERSETRANSFORM SELF YT TRANSFORM THE GIVEN INDICATOR MATRIX INTO LABEL SETS  
SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR  
TRANSFORM SELF Y TRANSFORM THE GIVEN LABEL SETS  
INIT SELFCLASSESNONE SPARSEOUTPUTFALSE  
FITSELFY  
FIT THE LABEL SETS BINARIZER STORING CLASSES  
PARAMETERS  
YITERABLE OF ITERABLES A SET OF LABELS ANY ORDERABLE AND HASHABLE OBJECT FOR EACH SAMPLE  
IF THECLASSES PARAMETER IS SET YWILL NOT BE ITERATED  
RETURNS  
SELF RETURNS THIS MULTILABELBINARIZER INSTANCE  
FITTRANSFORM SELFY  
FIT THE LABEL SETS BINARIZER AND TRANSFORM THE GIVEN LABEL SETS  
PARAMETERS  
YITERABLE OF ITERABLES A SET OF LABELS ANY ORDERABLE AND HASHABLE OBJECT FOR EACH SAMPLE  
IF THECLASSES PARAMETER IS SET YWILL NOT BE ITERATED  
RETURNS  
YINDICATOR ARRAY OR CSR MATRIX SHAPE NSAMPLES NCLASSES A MATRIX SUCH THAT  
YINDICATORI J 1 IFFCLASSESJ IS INYI AND 0 OTHERWISE  
GETPARAMS SELFDEEPTTRUE  
GET PARAMETERS FOR THIS ESTIMATOR  
PARAMETERS  
DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBOBJECTS THAT ARE ESTIMATORS  
RETURNS  
PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES  
INVERSETRANSFORM SELFYT  
TRANSFORM THE GIVEN INDICATOR MATRIX INTO LABEL SETS  
PARAMETERS  
634SKLEARNPREPROCESSING PREPROCESSING AND NORMALIZATION 2247

SCIKITLEARN USER GUIDE RELEASE 0213

Y TARRAY OR SPARSE MATRIX OF SHAPE NSAMPLES NCLASSES A MATRIX CONTAINING ONLY 1S ANDS  
0S

RETURNS

YLIST OF TUPLES THE SET OF LABELS FOR EACH SAMPLE SUCH THAT YI CONSISTS OF  
CLASSESJ FOR EACHYTI J 1

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE  
PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT  
OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELFY

TRANSFORM THE GIVEN LABEL SETS

PARAMETERS

YITERABLE OF ITERABLES A SET OF LABELS ANY ORDERABLE AND HASHABLE OBJECT FOR EACH SAMPLE  
IF THECLASSES PARAMETER IS SET YWILL NOT BE ITERATED

RETURNS

YINDICATOR ARRAY OR CSR MATRIX SHAPE NSAMPLES NCLASSES A MATRIX SUCH THAT  
YINDICATORI J 1 IFFCLASSESJ IS INYI AND 0 OTHERWISE

6348SKLEARNPREPROCESSING MAXABSSCALER

CLASSSSKLEARNPREPROCESSING MAXABSSCALER COPYTRUE

SCALE EACH FEATURE BY ITS MAXIMUM ABSOLUTE VALUE

THIS ESTIMATOR SCALES AND TRANSLATES EACH FEATURE INDIVIDUALLY SUCH THAT THE MAXIMAL ABSOLUTE VALUE OF EACH  
FEATURE IN THE TRAINING SET WILL BE 10 IT DOES NOT SHIFTCENTER THE DATA AND THUS DOES NOT DESTROY ANY SPARSITY

THIS SCALER CAN ALSO BE APPLIED TO SPARSE CSR OR CSC MATRICES

NEW IN VERSION 017

PARAMETERS

COPY BOOLEAN OPTIONAL DEFAULT IS TRUE SET TO FALSE TO PERFORM INPLACE SCALING AND AVOID A  
COPY IF THE INPUT IS ALREADY A NUMPY ARRAY

ATTRIBUTES

SCALE NDARRAY SHAPE NFEATURES PER FEATURE RELATIVE SCALING OF THE DATA

NEW IN VERSION 017 SCALE ATTRIBUTE

MAXABS NDARRAY SHAPE NFEATURES PER FEATURE MAXIMUM ABSOLUTE VALUE

NSAMPLESSEEN INT THE NUMBER OF SAMPLES PROCESSED BY THE ESTIMATOR WILL BE RESET ON  
NEW CALLS TO FIT BUT INCREMENTS ACROSS PARTIALFIT CALLS

SEE ALSO

MAXABSSCALE EQUIVALENT FUNCTION WITHOUT THE ESTIMATOR API

2248 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

NOTES

NANS ARE TREATED AS MISSING VALUES DISREGARDED IN FIT AND MAINTAINED IN TRANSFORM  
FOR A COMPARISON OF THE DIFFERENT SCALERS TRANSFORMERS AND NORMALIZERS SEE EXAM  
PLESPREPROCESSINGPLOTALLSCALINGPY

EXAMPLES

FROM SKLEARNPREPROCESSING IMPORT MAXABSSCALER

X = [[1, 1, 2,

2, 0, 0,

0, 1, 1]

TRANSFORMER = MAXABSSCALER(FIT\_X,

TRANSFORMER,

MAXABSSCALER\_COPY=True,

TRANSFORMER\_TRANSFORM\_X,

ARRAY\_05\_1\_1,

[1, 0, 0,

0, 1, 0.5]

METHODS

fit(X, y) COMPUTE THE MAXIMUM ABSOLUTE VALUE TO BE USED FOR

LATER SCALING

fit\_transform(X, y) FIT TO DATA THEN TRANSFORM IT

get\_params(deep=True) GET PARAMETERS FOR THIS ESTIMATOR

inverse\_transform(X) SCALE BACK THE DATA TO THE ORIGINAL REPRESENTATION

partial\_fit(X, y) ONLINE COMPUTATION OF MAX ABSOLUTE VALUE OF X FOR LATER

SCALING

set\_params(\*\*params) SET THE PARAMETERS OF THIS ESTIMATOR

transform(X) SCALE THE DATA

init(copy=True)

fit(X, y=None)

COMPUTE THE MAXIMUM ABSOLUTE VALUE TO BE USED FOR LATER SCALING

PARAMETERS

X: ARRAY-LIKE SPARSE MATRIX SHAPE (N\_SAMPLES, N\_FEATURES) THE DATA USED TO COMPUTE THE

PER-FEATURE MINIMUM AND MAXIMUM USED FOR LATER SCALING ALONG THE FEATURES AXIS

fit\_transform(X, y=None, fit\_params=None)

FIT TO DATA THEN TRANSFORM IT

fit\_transformer(X, y) WITH OPTIONAL PARAMETERS fit\_params AND RETURNS A TRANSFORMED VERSION OF X

PARAMETERS

X: NUMPY ARRAY OF SHAPE (N\_SAMPLES, N\_FEATURES) TRAINING SET

Y: NUMPY ARRAY OF SHAPE (N\_SAMPLES,) TARGET VALUES

RETURNS

634 SKLEARNPREPROCESSING PREPROCESSING AND NORMALIZATION 2249

SCIKITLEARN USER GUIDE RELEASE 0213

XNEW NUMPY ARRAY OF SHAPE NSAMPLES NFEATURESNEW TRANSFORMED ARRAY

GETPARAMS SELFDEEPTTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

INVERSETRANSFORM SELF

SCALE BACK THE DATA TO THE ORIGINAL REPRESENTATION

PARAMETERS

XARRAYLIKE SPARSE MATRIX THE DATA THAT SHOULD BE TRANSFORMED BACK

PARTIALFIT SELFXYNONE

ONLINE COMPUTATION OF MAX ABSOLUTE VALUE OF X FOR LATER SCALING ALL OF X IS PROCESSED AS A SINGLE BATCH

THIS IS INTENDED FOR CASES WHEN FIT IS NOT FEASIBLE DUE TO VERY LARGE NUMBER OF NSAMPLES OR BECAUSE X

IS READ FROM A CONTINUOUS STREAM

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE DATA USED TO COMPUTE THE

MEAN AND STANDARD DEVIATION USED FOR LATER SCALING ALONG THE FEATURES AXIS

YIGNORED

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE

PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT

OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELF

SCALE THE DATA

PARAMETERS

XARRAYLIKE SPARSE MATRIX THE DATA THAT SHOULD BE SCALED

EXAMPLES USING SKLEARNPREPROCESSINGMAXABSSCALER

- COMPARE THE EFFECT OF DIFFERENT SCALERS ON DATA WITH OUTLIERS

6349SKLEARNPREPROCESSING MINMAXSCALER

CLASSSKLEARNPREPROCESSING MINMAXSCALER FEATURERANGE0 1COPYTRUE

TRANSFORMS FEATURES BY SCALING EACH FEATURE TO A GIVEN RANGE

THIS ESTIMATOR SCALES AND TRANSLATES EACH FEATURE INDIVIDUALLY SUCH THAT IT IS IN THE GIVEN RANGE ON THE TRAINING SET

EG BETWEEN ZERO AND ONE

2250 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

THE TRANSFORMATION IS GIVEN BY

$$XSTD = \frac{X - XMINAXISO}{XMAXAXISO - XMINAXISO}$$

XSCALED = XSTD \* MAX\_MIN

WHERE MIN = MAX - FEATURERANGE

THE TRANSFORMATION IS CALCULATED AS

$$XSCALED = \frac{SCALE \times (X - XMINAXISO)}{XMAXAXISO - XMINAXISO}$$

WHERE SCALE = MAX - MIN

THIS TRANSFORMATION IS OFTEN USED AS AN ALTERNATIVE TO ZERO MEAN UNIT VARIANCE SCALING

READ MORE IN THE USER GUIDE

PARAMETERS

FEATURERANGE = TUPLE (MIN, MAX) DEFAULT 0 1 DESIRED RANGE OF TRANSFORMED DATA

COPY = BOOLEAN OPTIONAL DEFAULT TRUE SET TO FALSE TO PERFORM INPLACE ROW NORMALIZATION AND AVOID A COPY IF THE INPUT IS ALREADY A NUMPY ARRAY

ATTRIBUTES

MIN = NDARRAY SHAPE (NFEATURES) PER FEATURE ADJUSTMENT FOR MINIMUM EQUIVALENT TO MIN

XMINAXISO = SELFSCALE

SCALE = NDARRAY SHAPE (NFEATURES) PER FEATURE RELATIVE SCALING OF THE DATA EQUIVALENT TO MAX - MIN

XMAXAXISO = XMINAXISO

NEW IN VERSION 017: SCALE ATTRIBUTE

DATAMIN = NDARRAY SHAPE (NFEATURES) PER FEATURE MINIMUM SEEN IN THE DATA

NEW IN VERSION 017: DATAMIN

DATAMAX = NDARRAY SHAPE (NFEATURES) PER FEATURE MAXIMUM SEEN IN THE DATA

NEW IN VERSION 017: DATAMAX

DATARANGE = NDARRAY SHAPE (NFEATURES) PER FEATURE RANGE DATAMAX - DATAMIN

SEEN IN THE DATA

NEW IN VERSION 017: DATARANGE

SEE ALSO

MINMAXSCALE EQUIVALENT FUNCTION WITHOUT THE ESTIMATOR API

NOTES

NANS ARE TREATED AS MISSING VALUES DISREGARDED IN FIT AND MAINTAINED IN TRANSFORM

FOR A COMPARISON OF THE DIFFERENT SCALERS TRANSFORMERS AND NORMALIZERS SEE EXAMPLES

PREPROCESSINGPLOTALLSCALINGPY

EXAMPLES

634SKLEARNPREPROCESSING PREPROCESSING AND NORMALIZATION 2251

SCIKITLEARN USER GUIDE RELEASE 0213

FROM SKLEARNPREPROCESSING IMPORT MINMAXSCALER

DATA 1 2 05 6 0 10 1 18

SCALER MINMAXSCALER

PRINTSCALERFITDATA

MINMAXSCALERCOPYTRUE FEATURERANGE0 1

PRINTSCALERDATAMAX

1 18

PRINTSCALERTRANSFORMDATA

0 0

025 025

05 05

1 1

PRINTSCALERTRANSFORM2 2

15 0

METHODS

FITSELF X Y COMPUTE THE MINIMUM AND MAXIMUM TO BE USED FOR

LATER SCALING

FITTRANSFORM SELF X Y FIT TO DATA THEN TRANSFORM IT

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

INVERSETRANSFORM SELF X UNDO THE SCALING OF X ACCORDING TO FEATURERANGE

PARTIALFIT SELF X Y ONLINE COMPUTATION OF MIN AND MAX ON X FOR LATER SCAL

ING

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF X SCALING FEATURES OF X ACCORDING TO FEATURERANGE

INIT SELFFEATURERANGE0 1COPYTRUE

FITSELFXYNONE

COMPUTE THE MINIMUM AND MAXIMUM TO BE USED FOR LATER SCALING

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES THE DATA USED TO COMPUTE THE PERFEATURE MIN

IMUM AND MAXIMUM USED FOR LATER SCALING ALONG THE FEATURES AXIS

FITTRANSFORM SELFXYNONE FITPARAMS

FIT TO DATA THEN TRANSFORM IT

FITS TRANSFORMER TO X AND Y WITH OPTIONAL PARAMETERS FITPARAMS AND RETURNS A TRANSFORMED VERSION OF X

PARAMETERS

XNUMPY ARRAY OF SHAPE NSAMPLES NFEATURES TRAINING SET

YNUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES

RETURNS

XNEW NUMPY ARRAY OF SHAPE NSAMPLES NFEATURESNEW TRANSFORMED ARRAY

GETPARAMS SELFDEEPTTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

2252 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

INVERSETRANSFORM SELF

UNDO THE SCALING OF X ACCORDING TO FEATURERANGE

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES INPUT DATA THAT WILL BE TRANSFORMED IT CANNOT BE SPARSE

PARTIALFIT SELFXYNONE

ONLINE COMPUTATION OF MIN AND MAX ON X FOR LATER SCALING ALL OF X IS PROCESSED AS A SINGLE BATCH THIS IS INTENDED FOR CASES WHEN FIT IS NOT FEASIBLE DUE TO VERY LARGE NUMBER OF NSAMPLES OR BECAUSE X IS READ FROM A CONTINUOUS STREAM

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES THE DATA USED TO COMPUTE THE MEAN AND STAN DARD DEVIATION USED FOR LATER SCALING ALONG THE FEATURES AXIS

YIGNORED

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELF

SCALING FEATURES OF X ACCORDING TO FEATURERANGE

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES INPUT DATA THAT WILL BE TRANSFORMED

EXAMPLES USING SKLEARNPREPROCESSINGMINMAXSCALER

- COMPARE STOCHASTIC LEARNING STRATEGIES FOR MLPCLASSIFIER
- COMPARE THE EFFECT OF DIFFERENT SCALERS ON DATA WITH OUTLIERS

63410SKLEARNPREPROCESSING NORMALIZER

CLASSSSKLEARNPREPROCESSING NORMALIZER NORM‘L2’ COPYTRUE

NORMALIZE SAMPLES INDIVIDUALLY TO UNIT NORM

EACH SAMPLE IE EACH ROW OF THE DATA MATRIX WITH AT LEAST ONE NON ZERO COMPONENT IS RESCALED INDEPENDENTLY OF OTHER SAMPLES SO THAT ITS NORM L1 OR L2 EQUALS ONE

THIS TRANSFORMER IS ABLE TO WORK BOTH WITH DENSE NUMPY ARRAYS AND SCIPYSPARSE MATRIX USE CSR FORMAT IF YOU WANT TO AVOID THE BURDEN OF A COPY CONVERSION

6345SKLEARNPREPROCESSING PREPROCESSING AND NORMALIZATION 2253

SCIKITLEARN USER GUIDE RELEASE 0213

SCALING INPUTS TO UNIT NORMS IS A COMMON OPERATION FOR TEXT CLASSIFICATION OR CLUSTERING FOR INSTANCE FOR INSTANCE THE DOT PRODUCT OF TWO L2NORMALIZED TFIDF VECTORS IS THE COSINE SIMILARITY OF THE VECTORS AND IS THE BASE SIMILARITY METRIC FOR THE VECTOR SPACE MODEL COMMONLY USED BY THE INFORMATION RETRIEVAL COMMUNITY

READ MORE IN THE USER GUIDE

PARAMETERS

NORM ‘L1’ ‘L2’ OR ‘MAX’ OPTIONAL ‘L2’ BY DEFAULT THE NORM TO USE TO NORMALIZE EACH NON ZERO SAMPLE

COPY BOOLEAN OPTIONAL DEFAULT TRUE SET TO FALSE TO PERFORM INPLACE ROW NORMALIZATION AND AVOID A COPY IF THE INPUT IS ALREADY A NUMPY ARRAY OR A SCIPYSPARSE CSR MATRIX

SEE ALSO

NORMALIZE EQUIVALENT FUNCTION WITHOUT THE ESTIMATOR API

NOTES

THIS ESTIMATOR IS STATELESS BESIDES CONSTRUCTOR PARAMETERS THE FIT METHOD DOES NOTHING BUT IS USEFUL WHEN USED IN A PIPELINE

FOR A COMPARISON OF THE DIFFERENT SCALERS TRANSFORMERS AND NORMALIZERS SEE EXAM PLES

PREPROCESSINGPLOTALLSCALINGPY

EXAMPLES

```
FROM SKLEARNPREPROCESSING IMPORT NORMALIZER
X = [[4, 1, 2, 2],
     [1, 3, 9, 3],
     [5, 7, 5, 1]]
TRANSFORMER = NORMALIZER(FITX = FIT DOES NOTHING)
TRANSFORMER = NORMALIZER(COPY=True, NORM=L2)
TRANSFORMER.TRANSFORM(X)
ARRAY[[0.8, 0.2, 0.4, 0.4],
       [0.1, 0.3, 0.9, 0.3],
       [0.5, 0.7, 0.5, 0.1]]
```

METHODS

FITSELF X Y DO NOTHING AND RETURN THE ESTIMATOR UNCHANGED

FITTRANSFORM SELF X Y FIT TO DATA THEN TRANSFORM IT

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF X COPY SCALE EACH NON ZERO ROW OF X TO UNIT NORM

INIT SELFNORM‘L2’ COPYTRUE

FITSELFXYNONE

DO NOTHING AND RETURN THE ESTIMATOR UNCHANGED

THIS METHOD IS JUST THERE TO IMPLEMENT THE USUAL API AND HENCE WORK IN PIPELINES

2254 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

PARAMETERS

XARRAYLIKE

FITTRANSFORM SELFXYNONE FITPARAMS

FIT TO DATA THEN TRANSFORM IT

FITS TRANSFORMER TO X AND Y WITH OPTIONAL PARAMETERS FITPARAMS AND RETURNS A TRANSFORMED VERSION OF X

PARAMETERS

XNUMPY ARRAY OF SHAPE NSAMPLES NFEATURES TRAINING SET

YNUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES

RETURNS

XNEW NUMPY ARRAY OF SHAPE NSAMPLES NFEATURESNEW TRANSFORMED ARRAY

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE

PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT

OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELFXCOPYNONE

SCALE EACH NON ZERO ROW OF X TO UNIT NORM

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE DATA TO NORMALIZE ROW

BY ROW SCIPYSPARSE MATRICES SHOULD BE IN CSR FORMAT TO AVOID AN UNNECESSARY COPY

COPY BOOL OPTIONAL DEFAULT NONE COPY THE INPUT X OR NOT

EXAMPLES USING SKLEARNPREPROCESSINGNORMALIZER

- COMPARE THE EFFECT OF DIFFERENT SCALERS ON DATA WITH OUTLIERS
- CLUSTERING TEXT DOCUMENTS USING KMEANS

634SKLEARNPREPROCESSING PREPROCESSING AND NORMALIZATION 2255

SCIKITLEARN USER GUIDE RELEASE 0213

63411SKLEARNPREPROCESSING ONEHOTENCODER

CLASSSSKLEARNPREPROCESSING ONEHOTENCODER NVALUESNONE CATEGORICALFEATURESNONE

CATEGORIESNONE DROPNONE SPARSETRUE

DTYPECLASS 'NUMPYFLOAT64' HAN

DLEUNKNOWN'ERROR'

ENCODE CATEGORICAL INTEGER FEATURES AS A ONEHOT NUMERIC ARRAY

THE INPUT TO THIS TRANSFORMER SHOULD BE AN ARRAYLIKE OF INTEGERS OR STRINGS DENOTING THE VALUES TAKEN ON BY CATEGORICAL DISCRETE FEATURES THE FEATURES ARE ENCODED USING A ONEHOT AKA 'ONEOFK' OR 'DUMMY' ENCODING SCHEME THIS CREATES A BINARY COLUMN FOR EACH CATEGORY AND RETURNS A SPARSE MATRIX OR DENSE ARRAY BY DEFAULT THE ENCODER DERIVES THE CATEGORIES BASED ON THE UNIQUE VALUES IN EACH FEATURE ALTERNATIVELY YOU CAN ALSO SPECIFY THE CATEGORIES MANUALLY THE ONEHOTENCODER PREVIOUSLY ASSUMED THAT THE INPUT FEATURES TAKE ON VALUES IN THE RANGE 0 MAXVALUES THIS BEHAVIOUR IS DEPRECATED

THIS ENCODING IS NEEDED FOR FEEDING CATEGORICAL DATA TO MANY SCIKITLEARN ESTIMATORS NOTABLY LINEAR MODELS AND SVMS WITH THE STANDARD KERNELS

NOTE A ONEHOT ENCODING OF Y LABELS SHOULD USE A LABELBINARIZER INSTEAD

READ MORE IN THE USER GUIDE

PARAMETERS

CATEGORIES 'AUTO' OR A LIST OF LISTSARRAYS OF VALUES DEFAULT'AUTO' CATEGORIES UNIQUE VALUES PER FEATURE

- 'AUTO' DETERMINE CATEGORIES AUTOMATICALLY FROM THE TRAINING DATA
- LIST CATEGORIESI HOLDS THE CATEGORIES EXPECTED IN THE ITH COLUMN THE PASSED CATEGORIES SHOULD NOT MIX STRINGS AND NUMERIC VALUES WITHIN A SINGLE FEATURE AND SHOULD BE SORTED IN CASE OF NUMERIC VALUES

THE USED CATEGORIES CAN BE FOUND IN THE CATEGORIES ATTRIBUTE

DROP 'FIRST' OR A LISTARRAY OF SHAPE NFEATURES DEFAULTNONE SPECIFIES A METHODOLOGY TO USE TO DROP ONE OF THE CATEGORIES PER FEATURE THIS IS USEFUL IN SITUATIONS WHERE PERFECTLY COLLINEAR FEATURES CAUSE PROBLEMS SUCH AS WHEN FEEDING THE RESULTING DATA INTO A NEURAL NETWORK OR AN UNREGULARIZED REGRESSION

- NONE RETAIN ALL FEATURES THE DEFAULT
- 'FIRST' DROP THE FIRST CATEGORY IN EACH FEATURE IF ONLY ONE CATEGORY IS PRESENT THE FEATURE WILL BE DROPPED ENTIRELY
- ARRAY DROPI IS THE CATEGORY IN FEATURE X I THAT SHOULD BE DROPPED

SPARSE BOOLEAN DEFAULTTRUE WILL RETURN SPARSE MATRIX IF SET TRUE ELSE WILL RETURN AN ARRAY

DTYPE NUMBER TYPE DEFAULTNPFLOAT DESIRED DTYPE OF OUTPUT

HANDLEUNKNOWN 'ERROR' OR 'IGNORE' DEFAULT'ERROR' WHETHER TO RAISE AN ERROR OR IGNORE IF AN UNKNOWN CATEGORICAL FEATURE IS PRESENT DURING TRANSFORM DEFAULT IS TO RAISE WHEN THIS PARAMETER IS SET TO 'IGNORE' AND AN UNKNOWN CATEGORY IS ENCOUNTERED DURING TRANSFORM THE RESULTING ONEHOT ENCODED COLUMNS FOR THIS FEATURE WILL BE ALL ZEROS IN THE INVERSE TRANSFORM AN UNKNOWN CATEGORY WILL BE DENOTED AS NONE

NVALUES 'AUTO' INT OR ARRAY OF INTS DEFAULT'AUTO' NUMBER OF VALUES PER FEATURE

- 'AUTO' DETERMINE VALUE RANGE FROM TRAINING DATA
- INTNUMBER OF CATEGORICAL VALUES PER FEATURE EACH FEATURE VALUE SHOULD BE IN RANGENVVALUES

2256 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

•ARRAY NVALUESI IS THE NUMBER OF CATEGORICAL VALUES IN X I EACH FEATURE VALUE SHOULD BE IN RANGENVVALUESI

DEPRECATED SINCE VERSION 020 THE NVALUES KEYWORD WAS DEPRECATED IN VERSION 020 AND WILL BE REMOVED IN 022 USE CATEGORIES INSTEAD

CATEGORICALFEATURES 'ALL' OR ARRAY OF INDICES OR MASK DEFAULT'ALL' SPECIFY WHAT FEATURES ARE TREATED AS CATEGORICAL

- 'ALL' ALL FEATURES ARE TREATED AS CATEGORICAL
- ARRAY OF INDICES ARRAY OF CATEGORICAL FEATURE INDICES
- MASK ARRAY OF LENGTH NFEATURES AND WITH DTYPEBOOL

NONCATEGORICAL FEATURES ARE ALWAYS STACKED TO THE RIGHT OF THE MATRIX

DEPRECATED SINCE VERSION 020 THE CATEGORICALFEATURES KEYWORD WAS DEPRECATED IN VERSION 020 AND WILL BE REMOVED IN 022 YOU CAN USE THE COLUMNTRANSFORMER INSTEAD

ATTRIBUTES

CATEGORIES LIST OF ARRAYS THE CATEGORIES OF EACH FEATURE DETERMINED DURING FITTING IN ORDER OF THE FEATURES IN X AND CORRESPONDING WITH THE OUTPUT OF TRANSFORM THIS INCLUDES THE CATEGORY SPECIFIED IN DROP IF ANY

DROPIDX ARRAY OF SHAPE NFEATURES DROPIDX I IS THE INDEX IN CATEGORIES I

OF THE CATEGORY TO BE DROPPED FOR EACH FEATURE NONE IF ALL THE TRANSFORMED FEATURES WILL BE RETAINED

ACTIVEFEATURES ARRAY INDICES FOR ACTIVE FEATURES MEANING VALUES THAT ACTUALLY OCCUR IN THE TRAINING SET ONLY AVAILABLE WHEN NVALUES IS AUTO

DEPRECATED SINCE VERSION 020 THE ACTIVEFEATURES ATTRIBUTE WAS DEPRECATED IN VER SION 020 AND WILL BE REMOVED IN 022

FEATUREINDICES ARRAY OF SHAPE NFEATURES INDICES TO FEATURE RANGES FEATURE

I IN THE ORIGINAL DATA IS MAPPED TO FEATURES FROM FEATUREINDICES I TO FEATUREINDICES I1 AND THEN POTENTIALLY MASKED BY ACTIVEFEATURES AF

TERWARDS

DEPRECATED SINCE VERSION 020 THE FEATUREINDICES ATTRIBUTE WAS DEPRECATED IN VER SION 020 AND WILL BE REMOVED IN 022

NVALUES ARRAY OF SHAPE NFEATURES MAXIMUM NUMBER OF VALUES PER FEATURE

DEPRECATED SINCE VERSION 020 THE NVALUES ATTRIBUTE WAS DEPRECATED IN VERSION 020 AND WILL BE REMOVED IN 022

SEE ALSO SKLEARNPREPROCESSINGORDINALENCODER PERFORMS AN ORDINAL INTEGER ENCODING OF THE CATEGORICAL FEATURES

SKLEARNFEATUREEXTRACTIONDICTVECTORIZER PERFORMS A ONEHOT ENCODING OF DICTIONARY ITEMS ALSO HANDLES STRINGVALUED FEATURES

SKLEARNFEATUREEXTRACTIONFEATUREHASHER PERFORMS AN APPROXIMATE ONEHOT ENCODING OF DIC TIONARY ITEMS OR STRINGS

SKLEARNPREPROCESSINGLABELBINARIZER BINARIZES LABELS IN A ONEVSALL FASHION 634SKLEARNPREPROCESSING PREPROCESSING AND NORMALIZATION 2257

SCIKITLEARN USER GUIDE RELEASE 0213  
SKLEARNPREPROCESSINGMULTILABELBINARIZER TRANSFORMS BETWEEN ITERABLE OF ITERABLES AND A  
MULTILABEL FORMAT EG A SAMPLES X CLASSES BINARY MATRIX INDICATING THE PRESENCE OF A CLASS LABEL  
EXAMPLES  
GIVEN A DATASET WITH TWO FEATURES WE LET THE ENCODER FIND THE UNIQUE VALUES PER FEATURE AND TRANSFORM THE DATA  
TO A BINARY ONEHOT ENCODING  
FROM SKLEARNPREPROCESSING IMPORT ONEHOTENCODER  
ENC ONEHOTENCODERHANDLEUNKNOWNIGNORE  
X MALE 1 FEMALE 3 FEMALE 2  
ENCFITX

ONEHOTENCODERCATEGORICALFEATURESNONE CATEGORIESNONE DROPNONE  
DTYPE NUMPYFLOAT64 HANDLEUNKNOWNIGNORE  
NVALUESNONE SPARSETRUE  
ENCCATEGORIES  
ARRAYFEMALE MALE DTYPEOBJECT ARRAY1 2 3 DTYPEOBJECT  
ENCTRANSFORMFEMALE 1 MALE 4TOARRAY  
ARRAY1 0 1 0 0  
0 1 0 0 0  
ENCINVERSETRANSFORM0 1 1 0 0 0 0 0 1 0  
ARRAYMALE 1  
NONE 2 DTYPEOBJECT  
ENCGETFEATURENAMES  
ARRAYX0FEMALE X0MALE X11 X12 X13 DTYPEOBJECT  
DROPENC ONEHOTENCODERDROPFIRSTFITX  
DROPENCCATEGORIES  
ARRAYFEMALE MALE DTYPEOBJECT ARRAY1 2 3 DTYPEOBJECT  
DROPENCTRANSFORMFEMALE 1 MALE 2TOARRAY  
ARRAY0 0 0  
1 1 0  
METHODS  
FITSELF X Y FIT ONEHOTENCODER TO X  
FITTRANSFORM SELF X Y FIT ONEHOTENCODER TO X THEN TRANSFORM X  
GETFEATURENAMES SELF INPUTFEATURES RETURN FEATURE NAMES FOR OUTPUT FEATURES  
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR  
INVERSETRANSFORM SELF X CONVERT THE BACK DATA TO THE ORIGINAL REPRESENTATION  
SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR  
TRANSFORM SELF X TRANSFORM X USING ONEHOT ENCODING  
INIT SELFNVALUESNONE CATEGORICALFEATURESNONE CATEGORIESNONE DROPNONE  
SPARSETRUE DTYPECLASS 'NUMPYFLOAT64' HANDLEUNKNOWN'ERROR'  
FITSELFXYNONE  
FIT ONEHOTENCODER TO X  
PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES THE DATA TO DETERMINE THE CATEGORIES OF EACH  
2258 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

FEATURE

RETURNS

SELF

FITTRANSFORM SELFXYNONE

FIT ONEHOTENCODER TO X THEN TRANSFORM X

EQUIVALENT TO FITXTRANSFORMX BUT MORE CONVENIENT

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES THE DATA TO ENCODE

RETURNS

XOUT SPARSE MATRIX IF SPARSETRUE ELSE A 2D ARRAY TRANSFORMED INPUT

GETFEATURENAMES SELFINPUTFEATURESNONE

RETURN FEATURE NAMES FOR OUTPUT FEATURES

PARAMETERS

INPUTFEATURES LIST OF STRING LENGTH NFEATURES OPTIONAL STRING NAMES FOR INPUT FEATURES IF AVAILABLE BY DEFAULT "X0" "X1" "XNFEATURES" IS USED

RETURNS

OUTPUTFEATURENAMES ARRAY OF STRING LENGTH NOUTPUTFEATURES

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

INVERSETRANSFORM SELF

CONVERT THE BACK DATA TO THE ORIGINAL REPRESENTATION

IN CASE UNKNOWN CATEGORIES ARE ENCOUNTERED ALL ZEROS IN THE ONEHOT ENCODING NONE IS USED TO REPRESENT THIS CATEGORY

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NENCODEDFEATURES THE TRANSFORMED DATA

RETURNS

XTR ARRAYLIKE SHAPE NSAMPLES NFEATURES INVERSE TRANSFORMED ARRAY

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

634SKLEARNPREPROCESSING PREPROCESSING AND NORMALIZATION 2259

SCIKITLEARN USER GUIDE RELEASE 0213

TRANSFORM SELF

TRANSFORM X USING ONEHOT ENCODING

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES THE DATA TO ENCODE

RETURNS

XOUT SPARSE MATRIX IF SPARSETRUE ELSE A 2D ARRAY TRANSFORMED INPUT

EXAMPLES USING SKLEARNPREPROCESSINGONEHOTENCODER

- COLUMN TRANSFORMER WITH MIXED TYPES
- FEATURE TRANSFORMATIONS WITH ENSEMBLES OF TREES

63412SKLEARNPREPROCESSING ORDINALENCODER

CLASSSKLEARNPREPROCESSING ORDINALENCODER CATEGORIES'AUTO' DTYPECLASS

'NUMPYFLOAT64'

ENCODE CATEGORICAL FEATURES AS AN INTEGER ARRAY

THE INPUT TO THIS TRANSFORMER SHOULD BE AN ARRAYLIKE OF INTEGERS OR STRINGS DENOTING THE VALUES TAKEN ON BY CATEGORICAL DISCRETE FEATURES THE FEATURES ARE CONVERTED TO ORDINAL INTEGERS THIS RESULTS IN A SINGLE COLUMN OF INTEGERS 0 TO NCATEGORIES - 1 PER FEATURE

READ MORE IN THE USER GUIDE

PARAMETERS

CATEGORIES 'AUTO' OR A LIST OF LISTSARRAYS OF VALUES CATEGORIES UNIQUE VALUES PER FEATURE

- 'AUTO' DETERMINE CATEGORIES AUTOMATICALLY FROM THE TRAINING DATA
- LIST CATEGORIESI HOLDS THE CATEGORIES EXPECTED IN THE ITH COLUMN THE PASSED

CATEGORIES SHOULD NOT MIX STRINGS AND NUMERIC VALUES AND SHOULD BE SORTED IN CASE OF NUMERIC VALUES

THE USED CATEGORIES CAN BE FOUND IN THE CATEGORIES ATTRIBUTE

DTYPE NUMBER TYPE DEFAULT NPFLOAT64 DESIRED DTYPE OF OUTPUT

ATTRIBUTES

CATEGORIES LIST OF ARRAYS THE CATEGORIES OF EACH FEATURE DETERMINED DURING FITTING IN ORDER OF THE FEATURES IN X AND CORRESPONDING WITH THE OUTPUT OF TRANSFORM

SEE ALSO

SKLEARNPREPROCESSINGONEHOTENCODER PERFORMS A ONEHOT ENCODING OF CATEGORICAL FEATURES

SKLEARNPREPROCESSINGLABELENCODER ENCODES TARGET LABELS WITH VALUES BETWEEN 0 AND NCLASSES - 1

EXAMPLES

GIVEN A DATASET WITH TWO FEATURES WE LET THE ENCODER FIND THE UNIQUE VALUES PER FEATURE AND TRANSFORM THE DATA TO AN ORDINAL ENCODING

2260 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213  
FROM SKLEARNPREPROCESSING IMPORT ORDINALENCODER  
ENC ORDINALENCODER  
X MALE 1 FEMALE 3 FEMALE 2  
ENCFITX

ORDINALENCODERCATEGORIESAUTO DTYPE NUMPYFLOAT64  
ENCCATEGORIES  
ARRAYFEMALE MALE DTYPEOBJECT ARRAY1 2 3 DTYPEOBJECT  
ENCTRANSFORMFEMALE 3 MALE 1  
ARRAY0 2  
1 0  
ENCINVERSETRANSFORM1 0 0 1  
ARRAYMALE 1  
FEMALE 2 DTYPEOBJECT  
METHODS  
FITSELF X Y FIT THE ORDINALENCODER TO X  
FITTRANSFORM SELF X Y FIT TO DATA THEN TRANSFORM IT  
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR  
INVERSETRANSFORM SELF X CONVERT THE DATA BACK TO THE ORIGINAL REPRESENTATION  
SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR  
TRANSFORM SELF X TRANSFORM X TO ORDINAL CODES  
INIT SELF CATEGORIES' AUTO' DTYPECLASS 'NUMPYFLOAT64'  
FITSELFXYNONE  
FIT THE ORDINALENCODER TO X  
PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES THE DATA TO DETERMINE THE CATEGORIES OF EACH  
FEATURE  
RETURNS  
SELF  
FITTRANSFORM SELFXYNONE FITPARAMS  
FIT TO DATA THEN TRANSFORM IT  
FITS TRANSFORMER TO X AND Y WITH OPTIONAL PARAMETERS FITPARAMS AND RETURNS A TRANSFORMED VERSION OF X  
PARAMETERS  
XNUMPY ARRAY OF SHAPE NSAMPLES NFEATURES TRAINING SET  
YNUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES  
RETURNS  
XNEW NUMPY ARRAY OF SHAPE NSAMPLES NFEATURESNEW TRANSFORMED ARRAY  
GETPARAMS SELFDEEPTURE  
GET PARAMETERS FOR THIS ESTIMATOR  
PARAMETERS  
634SKLEARNPREPROCESSING PREPROCESSING AND NORMALIZATION 2261

SCIKITLEARN USER GUIDE RELEASE 0213

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

INVERSETRANSFORM SELF

CONVERT THE DATA BACK TO THE ORIGINAL REPRESENTATION

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NENCODEDFEATURES THE TRANSFORMED DATA

RETURNS

XTR ARRAYLIKE SHAPE NSAMPLES NFEATURES INVERSE TRANSFORMED ARRAY

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELF

TRANSFORM X TO ORDINAL CODES

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES THE DATA TO ENCODE

RETURNS

XOUT SPARSE MATRIX OR A 2D ARRAY TRANSFORMED INPUT

63413SKLEARNPREPROCESSING POLYNOMIALFEATURES

CLASSSKLEARNPREPROCESSING POLYNOMIALFEATURES DEGREE2 INTERACTIONONLYFALSE IN CLUDEBIASTRUE ORDER'C'

GENERATE POLYNOMIAL AND INTERACTION FEATURES

GENERATE A NEW FEATURE MATRIX CONSISTING OF ALL POLYNOMIAL COMBINATIONS OF THE FEATURES WITH DEGREE LESS THAN OR EQUAL TO THE SPECIFIED DEGREE FOR EXAMPLE IF AN INPUT SAMPLE IS TWO DIMENSIONAL AND OF THE FORM A B THE DEGREE2 POLYNOMIAL FEATURES ARE 1 A B A2 AB B2

PARAMETERS

DEGREE INTEGER THE DEGREE OF THE POLYNOMIAL FEATURES DEFAULT 2

INTERACTIONONLY BOOLEAN DEFAULT FALSE IF TRUE ONLY INTERACTION FEATURES ARE PRODUCED FEATURES THAT ARE PRODUCTS OF AT MOST DEGREE DISTINCT INPUT FEATURES SO NOT X12 X0X23 ETC

INCLUDEBIAS BOOLEAN IF TRUE DEFAULT THEN INCLUDE A BIAS COLUMN THE FEATURE IN WHICH ALL POLYNOMIAL POWERS ARE ZERO IE A COLUMN OF ONES ACTS AS AN INTERCEPT TERM IN A LINEAR MODEL

2262 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

ORDER STR IN ‘C’ ‘F’ DEFAULT ‘C’ ORDER OF OUTPUT ARRAY IN THE DENSE CASE ‘F’ ORDER IS FASTER TO COMPUTE BUT MAY SLOW DOWN SUBSEQUENT ESTIMATORS

NEW IN VERSION 021

ATTRIBUTES

POWERS ARRAY SHAPE NOUTPUTFEATURES NINPUTFEATURES POWERSI J IS THE EXPONENT OF THE JTH INPUT IN THE ITH OUTPUT

NINPUTFEATURES INT THE TOTAL NUMBER OF INPUT FEATURES

NOUTPUTFEATURES INT THE TOTAL NUMBER OF POLYNOMIAL OUTPUT FEATURES THE NUMBER OF OUTPUT FEATURES IS COMPUTED BY ITERATING OVER ALL SUITABLY SIZED COMBINATIONS OF INPUT FEATURES

NOTES

BE AWARE THAT THE NUMBER OF FEATURES IN THE OUTPUT ARRAY SCALES POLYNOMIALLY IN THE NUMBER OF FEATURES OF THE INPUT ARRAY AND EXPONENTIALLY IN THE DEGREE HIGH DEGREES CAN CAUSE OVERFITTING

SEEEXAMPLESLINEARMODELPLOTPOLYNOMIALINTERPOLATIONPY

EXAMPLES

```
import numpy as np
from sklearn.preprocessing import PolynomialFeatures
X = np.arange(6).reshape(3, 2)
X
array([[0, 1],
       [2, 3],
       [4, 5]])
poly = PolynomialFeatures(2)
poly.fit_transform(X)
array([[1, 0, 1, 0, 0, 1],
       [1, 2, 3, 4, 6, 9],
       [1, 4, 5, 16, 20, 25]])
poly = PolynomialFeatures(interaction_only=True)
poly.fit_transform(X)
array([[1, 0, 1, 0],
       [1, 2, 3, 6],
       [1, 4, 5, 20]])
```

METHODS

fitself X Y COMPUTE NUMBER OF OUTPUT FEATURES

fittransform self X Y FIT TO DATA THEN TRANSFORM IT

getfeaturenames self INPUTFEATURES RETURN FEATURE NAMES FOR OUTPUT FEATURES

getparams self DEEP GET PARAMETERS FOR THIS ESTIMATOR

setparams self PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

transform self X TRANSFORM DATA TO POLYNOMIAL FEATURES

init selfDEGREE2 INTERACTIONONLYFALSE INCLUDEBIASTRUE ORDER‘C’

fitselfFXNONE

634SKLEARNPREPROCESSING PREPROCESSING AND NORMALIZATION 2263

SCIKITLEARN USER GUIDE RELEASE 0213  
COMPUTE NUMBER OF OUTPUT FEATURES  
PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES THE DATA  
RETURNS  
SELF INSTANCE  
FITTRANSFORM SELFXYNONE FITPARAMS  
FIT TO DATA THEN TRANSFORM IT  
FITS TRANSFORMER TO X AND Y WITH OPTIONAL PARAMETERS FITPARAMS AND RETURNS A TRANSFORMED VERSION OF X  
PARAMETERS  
XNUMPY ARRAY OF SHAPE NSAMPLES NFEATURES TRAINING SET  
YNUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES  
RETURNS  
XNEW NUMPY ARRAY OF SHAPE NSAMPLES NFEATURESNEW TRANSFORMED ARRAY  
GETFEATURENAMES SELFINPUTFEATURESNONE  
RETURN FEATURE NAMES FOR OUTPUT FEATURES  
PARAMETERS  
INPUTFEATURES LIST OF STRING LENGTH NFEATURES OPTIONAL STRING NAMES FOR INPUT FEATURES IF  
AVAILABLE BY DEFAULT “X0” “X1” “XNFEATURES” IS USED  
RETURNS  
OUTPUTFEATURENAMES LIST OF STRING LENGTH NOUTPUTFEATURES  
GETPARAMS SELFDEEPTURE  
GET PARAMETERS FOR THIS ESTIMATOR  
PARAMETERS  
DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBOBJECTS THAT ARE ESTIMATORS  
RETURNS  
PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES  
SETPARAMS SELFPARAMS  
SET THE PARAMETERS OF THIS ESTIMATOR  
THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE  
PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT  
OF A NESTED OBJECT  
RETURNS  
SELF  
TRANSFORM SELF  
TRANSFORM DATA TO POLYNOMIAL FEATURES  
PARAMETERS  
XARRAYLIKE OR CSR CSC SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE DATA TO TRANS  
FORM ROW BY ROW  
2264 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

PREFER CSR OVER CSC FOR SPARSE INPUT FOR SPEED BUT CSC IS REQUIRED IF THE DEGREE IS 4 OR HIGHER IF THE DEGREE IS LESS THAN 4 AND THE INPUT FORMAT IS CSC IT WILL BE CONVERTED TO CSR HAVE ITS POLYNOMIAL FEATURES GENERATED THEN CONVERTED BACK TO CSC IF THE DEGREE IS 2 OR 3 THE METHOD DESCRIBED IN “LEVERAGING SPARSITY TO SPEED UP POLYNOMIAL FEATURE EXPANSIONS OF CSR MATRICES USING KSIMPLEX NUMBERS” BY ANDREW NYSTROM AND JOHN HUGHES IS USED WHICH IS MUCH FASTER THAN THE METHOD USED ON CSC INPUT FOR THIS REASON A CSC INPUT WILL BE CONVERTED TO CSR AND THE OUTPUT WILL BE CONVERTED BACK TO CSC PRIOR TO BEING RETURNED HENCE THE PREFERENCE OF CSR

RETURNS

XPNPNDARRAY OR CSRCSC SPARSE MATRIX SHAPE NSAMPLES NP THE MATRIX OF FEATURES WHERE NP IS THE NUMBER OF POLYNOMIAL FEATURES GENERATED FROM THE COMBINATION OF INPUTS EXAMPLES USING SKLEARNPREPROCESSINGPOLYNOMIALFEATURES

- POLYNOMIAL INTERPOLATION
- ROBUST LINEAR ESTIMATOR FITTING
- UNDERFITTING VS OVERFITTING

63414SKLEARNPREPROCESSING POWERTRANSFORMER

CLASSSKLEARNPREPROCESSING POWERTRANSFORMER METHOD‘YEOJOHNSON’ STANDARDIZETRUE COPYTRUE

APPLY A POWER TRANSFORM FEATUREWISE TO MAKE DATA MORE GAUSSIANLIKE POWER TRANSFORMS ARE A FAMILY OF PARAMETRIC MONOTONIC TRANSFORMATIONS THAT ARE APPLIED TO MAKE DATA MORE GAUSSIANLIKE THIS IS USEFUL FOR MODELING ISSUES RELATED TO HETEROSCEDASTICITY NONCONSTANT VARIANCE OR OTHER SITUATIONS WHERE NORMALITY IS DESIRED CURRENTLY POWERTRANSFORMER SUPPORTS THE BOXCOX TRANSFORM AND THE YEOJOHNSON TRANSFORM THE OPTIMAL PARAMETER FOR STABILIZING VARIANCE AND MINIMIZING SKEWNESS IS ESTIMATED THROUGH MAXIMUM LIKELIHOOD BOXCOX REQUIRES INPUT DATA TO BE STRICTLY POSITIVE WHILE YEOJOHNSON SUPPORTS BOTH POSITIVE OR NEGATIVE DATA BY DEFAULT ZEROMEAN UNITVARIANCE NORMALIZATION IS APPLIED TO THE TRANSFORMED DATA READ MORE IN THE USER GUIDE

PARAMETERS

METHOD STR DEFAULT‘YEOJOHNSON’ THE POWER TRANSFORM METHOD AVAILABLE METHODS ARE

- ‘YEOJOHNSON’ RF3E1504535DE1 WORKS WITH POSITIVE AND NEGATIVE VALUES
- ‘BOXCOX’ RF3E1504535DE2 ONLY WORKS WITH STRICTLY POSITIVE VALUES

STANDARDIZE BOOLEAN DEFAULTTRUE SET TO TRUE TO APPLY ZEROMEAN UNITVARIANCE NORMALIZATION TO THE TRANSFORMED OUTPUT

COPY BOOLEAN OPTIONAL DEFAULTTRUE SET TO FALSE TO PERFORM INPLACE COMPUTATION DURING TRANSFORMATION

ATTRIBUTES

LAMBDAS ARRAY OF FLOAT SHAPE NFEATURES THE PARAMETERS OF THE POWER TRANSFORMATION FOR THE SELECTED FEATURES

SEE ALSO

634SKLEARNPREPROCESSING PREPROCESSING AND NORMALIZATION 2265

SCIKITLEARN USER GUIDE RELEASE 0213

POWERTRANSFORM EQUIVALENT FUNCTION WITHOUT THE ESTIMATOR API

QUANTILETRANSFORMER MAPS DATA TO A STANDARD NORMAL DISTRIBUTION WITH THE PARAMETER OUTPUTDISTRIBUTIONNORMAL

NOTES

NANS ARE TREATED AS MISSING VALUES DISREGARDED IN FIT AND MAINTAINED IN TRANSFORM

FOR A COMPARISON OF THE DIFFERENT SCALERS TRANSFORMERS AND NORMALIZERS SEE EXAMP  
PLESPREPROCESSINGPLOTALLSCALINGPY

REFERENCES

RF3E1504535DE1 RF3E1504535DE2

EXAMPLES

```
import numpy as np
from sklearn.preprocessing import PowerTransformer
pt = PowerTransformer()
data = [1, 2, 3, 2, 4, 5]
print(pt.fit(data))
power_transformer = copy.deepcopy(pt)
power_transformer.method = 'yeojohnson'
power_transformer.standardize = True
print(pt.lambda_)
1386 3100
print(pt.transform(data))
1316 0707
0209 0707
1106 1414
```

METHODS

fit(self, X, y=None) Estimate the optimal parameter lambda for each feature

fit\_transform(self, X, y=None) Estimate the optimal parameter lambda for each feature

get\_params(self, deep=True) Get parameters for this estimator

inverse\_transform(self, X) Apply the inverse power transformation using the fitted lambda

set\_params(self, \*\*kwargs) Set the parameters of this estimator

transform(self, X) Apply the power transform to each feature using the fitted lambda

init(self, method='yeojohnson', standardize=True, copy=True) Initialize the estimator

fit(self, X, y=None) Estimate the optimal parameter lambda for each feature

The optimal lambda parameter for minimizing skewness is estimated on each feature independently using maximum likelihood

PARAMETERS

2266 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

XARRAYLIKE SHAPE NSAMPLES NFEATURES THE DATA USED TO ESTIMATE THE OPTIMAL TRANSFORMATION PARAMETERS

YIGNORED

RETURNS

SELF OBJECT

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

INVERSETRANSFORM SELF

APPLY THE INVERSE POWER TRANSFORMATION USING THE FITTED LAMBDA

THE INVERSE OF THE BOXCOX TRANSFORMATION IS GIVEN BY

IF LAMBDA 0

X EXPTRANS

ELSE

X XTRANS LAMBDA 1 LAMBDA

THE INVERSE OF THE YEOJOHNSON TRANSFORMATION IS GIVEN BY

IFX 0 AND LAMBDA 0

X EXPTRANS 1

ELIFX 0 AND LAMBDA 0

X XTRANS LAMBDA 1 LAMBDA 1

ELIFX 0 AND LAMBDA 2

X 1 2 LAMBDA XTRANS 1 1 2 LAMBDA

ELIFX 0 AND LAMBDA 2

X 1 EXPTRANS

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES THE TRANSFORMED DATA

RETURNS

XARRAYLIKE SHAPE NSAMPLES NFEATURES THE ORIGINAL DATA

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELF

APPLY THE POWER TRANSFORM TO EACH FEATURE USING THE FITTED LAMBDA

634SKLEARNPREPROCESSING PREPROCESSING AND NORMALIZATION 2267

SCIKITLEARN USER GUIDE RELEASE 0213

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES THE DATA TO BE TRANSFORMED USING A POWER TRANSFORMATION

RETURNS

XTRANS ARRAYLIKE SHAPE NSAMPLES NFEATURES THE TRANSFORMED DATA

EXAMPLES USING SKLEARNPREPROCESSINGPOWERTRANSFORMER

- MAP DATA TO A NORMAL DISTRIBUTION
- COMPARE THE EFFECT OF DIFFERENT SCALERS ON DATA WITH OUTLIERS

63415SKLEARNPREPROCESSING QUANTILETRANSFORMER

CLASSSSKLEARNPREPROCESSING QUANTILETRANSFORMER NQUANTILES1000 OUT

PUTDISTRIBUTION‘UNIFORM’ IG

NOREIMPLICITZEROSFALSE SUBSAM

PLE1000000 RANDOMSTATENONE

COPYTRUE

TRANSFORM FEATURES USING QUANTILES INFORMATION

THIS METHOD TRANSFORMS THE FEATURES TO FOLLOW A UNIFORM OR A NORMAL DISTRIBUTION THEREFORE FOR A GIVEN FEATURE THIS TRANSFORMATION TENDS TO SPREAD OUT THE MOST FREQUENT VALUES IT ALSO REDUCES THE IMPACT OF MARGINAL OUTLIERS THIS IS THEREFORE A ROBUST PREPROCESSING SCHEME

THE TRANSFORMATION IS APPLIED ON EACH FEATURE INDEPENDENTLY FIRST AN ESTIMATE OF THE CUMULATIVE DISTRIBUTION FUNCTION OF A FEATURE IS USED TO MAP THE ORIGINAL VALUES TO A UNIFORM DISTRIBUTION THE OBTAINED VALUES ARE THEN MAPPED TO THE DESIRED OUTPUT DISTRIBUTION USING THE ASSOCIATED QUANTILE FUNCTION FEATURES VALUES OF NEWUNSEEN DATA THAT FALL BELOW OR ABOVE THE FITTED RANGE WILL BE MAPPED TO THE BOUNDS OF THE OUTPUT DISTRIBUTION NOTE THAT THIS TRANSFORM IS NONLINEAR IT MAY DISTORT LINEAR CORRELATIONS BETWEEN VARIABLES MEASURED AT THE SAME SCALE BUT RENDERS VARIABLES MEASURED AT DIFFERENT SCALES MORE DIRECTLY COMPARABLE

READ MORE IN THE USER GUIDE

PARAMETERS

NQUANTILES INT OPTIONAL DEFAULT1000 OR NSAMPLES NUMBER OF QUANTILES TO BE COMPUTED IT CORRESPONDS TO THE NUMBER OF LANDMARKS USED TO DISCRETIZE THE CUMULATIVE DISTRIBUTION FUNCTION IF NQUANTILES IS LARGER THAN THE NUMBER OF SAMPLES NQUANTILES IS SET TO THE NUMBER OF SAMPLES AS A LARGER NUMBER OF QUANTILES DOES NOT GIVE A BETTER APPROXIMATION OF THE CUMULATIVE DISTRIBUTION FUNCTION ESTIMATOR

OUTPUTDISTRIBUTION STR OPTIONAL DEFAULT‘UNIFORM’ MARGINAL DISTRIBUTION FOR THE TRANSFORMED DATA THE CHOICES ARE ‘UNIFORM’ DEFAULT OR ‘NORMAL’

IGNOREIMPLICITZEROS BOOL OPTIONAL DEFAULTFALSE ONLY APPLIES TO SPARSE MATRICES IF TRUE THE SPARSE ENTRIES OF THE MATRIX ARE DISCARDED TO COMPUTE THE QUANTILE STATISTICS IF FALSE THESE ENTRIES ARE TREATED AS ZEROS

SUBSAMPLE INT OPTIONAL DEFAULT1E5 MAXIMUM NUMBER OF SAMPLES USED TO ESTIMATE THE QUANTILES FOR COMPUTATIONAL EFFICIENCY NOTE THAT THE SUBSAMPLING PROCEDURE MAY DIFFER FOR VALUEIDENTICAL SPARSE AND DENSE MATRICES

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RAN

DOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE

RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS

2268 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

THE RANDOMSTATE INSTANCE USED BY NPRANDOM NOTE THAT THIS IS USED BY SUBSAMPLING AND SMOOTHING NOISE

COPY BOOLEAN OPTIONAL DEFAULTTRUE SET TO FALSE TO PERFORM INPLACE TRANSFORMATION AND AVOID A COPY IF THE INPUT IS ALREADY A NUMPY ARRAY

ATTRIBUTES

NQUANTILES INTEGER THE ACTUAL NUMBER OF QUANTILES USED TO DISCRETIZE THE CUMULATIVE DISTRIBUTION FUNCTION

QUANTILES NDARRAY SHAPE NQUANTILES NFEATURES THE VALUES CORRESPONDING THE QUANTILES OF REFERENCE

REFERENCES NDARRAY SHAPENQUANTILES QUANTILES OF REFERENCES

SEE ALSO

QUANTILETRANSFORM EQUIVALENT FUNCTION WITHOUT THE ESTIMATOR API

POWERTRANSFORMER PERFORM MAPPING TO A NORMAL DISTRIBUTION USING A POWER TRANSFORM

STANDARDSCALER PERFORM STANDARDIZATION THAT IS FASTER BUT LESS ROBUST TO OUTLIERS

ROBUSTSCALER PERFORM ROBUST STANDARDIZATION THAT REMOVES THE INFLUENCE OF OUTLIERS BUT DOES NOT PUT OUTLIERS AND INLIERS ON THE SAME SCALE

NOTES

NANS ARE TREATED AS MISSING VALUES DISREGARDED IN FIT AND MAINTAINED IN TRANSFORM

FOR A COMPARISON OF THE DIFFERENT SCALERS TRANSFORMERS AND NORMALIZERS SEE EXAMPLESPREPROCESSINGPLOTALLSCALINGPY

EXAMPLES

```
import numpy as np
from sklearn.preprocessing import QuantileTransformer
rng = np.random.RandomState(0)
X = np.sort(rng.normal(loc=0.5, scale=0.25, size=(25, 1)), axis=0)
qt = QuantileTransformer(n_quantiles=10, random_state=0)
qt.fit_transform(X)
```

ARRAY

METHODS

fit(self, X, y) COMPUTE THE QUANTILES USED FOR TRANSFORMING

fit\_transform(self, X, y) FIT TO DATA THEN TRANSFORM IT

get\_params(self, deep=True) GET PARAMETERS FOR THIS ESTIMATOR

inverse\_transform(self, X) BACKPROJECTION TO THE ORIGINAL SPACE

set\_params(self, \*\*params) SET THE PARAMETERS OF THIS ESTIMATOR

transform(self, X) FEATUREWISE TRANSFORMATION OF THE DATA

634SKLEARNPREPROCESSING PREPROCESSING AND NORMALIZATION 2269

SCIKITLEARN USER GUIDE RELEASE 0213

INIT SELF NQUANTILES 1000 OUTPUT DISTRIBUTION 'UNIFORM' IGNORE IMPLICIT ZEROS FALSE SUB  
SAMPLE 100000 RANDOM STATE NONE COPY TRUE

FIT SELF X NONE

COMPUTE THE QUANTILES USED FOR TRANSFORMING  
PARAMETERS

X NDARRAY OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE DATA USED TO SCALE  
ALONG THE FEATURES AXIS IF A SPARSE MATRIX IS PROVIDED IT WILL BE CONVERTED INTO  
A SPARSE CSC MATRIX ADDITIONALLY THE SPARSE MATRIX NEEDS TO BE NONNEGATIVE IF  
IGNORE IMPLICIT ZEROS IS FALSE

RETURNS  
SELF OBJECT

FIT TRANSFORM SELF X NONE FIT PARAMS

FIT TO DATA THEN TRANSFORM IT

FITS TRANSFORMER TO X AND Y WITH OPTIONAL PARAMETERS FIT PARAMS AND RETURNS A TRANSFORMED VERSION OF X  
PARAMETERS

X NUMPY ARRAY OF SHAPE NSAMPLES NFEATURES TRAINING SET  
Y NUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES

RETURNS  
X NEW NUMPY ARRAY OF SHAPE NSAMPLES NFEATURES NEW TRANSFORMED ARRAY

GET PARAMS SELF DEEP TRUE

GET PARAMETERS FOR THIS ESTIMATOR  
PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBJECTS THAT ARE ESTIMATORS

RETURNS  
PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

INVERSE TRANSFORM SELF X

BACK PROJECTION TO THE ORIGINAL SPACE  
PARAMETERS

X NDARRAY OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE DATA USED TO SCALE  
ALONG THE FEATURES AXIS IF A SPARSE MATRIX IS PROVIDED IT WILL BE CONVERTED INTO  
A SPARSE CSC MATRIX ADDITIONALLY THE SPARSE MATRIX NEEDS TO BE NONNEGATIVE IF  
IGNORE IMPLICIT ZEROS IS FALSE

RETURNS  
X NDARRAY OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE PROJECTED DATA

SET PARAMS SELF PARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE  
PARAMETERS OF THE FORM COMPONENT PARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT  
OF A NESTED OBJECT

RETURNS

2270 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

SELF

TRANSFORM SELF

FEATUREWISE TRANSFORMATION OF THE DATA

PARAMETERS

XNDARRAY OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE DATA USED TO SCALE  
ALONG THE FEATURES AXIS IF A SPARSE MATRIX IS PROVIDED IT WILL BE CONVERTED INTO  
A SPARSECSCMATRIX ADDITIONALLY THE SPARSE MATRIX NEEDS TO BE NONNEGATIVE IF  
IGNOREIMPLICITZEROS IS FALSE

RETURNS

XTNDARRAY OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE PROJECTED DATA

EXAMPLES USING SKLEARNPREPROCESSINGQUANTILETRANSFORMER

- EFFECT OF TRANSFORMING THE TARGETS IN REGRESSION MODEL
- MAP DATA TO A NORMAL DISTRIBUTION
- COMPARE THE EFFECT OF DIFFERENT SCALERS ON DATA WITH OUTLIERS

63416SKLEARNPREPROCESSING ROBUSTSCALER

CLASSSSKLEARNPREPROCESSING ROBUSTSCALER WITHCENTERINGTRUE WITHSCALINGTRUE QUAN  
TILERANGE250 750 COPYTRUE

SCALE FEATURES USING STATISTICS THAT ARE ROBUST TO OUTLIERS

THIS SCALER REMOVES THE MEDIAN AND SCALES THE DATA ACCORDING TO THE QUANTILE RANGE DEFAULTS TO IQR INTERQUARTILE  
RANGE THE IQR IS THE RANGE BETWEEN THE 1ST QUANTILE 25TH QUANTILE AND THE 3RD QUANTILE 75TH QUANTILE  
CENTERING AND SCALING HAPPEN INDEPENDENTLY ON EACH FEATURE BY COMPUTING THE RELEVANT STATISTICS ON THE SAMPLES  
IN THE TRAINING SET MEDIAN AND INTERQUARTILE RANGE ARE THEN STORED TO BE USED ON LATER DATA USING THE TRANSFORM  
METHOD

STANDARDIZATION OF A DATASET IS A COMMON REQUIREMENT FOR MANY MACHINE LEARNING ESTIMATORS TYPICALLY THIS IS  
DONE BY REMOVING THE MEAN AND SCALING TO UNIT VARIANCE HOWEVER OUTLIERS CAN OFTEN INFLUENCE THE SAMPLE MEAN  
VARIANCE IN A NEGATIVE WAY IN SUCH CASES THE MEDIAN AND THE INTERQUARTILE RANGE OFTEN GIVE BETTER RESULTS

NEW IN VERSION 017

READ MORE IN THE USER GUIDE

PARAMETERS

WITHCENTERING BOOLEAN TRUE BY DEFAULT IF TRUE CENTER THE DATA BEFORE SCALING THIS WILL  
CAUSETRANSFORM TO RAISE AN EXCEPTION WHEN ATTEMPTED ON SPARSE MATRICES BECAUSE CEN  
TERING THEM ENTAILS BUILDING A DENSE MATRIX WHICH IN COMMON USE CASES IS LIKELY TO BE TOO  
LARGE TO FIT IN MEMORY

WITHSCALING BOOLEAN TRUE BY DEFAULT IF TRUE SCALE THE DATA TO INTERQUARTILE RANGE

QUANTILERANGE TUPLE QMIN QMAX 00 QMIN QMAX 1000 DEFAULT 250 750

1ST QUANTILE 3RD QUANTILE IQR QUANTILE RANGE USED TO CALCULATE SCALE

NEW IN VERSION 018

COPY BOOLEAN OPTIONAL DEFAULT IS TRUE IF FALSE TRY TO AVOID A COPY AND DO INPLACE SCALING  
INSTEAD THIS IS NOT GUARANTEED TO ALWAYS WORK INPLACE EG IF THE DATA IS NOT A NUMPY ARRAY  
OR SCIPYSPARSE CSR MATRIX A COPY MAY STILL BE RETURNED

6345SKLEARNPREPROCESSING PREPROCESSING AND NORMALIZATION 2271

SCIKITLEARN USER GUIDE RELEASE 0213

ATTRIBUTES

CENTER ARRAY OF FLOATS THE MEDIAN VALUE FOR EACH FEATURE IN THE TRAINING SET

SCALE ARRAY OF FLOATS THE SCALED INTERQUARTILE RANGE FOR EACH FEATURE IN THE TRAINING SET

NEW IN VERSION 017 SCALE ATTRIBUTE

SEE ALSO

ROBUSTSCALE EQUIVALENT FUNCTION WITHOUT THE ESTIMATOR API

SKLEARNDECOMPOSITIONPCA FURTHER REMOVES THE LINEAR CORRELATION ACROSS FEATURES WITH 'WHITENTRUE'

NOTES

FOR A COMPARISON OF THE DIFFERENT SCALERS TRANSFORMERS AND NORMALIZERS SEE EXAM

PLESPREPROCESSINGPLOTALLSCALINGPY

HTTPSENWIKIPEDIAORGWIKIMEDIAN HTTPSENWIKIPEDIAORGWIKIINTERQUANTILERANGE

EXAMPLES

FROM SKLEARNPREPROCESSING IMPORT ROBUSTSCALER

X 1 2 2

2 1 3

4 1 2

TRANSFORMER ROBUSTSCALERFITX

TRANSFORMER

ROBUSTSCALERCOPYTRUE QUANTILERANGE250 750 WITHCENTERINGTRUE

WITHSCALINGTRUE

TRANSFORMERTRANSFORMX

ARRAY 0 2 0

1 0 04

1 0 16

METHODS

FITSELF X Y COMPUTE THE MEDIAN AND QUANTILES TO BE USED FOR SCAL

ING

FITTRANSFORM SELF X Y FIT TO DATA THEN TRANSFORM IT

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

INVERSETRANSFORM SELF X SCALE BACK THE DATA TO THE ORIGINAL REPRESENTATION

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF X CENTER AND SCALE THE DATA

INIT SELFWITHCENTERINGTRUE WITHSCALINGTRUE QUANTILERANGE250 750 COPYTRUE

FITSELFXYNONE

COMPUTE THE MEDIAN AND QUANTILES TO BE USED FOR SCALING

PARAMETERS

2272 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

XARRAYLIKE SHAPE NSAMPLES NFEATURES THE DATA USED TO COMPUTE THE MEDIAN AND QUANTILES USED FOR LATER SCALING ALONG THE FEATURES AXIS

FITTRANSFORM SELFXYNONE FITPARAMS

FIT TO DATA THEN TRANSFORM IT

FITS TRANSFORMER TO X AND Y WITH OPTIONAL PARAMETERS FITPARAMS AND RETURNS A TRANSFORMED VERSION OF X

PARAMETERS

XNUMPY ARRAY OF SHAPE NSAMPLES NFEATURES TRAINING SET

YNUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES

RETURNS

XNEW NUMPY ARRAY OF SHAPE NSAMPLES NFEATURESNEW TRANSFORMED ARRAY

GETPARAMS SELFDEEPTTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

INVERSETRANSFORM SELF

SCALE BACK THE DATA TO THE ORIGINAL REPRESENTATION

PARAMETERS

XARRAYLIKE THE DATA USED TO SCALE ALONG THE SPECIFIED AXIS

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE

PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT

OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELF

CENTER AND SCALE THE DATA

PARAMETERS

XARRAYLIKE SPARSE MATRIX THE DATA USED TO SCALE ALONG THE SPECIFIED AXIS

EXAMPLES USING SKLEARNPREPROCESSINGROBUSTSCALER

- COMPARE THE EFFECT OF DIFFERENT SCALERS ON DATA WITH OUTLIERS

63417SKLEARNPREPROCESSING STANDARDSCALER

CLASSSSKLEARNPREPROCESSING STANDARDSCALER COPYTRUE WITHMEANTRUE WITHSTDTRUE

STANDARDIZE FEATURES BY REMOVING THE MEAN AND SCALING TO UNIT VARIANCE

634SKLEARNPREPROCESSING PREPROCESSING AND NORMALIZATION 2273

SCIKITLEARN USER GUIDE RELEASE 0213

THE STANDARD SCORE OF A SAMPLE  $X_i$  IS CALCULATED AS

$$Z = \frac{X_i - \mu}{\sigma}$$

WHERE  $\mu$  IS THE MEAN OF THE TRAINING SAMPLES OR ZERO IF `WITHMEAN=False` AND  $\sigma$  IS THE STANDARD DEVIATION OF THE TRAINING SAMPLES OR ONE IF `WITHSTD=False`

CENTERING AND SCALING HAPPEN INDEPENDENTLY ON EACH FEATURE BY COMPUTING THE RELEVANT STATISTICS ON THE SAMPLES IN THE TRAINING SET MEAN AND STANDARD DEVIATION ARE THEN STORED TO BE USED ON LATER DATA USING THE TRANSFORM METHOD

STANDARDIZATION OF A DATASET IS A COMMON REQUIREMENT FOR MANY MACHINE LEARNING ESTIMATORS THEY MIGHT BEHAVE BADLY IF THE INDIVIDUAL FEATURES DO NOT MORE OR LESS LOOK LIKE STANDARD NORMALLY DISTRIBUTED DATA EG GAUSSIAN WITH 0 MEAN AND UNIT VARIANCE

FOR INSTANCE MANY ELEMENTS USED IN THE OBJECTIVE FUNCTION OF A LEARNING ALGORITHM SUCH AS THE RBF KERNEL OF SUPPORT VECTOR MACHINES OR THE L1 AND L2 REGULARIZERS OF LINEAR MODELS ASSUME THAT ALL FEATURES ARE CENTERED AROUND 0 AND HAVE VARIANCE IN THE SAME ORDER IF A FEATURE HAS A VARIANCE THAT IS ORDERS OF MAGNITUDE LARGER THAN OTHERS IT MIGHT DOMINATE THE OBJECTIVE FUNCTION AND MAKE THE ESTIMATOR UNABLE TO LEARN FROM OTHER FEATURES CORRECTLY AS EXPECTED

THIS SCALER CAN ALSO BE APPLIED TO SPARSE CSR OR CSC MATRICES BY PASSING `WITHMEAN=False` TO AVOID BREAKING THE SPARSITY STRUCTURE OF THE DATA

READ MORE IN THE USER GUIDE

PARAMETERS

`COPY` BOOLEAN OPTIONAL DEFAULT `True` IF `False` TRY TO AVOID A COPY AND DO INPLACE SCALING INSTEAD THIS IS NOT GUARANTEED TO ALWAYS WORK INPLACE EG IF THE DATA IS NOT A NUMPY ARRAY OR `ScipySparse` CSR MATRIX A COPY MAY STILL BE RETURNED

`WITHMEAN` BOOLEAN `True` BY DEFAULT IF `True` CENTER THE DATA BEFORE SCALING THIS DOES NOT WORK AND WILL RAISE AN EXCEPTION WHEN ATTEMPTED ON SPARSE MATRICES BECAUSE CENTERING THEM ENTAILS BUILDING A DENSE MATRIX WHICH IN COMMON USE CASES IS LIKELY TO BE TOO LARGE TO FIT IN MEMORY

`WITHSTD` BOOLEAN `True` BY DEFAULT IF `True` SCALE THE DATA TO UNIT VARIANCE OR EQUIVALENTLY UNIT STANDARD DEVIATION

ATTRIBUTES

`SCALE` NDARRAY OR `None` SHAPE `n_features` PER FEATURE RELATIVE SCALING OF THE DATA THIS IS CALCULATED USING `Np.sqrt(var)` EQUAL TO `None` WHEN `WITHSTD=False`

NEW IN VERSION 0.17 `SCALE`

`MEAN` NDARRAY OR `None` SHAPE `n_features` THE MEAN VALUE FOR EACH FEATURE IN THE TRAINING SET EQUAL TO `None` WHEN `WITHMEAN=False`

`VAR` NDARRAY OR `None` SHAPE `n_features` THE VARIANCE FOR EACH FEATURE IN THE TRAINING SET USED TO COMPUTE `SCALE` EQUAL TO `None` WHEN `WITHSTD=False`

`NSAMPLESSEEN` INT OR ARRAY SHAPE `n_features` THE NUMBER OF SAMPLES PROCESSED BY THE ESTIMATOR FOR EACH FEATURE IF THERE ARE NOT MISSING SAMPLES THE `NSAMPLESSEEN` WILL BE AN INTEGER OTHERWISE IT WILL BE AN ARRAY WILL BE RESET ON NEW CALLS TO `fit` BUT INCREMENTS ACROSS `partial_fit` CALLS

SEE ALSO

`SCALE` EQUIVALENT FUNCTION WITHOUT THE ESTIMATOR API

2274 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

SKLEARNDECOMPOSITIONPCA FURTHER REMOVES THE LINEAR CORRELATION ACROSS FEATURES WITH 'WHITENTRUE'

NOTES

NANS ARE TREATED AS MISSING VALUES DISREGARDED IN FIT AND MAINTAINED IN TRANSFORM

WE USE A BIASED ESTIMATOR FOR THE STANDARD DEVIATION EQUIVALENT TO NUMPYSTD DDOF0 NOTE THAT THE CHOICE OFDDOF IS UNLIKELY TO AFFECT MODEL PERFORMANCE

FOR A COMPARISON OF THE DIFFERENT SCALERS TRANSFORMERS AND NORMALIZERS SEE EXAM PLES

PREPROCESSINGPLOTALLSCALINGPY

EXAMPLES

```
FROM SKLEARNPREPROCESSING IMPORT STANDARDSCALER
DATA 0 0 0 0 1 1 1 1
SCALER STANDARDSCALER
PRINTSCALERFITDATA
STANDARDSCALERCOPYTRUE WITHMEANTRUE WITHSTDTRUE
PRINTSCALERMEAN
05 05
PRINTSCALERTRANSFORMDATA
1 1
1 1
1 1
1 1
PRINTSCALERTRANSFORM2 2
3 3
```

METHODS

FITSELF X Y COMPUTE THE MEAN AND STD TO BE USED FOR LATER SCALING

FITTRANSFORM SELF X Y FIT TO DATA THEN TRANSFORM IT

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

INVERSETRANSFORM SELF X COPY SCALE BACK THE DATA TO THE ORIGINAL REPRESENTATION

PARTIALFIT SELF X Y ONLINE COMPUTATION OF MEAN AND STD ON X FOR LATER SCALING

ING

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF X COPY PERFORM STANDARDIZATION BY CENTERING AND SCALING

INIT SELF COPYTRUE WITHMEANTRUE WITHSTDTRUE

FITSELFXYNONE

COMPUTE THE MEAN AND STD TO BE USED FOR LATER SCALING

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE DATA USED TO COMPUTE THE MEAN AND STANDARD DEVIATION USED FOR LATER SCALING ALONG THE FEATURES AXIS

YIGNORED

634SKLEARNPREPROCESSING PREPROCESSING AND NORMALIZATION 2275

SCIKITLEARN USER GUIDE RELEASE 0213

FITTRANSFORM SELFXYNONE FITPARAMS

FIT TO DATA THEN TRANSFORM IT

FITS TRANSFORMER TO X AND Y WITH OPTIONAL PARAMETERS FITPARAMS AND RETURNS A TRANSFORMED VERSION OF X

PARAMETERS

XNUMPY ARRAY OF SHAPE NSAMPLES NFEATURES TRAINING SET

YNUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES

RETURNS

XNEW NUMPY ARRAY OF SHAPE NSAMPLES NFEATURESNEW TRANSFORMED ARRAY

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

INVERSETRANSFORM SELFXCOPYNONE

SCALE BACK THE DATA TO THE ORIGINAL REPRESENTATION

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES THE DATA USED TO SCALE ALONG THE FEATURES AXIS

COPY BOOL OPTIONAL DEFAULT NONE COPY THE INPUT X OR NOT

RETURNS

XTR ARRAYLIKE SHAPE NSAMPLES NFEATURES TRANSFORMED ARRAY

PARTIALFIT SELFXYNONE

ONLINE COMPUTATION OF MEAN AND STD ON X FOR LATER SCALING ALL OF X IS PROCESSED AS A SINGLE BATCH THIS IS

INTENDED FOR CASES WHEN FIT IS NOT FEASIBLE DUE TO VERY LARGE NUMBER OF NSAMPLES OR BECAUSE X IS READ

FROM A CONTINUOUS STREAM

THE ALGORITHM FOR INCREMENTAL MEAN AND STD IS GIVEN IN EQUATION 15AB IN CHAN TONY F GENE H GOLUB

AND RANDALL J LEVEQUE “ALGORITHMS FOR COMPUTING THE SAMPLE VARIANCE ANALYSIS AND RECOMMENDATIONS”

THE AMERICAN STATISTICIAN 373 1983 242247

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE DATA USED TO COMPUTE THE

MEAN AND STANDARD DEVIATION USED FOR LATER SCALING ALONG THE FEATURES AXIS

YIGNORED

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE

PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT

OF A NESTED OBJECT

RETURNS

SELF

2276 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

TRANSFORM SELFXCOPYNONE

PERFORM STANDARDIZATION BY CENTERING AND SCALING

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES THE DATA USED TO SCALE ALONG THE FEATURES AXIS

COPY BOOL OPTIONAL DEFAULT NONE COPY THE INPUT X OR NOT

EXAMPLES USING SKLEARNPREPROCESSINGSTANDARDSCALER

- PREDICTION LATENCY
- CLASSIFIER COMPARISON
- DEMO OF DBSCAN CLUSTERING ALGORITHM
- COMPARING DIFFERENT HIERARCHICAL LINKAGE METHODS ON TOY DATASETS
- COMPARING DIFFERENT CLUSTERING ALGORITHMS ON TOY DATASETS
- COLUMN TRANSFORMER WITH MIXED TYPES
- MNIST CLASSIFICATION USING MULTINOMIAL LOGISTIC L1
- L1 PENALTY AND SPARSITY IN LOGISTIC REGRESSION
- COMPARING NEAREST NEIGHBORS WITH AND WITHOUT NEIGHBORHOOD COMPONENTS ANALYSIS
- DIMENSIONALITY REDUCTION WITH NEIGHBORHOOD COMPONENTS ANALYSIS
- VARYING REGULARIZATION IN MULTILAYER PERCEPTRON
- IMPORTANCE OF FEATURE SCALING
- FEATURE DISCRETIZATION
- COMPARE THE EFFECT OF DIFFERENT SCALERS ON DATA WITH OUTLIERS
- SVMANOVA SVM WITH UNIVARIATE FEATURE SELECTION
- RBF SVM PARAMETERS

PREPROCESSINGADDDUMMYFEATURE X VALUE AUGMENT DATASET WITH AN ADDITIONAL DUMMY FEATURE

PREPROCESSINGBINARIZE X THRESHOLD COPY BOOLEAN THRESHOLDING OF ARRAYLIKE OR SCIPYSPARSE MATRIX

PREPROCESSINGLABELBINARIZE Y CLASSES BINARIZE LABELS IN A ONEVSALL FASHION

PREPROCESSINGMAXABSSCALE X AXIS COPY SCALE EACH FEATURE TO THE 1 1 RANGE WITHOUT BREAKING THE SPARSITY

PREPROCESSINGMINMAXSCALE X TRANSFORMS FEATURES BY SCALING EACH FEATURE TO A GIVEN RANGE

PREPROCESSINGNORMALIZE X NORM AXIS SCALE INPUT VECTORS INDIVIDUALLY TO UNIT NORM VECTOR LENGTH

PREPROCESSINGQUANTILETRANSFORM X AXIS TRANSFORM FEATURES USING QUANTILES INFORMATION

PREPROCESSINGROBUSTSCALE X AXIS STANDARDIZE A DATASET ALONG ANY AXIS

PREPROCESSINGSSCALE X AXIS WITHMEAN STANDARDIZE A DATASET ALONG ANY AXIS

PREPROCESSINGPOWERTRANSFORM X METHOD POWER TRANSFORMS ARE A FAMILY OF PARAMETRIC MONO TONIC TRANSFORMATIONS THAT ARE APPLIED TO MAKE DATA MORE GAUSSIANLIKE

634SKLEARNPREPROCESSING PREPROCESSING AND NORMALIZATION 2277

SCIKITLEARN USER GUIDE RELEASE 0213

63418SKLEARNPREPROCESSING ADDDUMMYFEATURE

SKLEARNPREPROCESSING ADDDUMMYFEATURE XVALUE10

AUGMENT DATASET WITH AN ADDITIONAL DUMMY FEATURE

THIS IS USEFUL FOR FITTING AN INTERCEPT TERM WITH IMPLEMENTATIONS WHICH CANNOT OTHERWISE FIT IT DIRECTLY

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES DATA

VALUE FLOAT VALUE TO USE FOR THE DUMMY FEATURE

RETURNS

XARRAY SPARSE MATRIX SHAPE NSAMPLES NFEATURES 1 SAME DATA WITH DUMMY FEATURE

ADDED AS FIRST COLUMN

EXAMPLES

FROM SKLEARNPREPROCESSING IMPORT ADDDUMMYFEATURE

ADDDUMMYFEATURE0 1 1 0

ARRAY1 0 1

1 1 0

63419SKLEARNPREPROCESSING BINARIZE

SKLEARNPREPROCESSING BINARIZE XTHRESHOLD00 COPYTRUE

BOOLEAN THRESHOLDING OF ARRAYLIKE OR SCIPYSPARSE MATRIX

READ MORE IN THE USER GUIDE

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE DATA TO BINARIZE ELEMENT BY

ELEMENT SCIPYSPARSE MATRICES SHOULD BE IN CSR OR CSC FORMAT TO AVOID AN UNNECESSARY

COPY

THRESHOLD FLOAT OPTIONAL 00 BY DEFAULT FEATURE VALUES BELOW OR EQUAL TO THIS ARE REPLACED BY

0 ABOVE IT BY 1 THRESHOLD MAY NOT BE LESS THAN 0 FOR OPERATIONS ON SPARSE MATRICES

COPY BOOLEAN OPTIONAL DEFAULT TRUE SET TO FALSE TO PERFORM INPLACE BINARIZATION AND AVOID A

COPY IF THE INPUT IS ALREADY A NUMPY ARRAY OR A SCIPYSPARSE CSR CSC MATRIX AND IF AXIS

IS 1

SEE ALSO

BINARIZER PERFORMS BINARIZATION USING THE TRANSFORMER API EG AS PART OF A PREPROCESSING SKLEARN

PIPELINEPIPELINE

63420SKLEARNPREPROCESSING LABELBINARIZE

SKLEARNPREPROCESSING LABELBINARIZE Y CLASSES NEGLABEL0 POSLABEL1

SPARSEOUTPUTFALSE

BINARIZE LABELS IN A ONEVSALL FASHION

2278 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

SEVERAL REGRESSION AND BINARY CLASSIFICATION ALGORITHMS ARE AVAILABLE IN SCIKITLEARN A SIMPLE WAY TO EXTEND THESE ALGORITHMS TO THE MULTICLASS CLASSIFICATION CASE IS TO USE THE SOCALLED ONEVSALL SCHEME

THIS FUNCTION MAKES IT POSSIBLE TO COMPUTE THIS TRANSFORMATION FOR A FIXED SET OF CLASS LABELS KNOWN AHEAD OF TIME

PARAMETERS

YARRAYLIKE SEQUENCE OF INTEGER LABELS OR MULTILABEL DATA TO ENCODE

CLASSES ARRAYLIKE OF SHAPE NCLASSES UNIQUELY HOLDS THE LABEL FOR EACH CLASS

NEGLABEL INT DEFAULT 0 VALUE WITH WHICH NEGATIVE LABELS MUST BE ENCODED

POSLABEL INT DEFAULT 1 VALUE WITH WHICH POSITIVE LABELS MUST BE ENCODED

SPARSEOUTPUT BOOLEAN DEFAULT FALSE SET TO TRUE IF OUTPUT BINARY ARRAY IS DESIRED IN CSR

SPARSE FORMAT

RETURNS

YNUMPY ARRAY OR CSR MATRIX OF SHAPE NSAMPLES NCLASSES SHAPE WILL BE NSAMPLES 1

FOR BINARY PROBLEMS

SEE ALSO

LABELBINARIZER CLASS USED TO WRAP THE FUNCTIONALITY OF LABELBINARIZE AND ALLOW FOR FITTING TO CLASSES INDEPENDENTLY OF THE TRANSFORM OPERATION

EXAMPLES

FROM SKLEARNPREPROCESSING IMPORT LABELBINARIZE

LABELBINARIZE1 6 CLASSES1 2 4 6

ARRAY1 0 0 0

0 0 0 1

THE CLASS ORDERING IS PRESERVED

LABELBINARIZE1 6 CLASSES1 6 4 2

ARRAY1 0 0 0

0 1 0 0

BINARY TARGETS TRANSFORM TO A COLUMN VECTOR

LABELBINARIZEYES NO NO YES CLASSESNO YES

ARRAY1

0

0

1

EXAMPLES USING SKLEARNPREPROCESSINGLABELBINARIZE

- RECEIVER OPERATING CHARACTERISTIC ROC
- PRECISIONRECALL

634SKLEARNPREPROCESSING PREPROCESSING AND NORMALIZATION 2279

SCIKITLEARN USER GUIDE RELEASE 0213

63421SKLEARNPREPROCESSING MAXABSSCALE

SKLEARNPREPROCESSING MAXABSSCALE XAXIS0 COPYTRUE

SCALE EACH FEATURE TO THE 1 1 RANGE WITHOUT BREAKING THE SPARSITY

THIS ESTIMATOR SCALES EACH FEATURE INDIVIDUALLY SUCH THAT THE MAXIMAL ABSOLUTE VALUE OF EACH FEATURE IN THE TRAINING SET WILL BE 10

THIS SCALER CAN ALSO BE APPLIED TO SPARSE CSR OR CSC MATRICES

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES THE DATA

AXIS INT 0 BY DEFAULT AXIS USED TO SCALE ALONG IF 0 INDEPENDENTLY SCALE EACH FEATURE OTHER WISE IF 1 SCALE EACH SAMPLE

COPY BOOLEAN OPTIONAL DEFAULT IS TRUE SET TO FALSE TO PERFORM INPLACE SCALING AND AVOID A COPY IF THE INPUT IS ALREADY A NUMPY ARRAY

SEE ALSO

MAXABSSCALER PERFORMS SCALING TO THE 1 1 RANGE USING THE“TRANSFORMER” API EG AS PART OF A PREPROCESSINGSKLEARNPIPELINEPIPELINE

NOTES

NANS ARE TREATED AS MISSING VALUES DISREGARDED TO COMPUTE THE STATISTICS AND MAINTAINED DURING THE DATA TRANSFORMATION

FOR A COMPARISON OF THE DIFFERENT SCALERS TRANSFORMERS AND NORMALIZERS SEE EXAMPLESPREPROCESSINGPLOTALLSCALINGPY

63422SKLEARNPREPROCESSING MINMAXSCALE

SKLEARNPREPROCESSING MINMAXSCALE XFEATURERANGE0 1AXIS0 COPYTRUE

TRANSFORMS FEATURES BY SCALING EACH FEATURE TO A GIVEN RANGE

THIS ESTIMATOR SCALES AND TRANSLATES EACH FEATURE INDIVIDUALLY SUCH THAT IT IS IN THE GIVEN RANGE ON THE TRAINING SET IE BETWEEN ZERO AND ONE

THE TRANSFORMATION IS GIVEN BY WHEN AXIS0

$$X_{STD} = \frac{X - X_{MINAXIS0}}{X_{MAXAXIS0} - X_{MINAXIS0}}$$

$$X_{SCALED} = X_{STD} \cdot MAX\_MIN\_MIN$$

WHERE MIN MAX FEATURERANGE

THE TRANSFORMATION IS CALCULATED AS WHEN AXIS0

$$X_{SCALED} = SCALE \cdot X - MIN\_XMINAXIS0 \cdot SCALE$$

WHERE SCALE MAX MIN XMAXAXIS0 XMINAXIS0

THIS TRANSFORMATION IS OFTEN USED AS AN ALTERNATIVE TO ZERO MEAN UNIT VARIANCE SCALING

READ MORE IN THE USER GUIDE

NEW IN VERSION 017 MINMAXSCALE FUNCTION INTERFACE TO SKLEARNPREPROCESSINGMINMAXSCALER

PARAMETERS

2280 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

XARRAYLIKE SHAPE NSAMPLES NFEATURES THE DATA

FEATURERANGE TUPLE MIN MAX DEFAULT0 1 DESIRED RANGE OF TRANSFORMED DATA

AXIS INT 0 BY DEFAULT AXIS USED TO SCALE ALONG IF 0 INDEPENDENTLY SCALE EACH FEATURE OTHERWISE IF 1 SCALE EACH SAMPLE

COPY BOOLEAN OPTIONAL DEFAULT IS TRUE SET TO FALSE TO PERFORM INPLACE SCALING AND AVOID A COPY IF THE INPUT IS ALREADY A NUMPY ARRAY

SEE ALSO

MINMAXSCALER PERFORMS SCALING TO A GIVEN RANGE USING THE“TRANSFORMER” API EG AS PART OF A PREPROCESSINGSKLEARNPIPELINEPIPELINE

NOTES

FOR A COMPARISON OF THE DIFFERENT SCALERS TRANSFORMERS AND NORMALIZERS SEE EXAMPLESPREPROCESSINGPLOTALLSCALINGPY

EXAMPLES USING SKLEARNPREPROCESSINGMINMAXSCALE

- COMPARE THE EFFECT OF DIFFERENT SCALERS ON DATA WITH OUTLIERS

63423SKLEARNPREPROCESSING NORMALIZE

SKLEARNPREPROCESSING NORMALIZE XNORM‘L2’ AXIS1 COPYTRUE RETURNNORMFALSE

SCALE INPUT VECTORS INDIVIDUALLY TO UNIT NORM VECTOR LENGTH

READ MORE IN THE USER GUIDE

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE DATA TO NORMALIZE ELEMENT BY ELEMENT SCIPYSPARSE MATRICES SHOULD BE IN CSR FORMAT TO AVOID AN UNNECESSARY COPY

NORM ‘L1’ ‘L2’ OR ‘MAX’ OPTIONAL ‘L2’ BY DEFAULT THE NORM TO USE TO NORMALIZE EACH NONZERO SAMPLE OR EACH NONZERO FEATURE IF AXIS IS 0

AXIS 0 OR 1 OPTIONAL 1 BY DEFAULT AXIS USED TO NORMALIZE THE DATA ALONG IF 1 INDEPENDENTLY NORMALIZE EACH SAMPLE OTHERWISE IF 0 NORMALIZE EACH FEATURE

COPY BOOLEAN OPTIONAL DEFAULT TRUE SET TO FALSE TO PERFORM INPLACE ROW NORMALIZATION AND AVOID A COPY IF THE INPUT IS ALREADY A NUMPY ARRAY OR A SCIPYSPARSE CSR MATRIX AND IF AXIS IS 1

RETURNNORM BOOLEAN DEFAULT FALSE WHETHER TO RETURN THE COMPUTED NORMS

RETURNS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES NORMALIZED INPUT X

NORMS ARRAY SHAPE NSAMPLES IF AXIS1 ELSE NFEATURES AN ARRAY OF NORMS ALONG GIVEN AXIS FOR X WHEN X IS SPARSE A NOTIMPLEMENTEDERROR WILL BE RAISED FOR NORM ‘L1’ OR ‘L2’

SEE ALSO

NORMALIZER PERFORMS NORMALIZATION USING THE TRANSFORMER API EG AS PART OF A PREPROCESSINGSKLEARNPIPELINEPIPELINE

6345SKLEARNPREPROCESSING PREPROCESSING AND NORMALIZATION 2281

NOTES

FOR A COMPARISON OF THE DIFFERENT SCALERS TRANSFORMERS AND NORMALIZERS SEE EXAM

PLESPREPROCESSINGPLOTALLSCALINGPY

63424SKLEARNPREPROCESSING QUANTILETRANSFORM

SKLEARNPREPROCESSING QUANTILETRANSFORM X AXIS0 NQUANTILES1000 OUT

PUTDISTRIBUTION'UNIFORM' IG

NOREIMPLICITZEROSFALSE SUBSAMPLE100000

RANDOMSTATENONE COPY'WARN'

TRANSFORM FEATURES USING QUANTILES INFORMATION

THIS METHOD TRANSFORMS THE FEATURES TO FOLLOW A UNIFORM OR A NORMAL DISTRIBUTION THEREFORE FOR A GIVEN FEATURE

THIS TRANSFORMATION TENDS TO SPREAD OUT THE MOST FREQUENT VALUES IT ALSO REDUCES THE IMPACT OF MARGINAL OUTLIERS

THIS IS THEREFORE A ROBUST PREPROCESSING SCHEME

THE TRANSFORMATION IS APPLIED ON EACH FEATURE INDEPENDENTLY FIRST AN ESTIMATE OF THE CUMULATIVE DISTRIBUTION

FUNCTION OF A FEATURE IS USED TO MAP THE ORIGINAL VALUES TO A UNIFORM DISTRIBUTION THE OBTAINED VALUES ARE THEN

MAPPED TO THE DESIRED OUTPUT DISTRIBUTION USING THE ASSOCIATED QUANTILE FUNCTION FEATURES VALUES OF NEWUNSEEN

DATA THAT FALL BELOW OR ABOVE THE FITTED RANGE WILL BE MAPPED TO THE BOUNDS OF THE OUTPUT DISTRIBUTION NOTE THAT

THIS TRANSFORM IS NONLINEAR IT MAY DISTORT LINEAR CORRELATIONS BETWEEN VARIABLES MEASURED AT THE SAME SCALE BUT

RENDERS VARIABLES MEASURED AT DIFFERENT SCALES MORE DIRECTLY COMPARABLE

READ MORE IN THE USER GUIDE

PARAMETERS

XARRAYLIKE SPARSE MATRIX THE DATA TO TRANSFORM

AXIS INT DEFAULT0 AXIS USED TO COMPUTE THE MEANS AND STANDARD DEVIATIONS ALONG IF 0

TRANSFORM EACH FEATURE OTHERWISE IF 1 TRANSFORM EACH SAMPLE

NQUANTILES INT OPTIONAL DEFAULT1000 OR NSAMPLES NUMBER OF QUANTILES TO BE COMPUTED

IT CORRESPONDS TO THE NUMBER OF LANDMARKS USED TO DISCRETIZE THE CUMULATIVE DISTRIBUTION

FUNCTION IF NQUANTILES IS LARGER THAN THE NUMBER OF SAMPLES NQUANTILES IS SET TO THE

NUMBER OF SAMPLES AS A LARGER NUMBER OF QUANTILES DOES NOT GIVE A BETTER APPROXIMATION OF

THE CUMULATIVE DISTRIBUTION FUNCTION ESTIMATOR

OUTPUTDISTRIBUTION STR OPTIONAL DEFAULT'UNIFORM' MARGINAL DISTRIBUTION FOR THE TRANS

FORMED DATA THE CHOICES ARE 'UNIFORM' DEFAULT OR 'NORMAL'

IGNOREIMPLICITZEROS BOOL OPTIONAL DEFAULTFALSE ONLY APPLIES TO SPARSE MATRICES IF TRUE

THE SPARSE ENTRIES OF THE MATRIX ARE DISCARDED TO COMPUTE THE QUANTILE STATISTICS IF FALSE

THESE ENTRIES ARE TREATED AS ZEROS

SUBSAMPLE INT OPTIONAL DEFAULT1E5 MAXIMUM NUMBER OF SAMPLES USED TO ESTIMATE THE

QUANTILES FOR COMPUTATIONAL EFFICIENCY NOTE THAT THE SUBSAMPLING PROCEDURE MAY DIFFER

FOR VALUEIDENTICAL SPARSE AND DENSE MATRICES

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RAN

DOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE

RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS

THE RANDOMSTATE INSTANCE USED BY NPRANDOM NOTE THAT THIS IS USED BY SUBSAMPLING AND

SMOOTHING NOISE

COPY BOOLEAN OPTIONAL DEFAULT"WARN" SET TO FALSE TO PERFORM INPLACE TRANSFORMATION AND

AVOID A COPY IF THE INPUT IS ALREADY A NUMPY ARRAY IF TRUE A COPY OF XIS TRANSFORMED

LEAVING THE ORIGINAL XUNCHANGED



SCIKITLEARN USER GUIDE RELEASE 0213

DEPRECATED SINCE VERSION 021 THE DEFAULT VALUE OF PARAMETER COPY WILL BE CHANGED FROM FALSE TO TRUE IN 023 THE CURRENT DEFAULT OF FALSE IS BEING CHANGED TO MAKE IT MORE CONSISTENT WITH THE DEFAULT COPY VALUES OF OTHER FUNCTIONS IN SKLEARNPREPROCESSING DATA FURTHERMORE THE CURRENT DEFAULT OF FALSE MAY HAVE UNEXPECTED SIDE EFFECTS BY MODIFYING THE VALUE OF XINPLACE

RETURNS

XTNDARRAY OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE TRANSFORMED DATA

SEE ALSO

QUANTILETRANSFORMER PERFORMS QUANTILEBASED SCALING USING THE TRANSFORMER API EG AS PART OF A PREPROCESSING SKLEARNPIPELINEPIPELINE

POWERTRANSFORM MAPS DATA TO A NORMAL DISTRIBUTION USING A POWER TRANSFORMATION

SCALE PERFORMS STANDARDIZATION THAT IS FASTER BUT LESS ROBUST TO OUTLIERS

ROBUSTSCALE PERFORMS ROBUST STANDARDIZATION THAT REMOVES THE INFLUENCE OF OUTLIERS BUT DOES NOT PUT OUT LIERS AND INLIERS ON THE SAME SCALE

NOTES

NANS ARE TREATED AS MISSING VALUES DISREGARDED IN FIT AND MAINTAINED IN TRANSFORM

FOR A COMPARISON OF THE DIFFERENT SCALERS TRANSFORMERS AND NORMALIZERS SEE EXAM

PLESPREPROCESSINGPLOTALLSCALINGPY

EXAMPLES

```
import numpy as np
from sklearn.preprocessing import QuantileTransformer
rng = np.random.RandomState(0)
X = np.sort(rng.normal(loc=0.5, scale=0.25, size=(25, 1)), axis=0)
quantile_transform(X, quantiles=10, random_state=0, copy=True)
```

ARRAY

EXAMPLES USING SKLEARNPREPROCESSINGQUANTILETRANSFORM

•EFFECT OF TRANSFORMING THE TARGETS IN REGRESSION MODEL

63425SKLEARNPREPROCESSING ROBUSTSCALE

SKLEARNPREPROCESSING ROBUSTSCALE XAXIS0 WITHCENTERINGTRUE WITHSCALINGTRUE

QUANTILERANGE250 750 COPYTRUE

STANDARDIZE A DATASET ALONG ANY AXIS

CENTER TO THE MEDIAN AND COMPONENT WISE SCALE ACCORDING TO THE INTERQUARTILE RANGE

READ MORE IN THE USER GUIDE

PARAMETERS

XARRAYLIKE THE DATA TO CENTER AND SCALE

634SKLEARNPREPROCESSING PREPROCESSING AND NORMALIZATION 2283

SCIKITLEARN USER GUIDE RELEASE 0213

AXIS INT 0 BY DEFAULT AXIS USED TO COMPUTE THE MEDIANS AND IQR ALONG IF 0 INDEPENDENTLY  
SCALE EACH FEATURE OTHERWISE IF 1 SCALE EACH SAMPLE  
WITHCENTERING BOOLEAN TRUE BY DEFAULT IF TRUE CENTER THE DATA BEFORE SCALING  
WITHSCALING BOOLEAN TRUE BY DEFAULT IF TRUE SCALE THE DATA TO UNIT VARIANCE OR EQUIVALENTLY  
UNIT STANDARD DEVIATION  
QUANTILERANGE TUPLE QMIN QMAX 00 QMIN QMAX 1000 DEFAULT 250 750  
1ST QUANTILE 3RD QUANTILE IQR QUANTILE RANGE USED TO CALCULATE SCALE  
NEW IN VERSION 018

COPY BOOLEAN OPTIONAL DEFAULT IS TRUE SET TO FALSE TO PERFORM INPLACE ROW NORMALIZATION AND  
AVOID A COPY IF THE INPUT IS ALREADY A NUMPY ARRAY OR A SCIPYSPARSE CSR MATRIX AND IF AXIS  
IS 1

SEE ALSO  
ROBUSTSCALER PERFORMS CENTERING AND SCALING USING THE TRANSFORMER API EG AS PART OF A PREPROCESS  
INGSKLEARNPIPELINEPIPELINE

NOTES  
THIS IMPLEMENTATION WILL REFUSE TO CENTER SCIPYSPARSE MATRICES SINCE IT WOULD MAKE THEM NONSPARSE AND WOULD  
POTENTIALLY CRASH THE PROGRAM WITH MEMORY EXHAUSTION PROBLEMS  
INSTEAD THE CALLER IS EXPECTED TO EITHER SET EXPLICITLY WITHCENTERINGFALSE IN THAT CASE ONLY VARIANCE  
SCALING WILL BE PERFORMED ON THE FEATURES OF THE CSR MATRIX OR TO CALL XTOARRAY IF HESHE EXPECTS THE  
MATERIALIZED DENSE ARRAY TO FIT IN MEMORY  
TO AVOID MEMORY COPY THE CALLER SHOULD PASS A CSR MATRIX  
FOR A COMPARISON OF THE DIFFERENT SCALERS TRANSFORMERS AND NORMALIZERS SEE EXAM  
PLESPREPROCESSINGPLOTALLSCALINGPY  
63426SKLEARNPREPROCESSING SCALE  
SKLEARNPREPROCESSING SCALEXAXIS0 WITHMEANTRUE WITHSTDTRUE COPYTRUE  
STANDARDIZE A DATASET ALONG ANY AXIS  
CENTER TO THE MEAN AND COMPONENT WISE SCALE TO UNIT VARIANCE  
READ MORE IN THE USER GUIDE

PARAMETERS  
XARRAYLIKE SPARSE MATRIX THE DATA TO CENTER AND SCALE  
AXIS INT 0 BY DEFAULT AXIS USED TO COMPUTE THE MEANS AND STANDARD DEVIATIONS ALONG IF 0  
INDEPENDENTLY STANDARDIZE EACH FEATURE OTHERWISE IF 1 STANDARDIZE EACH SAMPLE  
WITHMEAN BOOLEAN TRUE BY DEFAULT IF TRUE CENTER THE DATA BEFORE SCALING  
WITHSTD BOOLEAN TRUE BY DEFAULT IF TRUE SCALE THE DATA TO UNIT VARIANCE OR EQUIVALENTLY UNIT  
STANDARD DEVIATION  
COPY BOOLEAN OPTIONAL DEFAULT TRUE SET TO FALSE TO PERFORM INPLACE ROW NORMALIZATION AND  
AVOID A COPY IF THE INPUT IS ALREADY A NUMPY ARRAY OR A SCIPYSPARSE CSC MATRIX AND IF AXIS  
IS 1  
2284 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

SEE ALSO

STANDARDSCALER PERFORMS SCALING TO UNIT VARIANCE USING THE“TRANSFORMER” API EG AS PART OF A PREPRO  
CESSINGSKLEARNPIPELINEPIPELINE

NOTES

THIS IMPLEMENTATION WILL REFUSE TO CENTER SCIPYSPARSE MATRICES SINCE IT WOULD MAKE THEM NONSPARSE AND WOULD  
POTENTIALLY CRASH THE PROGRAM WITH MEMORY EXHAUSTION PROBLEMS

INSTEAD THE CALLER IS EXPECTED TO EITHER SET EXPLICITLY WITHMEANFALSE IN THAT CASE ONLY VARIANCE SCALING  
WILL BE PERFORMED ON THE FEATURES OF THE CSC MATRIX OR TO CALL XTOARRAY IF HESHE EXPECTS THE MATERIALIZED  
DENSE ARRAY TO FIT IN MEMORY

TO AVOID MEMORY COPY THE CALLER SHOULD PASS A CSC MATRIX

NANS ARE TREATED AS MISSING VALUES DISREGARDED TO COMPUTE THE STATISTICS AND MAINTAINED DURING THE DATA TRANS  
FORMATION

WE USE A BIASED ESTIMATOR FOR THE STANDARD DEVIATION EQUIVALENT TO NUMPYSTD<sub>X</sub> DDOF0 NOTE THAT THE  
CHOICE OFDDOF IS UNLIKELY TO AFFECT MODEL PERFORMANCE

FOR A COMPARISON OF THE DIFFERENT SCALERS TRANSFORMERS AND NORMALIZERS SEE EXAM  
PLESPREPROCESSINGPLOTALLSCALINGPY

EXAMPLES USING SKLEARNPREPROCESSINGSCALE

•A DEMO OF KMEANS CLUSTERING ON THE HANDWRITTEN DIGITS DATA

63427SKLEARNPREPROCESSING POWERTRANSFORM

SKLEARNPREPROCESSING POWERTRANSFORM XMETHOD‘WARN’ STANDARDIZETRUE COPYTRUE

POWER TRANSFORMS ARE A FAMILY OF PARAMETRIC MONOTONIC TRANSFORMATIONS THAT ARE APPLIED TO MAKE DATA MORE  
GAUSSIANLIKE THIS IS USEFUL FOR MODELING ISSUES RELATED TO HETEROSCEDASTICITY NONCONSTANT VARIANCE OR OTHER  
SITUATIONS WHERE NORMALITY IS DESIRED

CURRENTLY POWERTRANSFORM SUPPORTS THE BOXCOX TRANSFORM AND THE YEOJOHNSON TRANSFORM THE OPTIMAL PA  
RAMETER FOR STABILIZING VARIANCE AND MINIMIZING SKEWNESS IS ESTIMATED THROUGH MAXIMUM LIKELIHOOD

BOXCOX REQUIRES INPUT DATA TO BE STRICTLY POSITIVE WHILE YEOJOHNSON SUPPORTS BOTH POSITIVE OR NEGATIVE DATA  
BY DEFAULT ZEROMEAN UNITVARIANCE NORMALIZATION IS APPLIED TO THE TRANSFORMED DATA

READ MORE IN THE USER GUIDE

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES THE DATA TO BE TRANSFORMED USING A POWER TRANS  
FORMATION

METHOD STR THE POWER TRANSFORM METHOD AVAILABLE METHODS ARE

• ‘YEOJOHNSON’ 1 WORKS WITH POSITIVE AND NEGATIVE VALUES

• ‘BOXCOX’ 2 ONLY WORKS WITH STRICTLY POSITIVE VALUES

THE DEFAULT METHOD WILL BE CHANGED FROM ‘BOXCOX’ TO ‘YEOJOHNSON’ IN VERSION 023 TO

SUPPRESS THE FUTUREWARNING EXPLICITLY SET THE PARAMETER

6345SKLEARNPREPROCESSING PREPROCESSING AND NORMALIZATION 2285

SCIKITLEARN USER GUIDE RELEASE 0213

STANDARDIZE BOOLEAN DEFAULTTRUE SET TO TRUE TO APPLY ZERO MEAN UNIT VARIANCE NORMALIZATION TO THE TRANSFORMED OUTPUT

COPY BOOLEAN OPTIONAL DEFAULTTRUE SET TO FALSE TO PERFORM INPLACE COMPUTATION DURING TRANSFORMATION

RETURNS

XTRANS ARRAYLIKE SHAPE NSAMPLES NFEATURES THE TRANSFORMED DATA

SEE ALSO

POWERTRANSFORMER EQUIVALENT TRANSFORMATION WITH THE TRANSFORMER API EG AS PART OF A PREPROCESSING PIPELINE

QUANTILETRANSFORM MAPS DATA TO A STANDARD NORMAL DISTRIBUTION WITH THE PARAMETER OUTPUTDISTRIBUTIONNORMAL

NOTES

NANS ARE TREATED AS MISSING VALUES DISREGARDED IN FIT AND MAINTAINED IN TRANSFORM

FOR A COMPARISON OF THE DIFFERENT SCALERS TRANSFORMERS AND NORMALIZERS SEE EXAMPLESPREPROCESSINGPLOTALLSCALINGPY

REFERENCES

12

EXAMPLES

```
import numpy as np
from sklearn.preprocessing import PowerTransformer
data = [1, 2, 3, 2, 4, 5]
print(PowerTransformer(data).method)
# 0.707
# 0.707
# 1.414
```

635SKLEARNRANDOMPROJECTION RANDOM PROJECTION

RANDOM PROJECTION TRANSFORMERS

RANDOM PROJECTIONS ARE A SIMPLE AND COMPUTATIONALLY EFFICIENT WAY TO REDUCE THE DIMENSIONALITY OF THE DATA BY TRADING A CONTROLLED AMOUNT OF ACCURACY AS ADDITIONAL VARIANCE FOR FASTER PROCESSING TIMES AND SMALLER MODEL SIZES

THE DIMENSIONS AND DISTRIBUTION OF RANDOM PROJECTIONS MATRICES ARE CONTROLLED SO AS TO PRESERVE THE PAIRWISE DISTANCES BETWEEN ANY TWO SAMPLES OF THE DATASET

THE MAIN THEORETICAL RESULT BEHIND THE EFFICIENCY OF RANDOM PROJECTION IS THE JOHNSON-LINDENSTRAUSS LEMMA QUOTING WIKIPEDIA

2286 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

IN MATHEMATICS THE JOHNSONLINDENSTRAUSS LEMMA IS A RESULT CONCERNING LOWDISTORTION EMBEDDINGS OF POINTS FROM HIGHDIMENSIONAL INTO LOWDIMENSIONAL EUCLIDEAN SPACE THE LEMMA STATES THAT A SMALL SET OF POINTS IN A HIGHDIMENSIONAL SPACE CAN BE EMBEDDED INTO A SPACE OF MUCH LOWER DIMENSION IN SUCH A WAY THAT DISTANCES BETWEEN THE POINTS ARE NEARLY PRESERVED THE MAP USED FOR THE EMBEDDING IS AT LEAST LIPSCHITZ AND CAN EVEN BE TAKEN TO BE AN ORTHOGONAL PROJECTION

USER GUIDE SEE THE RANDOM PROJECTION SECTION FOR FURTHER DETAILS

RANDOMPROJECTION

GAUSSIANRANDOMPROJECTION REDUCE DIMENSIONALITY THROUGH GAUSSIAN RANDOM PROJECTION

RANDOMPROJECTION

SPARSERANDOMPROJECTION REDUCE DIMENSIONALITY THROUGH SPARSE RANDOM PROJECTION

6351SKLEARNRANDOMPROJECTION GAUSSIANRANDOMPROJECTION

CLASSSSKLEARNRANDOMPROJECTION GAUSSIANRANDOMPROJECTION NCOMPONENTS'AUTO'

EPS01 RAN

DOMSTATENONE

REDUCE DIMENSIONALITY THROUGH GAUSSIAN RANDOM PROJECTION

THE COMPONENTS OF THE RANDOM MATRIX ARE DRAWN FROM N0 1 NCOMPONENTS

READ MORE IN THE USER GUIDE

PARAMETERS

NCOMPONENTS INT OR 'AUTO' OPTIONAL DEFAULT 'AUTO' DIMENSIONALITY OF THE TARGET PROJECTION SPACE

NCOMPONENTS CAN BE AUTOMATICALLY ADJUSTED ACCORDING TO THE NUMBER OF SAMPLES IN THE DATASET AND THE BOUND GIVEN BY THE JOHNSONLINDENSTRAUSS LEMMA IN THAT CASE THE QUALITY OF THE EMBEDDING IS CONTROLLED BY THE EPS PARAMETER

IT SHOULD BE NOTED THAT JOHNSONLINDENSTRAUSS LEMMA CAN YIELD VERY CONSERVATIVE ESTIMATED OF THE REQUIRED NUMBER OF COMPONENTS AS IT MAKES NO ASSUMPTION ON THE STRUCTURE OF THE DATASET

EPS STRICTLY POSITIVE FLOAT OPTIONAL DEFAULT01 PARAMETER TO CONTROL THE QUALITY OF THE EMBEDDING ACCORDING TO THE JOHNSONLINDENSTRAUSS LEMMA WHEN NCOMPONENTS IS SET TO 'AUTO'

SMALLER VALUES LEAD TO BETTER EMBEDDING AND HIGHER NUMBER OF DIMENSIONS NCOMPONENTS IN THE TARGET PROJECTION SPACE

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE CONTROL THE PSEUDO RANDOM NUMBER GENERATOR USED TO GENERATE THE MATRIX AT FIT TIME IF INT RANDOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NPRANDOM

ATTRIBUTES

NCOMPONENT INT CONCRETE NUMBER OF COMPONENTS COMPUTED WHEN NCOMPONENTS" AUTO"

COMPONENTS NUMPY ARRAY OF SHAPE NCOMPONENTS NFEATURES RANDOM MATRIX USED FOR THE PROJECTION

SEE ALSO

SPARSERANDOMPROJECTION

6355SKLEARNRANDOMPROJECTION RANDOM PROJECTION 2287

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES

```
import numpy as np
from sklearn.random_projection import GaussianRandomProjection
rng = np.random.RandomState(42)
X = rng.rand(100, 10000)
transformer = GaussianRandomProjection(random_state=rng)
X_new = transformer.fit_transform(X)
X_new.shape
```

100 3947

METHODS

`fit` `X` `Y` GENERATE A SPARSE RANDOM PROJECTION MATRIX

`fit_transform` `self` `X` `Y` FIT TO DATA THEN TRANSFORM IT

`get_params` `self` DEEP GET PARAMETERS FOR THIS ESTIMATOR

`set_params` `self` `params` SET THE PARAMETERS OF THIS ESTIMATOR

`transform` `self` `X` PROJECT THE DATA BY USING MATRIX PRODUCT WITH THE RANDOM MATRIX

`init` `self` `n_components` `'auto'` `eps` `0.1` `random_state` `None`

`fit` `self` `X` `Y` `None`

GENERATE A SPARSE RANDOM PROJECTION MATRIX

PARAMETERS

`X` NUMPY ARRAY OR SCIPYSPARSE OF SHAPE `n_samples` `n_features` TRAINING SET ONLY THE SHAPE IS USED TO FIND OPTIMAL RANDOM MATRIX DIMENSIONS BASED ON THE THEORY REFERENCED IN THE AFORE MENTIONED PAPERS

`Y` IGNORED

RETURNS

`self`

`fit_transform` `self` `X` `Y` `None` `fit_params`

FIT TO DATA THEN TRANSFORM IT

FITS TRANSFORMER TO `X` AND `Y` WITH OPTIONAL PARAMETERS `fit_params` AND RETURNS A TRANSFORMED VERSION OF `X`

PARAMETERS

`X` NUMPY ARRAY OF SHAPE `n_samples` `n_features` TRAINING SET

`Y` NUMPY ARRAY OF SHAPE `n_samples` TARGET VALUES

RETURNS

`X_new` NUMPY ARRAY OF SHAPE `n_samples` `n_features_new` NEW TRANSFORMED ARRAY

`get_params` `self` `deep` `True`

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

`deep` BOOLEAN OPTIONAL IF `True` WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

2288 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

TRANSFORM SELF

PROJECT THE DATA BY USING MATRIX PRODUCT WITH THE RANDOM MATRIX

PARAMETERS

XNUMPY ARRAY OR SCIPYSPARSE OF SHAPE NSAMPLES NFEATURES THE INPUT DATA TO PROJECT INTO A SMALLER DIMENSIONAL SPACE

RETURNS

XNEW NUMPY ARRAY OR SCIPY SPARSE OF SHAPE NSAMPLES NCOMPONENTS PROJECTED ARRAY

6352SKLEARNRANDOMPROJECTION SPARSERANDOMPROJECTION

CLASSSSKLEARNRANDOMPROJECTION SPARSERANDOMPROJECTION NCOMPONENTS’AUTO’

DENSITY’AUTO’ EPS01

DENSEOUTPUTFALSE RAN

DOMSTATENONE

REDUCE DIMENSIONALITY THROUGH SPARSE RANDOM PROJECTION

SPARSE RANDOM MATRIX IS AN ALTERNATIVE TO DENSE RANDOM PROJECTION MATRIX THAT GUARANTEES SIMILAR EMBEDDING QUALITY WHILE BEING MUCH MORE MEMORY EFFICIENT AND ALLOWING FASTER COMPUTATION OF THE PROJECTED DATA

IF WE NOTES 1 DENSITY THE COMPONENTS OF THE RANDOM MATRIX ARE DRAWN FROM

- SQRTS SQRTNCOMPONENTS WITH PROBABILITY 1 2S
- 0 WITH PROBABILITY 1 1 S
- SQRTS SQRTNCOMPONENTS WITH PROBABILITY 1 2S

READ MORE IN THE USER GUIDE

PARAMETERS

NCOMPONENTS INT OR ‘AUTO’ OPTIONAL DEFAULT ‘AUTO’ DIMENSIONALITY OF THE TARGET PROJECTION SPACE

NCOMPONENTS CAN BE AUTOMATICALLY ADJUSTED ACCORDING TO THE NUMBER OF SAMPLES IN THE DATASET AND THE BOUND GIVEN BY THE JOHNSONLINDENSTRAUSS LEMMA IN THAT CASE THE QUALITY OF THE EMBEDDING IS CONTROLLED BY THE EPS PARAMETER

IT SHOULD BE NOTED THAT JOHNSONLINDENSTRAUSS LEMMA CAN YIELD VERY CONSERVATIVE ESTIMATED OF THE REQUIRED NUMBER OF COMPONENTS AS IT MAKES NO ASSUMPTION ON THE STRUCTURE OF THE DATASET

DENSITY FLOAT IN RANGE 0 1 OPTIONAL DEFAULT’AUTO’ RATIO OF NONZERO COMPONENT IN THE RANDOM PROJECTION MATRIX

6355SKLEARNRANDOMPROJECTION RANDOM PROJECTION 2289

SCIKITLEARN USER GUIDE RELEASE 0213

IF DENSITY ‘AUTO’ THE VALUE IS SET TO THE MINIMUM DENSITY AS RECOMMENDED BY PING LI ET AL 1 SQRTNFEATURES

USE DENSITY 1 30 IF YOU WANT TO REPRODUCE THE RESULTS FROM ACHLIOPTAS 2001

EPS STRICTLY POSITIVE FLOAT OPTIONAL DEFAULT01 PARAMETER TO CONTROL THE QUALITY OF THE EM BEDDING ACCORDING TO THE JOHNSONLINDENSTRAUSS LEMMA WHEN NCOMPONENTS IS SET TO ‘AUTO’ SMALLER VALUES LEAD TO BETTER EMBEDDING AND HIGHER NUMBER OF DIMENSIONS NCOMPONENTS IN THE TARGET PROJECTION SPACE

DENSEOUTPUT BOOLEAN OPTIONAL DEFAULTFALSE IF TRUE ENSURE THAT THE OUTPUT OF THE RANDOM PROJECTION IS A DENSE NUMPY ARRAY EVEN IF THE INPUT AND RANDOM PROJECTION MATRIX ARE BOTH SPARSE IN PRACTICE IF THE NUMBER OF COMPONENTS IS SMALL THE NUMBER OF ZERO COMPONENTS IN THE PROJECTED DATA WILL BE VERY SMALL AND IT WILL BE MORE CPU AND MEMORY EFFICIENT TO USE A DENSE REPRESENTATION

IF FALSE THE PROJECTED DATA USES A SPARSE REPRESENTATION IF THE INPUT IS SPARSE

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE CONTROL THE PSEUDO RANDOM NUMBER GENERATOR USED TO GENERATE THE MATRIX AT FIT TIME IF INT RANDOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NPRANDOM

ATTRIBUTES

NCOMPONENT INT CONCRETE NUMBER OF COMPONENTS COMPUTED WHEN NCOMPONENTS”AUTO”

COMPONENTS CSR MATRIX WITH SHAPE NCOMPONENTS NFEATURES RANDOM MATRIX USED FOR THE PROJECTION

DENSITY FLOAT IN RANGE 00 10 CONCRETE DENSITY COMPUTED FROM WHEN DENSITY “AUTO”

SEE ALSO

GAUSSIANRANDOMPROJECTION

REFERENCES

R0FECF191E4B81 R0FECF191E4B82

EXAMPLES

```
import numpy as np
from sklearn.randomprojection import sparserandomprojection
rng = np.random.RandomState(42)
X = rng.randn(100, 10000)
transformer = sparserandomprojection.RandomState(rng)
X_new = transformer.fit_transform(X)
X_new.shape
(100, 3947)
# VERY FEW COMPONENTS ARE NONZERO
np.mean(transformer.components_ > 0)
0.0100
```

2290 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

METHODS

FITSELF X Y GENERATE A SPARSE RANDOM PROJECTION MATRIX

FITTRANSFORM SELF X Y FIT TO DATA THEN TRANSFORM IT

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF X PROJECT THE DATA BY USING MATRIX PRODUCT WITH THE RANDOM MATRIX

INIT SELFNCOMPONENTS'AUTO' DENSITY'AUTO' EPS01 DENSEOUTPUTFALSE RANDOM

DOMSTATENONE

FITSELFXYNONE

GENERATE A SPARSE RANDOM PROJECTION MATRIX

PARAMETERS

XNUMPY ARRAY OR SCIPYSPARSE OF SHAPE NSAMPLES NFEATURES TRAINING SET ONLY THE SHAPE IS USED TO FIND OPTIMAL RANDOM MATRIX DIMENSIONS BASED ON THE THEORY REFERENCED IN THE AFORE MENTIONED PAPERS

YIGNORED

RETURNS

SELF

FITTRANSFORM SELFXYNONE FITPARAMS

FIT TO DATA THEN TRANSFORM IT

FITS TRANSFORMER TO X AND Y WITH OPTIONAL PARAMETERS FITPARAMS AND RETURNS A TRANSFORMED VERSION OF X

PARAMETERS

XNUMPY ARRAY OF SHAPE NSAMPLES NFEATURES TRAINING SET

YNUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES

RETURNS

XNEW NUMPY ARRAY OF SHAPE NSAMPLES NFEATURESNEW TRANSFORMED ARRAY

GETPARAMS SELFDEEPTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

635SKLEARNRANDOMPROJECTION RANDOM PROJECTION 2291

SCIKITLEARN USER GUIDE RELEASE 0213

TRANSFORM SELF

PROJECT THE DATA BY USING MATRIX PRODUCT WITH THE RANDOM MATRIX

PARAMETERS

XNUMPY ARRAY OR SCIPYSPARSE OF SHAPE NSAMPLES NFEATURES THE INPUT DATA TO PROJECT INTO A SMALLER DIMENSIONAL SPACE

RETURNS

XNEW NUMPY ARRAY OR SCIPY SPARSE OF SHAPE NSAMPLES NCOMPONENTS PROJECTED ARRAY

EXAMPLES USING SKLEARNRANDOMPROJECTIONS

- THE JOHNSONLINDENSTRAUSS BOUND FOR EMBEDDING WITH RANDOM PROJECTIONS
- MANIFOLD LEARNING ON HANDWRITTEN DIGITS LOCALLY LINEAR EMBEDDING ISOMAP

RANDOMPROJECTION

JOHNSONLINDENSTRAUSSMINDIM FIND A 'SAFE' NUMBER OF COMPONENTS TO RANDOMLY PROJECT TO

6353SKLEARNRANDOMPROJECTION JOHNSONLINDENSTRAUSSMINDIM

SKLEARNRANDOMPROJECTION JOHNSONLINDENSTRAUSSMINDIM NSAMPLES EPS01

FIND A 'SAFE' NUMBER OF COMPONENTS TO RANDOMLY PROJECT TO

THE DISTORTION INTRODUCED BY A RANDOM PROJECTION ONLY CHANGES THE DISTANCE BETWEEN TWO POINTS BY A FACTOR 1 EPS IN AN EUCLIDEAN SPACE WITH GOOD PROBABILITY THE PROJECTION IS AN EMBEDDING AS DEFINED BY

1 EPS U V2 PU V2 1 EPS U V2

WHERE U AND V ARE ANY ROWS TAKEN FROM A DATASET OF SHAPE NSAMPLES NFEATURES EPS IS IN 0 1 AND P IS A PROJECTION BY A RANDOM GAUSSIAN N0 1 MATRIX WITH SHAPE NCOMPONENTS NFEATURES OR A SPARSE ACHLIPTAS MATRIX

THE MINIMUM NUMBER OF COMPONENTS TO GUARANTEE THE EMBEDDING IS GIVEN BY

NCOMPONENTS 4 LOGNSAMPLES EPS2 2 EPS3 3

NOTE THAT THE NUMBER OF DIMENSIONS IS INDEPENDENT OF THE ORIGINAL NUMBER OF FEATURES BUT INSTEAD DEPENDS ON THE SIZE OF THE DATASET THE LARGER THE DATASET THE HIGHER IS THE MINIMAL DIMENSIONALITY OF AN EMBEDDING

READ MORE IN THE USER GUIDE

PARAMETERS

NSAMPLES INT OR NUMPY ARRAY OF INT GREATER THAN 0 NUMBER OF SAMPLES IF AN ARRAY IS GIVEN

IT WILL COMPUTE A SAFE NUMBER OF COMPONENTS ARRAYWISE

EPS FLOAT OR NUMPY ARRAY OF FLOAT IN 01 OPTIONAL DEFAULT01 MAXIMUM DISTORTION RATE AS DEFINED BY THE JOHNSONLINDENSTRAUSS LEMMA IF AN ARRAY IS GIVEN IT WILL COMPUTE A SAFE

NUMBER OF COMPONENTS ARRAYWISE

RETURNS

NCOMPONENTS INT OR NUMPY ARRAY OF INT THE MINIMAL NUMBER OF COMPONENTS TO GUARANTEE WITH GOOD PROBABILITY AN EMBEDDING WITH NSAMPLES

2292 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

REFERENCES

12

EXAMPLES

JOHNSONLINDENSTRAUSSMINDIM1E6 EPS05

663

JOHNSONLINDENSTRAUSSMINDIM1E6 EPS05 01 001

ARRAY 663 11841 1112658

JOHNSONLINDENSTRAUSSMINDIM1E4 1E5 1E6 EPS01

ARRAY 7894 9868 11841

EXAMPLES USING SKLEARNRANDOMPROJECTIONJOHNSONLINDENSTRAUSSMINDIM

•THE JOHNSONLINDENSTRAUSS BOUND FOR EMBEDDING WITH RANDOM PROJECTIONS

636SKLEARNSEMISUPERVISED SEMISUPERVISED LEARNING

THESKLEARNSEMISUPERVISED MODULE IMPLEMENTS SEMISUPERVISED LEARNING ALGORITHMS THESE ALGORITHMS

UTILIZED SMALL AMOUNTS OF LABELED DATA AND LARGE AMOUNTS OF UNLABELED DATA FOR CLASSIFICATION TASKS THIS MODULE

INCLUDES LABEL PROPAGATION

USER GUIDE SEE THE SEMISUPERVISED SECTION FOR FURTHER DETAILS

SEMISUPERVISEDLABELPROPAGATION KERNEL

    LABEL PROPAGATION CLASSIFIER

SEMISUPERVISEDLABELSPREADING KERNEL

    LABELSPREADING MODEL FOR SEMISUPERVISED LEARNING

6361SKLEARNSEMISUPERVISED LABELPROPAGATION

CLASSSSKLEARNSEMISUPERVISED LABELPROPAGATION KERNEL'RBF' GAMMA20 NNEIGHBORS7

MAXITER1000 TOL0001 NJOBSNONE

LABEL PROPAGATION CLASSIFIER

READ MORE IN THE USER GUIDE

PARAMETERS

KERNEL 'KNN' 'RBF' CALLABLE STRING IDENTIFIER FOR KERNEL FUNCTION TO USE OR THE KERNEL FUNCTION

ITSELF ONLY 'RBF' AND 'KNN' STRINGS ARE VALID INPUTS THE FUNCTION PASSED SHOULD TAKE TWO

INPUTS EACH OF SHAPE NSAMPLES NFEATURES AND RETURN A NSAMPLES NSAMPLES SHAPED

WEIGHT MATRIX

GAMMA FLOAT PARAMETER FOR RBF KERNEL

NNEIGHBORS INTEGER 0 PARAMETER FOR KNN KERNEL

636SKLEARNSEMISUPERVISED SEMISUPERVISED LEARNING 2293

SCIKITLEARN USER GUIDE RELEASE 0213

MAXITER INTEGER CHANGE MAXIMUM NUMBER OF ITERATIONS ALLOWED

TOLFLOAT CONVERGENCE TOLERANCE THRESHOLD TO CONSIDER THE SYSTEM AT STEADY STATE

NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF PARALLEL JOBS TO RUN NONE MEANS

1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL PROCESSORS SEE

GLOSSARY FOR MORE DETAILS

ATTRIBUTES

XARRAY SHAPE NSAMPLES NFEATURES INPUT ARRAY

CLASSES ARRAY SHAPE NCLASSES THE DISTINCT LABELS USED IN CLASSIFYING INSTANCES

LABELDISTRIBUTIONS ARRAY SHAPE NSAMPLES NCLASSES CATEGORICAL DISTRIBUTION FOR EACH

ITEM

TRANSDUCTION ARRAY SHAPE NSAMPLES LABEL ASSIGNED TO EACH ITEM VIA THE TRANSDUCTION

NITER INT NUMBER OF ITERATIONS RUN

SEE ALSO

LABELSPREADING ALTERNATE LABEL PROPAGATION STRATEGY MORE ROBUST TO NOISE

REFERENCES

XIAOJIN ZHU AND ZOUBIN GHAHRAMANI LEARNING FROM LABELED AND UNLABELED DATA WITH LABEL PROPAGATION TECH

NICAL REPORT CMUCALD02107 CARNEGIE MELLON UNIVERSITY 2002 HTTPPAGESCSWISCEDUJERRYZHUPUB

CMUCALD02107PDF

EXAMPLES

IMPORT NUMPY AS NP

FROM SKLEARN IMPORT DATASETS

FROM SKLEARNSEMISSUPERVISED IMPORT LABELPROPAGATION

LABELPROPMODEL LABELPROPAGATION

IRIS DATASETSLOADIRIS

RNG NPRANDOMRANDOMSTATE42

RANDOMUNLABELEDPOINTS RNGRANDLENIRISTARGET 03

LABELS NPCOPYIRISTARGET

LABELSRANDOMUNLABELEDPOINTS 1

LABELPROPMODELFITIRISDATA LABELS

LABELPROPAGATION

METHODS

FITSELF X Y

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PREDICT SELF X PERFORMS INDUCTIVE INFERENCE ACROSS THE MODEL

PREDICTPROBA SELF X PREDICT PROBABILITY FOR EACH POSSIBLE OUTCOME

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND

LABELS

CONTINUED ON NEXT PAGE

2294 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

TABLE 6270 – CONTINUED FROM PREVIOUS PAGE

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELFKERNEL'RBF' GAMMA20 NNEIGHBORS7 MAXITER1000 TOL0001 NJOBSNONE

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELF

PERFORMS INDUCTIVE INFERENCE ACROSS THE MODEL

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES

RETURNS

YARRAYLIKE SHAPE NSAMPLES PREDICTIONS FOR INPUT DATA

PREDICTPROBA SELF

PREDICT PROBABILITY FOR EACH POSSIBLE OUTCOME

COMPUTE THE PROBABILITY ESTIMATES FOR EACH SINGLE SAMPLE IN X AND EACH POSSIBLE OUTCOME SEEN DURING

TRAINING CATEGORICAL DISTRIBUTION

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES

RETURNS

PROBABILITIES ARRAY SHAPE NSAMPLES NCLASSES NORMALIZED PROBABILITY DISTRIBUTIONS

ACROSS CLASS LABELS

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS

IN MULTILABEL CLASSIFICATION THIS IS THE SUBSET ACCURACY WHICH IS A HARSH METRIC SINCE YOU REQUIRE FOR EACH

SAMPLE THAT EACH LABEL SET BE CORRECTLY PREDICTED

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE LABELS FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT MEAN ACCURACY OF SELF PREDICTX WRT Y

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE

PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT

OF A NESTED OBJECT

6365KLEARNSEMISSUPERVISED SEMISUPERVISED LEARNING 2295

SCIKITLEARN USER GUIDE RELEASE 0213

RETURNS

SELF

6362SKLEARNSEMISUPERVISED LABELSPREADING

CLASSSKLEARNSEMISUPERVISED LABELSPREADING KERNEL'RBF' GAMMA20 NNEIGHBORS7

ALPHA02 MAXITER30 TOL0001

NJOBSNONE

LABELSPREADING MODEL FOR SEMISUPERVISED LEARNING

THIS MODEL IS SIMILAR TO THE BASIC LABEL PROPAGATION ALGORITHM BUT USES AFFINITY MATRIX BASED ON THE NORMALIZED GRAPH LAPLACIAN AND SOFT CLAMPING ACROSS THE LABELS

READ MORE IN THE USER GUIDE

PARAMETERS

KERNEL 'KNN' 'RBF' CALLABLE STRING IDENTIFIER FOR KERNEL FUNCTION TO USE OR THE KERNEL FUNCTION ITSELF ONLY 'RBF' AND 'KNN' STRINGS ARE VALID INPUTS THE FUNCTION PASSED SHOULD TAKE TWO INPUTS EACH OF SHAPE NSAMPLES NFEATURES AND RETURN A NSAMPLES NSAMPLES SHAPED WEIGHT MATRIX

GAMMA FLOAT PARAMETER FOR RBF KERNEL

NNEIGHBORS INTEGER 0 PARAMETER FOR KNN KERNEL

ALPHA FLOAT CLAMPING FACTOR A VALUE IN 0 1 THAT SPECIFIES THE RELATIVE AMOUNT THAT AN INSTANCE SHOULD ADOPT THE INFORMATION FROM ITS NEIGHBORS AS OPPOSED TO ITS INITIAL LABEL ALPHA0 MEANS KEEPING THE INITIAL LABEL INFORMATION ALPHA1 MEANS REPLACING ALL INITIAL INFORMATION

MAXITER INTEGER MAXIMUM NUMBER OF ITERATIONS ALLOWED

TOLFLOAT CONVERGENCE TOLERANCE THRESHOLD TO CONSIDER THE SYSTEM AT STEADY STATE

NJOBS INT OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF PARALLEL JOBS TO RUN NONE MEANS 1 UNLESS IN A JOBLIBPARALLELBACKEND CONTEXT1MEANS USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

ATTRIBUTES

XARRAY SHAPE NSAMPLES NFEATURES INPUT ARRAY

CLASSES ARRAY SHAPE NCLASSES THE DISTINCT LABELS USED IN CLASSIFYING INSTANCES

LABELDISTRIBUTIONS ARRAY SHAPE NSAMPLES NCLASSES CATEGORICAL DISTRIBUTION FOR EACH ITEM

TRANSDUCTION ARRAY SHAPE NSAMPLES LABEL ASSIGNED TO EACH ITEM VIA THE TRANSDUCTION

NITER INT NUMBER OF ITERATIONS RUN

SEE ALSO

LABELPROPAGATION UNREGULARIZED GRAPH BASED SEMISUPERVISED LEARNING

REFERENCES

DENGYONG ZHOU OLIVIER BOUSQUET THOMAS NAVIN LAL JASON WESTON BERNHARD SCHOELKOPF LEARNING WITH LOCAL AND GLOBAL CONSISTENCY 2004 HTTPCITESEERISTPSUEDUVIEWDOCSUMMARYDOI10111153219 2296 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

EXAMPLES

```
import numpy as np
from sklearn import datasets
from sklearn.semi_supervised import LabelSpreading
iris = datasets.load_iris()
rng = np.random.RandomState(42)
random_unlabeled_points = rng.rand(len(iris.target)) < 0.3
labels = np.copy(iris.target)
label_prop_model = LabelSpreading()
label_prop_model.fit(iris.data, labels)
```

LabelSpreading

Methods

```
fit(self, X, Y) Fit a semi-supervised label propagation model based
get_params(self, deep=True) Get parameters for this estimator
predict(self, X) Performs inductive inference across the model
predict_proba(self, X) Predict probability for each possible outcome
score(self, X, Y, sample_weight) Returns the mean accuracy on the given test data and
labels
set_params(self, **kwargs) Set the parameters of this estimator
init(self, kernel='rbf', gamma=20, n_neighbors=7, alpha=0.2, max_iter=30, tol=0.001,
n_jobs=None)
fit(self, X, Y) Fit a semi-supervised label propagation model based
all the input data is provided matrix X labeled and unlabeled and corresponding label matrix Y with a
dedicated marker value for unlabeled samples
parameters
X: array-like shape (n_samples, n_features) A (n_samples by n_samples) size matrix
will be created from this
Y: array-like shape (n_samples, n_labeled_samples) Unlabeled points are marked as 1
all unlabeled samples will be transductively assigned labels
returns
self Returns an instance of self
get_params(self, deep=True) Get parameters for this estimator
parameters
deep: boolean optional if True will return the parameters for this estimator and contained
subobjects that are estimators
returns
params Mapping of string to any parameter names mapped to their values
636SKLEARNSEMI-SUPERVISED SEMI-SUPERVISED LEARNING 2297
```

SCIKITLEARN USER GUIDE RELEASE 0213

PREDICTSELF  
PERFORMS INDUCTIVE INFERENCE ACROSS THE MODEL

PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES

RETURNS  
YARRAYLIKE SHAPE NSAMPLES PREDICTIONS FOR INPUT DATA

PREDICTPROBA SELF  
PREDICT PROBABILITY FOR EACH POSSIBLE OUTCOME  
COMPUTE THE PROBABILITY ESTIMATES FOR EACH SINGLE SAMPLE IN X AND EACH POSSIBLE OUTCOME SEEN DURING TRAINING CATEGORICAL DISTRIBUTION

PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES

RETURNS  
PROBABILITIES ARRAY SHAPE NSAMPLES NCLASSES NORMALIZED PROBABILITY DISTRIBUTIONS  
ACROSS CLASS LABELS

SCORESELFXYSAMPLEWEIGHTNONE  
RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS  
IN MULTILABEL CLASSIFICATION THIS IS THE SUBSET ACCURACY WHICH IS A HARSH METRIC SINCE YOU REQUIRE FOR EACH SAMPLE THAT EACH LABEL SET BE CORRECTLY PREDICTED

PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES  
YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE LABELS FOR X  
SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS  
SCORE FLOAT MEAN ACCURACY OF SELF PREDICTX WRT Y

SETPARAMS SELF  
SET THE PARAMETERS OF THIS ESTIMATOR  
THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS  
SELF

EXAMPLES USING SKLEARNSEMISSUPERVISED LABELSPREADING

- DECISION BOUNDARY OF LABEL PROPAGATION VERSUS SVM ON THE IRIS DATASET
- LABEL PROPAGATION LEARNING A COMPLEX STRUCTURE
- LABEL PROPAGATION DIGITS DEMONSTRATING PERFORMANCE
- LABEL PROPAGATION DIGITS ACTIVE LEARNING

2298 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

6375SKLEARN SVM SUPPORT VECTOR MACHINES

THE SKLEARN SVM MODULE INCLUDES SUPPORT VECTOR MACHINE ALGORITHMS

USER GUIDE SEE THE SUPPORT VECTOR MACHINES SECTION FOR FURTHER DETAILS

6371 ESTIMATORS

SVMLINEARSVC PENALTY LOSS DUAL TOL C LINEAR SUPPORT VECTOR CLASSIFICATION

SVMLINEARSVR EPSILON TOL C LOSS LINEAR SUPPORT VECTOR REGRESSION

SVMNUSVC NU KERNEL DEGREE GAMMA NUSUPPORT VECTOR CLASSIFICATION

SVMNUSVR NU C KERNEL DEGREE GAMMA NU SUPPORT VECTOR REGRESSION

SVMONECLASS SVM KERNEL DEGREE GAMMA UNSUPERVISED OUTLIER DETECTION

SVMSVC C KERNEL DEGREE GAMMA COEF0 CSUPPORT VECTOR CLASSIFICATION

SVMSVR KERNEL DEGREE GAMMA COEF0 TOL EPSILONSUPPORT VECTOR REGRESSION

SKLEARN SVM LINEARSVC

CLASS SKLEARN SVM LINEARSVC PENALTY 'L2' LOSS 'SQUARED HINGE' DUAL TRUE TOL 0.0001

C10 MULTICLASS 'OVR' FIT INTERCEPT TRUE INTERCEPT SCALING 1

CLASS WEIGHT NONE VERBOSE 0 RANDOM STATE NONE MAX ITER 1000

LINEAR SUPPORT VECTOR CLASSIFICATION

SIMILAR TO SVC WITH PARAMETER KERNEL 'LINEAR' BUT IMPLEMENTED IN TERMS OF LIBLINEAR RATHER THAN LIBSVM SO IT HAS MORE FLEXIBILITY IN THE CHOICE OF PENALTIES AND LOSS FUNCTIONS AND SHOULD SCALE BETTER TO LARGE NUMBERS OF SAMPLES

THIS CLASS SUPPORTS BOTH DENSE AND SPARSE INPUT AND THE MULTICLASS SUPPORT IS HANDLED ACCORDING TO A ONE VS THE REST SCHEME

READ MORE IN THE USER GUIDE

PARAMETERS

PENALTY STRING 'L1' OR 'L2' DEFAULT 'L2' SPECIFIES THE NORM USED IN THE PENALIZATION THE 'L2' PENALTY IS THE STANDARD USED IN SVC THE 'L1' LEADS TO COEF VECTORS THAT ARE SPARSE

LOSS STRING 'HINGE' OR 'SQUARED HINGE' DEFAULT 'SQUARED HINGE' SPECIFIES THE LOSS FUNCTION 'HINGE' IS THE STANDARD SVM LOSS USED EG BY THE SVC CLASS WHILE 'SQUARED HINGE' IS THE SQUARE OF THE HINGE LOSS

DUAL BOOL DEFAULT TRUE SELECT THE ALGORITHM TO EITHER SOLVE THE DUAL OR PRIMAL OPTIMIZATION PROBLEM PREFER DUAL FALSE WHEN NSAMPLES NFEATURES

TOL FLOAT OPTIONAL DEFAULT 1E-4 TOLERANCE FOR STOPPING CRITERIA

CFLOAT OPTIONAL DEFAULT 10 PENALTY PARAMETER C OF THE ERROR TERM

MULTICLASS STRING 'OVR' OR 'CRAMMERSINGER' DEFAULT 'OVR' DETERMINES THE MULTICLASS STRATEGY IF YCONTAINS MORE THAN TWO CLASSES OVR TRAINS NCLASSES ONE VS REST CLASSIFIERS WHILE CRAMMERSINGER OPTIMIZES A JOINT OBJECTIVE OVER ALL CLASSES WHILE CRAMMERSINGER IS INTERESTING FROM A THEORETICAL PERSPECTIVE AS IT IS CONSISTENT IT IS SELDOM USED IN PRACTICE AS IT RARELY LEADS TO BETTER ACCURACY AND IS MORE EXPENSIVE TO COMPUTE IF CRAMMERSINGER IS CHOSEN THE OPTIONS LOSS PENALTY AND DUAL WILL BE IGNORED

FIT INTERCEPT BOOLEAN OPTIONAL DEFAULT TRUE WHETHER TO CALCULATE THE INTERCEPT FOR THIS MODEL IF SET TO FALSE NO INTERCEPT WILL BE USED IN CALCULATIONS IE DATA IS EXPECTED TO

6375SKLEARN SVM SUPPORT VECTOR MACHINES 2299

BE ALREADY CENTERED

INTERCEPTSCALING FLOAT OPTIONAL DEFAULT1 WHEN SELFFITINTERCEPT IS TRUE INSTANCE VECTOR X BECOMESX SELFINTERCEPTSCALING IE A “SYNTHETIC” FEATURE WITH CONSTANT VALUE EQUALS TO INTERCEPTSCALING IS APPENDED TO THE INSTANCE VECTOR THE INTERCEPT BECOMES INTERCEPTSCALING SYNTHETIC FEATURE WEIGHT NOTE THE SYNTHETIC FEATURE WEIGHT IS SUBJECT TO L1L2 REGULARIZATION AS ALL OTHER FEATURES TO LESSEN THE EFFECT OF REGULARIZATION ON SYNTHETIC FEATURE WEIGHT AND THEREFORE ON THE INTERCEPT INTERCEPTSCALING HAS TO BE INCREASED CLASSWEIGHT DICT ‘BALANCED’ OPTIONAL SET THE PARAMETER C OF CLASS I TO CLASSWEIGHTI CFOR SVC IF NOT GIVEN ALL CLASSES ARE SUPPOSED TO HAVE WEIGHT ONE THE “BALANCED” MODE USES THE VALUES OF Y TO AUTOMATICALLY ADJUST WEIGHTS INVERSELY PROPORTIONAL TO CLASS FREQUENCIES IN THE INPUT DATA AS NSAMPLES NCLASSES NPBINCOUNTY

VERBOSE INT DEFAULT0 ENABLE VERBOSE OUTPUT NOTE THAT THIS SETTING TAKES ADVANTAGE OF A PER PROCESS RUNTIME SETTING IN LIBLINEAR THAT IF ENABLED MAY NOT WORK PROPERLY IN A MULTITHREADED CONTEXT

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE THE SEED OF THE PSEUDO RANDOM NUMBER GENERATOR TO USE WHEN SHUFFLING THE DATA FOR THE DUAL COORDINATE DESCENT IF DUALTRUE WHEN DUALFALSE THE UNDERLYING IMPLEMENTATION OF LINEARSVC IS NOT RANDOM AND RANDOMSTATE HAS NO EFFECT ON THE RESULTS IF INT RANDOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NPRANDOM

MAXITER INT DEFAULT1000 THE MAXIMUM NUMBER OF ITERATIONS TO BE RUN

ATTRIBUTES

COEF ARRAY SHAPE NFEATURES IF NCLASSES 2 ELSE NCLASSES NFEATURES WEIGHTS AS SIGNED TO THE FEATURES COEFFICIENTS IN THE PRIMAL PROBLEM THIS IS ONLY AVAILABLE IN THE CASE OF A LINEAR KERNEL

COEF IS A READONLY PROPERTY DERIVED FROM RAWCOEF THAT FOLLOWS THE INTERNAL MEMORY LAYOUT OF LIBLINEAR

INTERCEPT ARRAY SHAPE 1 IF NCLASSES 2 ELSE NCLASSES CONSTANTS IN DECISION FUNCTION

SEE ALSO

SVC IMPLEMENTATION OF SUPPORT VECTOR MACHINE CLASSIFIER USING LIBSVM THE KERNEL CAN BE NONLINEAR BUT ITS SMO ALGORITHM DOES NOT SCALE TO LARGE NUMBER OF SAMPLES AS LINEARSVC DOES FURTHERMORE SVC MULTICLASS MODE IS IMPLEMENTED USING ONE VS ONE SCHEME WHILE LINEARSVC USES ONE VS THE REST IT IS POSSIBLE TO IMPLEMENT ONE VS THE REST WITH SVC BY USING THE SKLEARNMULTICLASSONEVSRESTCLASSIFIER WRAPPER FINALLY SVC CAN FIT DENSE DATA WITHOUT MEMORY COPY IF THE INPUT IS CCONTIGUOUS SPARSE DATA WILL STILL INCUR MEMORY COPY THOUGH

SKLEARNLINEARMODELSGDCLASSIFIER SGDCLASSIFIER CAN OPTIMIZE THE SAME COST FUNCTION AS LINEARSVC BY ADJUSTING THE PENALTY AND LOSS PARAMETERS IN ADDITION IT REQUIRES LESS MEMORY ALLOWS INCREMENTAL ONLINE LEARNING AND IMPLEMENTS VARIOUS LOSS FUNCTIONS AND REGULARIZATION REGIMES

NOTES

THE UNDERLYING C IMPLEMENTATION USES A RANDOM NUMBER GENERATOR TO SELECT FEATURES WHEN FITTING THE MODEL IT IS THUS NOT UNCOMMON TO HAVE SLIGHTLY DIFFERENT RESULTS FOR THE SAME INPUT DATA IF THAT HAPPENS TRY WITH A SMALLERTOL PARAMETER

SCIKITLEARN USER GUIDE RELEASE 0213

THE UNDERLYING IMPLEMENTATION LIBLINEAR USES A SPARSE INTERNAL REPRESENTATION FOR THE DATA THAT WILL INCUR A MEMORY COPY

PREDICT OUTPUT MAY NOT MATCH THAT OF STANDALONE LIBLINEAR IN CERTAIN CASES SEE DIFFERENCES FROM LIBLINEAR IN THE NARRATIVE DOCUMENTATION

REFERENCES

LIBLINEAR A LIBRARY FOR LARGE LINEAR CLASSIFICATION

EXAMPLES

```
FROM SKLEARN SVM IMPORT LINEAR SVC
FROM SKLEARN DATASETS IMPORT MAKE CLASSIFICATION
X Y MAKE CLASSIFICATION NFEATURES 4 RANDOM STATE 0
CLF LINEAR SVC RANDOM STATE 0 TOL 1E5
CLF FIT X Y
LINEAR SVC C10 CLASSWEIGHT NONE DUAL TRUE FIT INTERCEPT TRUE
INTERCEPT SCALING 1 LOSS SQUARED HINGE MAX ITER 1000
MULTICLASS OVR PENALTY L2 RANDOM STATE 0 TOL 1E05 VERBOSE 0
PRINT CLF COEF
0085 0394 0498 0375
PRINT CLF INTERCEPT
0284
PRINT CLF PREDICT 0 0 0 0
1
```

METHODS

DECISION FUNCTION SELF X PREDICT CONFIDENCE SCORES FOR SAMPLES

DENSIFY SELF CONVERT COEFFICIENT MATRIX TO DENSE ARRAY FORMAT

FIT SELF X Y SAMPLEWEIGHT FIT THE MODEL ACCORDING TO THE GIVEN TRAINING DATA

GET PARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PREDICT SELF X PREDICT CLASS LABELS FOR SAMPLES IN X

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS

SET PARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

SPARSIFY SELF CONVERT COEFFICIENT MATRIX TO SPARSE FORMAT

```
INIT SELF PENALTY 'L2' LOSS 'SQUARED HINGE' DUAL TRUE TOL 0.0001 C10 MULTICLASS 'OVR'
FIT INTERCEPT TRUE INTERCEPT SCALING 1 CLASSWEIGHT NONE VERBOSE 0 RAN
DOM STATE NONE MAX ITER 1000
DECISION FUNCTION SELF X
PREDICT CONFIDENCE SCORES FOR SAMPLES
THE CONFIDENCE SCORE FOR A SAMPLE IS THE SIGNED DISTANCE OF THAT SAMPLE TO THE HYPERPLANE
PARAMETERS
X ARRAY LIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES SAMPLES
RETURNS
637 SKLEARN SVM SUPPORT VECTOR MACHINES 2301
```

SCIKITLEARN USER GUIDE RELEASE 0213

ARRAY SHAPENSAMPLES IF NCLASSES > 2 ELSE NSAMPLES NCLASSES CONFIDENCE  
SCORES PER SAMPLE CLASS COMBINATION IN THE BINARY CASE CONFIDENCE SCORE FOR  
SELFCLASSES1 WHERE 0 MEANS THIS CLASS WOULD BE PREDICTED

DENSIFYSELF  
CONVERT COEFFICIENT MATRIX TO DENSE ARRAY FORMAT

CONVERTS THE COEF MEMBER BACK TO A NUMPYNDARRAY THIS IS THE DEFAULT FORMAT OF COEF AND IS  
REQUIRED FOR FITTING SO CALLING THIS METHOD IS ONLY REQUIRED ON MODELS THAT HAVE PREVIOUSLY BEEN SPARSIFIED  
OTHERWISE IT IS A NOOP

RETURNS  
SELF ESTIMATOR

FITSELFXYSAMPLEWEIGHTNONE  
FIT THE MODEL ACCORDING TO THE GIVEN TRAINING DATA

PARAMETERS  
XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES TRAINING VECTOR WHERE  
NSAMPLES IN THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES  
YARRAYLIKE SHAPE NSAMPLES TARGET VECTOR RELATIVE TO X  
SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL ARRAY OF WEIGHTS THAT ARE AS  
SIGNED TO INDIVIDUAL SAMPLES IF NOT PROVIDED THEN EACH SAMPLE IS GIVEN UNIT WEIGHT

RETURNS  
SELF OBJECT

GETPARAMS SELFDEEPTURE  
GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS  
DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBOBJECTS THAT ARE ESTIMATORS

RETURNS  
PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELF  
PREDICT CLASS LABELS FOR SAMPLES IN X

PARAMETERS  
XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES SAMPLES

RETURNS  
CARRAY SHAPE NSAMPLES PREDICTED CLASS LABEL PER SAMPLE

SCORESELFXYSAMPLEWEIGHTNONE  
RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS  
IN MULTILABEL CLASSIFICATION THIS IS THE SUBSET ACCURACY WHICH IS A HARSH METRIC SINCE YOU REQUIRE FOR EACH  
SAMPLE THAT EACH LABEL SET BE CORRECTLY PREDICTED

PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES  
YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE LABELS FOR X

2302 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT MEAN ACCURACY OF SELF PREDICTX WRT Y

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

SPARSIFY SELF

CONVERT COEFFICIENT MATRIX TO SPARSE FORMAT

CONVERTS THE COEF MEMBER TO A SCIPYSPARSE MATRIX WHICH FOR L1REGULARIZED MODELS CAN BE MUCH MORE MEMORY AND STORAGEEFFICIENT THAN THE USUAL NUMPYNDARRAY REPRESENTATION

THEINTERCEPT MEMBER IS NOT CONVERTED

RETURNS

SELF ESTIMATOR

NOTES

FOR NONSPARSE MODELS IE WHEN THERE ARE NOT MANY ZEROS IN COEF THIS MAY ACTUALLY INCREASE MEMORY USAGE SO USE THIS METHOD WITH CARE A RULE OF THUMB IS THAT THE NUMBER OF ZERO ELEMENTS WHICH CAN BE COMPUTED WITH COEF 0SUM MUST BE MORE THAN 50 FOR THIS TO PROVIDE SIGNIFICANT BENEFITS

AFTER CALLING THIS METHOD FURTHER FITTING WITH THE PARTIALFIT METHOD IF ANY WILL NOT WORK UNTIL YOU CALL DENSIFY

EXAMPLES USING SKLEARN SVMLINEAR SVC

- EXPLICIT FEATURE MAP APPROXIMATION FOR RBF KERNELS
- COMPARISON OF CALIBRATION OF CLASSIFIERS
- PROBABILITY CALIBRATION CURVES
- SELECTING DIMENSIONALITY REDUCTION WITH PIPELINE AND GRIDSEARCH CV
- COLUMN TRANSFORMER WITH HETEROGENEOUS DATA SOURCES
- PIPELINE ANOVA SVM
- BALANCE MODEL COMPLEXITY AND CROSSVALIDATED SCORE
- PRECISIONRECALL
- FEATURE DISCRETIZATION
- PLOT DIFFERENT SVM CLASSIFIERS IN THE IRIS DATASET
- SCALING THE REGULARIZATION PARAMETER FOR SVCS
- CLASSIFICATION OF TEXT DOCUMENTS USING SPARSE FEATURES

637SKLEARN SVM SUPPORT VECTOR MACHINES 2303

SCIKITLEARN USER GUIDE RELEASE 0213

SKLEARN SVM LINEAR SVR

CLASS SKLEARN SVM LINEAR SVR EPSILON 0.001 C 1.0 LOSS 'EPSILON INSENSITIVE'

FIT INTERCEPT TRUE INTERCEPT SCALING 1.0 DUAL TRUE VERBOSE 0

RANDOM STATE NONE MAX ITER 1000

LINEAR SUPPORT VECTOR REGRESSION

SIMILAR TO SVR WITH PARAMETER KERNEL 'LINEAR' BUT IMPLEMENTED IN TERMS OF LIBLINEAR RATHER THAN LIBSVM SO IT HAS MORE FLEXIBILITY IN THE CHOICE OF PENALTIES AND LOSS FUNCTIONS AND SHOULD SCALE BETTER TO LARGE NUMBERS OF SAMPLES

THIS CLASS SUPPORTS BOTH DENSE AND SPARSE INPUT

READ MORE IN THE USER GUIDE

PARAMETERS

EPSILON FLOAT OPTIONAL DEFAULT 0.0 EPSILON PARAMETER IN THE EPSILON INSENSITIVE LOSS FUNCTION

NOTE THAT THE VALUE OF THIS PARAMETER DEPENDS ON THE SCALE OF THE TARGET VARIABLE Y IF UNSURE

SET EPSILON 0

TOL FLOAT OPTIONAL DEFAULT 1e-4 TOLERANCE FOR STOPPING CRITERIA

C FLOAT OPTIONAL DEFAULT 1.0 PENALTY PARAMETER C OF THE ERROR TERM THE PENALTY IS A SQUARED

L2 PENALTY THE BIGGER THIS PARAMETER THE LESS REGULARIZATION IS USED

LOSS STRING OPTIONAL DEFAULT 'EPSILON INSENSITIVE' SPECIFIES THE LOSS FUNCTION THE EPSILON

INSENSITIVE LOSS STANDARD SVR IS THE L1 LOSS WHILE THE SQUARED EPSILON INSENSITIVE LOSS

'SQUARED EPSILON INSENSITIVE' IS THE L2 LOSS

FIT INTERCEPT BOOLEAN OPTIONAL DEFAULT TRUE WHETHER TO CALCULATE THE INTERCEPT FOR THIS

MODEL IF SET TO FALSE NO INTERCEPT WILL BE USED IN CALCULATIONS IE DATA IS EXPECTED TO

BE ALREADY CENTERED

INTERCEPT SCALING FLOAT OPTIONAL DEFAULT 1 WHEN SELF FIT INTERCEPT IS TRUE INSTANCE VEC

TOR X BECOMES X SELF INTERCEPT SCALING IE A "SYNTHETIC" FEATURE WITH CONSTANT VALUE

EQUALS TO INTERCEPT SCALING IS APPENDED TO THE INSTANCE VECTOR THE INTERCEPT BECOMES IN

TERCEPT SCALING SYNTHETIC FEATURE WEIGHT NOTE THE SYNTHETIC FEATURE WEIGHT IS SUBJECT TO

L1 L2 REGULARIZATION AS ALL OTHER FEATURES TO LESSEN THE EFFECT OF REGULARIZATION ON SYNTHETIC

FEATURE WEIGHT AND THEREFORE ON THE INTERCEPT INTERCEPT SCALING HAS TO BE INCREASED

DUAL BOOL DEFAULT TRUE SELECT THE ALGORITHM TO EITHER SOLVE THE DUAL OR PRIMAL OPTIMIZATION

PROBLEM PREFER DUAL FALSE WHEN N SAMPLES N FEATURES

VERBOSE INT DEFAULT 0 ENABLE VERBOSE OUTPUT NOTE THAT THIS SETTING TAKES ADVANTAGE OF A PER

PROCESS RUNTIME SETTING IN LIBLINEAR THAT IF ENABLED MAY NOT WORK PROPERLY IN A MULTITHREADED

CONTEXT

RANDOM STATE INT RANDOM STATE INSTANCE OR NONE OPTIONAL DEFAULT NONE THE SEED OF THE

PSEUDO RANDOM NUMBER GENERATOR TO USE WHEN SHUFFLING THE DATA IF INT RANDOM STATE IS

THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOM STATE INSTANCE RANDOM STATE IS

THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOM STATE

INSTANCE USED BY N RANDOM

MAX ITER INT DEFAULT 1000 THE MAXIMUM NUMBER OF ITERATIONS TO BE RUN

ATTRIBUTES

COEF ARRAY SHAPE N FEATURES IF N CLASSES > 2 ELSE N CLASSES N FEATURES WEIGHTS AS

SIGNED TO THE FEATURES COEFFICIENTS IN THE PRIMAL PROBLEM THIS IS ONLY AVAILABLE IN THE CASE

OF A LINEAR KERNEL

2304 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

COEF IS A READONLY PROPERTY DERIVED FROM RAWCOEF THAT FOLLOWS THE INTERNAL MEMORY LAYOUT OF LIBLINEAR

INTERCEPT ARRAY SHAPE 1 IF NCLASSES 2 ELSE NCLASSES CONSTANTS IN DECISION FUNCTION

SEE ALSO

LINEARSVC IMPLEMENTATION OF SUPPORT VECTOR MACHINE CLASSIFIER USING THE SAME LIBRARY AS THIS CLASS LIBLINEAR

SVR IMPLEMENTATION OF SUPPORT VECTOR MACHINE REGRESSION USING LIBSVM THE KERNEL CAN BE NONLINEAR BUT ITS SMO ALGORITHM DOES NOT SCALE TO LARGE NUMBER OF SAMPLES AS LINEARSVC DOES

SKLEARNLINEARMODELSGDREGRESSOR SGDREGRESSOR CAN OPTIMIZE THE SAME COST FUNCTION AS LINEARSVR BY ADJUSTING THE PENALTY AND LOSS PARAMETERS IN ADDITION IT REQUIRES LESS MEMORY ALLOWS INCREMENTAL ONLINE LEARNING AND IMPLEMENTS VARIOUS LOSS FUNCTIONS AND REGULARIZATION REGIMES

EXAMPLES

```
FROM SKLEARN SVM IMPORT LINEARSVR
FROM SKLEARN DATASETS IMPORT MAKEREGRESSION
X Y MAKEREGRESSION NFEATURES 4 RANDOMSTATE 0
REGR LINEARSVR RANDOMSTATE 0 TOL 1E5
REGR FIT X Y
LINEARSVR C 10 DUAL TRUE EPSILON 00 FIT INTERCEPT TRUE
INTERCEPT SCALING 10 LOSS EPSILON INSENSITIVE MAX ITER 1000
RANDOMSTATE 0 TOL 1E05 VERBOSE 0
PRINT REGR COEF
1635 2691 4230 6047
PRINT REGR INTERCEPT
429
PRINT REGR PREDICT 0 0 0 0
429
```

METHODS

FIT SELF X Y SAMPLEWEIGHT FIT THE MODEL ACCORDING TO THE GIVEN TRAINING DATA

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PREDICT SELF X PREDICT USING THE LINEAR MODEL

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELF EPSILON 00 TOL 00001 C 10 LOSS 'EPSILON INSENSITIVE' FIT INTERCEPT TRUE INTERCEPT SCALING 10 DUAL TRUE VERBOSE 0 RANDOMSTATE NONE MAX ITER 1000

FIT SELF X Y SAMPLEWEIGHT NONE

FIT THE MODEL ACCORDING TO THE GIVEN TRAINING DATA

PARAMETERS

X ARRAY LIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES TRAINING VECTOR WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

Y ARRAY LIKE SHAPE NSAMPLES TARGET VECTOR RELATIVE TO X

637 SKLEARN SVM SUPPORT VECTOR MACHINES 2305

SCIKITLEARN USER GUIDE RELEASE 0213

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL ARRAY OF WEIGHTS THAT ARE AS SIGNED TO INDIVIDUAL SAMPLES IF NOT PROVIDED THEN EACH SAMPLE IS GIVEN UNIT WEIGHT

RETURNS

SELF OBJECT

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELF

PREDICT USING THE LINEAR MODEL

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES SAMPLES

RETURNS

CARRAY SHAPE NSAMPLES RETURNS PREDICTED VALUES

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{U}{V}$  WHERE U IS THE RESIDUAL SUM OF SQUARES YTRUE YPRED 2SUM AND V IS THE TOTAL SUM OF SQUARES YTRUE YTRUEMEAN 2SUM THE BEST POSSIBLE SCORE IS 10 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 00

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES FOR SOME ESTIMATORS THIS MAY BE A PRECOMPUTED KERNEL MATRIX INSTEAD SHAPE NSAMPLES NSAMPLESFITTED WHERE NSAMPLESFITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE VALUES FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT R2 OF SELF PREDICTX WRT Y

NOTES

THE R2 SCORE USED WHEN CALLING SCORE ON A REGRESSOR WILL USE MULTIOUTPUTUNIFORMAVERAGE FROM VERSION 023 TO KEEP CONSISTENT WITH METRICSR2SCORE THIS WILL INFLUENCE THE SCORE METHOD OF ALL THE MULTIOUTPUT REGRESSORS EXCEPT FOR MULTIOUTPUTMULTIOUTPUTREGRESSOR TO SPECIFY THE DEFAULT VALUE MANUALLY AND AVOID THE WARNING PLEASE EITHER CALL METRICSR2SCORE DIRECTLY OR MAKE A CUSTOM SCORER WITH METRICSMAKESCORER THE BUILTIN SCORER R2 USES MULTIOUTPUTUNIFORMAVERAGE

2306 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

SKLEARN SVM NUSVC

CLASS SKLEARN SVM NUSVC NU05 KERNEL 'RBF' DEGREE 3 GAMMA 'AUTO DEPRECATED'

COEF 0.00 SHRINKING TRUE PROBABILITY FALSE TOL 0.001

CACHE SIZE 200 CLASS WEIGHT NONE VERBOSE FALSE MAX ITER 1 DECI

SION FUNCTION SHAPE 'OVR' RANDOM STATE NONE

NUSUPPORT VECTOR CLASSIFICATION

SIMILAR TO SVC BUT USES A PARAMETER TO CONTROL THE NUMBER OF SUPPORT VECTORS

THE IMPLEMENTATION IS BASED ON LIBSVM

READ MORE IN THE USER GUIDE

PARAMETERS

NUFLOAT OPTIONAL DEFAULT 0.5 AN UPPER BOUND ON THE FRACTION OF TRAINING ERRORS AND A LOWER BOUND OF THE FRACTION OF SUPPORT VECTORS SHOULD BE IN THE INTERVAL 0 1

KERNEL STRING OPTIONAL DEFAULT 'RBF' SPECIFIES THE KERNEL TYPE TO BE USED IN THE ALGORITHM IT MUST BE ONE OF 'LINEAR' 'POLY' 'RBF' 'SIGMOID' 'PRECOMPUTED' OR A CALLABLE IF NONE IS GIVEN 'RBF' WILL BE USED IF A CALLABLE IS GIVEN IT IS USED TO PRECOMPUTE THE KERNEL MATRIX

DEGREE INT OPTIONAL DEFAULT 3 DEGREE OF THE POLYNOMIAL KERNEL FUNCTION 'POLY' IGNORED BY ALL OTHER KERNELS

GAMMA FLOAT OPTIONAL DEFAULT 'AUTO' KERNEL COEFFICIENT FOR 'RBF' 'POLY' AND 'SIGMOID' CURRENT DEFAULT IS 'AUTO' WHICH USES 1 / NFEATUES IF GAMMA SCALE IS PASSED THEN IT USES 1 / NFEATUES XVAR AS VALUE OF GAMMA THE CURRENT DEFAULT OF GAMMA 'AUTO' WILL CHANGE TO 'SCALE' IN VERSION 0.22 'AUTO DEPRECATED' A DEPRECATED VERSION OF 'AUTO' IS USED AS A DEFAULT INDICATING THAT NO EXPLICIT VALUE OF GAMMA WAS PASSED

COEF 0 FLOAT OPTIONAL DEFAULT 0.0 INDEPENDENT TERM IN KERNEL FUNCTION IT IS ONLY SIGNIFICANT IN 'POLY' AND 'SIGMOID'

SHRINKING BOOLEAN OPTIONAL DEFAULT TRUE WHETHER TO USE THE SHRINKING HEURISTIC

PROBABILITY BOOLEAN OPTIONAL DEFAULT FALSE WHETHER TO ENABLE PROBABILITY ESTIMATES THIS MUST BE ENABLED PRIOR TO CALLING FIT AND WILL SLOW DOWN THAT METHOD

TOL FLOAT OPTIONAL DEFAULT 1e-3 TOLERANCE FOR STOPPING CRITERION

CACHE SIZE FLOAT OPTIONAL SPECIFY THE SIZE OF THE KERNEL CACHE IN MB

CLASS WEIGHT DICT 'BALANCED' OPTIONAL SET THE PARAMETER C OF CLASS I TO CLASS WEIGHT I FOR SVC IF NOT GIVEN ALL CLASSES ARE SUPPOSED TO HAVE WEIGHT ONE THE "BALANCED" MODE USES THE VALUES OF Y TO AUTOMATICALLY ADJUST WEIGHTS INVERSELY PROPORTIONAL TO CLASS FREQUENCIES AS N SAMPLES N CLASSES N BINCOUNTY

637 SKLEARN SVM SUPPORT VECTOR MACHINES 2307

SCIKITLEARN USER GUIDE RELEASE 0213

VERBOSE BOOL DEFAULT FALSE ENABLE VERBOSE OUTPUT NOTE THAT THIS SETTING TAKES ADVANTAGE OF A PERPROCESS RUNTIME SETTING IN LIBSVM THAT IF ENABLED MAY NOT WORK PROPERLY IN A MULTITHREADED CONTEXT

MAXITER INT OPTIONAL DEFAULT1 HARD LIMIT ON ITERATIONS WITHIN SOLVER OR 1 FOR NO LIMIT

DECISIONFUNCTIONSHAPE 'OVO' 'OVR' DEFAULT'OVR' WHETHER TO RETURN A ONEVSREST 'OVR' DECISION FUNCTION OF SHAPE NSAMPLES NCLASSES AS ALL OTHER CLASSIFIERS OR THE ORIGINAL ONEVSONE 'OVO' DECISION FUNCTION OF LIBSVM WHICH HAS SHAPE NSAMPLES NCLASSES

NCLASSES 1 2

CHANGED IN VERSION 019 DECISIONFUNCTIONSHAPE IS 'OVR' BY DEFAULT

NEW IN VERSION 017 DECISIONFUNCTIONSHAPE'OVR' IS RECOMMENDED

CHANGED IN VERSION 017 DEPRECATED DECISIONFUNCTIONSHAPE'OVO' AND NONE

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE THE SEED OF THE PSEUDO RANDOM NUMBER GENERATOR USED WHEN SHUFFLING THE DATA FOR PROBABILITY ESTIMATES

IF INT RANDOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NPRANDOM

ATTRIBUTES

SUPPORT ARRAYLIKE SHAPE NSV INDICES OF SUPPORT VECTORS

SUPPORTVECTORS ARRAYLIKE SHAPE NSV NFEATURES SUPPORT VECTORS

NSUPPORT ARRAYLIKE DTYPEINT32 SHAPE NCLASS NUMBER OF SUPPORT VECTORS FOR EACH CLASS

DUALCOEF ARRAY SHAPE NCLASS1 NSV COEFFICIENTS OF THE SUPPORT VECTOR IN THE DECISION FUNCTION FOR MULTICLASS COEFFICIENT FOR ALL 1VS1 CLASSIFIERS THE LAYOUT OF THE COEFFICIENTS IN THE MULTICLASS CASE IS SOMEWHAT NONTRIVIAL SEE THE SECTION ABOUT MULTICLASS CLASSIFICATION IN THE SVM SECTION OF THE USER GUIDE FOR DETAILS

COEF ARRAY SHAPE NCLASS NCLASS1 2 NFEATURES WEIGHTS ASSIGNED TO THE FEATURES COEFFICIENTS IN THE PRIMAL PROBLEM THIS IS ONLY AVAILABLE IN THE CASE OF A LINEAR KERNEL

COEF IS READONLY PROPERTY DERIVED FROM DUALCOEF ANDSUPPORTVECTORS

INTERCEPT ARRAY SHAPE NCLASS NCLASS1 2 CONSTANTS IN DECISION FUNCTION

SEE ALSO

SVC SUPPORT VECTOR MACHINE FOR CLASSIFICATION USING LIBSVM

LINEARSVC SCALABLE LINEAR SUPPORT VECTOR MACHINE FOR CLASSIFICATION USING LIBLINEAR

NOTES

REFERENCES LIBSVM A LIBRARY FOR SUPPORT VECTOR MACHINES

EXAMPLES

IMPORT NUMPY AS NP

X NPARRAY1 1 2 1 1 1 2 1

Y NPARRAY1 1 2 2

FROM SKLEARN SVM IMPORT NUSVC

2308 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213  
CLF NUSVCGAMMASCALE  
CLFFITX Y  
NUSVCCACHESIZE200 CLASSWEIGHTNONE COEF000  
DECISIONFUNCTIONSHAPEOVR DEGREE3 GAMMASCALE KERNELRBF  
MAXITER1 NU05 PROBABILITYFALSE RANDOMSTATENONE  
SHRINKINGTRUE TOL0001 VERBOSEFALSE  
PRINTCLFPREDICT08 1  
1  
METHODS  
DECISIONFUNCTION SELF X EVALUATES THE DECISION FUNCTION FOR THE SAMPLES IN X  
FITSELF X Y SAMPLEWEIGHT FIT THE SVM MODEL ACCORDING TO THE GIVEN TRAINING DATA  
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR  
PREDICT SELF X PERFORM CLASSIFICATION ON SAMPLES IN X  
SCORE SELF X Y SAMPLEWEIGHT RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND  
LABELS  
SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR  
INIT SELFNU05 KERNEL'RBF' DEGREE3 GAMMA'AUTODEPRECATED' COEF000 SHRINK  
INGTRUE PROBABILITYFALSE TOL0001 CACHESIZE200 CLASSWEIGHTNONE VER  
BOSEFALSE MAXITER1 DECISIONFUNCTIONSHAPE'OVR' RANDOMSTATENONE  
DECISIONFUNCTION SELF X  
EVALUATES THE DECISION FUNCTION FOR THE SAMPLES IN X  
PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES  
RETURNS  
XARRAYLIKE SHAPE NSAMPLES NCLASSES NCLASSES1 2 RETURNS THE DECISION FUNC  
TION OF THE SAMPLE FOR EACH CLASS IN THE MODEL IF DECISIONFUNCTIONSHAPE'OVR' THE SHAPE  
IS NSAMPLES NCLASSES  
NOTES  
IF DECISIONFUNCTIONSHAPE'OVO' THE FUNCTION VALUES ARE PROPORTIONAL TO THE DISTANCE OF THE SAMPLES X TO  
THE SEPARATING HYPERPLANE IF THE EXACT DISTANCES ARE REQUIRED DIVIDE THE FUNCTION VALUES BY THE NORM OF  
THE WEIGHT VECTOR COEF SEE ALSO THIS QUESTION FOR FURTHER DETAILS IF DECISIONFUNCTIONSHAPE'OVR'  
THE DECISION FUNCTION IS A MONOTONIC TRANSFORMATION OF OVO DECISION FUNCTION  
FITSELFXYSAMPLEWEIGHTNONE  
FIT THE SVM MODEL ACCORDING TO THE GIVEN TRAINING DATA  
PARAMETERS  
XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES TRAINING VECTORS WHERE  
NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES FOR KER  
NEL"PRECOMPUTED" THE EXPECTED SHAPE OF X IS NSAMPLES NSAMPLES  
YARRAYLIKE SHAPE NSAMPLES TARGET VALUES CLASS LABELS IN CLASSIFICATION REAL NUMBERS  
IN REGRESSION  
637SKLEARN SVM SUPPORT VECTOR MACHINES 2309

SCIKITLEARN USER GUIDE RELEASE 0213

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES PERSAMPLE WEIGHTS RESCALE C PER SAMPLE  
HIGHER WEIGHTS FORCE THE CLASSIFIER TO PUT MORE EMPHASIS ON THESE POINTS

RETURNS  
SELF OBJECT

NOTES  
IF X AND Y ARE NOT CORDERED AND CONTIGUOUS ARRAYS OF NPFLOAT64 AND X IS NOT A SCIPYSPARSECSRMATRIX X  
ANDOR Y MAY BE COPIED  
IF X IS A DENSE ARRAY THEN THE OTHER METHODS WILL NOT SUPPORT SPARSE MATRICES AS INPUT

GETPARAMS SELFDEEPTURE  
GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS  
DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBOBJECTS THAT ARE ESTIMATORS

RETURNS  
PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELF  
PERFORM CLASSIFICATION ON SAMPLES IN X  
FOR AN ONECLASS MODEL 1 OR 1 IS RETURNED

PARAMETERS  
XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES FOR KERNEL"PRECOMPUTED"  
THE EXPECTED SHAPE OF X IS NSAMPLESTEST NSAMPLESTRAIN

RETURNS  
YPRED ARRAY SHAPE NSAMPLES CLASS LABELS FOR SAMPLES IN X

PREDICTLOGPROBA  
COMPUTE LOG PROBABILITIES OF POSSIBLE OUTCOMES FOR SAMPLES IN X  
THE MODEL NEED TO HAVE PROBABILITY INFORMATION COMPUTED AT TRAINING TIME FIT WITH ATTRIBUTE  
PROBABILITY SET TO TRUE

PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES FOR KERNEL"PRECOMPUTED" THE EXPECTED  
SHAPE OF X IS NSAMPLESTEST NSAMPLESTRAIN

RETURNS  
TARRAYLIKE SHAPE NSAMPLES NCLASSES RETURNS THE LOGPROBABILITIES OF THE SAMPLE FOR  
EACH CLASS IN THE MODEL THE COLUMNS CORRESPOND TO THE CLASSES IN SORTED ORDER AS THEY  
APPEAR IN THE ATTRIBUTE CLASSES

NOTES  
THE PROBABILITY MODEL IS CREATED USING CROSS VALIDATION SO THE RESULTS CAN BE SLIGHTLY DIFFERENT THAN THOSE  
OBTAINED BY PREDICT ALSO IT WILL PRODUCE MEANINGLESS RESULTS ON VERY SMALL DATASETS

2310 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

PREDICTPROBA

COMPUTE PROBABILITIES OF POSSIBLE OUTCOMES FOR SAMPLES IN X

THE MODEL NEED TO HAVE PROBABILITY INFORMATION COMPUTED AT TRAINING TIME FIT WITH ATTRIBUTE

PROBABILITY SET TO TRUE

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES FOR KERNEL"PRECOMPUTED" THE EXPECTED

SHAPE OF X IS NSAMPLESTEST NSAMPLESTRAIN

RETURNS

TARRAYLIKE SHAPE NSAMPLES NCLASSES RETURNS THE PROBABILITY OF THE SAMPLE FOR EACH

CLASS IN THE MODEL THE COLUMNS CORRESPOND TO THE CLASSES IN SORTED ORDER AS THEY APPEAR

IN THE ATTRIBUTE CLASSES

NOTES

THE PROBABILITY MODEL IS CREATED USING CROSS VALIDATION SO THE RESULTS CAN BE SLIGHTLY DIFFERENT THAN THOSE

OBTAINED BY PREDICT ALSO IT WILL PRODUCE MEANINGLESS RESULTS ON VERY SMALL DATASETS

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS

IN MULTILABEL CLASSIFICATION THIS IS THE SUBSET ACCURACY WHICH IS A HARSH METRIC SINCE YOU REQUIRE FOR EACH

SAMPLE THAT EACH LABEL SET BE CORRECTLY PREDICTED

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE LABELS FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT MEAN ACCURACY OF SELF PREDICTX WRT Y

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE

PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT

OF A NESTED OBJECT

RETURNS

SELF

EXAMPLES USING SKLEARN SVMNUSVC

•NONLINEAR SVM

637SKLEARN SVM SUPPORT VECTOR MACHINES 2311

SCIKITLEARN USER GUIDE RELEASE 0213

SKLEARN SVM NUSVR

CLASS SKLEARN SVM NUSVR NU05 C10 KERNEL 'RBF' DEGREE3 GAMMA 'AUTODEPRECATED'

COEF000 SHRINKING TRUE TOL0001 CACHESIZE200 VERBOSE FALSE

MAXITER1

NU SUPPORT VECTOR REGRESSION

SIMILAR TO NUSVC FOR REGRESSION USES A PARAMETER NU TO CONTROL THE NUMBER OF SUPPORT VECTORS HOWEVER UNLIKE NUSVC WHERE NU REPLACES C HERE NU REPLACES THE PARAMETER EPSILON OF EPSILON SVR

THE IMPLEMENTATION IS BASED ON LIBSVM

READ MORE IN THE USER GUIDE

PARAMETERS

NU FLOAT OPTIONAL AN UPPER BOUND ON THE FRACTION OF TRAINING ERRORS AND A LOWER BOUND ON THE FRACTION OF SUPPORT VECTORS SHOULD BE IN THE INTERVAL 0 1 BY DEFAULT 05 WILL BE TAKEN

CFLOAT OPTIONAL DEFAULT10 PENALTY PARAMETER C OF THE ERROR TERM

KERNEL STRING OPTIONAL DEFAULT 'RBF' SPECIFIES THE KERNEL TYPE TO BE USED IN THE ALGORITHM IT MUST BE ONE OF 'LINEAR' 'POLY' 'RBF' 'SIGMOID' 'PRECOMPUTED' OR A CALLABLE IF NONE IS GIVEN 'RBF' WILL BE USED IF A CALLABLE IS GIVEN IT IS USED TO PRECOMPUTE THE KERNEL MATRIX

DEGREE INT OPTIONAL DEFAULT3 DEGREE OF THE POLYNOMIAL KERNEL FUNCTION 'POLY' IGNORED BY ALL OTHER KERNELS

GAMMA FLOAT OPTIONAL DEFAULT 'AUTO' KERNEL COEFFICIENT FOR 'RBF' 'POLY' AND 'SIGMOID' CURRENT DEFAULT IS 'AUTO' WHICH USES 1 / NFEATUES IF GAMMA SCALE IS PASSED THEN IT USES 1 / NFEATUES XVAR AS VALUE OF GAMMA THE CURRENT DEFAULT OF GAMMA 'AUTO' WILL CHANGE TO 'SCALE' IN VERSION 022 'AUTODEPRECATED' A DEPRECATED VERSION OF 'AUTO' IS USED AS A DEFAULT INDICATING THAT NO EXPLICIT VALUE OF GAMMA WAS PASSED

COEF0 FLOAT OPTIONAL DEFAULT00 INDEPENDENT TERM IN KERNEL FUNCTION IT IS ONLY SIGNIFICANT IN 'POLY' AND 'SIGMOID'

SHRINKING BOOLEAN OPTIONAL DEFAULT TRUE WHETHER TO USE THE SHRINKING HEURISTIC

TOL FLOAT OPTIONAL DEFAULT1E3 TOLERANCE FOR STOPPING CRITERION

CACHESIZE FLOAT OPTIONAL SPECIFY THE SIZE OF THE KERNEL CACHE IN MB

VERBOSE BOOL DEFAULT FALSE ENABLE VERBOSE OUTPUT NOTE THAT THIS SETTING TAKES ADVANTAGE OF A PERPROCESS RUNTIME SETTING IN LIBSVM THAT IF ENABLED MAY NOT WORK PROPERLY IN A MULTITHREADED CONTEXT

MAXITER INT OPTIONAL DEFAULT1 HARD LIMIT ON ITERATIONS WITHIN SOLVER OR 1 FOR NO LIMIT

ATTRIBUTES

SUPPORT ARRAYLIKE SHAPE NSV INDICES OF SUPPORT VECTORS

SUPPORT VECTORS ARRAYLIKE SHAPE NSV NFEATUES SUPPORT VECTORS

DUALCOEF ARRAY SHAPE 1 NSV COEFFICIENTS OF THE SUPPORT VECTOR IN THE DECISION FUNCTION

COEF ARRAY SHAPE 1 NFEATUES WEIGHTS ASSIGNED TO THE FEATURES COEFFICIENTS IN THE PRIMAL PROBLEM THIS IS ONLY AVAILABLE IN THE CASE OF A LINEAR KERNEL

COEF IS READONLY PROPERTY DERIVED FROM DUALCOEF AND SUPPORT VECTORS

INTERCEPT ARRAY SHAPE 1 CONSTANTS IN DECISION FUNCTION

2312 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

SEE ALSO

NUSVC SUPPORT VECTOR MACHINE FOR CLASSIFICATION IMPLEMENTED WITH LIBSVM WITH A PARAMETER TO CONTROL THE NUMBER OF SUPPORT VECTORS

SVR EPSILON SUPPORT VECTOR MACHINE FOR REGRESSION IMPLEMENTED WITH LIBSVM

NOTES

REFERENCES LIBSVM A LIBRARY FOR SUPPORT VECTOR MACHINES

EXAMPLES

```
FROM SKLEARN SVM IMPORT NUSVR
IMPORT NUMPY AS NP
NSAMPLES NFEATURES 10 5
NPRANDOMSEED0
Y NPRANDOMRANDNNSAMPLES
X NPRANDOMRANDNNSAMPLES NFEATURES
CLF NUSVR GAMMA SCALE C10 NU01
CLFFITX Y
NUSVR C10 CACHESIZE200 COEF000 DEGREE3 GAMMA SCALE
KERNEL RBF MAXITER1 NU01 SHRINKING TRUE TOL0001
VERBOSE FALSE
```

METHODS

FITSELF X Y SAMPLEWEIGHT FIT THE SVM MODEL ACCORDING TO THE GIVEN TRAINING DATA

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PREDICT SELF X PERFORM REGRESSION ON SAMPLES IN X

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELF NU05 C10 KERNEL 'RBF' DEGREE3 GAMMA 'AUTODEPRECATED' COEF000 SHRINKING TRUE TOL0001 CACHESIZE200 VERBOSE FALSE MAXITER1

FITSELF X Y SAMPLEWEIGHT NONE

FIT THE SVM MODEL ACCORDING TO THE GIVEN TRAINING DATA

PARAMETERS

X ARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES TRAINING VECTORS WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES FOR KERNEL "PRECOMPUTED" THE EXPECTED SHAPE OF X IS NSAMPLES NSAMPLES

Y ARRAYLIKE SHAPE NSAMPLES TARGET VALUES CLASS LABELS IN CLASSIFICATION REAL NUMBERS IN REGRESSION

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES PER SAMPLE WEIGHTS RESCALE C PER SAMPLE HIGHER WEIGHTS FORCE THE CLASSIFIER TO PUT MORE EMPHASIS ON THESE POINTS

RETURNS

637 SKLEARN SVM SUPPORT VECTOR MACHINES 2313

SCIKITLEARN USER GUIDE RELEASE 0213

SELF OBJECT

NOTES

IF X AND Y ARE NOT CORDERED AND CONTIGUOUS ARRAYS OF NPFLOAT64 AND X IS NOT A SCIPYSPARSECSRMATRIX X ANDOR Y MAY BE COPIED

IF X IS A DENSE ARRAY THEN THE OTHER METHODS WILL NOT SUPPORT SPARSE MATRICES AS INPUT

GETPARAMS SELFDEEPTTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELF

PERFORM REGRESSION ON SAMPLES IN X

FOR AN ONECLASS MODEL 1 INLIER OR 1 OUTLIER IS RETURNED

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES FOR KERNEL"PRECOMPUTED"

THE EXPECTED SHAPE OF X IS NSAMPLESTEST NSAMPLESTRAIN

RETURNS

YPRED ARRAY SHAPE NSAMPLES

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{U}{V}$  WHERE U IS THE RESIDUAL SUM OF SQUARES  $Y_{TRUE} - Y_{PRED}$

$2SUM$  AND V IS THE TOTAL SUM OF SQUARES  $Y_{TRUE} - Y_{TRUEMEAN}$   $2SUM$  THE BEST POSSIBLE SCORE

IS 10 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS

PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 00

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES FOR SOME ESTIMATORS THIS MAY

BE A PRECOMPUTED KERNEL MATRIX INSTEAD SHAPE NSAMPLES NSAMPLESFITTED WHERE

NSAMPLESFITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE VALUES FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT R2 OF SELF PREDICTX WRT Y

NOTES

THE R2 SCORE USED WHEN CALLING SCORE ON A REGRESSOR WILL USE MULTIOUTPUTUNIFORMAVERAGE

FROM VERSION 023 TO KEEP CONSISTENT WITH METRICCSR2SCORE THIS WILL INFLUENCE THE SCORE

METHOD OF ALL THE MULTIOUTPUT REGRESSORS EXCEPT FOR MULTIOUTPUTMULTIOUTPUTREGRESSOR

2314 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

TO SPECIFY THE DEFAULT VALUE MANUALLY AND AVOID THE WARNING PLEASE EITHER CALL METRICSR2SCORE  
DIRECTLY OR MAKE A CUSTOM SCORER WITH METRICSMAKESCORER THE BUILTIN SCORER R2 USES  
MULTIOUTPUTUNIFORMAVERAGE

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE  
PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT  
OF A NESTED OBJECT

RETURNS  
SELF

EXAMPLES USING SKLEARN SVMNUSVR

- MODEL COMPLEXITY INFLUENCE

SKLEARN SVM ONECLASS SVM

CLASS SKLEARN SVM ONECLASS SVM KERNEL 'RBF' DEGREE 3 GAMMA 'AUTODEPRECATED' COEF 000  
TOL 0001 NU 05 SHRINKING TRUE CACHE SIZE 200 VERBOSE FALSE

MAX ITER 1 RANDOM STATE NONE

UNSUPERVISED OUTLIER DETECTION

ESTIMATE THE SUPPORT OF A HIGH DIMENSIONAL DISTRIBUTION

THE IMPLEMENTATION IS BASED ON LIBSVM

READ MORE IN THE USER GUIDE

PARAMETERS

KERNEL STRING OPTIONAL DEFAULT 'RBF' SPECIFIES THE KERNEL TYPE TO BE USED IN THE ALGORITHM  
IT MUST BE ONE OF 'LINEAR' 'POLY' 'RBF' 'SIGMOID' 'PRECOMPUTED' OR A CALLABLE IF NONE IS  
GIVEN 'RBF' WILL BE USED IF A CALLABLE IS GIVEN IT IS USED TO PRECOMPUTE THE KERNEL MATRIX

DEGREE INT OPTIONAL DEFAULT 3 DEGREE OF THE POLYNOMIAL KERNEL FUNCTION 'POLY' IGNORED  
BY ALL OTHER KERNELS

GAMMA FLOAT OPTIONAL DEFAULT 'AUTO' KERNEL COEFFICIENT FOR 'RBF' 'POLY' AND 'SIGMOID'  
CURRENT DEFAULT IS 'AUTO' WHICH USES 1 NFEATURES IF GAMMA SCALE IS PASSED THEN IT  
USES 1 NFEATURES XVAR AS VALUE OF GAMMA THE CURRENT DEFAULT OF GAMMA 'AUTO'  
WILL CHANGE TO 'SCALE' IN VERSION 022 'AUTODEPRECATED' A DEPRECATED VERSION OF 'AUTO' IS  
USED AS A DEFAULT INDICATING THAT NO EXPLICIT VALUE OF GAMMA WAS PASSED

COEF 0 FLOAT OPTIONAL DEFAULT 00 INDEPENDENT TERM IN KERNEL FUNCTION IT IS ONLY SIGNIFICANT  
IN 'POLY' AND 'SIGMOID'

TOL FLOAT OPTIONAL TOLERANCE FOR STOPPING CRITERION

NU FLOAT OPTIONAL AN UPPER BOUND ON THE FRACTION OF TRAINING ERRORS AND A LOWER BOUND OF THE  
FRACTION OF SUPPORT VECTORS SHOULD BE IN THE INTERVAL 0 1 BY DEFAULT 05 WILL BE TAKEN

SHRINKING BOOLEAN OPTIONAL WHETHER TO USE THE SHRINKING HEURISTIC

CACHE SIZE FLOAT OPTIONAL SPECIFY THE SIZE OF THE KERNEL CACHE IN MB

6375 SKLEARN SVM SUPPORT VECTOR MACHINES 2315

SCIKITLEARN USER GUIDE RELEASE 0213

VERBOSE BOOL DEFAULT FALSE ENABLE VERBOSE OUTPUT NOTE THAT THIS SETTING TAKES ADVANTAGE OF A PERPROCESS RUNTIME SETTING IN LIBSVM THAT IF ENABLED MAY NOT WORK PROPERLY IN A MULTITHREADED CONTEXT

MAXITER INT OPTIONAL DEFAULT1 HARD LIMIT ON ITERATIONS WITHIN SOLVER OR 1 FOR NO LIMIT

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IGNORED

DEPRECATED SINCE VERSION 020 RANDOMSTATE HAS BEEN DEPRECATED IN 020 AND WILL BE REMOVED IN 022

ATTRIBUTES

SUPPORT ARRAYLIKE SHAPE NSV INDICES OF SUPPORT VECTORS

SUPPORTVECTORS ARRAYLIKE SHAPE NSV NFEATURES SUPPORT VECTORS

DUALCOEF ARRAY SHAPE 1 NSV COEFFICIENTS OF THE SUPPORT VECTORS IN THE DECISION FUNCTION

COEF ARRAY SHAPE 1 NFEATURES WEIGHTS ASSIGNED TO THE FEATURES COEFFICIENTS IN THE PRIMAL PROBLEM THIS IS ONLY AVAILABLE IN THE CASE OF A LINEAR KERNEL

COEF IS READONLY PROPERTY DERIVED FROM DUALCOEF ANDSUPPORTVECTORS

INTERCEPT ARRAY SHAPE 1 CONSTANT IN THE DECISION FUNCTION

OFFSET FLOAT OFFSET USED TO DEFINE THE DECISION FUNCTION FROM THE RAW SCORES WE HAVE THE RELATION DECISIONFUNCTION SCORESAMPLES OFFSET THE OFFSET IS THE OPPOSITE OF INTERCEPT AND IS PROVIDED FOR CONSISTENCY WITH OTHER OUTLIER DETECTION ALGORITHMS

EXAMPLES

```
FROM SKLEARN SVM IMPORT ONECLASS SVM
X 0 044 045 046 1
CLF ONECLASS SVM GAMMA AUTO FIT X
CLF PREDICT X
ARRAY1 1 1 1 1
CLF SCORE SAMPLES X
ARRAY1 7798 20547 20556 20561 17332
```

METHODS

DECISIONFUNCTION SELF X SIGNED DISTANCE TO THE SEPARATING HYPERPLANE

FITSELF X Y SAMPLEWEIGHT DETECTS THE SOFT BOUNDARY OF THE SET OF SAMPLES X

FITPREDICT SELF X Y PERFORMS FIT ON X AND RETURNS LABELS FOR X

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PREDICT SELF X PERFORM CLASSIFICATION ON SAMPLES IN X

SCORESAMPLES SELF X RAW SCORING FUNCTION OF THE SAMPLES

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELF KERNEL'RBF' DEGREE3 GAMMA'AUTO DEPRECATED' COEF000 TOL0001 NU05

SHRINKINGTRUE CACHESIZE200 VERBOSEFALSE MAXITER1 RANDOMSTATENONE

DECISIONFUNCTION SELF X

SIGNED DISTANCE TO THE SEPARATING HYPERPLANE

2316 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

SIGNED DISTANCE IS POSITIVE FOR AN INLIER AND NEGATIVE FOR AN OUTLIER

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES

RETURNS

DEC ARRAYLIKE SHAPE NSAMPLES RETURNS THE DECISION FUNCTION OF THE SAMPLES

FITSELFXYNONE SAMPLEWEIGHTNONE PARAMS

DETECTS THE SOFT BOUNDARY OF THE SET OF SAMPLES X

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES SET OF SAMPLES WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES PERSAMPLE WEIGHTS RESCALE C PER SAMPLE

HIGHER WEIGHTS FORCE THE CLASSIFIER TO PUT MORE EMPHASIS ON THESE POINTS

YIGNORED NOT USED PRESENT FOR API CONSISTENCY BY CONVENTION

RETURNS

SELF OBJECT

NOTES

IF X IS NOT A CORDERED CONTIGUOUS ARRAY IT IS COPIED

FITPREDICT SELFXYNONE

PERFORMS FIT ON X AND RETURNS LABELS FOR X

RETURNS 1 FOR OUTLIERS AND 1 FOR INLIERS

PARAMETERS

XNDARRAY SHAPE NSAMPLES NFEATURES INPUT DATA

YIGNORED NOT USED PRESENT FOR API CONSISTENCY BY CONVENTION

RETURNS

YNDARRAY SHAPE NSAMPLES 1 FOR INLIERS 1 FOR OUTLIERS

GETPARAMS SELFDEEPTTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELFXY

PERFORM CLASSIFICATION ON SAMPLES IN X

FOR A ONECLASS MODEL 1 OR 1 IS RETURNED

PARAMETERS

637SKLEARN SVM SUPPORT VECTOR MACHINES 2317

SKITLEARN USER GUIDE RELEASE 0213  
 XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES FOR KERNEL"PRECOMPUTED"  
 THE EXPECTED SHAPE OF X IS NSAMPLESTEST NSAMPLESTRAIN  
 RETURNS  
 YPRED ARRAY SHAPE NSAMPLES CLASS LABELS FOR SAMPLES IN X  
 SCORESAMPLES SELF  
 RAW SCORING FUNCTION OF THE SAMPLES  
 PARAMETERS  
 XARRAYLIKE SHAPE NSAMPLES NFEATURES  
 RETURNS  
 SCORESAMPLES ARRAYLIKE SHAPE NSAMPLES RETURNS THE UNSHIFTED SCORING FUNCTION OF  
 THE SAMPLES  
 SETPARAMS SELFPARAMS  
 SET THE PARAMETERS OF THIS ESTIMATOR  
 THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE  
 PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT  
 OF A NESTED OBJECT  
 RETURNS  
 SELF  
 EXAMPLES USING SKLEARN SVM ONE CLASS SVM  
 •COMPARING ANOMALY DETECTION ALGORITHMS FOR OUTLIER DETECTION ON TOY DATASETS  
 •OUTLIER DETECTION ON A REAL DATA SET  
 •SPECIES DISTRIBUTION MODELING  
 •LIBSVM GUI  
 •ONECLASS SVM WITH NONLINEAR KERNEL RBF  
 SKLEARN SVM SVC  
 CLASS SKLEARN SVM SVCC10 KERNEL'RBF' DEGREE3 GAMMA'AUTODEPRECATED' COEF000 SHRINK  
 INGTRUE PROBABILITYFALSE TOL0001 CACHESIZE200 CLASSWEIGHTNONE  
 VERBOSEFALSE MAXITER1 DECISIONFUNCTIONSHAPE'OVR' RAN  
 DOMSTATENONE  
 CSUPPORT VECTOR CLASSIFICATION  
 THE IMPLEMENTATION IS BASED ON LIBSVM THE FIT TIME SCALES AT LEAST QUADRATICALLY WITH THE NUMBER OF SAMPLES  
 AND MAY BE IMPRACTICAL BEYOND TENS OF THOUSANDS OF SAMPLES FOR LARGE DATASETS CONSIDER USING SKLEARN  
 LINEAR MODEL LINEAR SVC OR SKLEARN LINEAR MODEL SGD CLASSIFIER INSTEAD POSSIBLY AFTER  
 ASKLEARN KERNEL APPROXIMATION NYSTROEM TRANSFORMER  
 THE MULTICLASS SUPPORT IS HANDLED ACCORDING TO A ONEVSONE SCHEME  
 FOR DETAILS ON THE PRECISE MATHEMATICAL FORMULATION OF THE PROVIDED KERNEL FUNCTIONS AND HOW GAMMA COEFO  
 AND DEGREE AFFECT EACH OTHER SEE THE CORRESPONDING SECTION IN THE NARRATIVE DOCUMENTATION KERNEL FUNCTIONS  
 READ MORE IN THE USER GUIDE  
 2318 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

PARAMETERS

CFLOAT OPTIONAL DEFAULT10 PENALTY PARAMETER C OF THE ERROR TERM

KERNEL STRING OPTIONAL DEFAULT'RBF' SPECIFIES THE KERNEL TYPE TO BE USED IN THE ALGORITHM IT MUST BE ONE OF 'LINEAR' 'POLY' 'RBF' 'SIGMOID' 'PRECOMPUTED' OR A CALLABLE IF NONE IS GIVEN 'RBF' WILL BE USED IF A CALLABLE IS GIVEN IT IS USED TO PRECOMPUTE THE KERNEL MATRIX FROM DATA MATRICES THAT MATRIX SHOULD BE AN ARRAY OF SHAPE NSAMPLES NSAMPLES

DEGREE INT OPTIONAL DEFAULT3 DEGREE OF THE POLYNOMIAL KERNEL FUNCTION 'POLY' IGNORED BY ALL OTHER KERNELS

GAMMA FLOAT OPTIONAL DEFAULT'AUTO' KERNEL COEFFICIENT FOR 'RBF' 'POLY' AND 'SIGMOID' CURRENT DEFAULT IS 'AUTO' WHICH USES 1 NFEATURES IF GAMMASCALE IS PASSED THEN IT USES 1 NFEATURES XVAR AS VALUE OF GAMMA THE CURRENT DEFAULT OF GAMMA 'AUTO' WILL CHANGE TO 'SCALE' IN VERSION 022 'AUTODEPRECATED' A DEPRECATED VERSION OF 'AUTO' IS USED AS A DEFAULT INDICATING THAT NO EXPLICIT VALUE OF GAMMA WAS PASSED

COEF0 FLOAT OPTIONAL DEFAULT00 INDEPENDENT TERM IN KERNEL FUNCTION IT IS ONLY SIGNIFICANT IN 'POLY' AND 'SIGMOID'

SHRINKING BOOLEAN OPTIONAL DEFAULTTRUE WHETHER TO USE THE SHRINKING HEURISTIC

PROBABILITY BOOLEAN OPTIONAL DEFAULTFALSE WHETHER TO ENABLE PROBABILITY ESTIMATES THIS MUST BE ENABLED PRIOR TO CALLING FIT AND WILL SLOW DOWN THAT METHOD

TOLFLOAT OPTIONAL DEFAULT1E3 TOLERANCE FOR STOPPING CRITERION

CACHESIZE FLOAT OPTIONAL SPECIFY THE SIZE OF THE KERNEL CACHE IN MB

CLASSWEIGHT DICT 'BALANCED' OPTIONAL SET THE PARAMETER C OF CLASS I TO CLASSWEIGHTIC FOR SVC IF NOT GIVEN ALL CLASSES ARE SUPPOSED TO HAVE WEIGHT ONE THE "BALANCED" MODE USES THE VALUES OF Y TO AUTOMATICALLY ADJUST WEIGHTS INVERSELY PROPORTIONAL TO CLASS FREQUENCIES IN THE INPUT DATA AS NSAMPLES NCLASSES NPBINCOUNTY

VERBOSE BOOL DEFAULT FALSE ENABLE VERBOSE OUTPUT NOTE THAT THIS SETTING TAKES ADVANTAGE OF A PERPROCESS RUNTIME SETTING IN LIBSVM THAT IF ENABLED MAY NOT WORK PROPERLY IN A MULTITHREADED CONTEXT

MAXITER INT OPTIONAL DEFAULT1 HARD LIMIT ON ITERATIONS WITHIN SOLVER OR 1 FOR NO LIMIT

DECISIONFUNCTIONSHAPE 'OVO' 'OVR' DEFAULT'OVR' WHETHER TO RETURN A ONEVSREST 'OVR' DECISION FUNCTION OF SHAPE NSAMPLES NCLASSES AS ALL OTHER CLASSIFIERS OR THE ORIGINAL ONEVSONE 'OVO' DECISION FUNCTION OF LIBSVM WHICH HAS SHAPE NSAMPLES NCLASSES NCLASSES 1 2 HOWEVER ONEVSONE 'OVO' IS ALWAYS USED AS MULTICLASS STRATEGY

CHANGED IN VERSION 019 DECISIONFUNCTIONSHAPE IS 'OVR' BY DEFAULT

NEW IN VERSION 017 DECISIONFUNCTIONSHAPE'OVR' IS RECOMMENDED

CHANGED IN VERSION 017 DEPRECATED DECISIONFUNCTIONSHAPE'OVO' AND NONE

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE THE SEED OF THE PSEUDO RANDOM NUMBER GENERATOR USED WHEN SHUFFLING THE DATA FOR PROBABILITY ESTIMATES IF INT RANDOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NPRANDOM

ATTRIBUTES

SUPPORT ARRAYLIKE SHAPE NSV INDICES OF SUPPORT VECTORS

637SKLEARN SVM SUPPORT VECTOR MACHINES 2319

SCIKITLEARN USER GUIDE RELEASE 0213

SUPPORTVECTORS ARRAYLIKE SHAPE NSV NFEATURES SUPPORT VECTORS

NSUPPORT ARRAYLIKE DTYPEINT32 SHAPE NCLASS NUMBER OF SUPPORT VECTORS FOR EACH CLASS

DUALCOEF ARRAY SHAPE NCLASS1 NSV COEFFICIENTS OF THE SUPPORT VECTOR IN THE DECISION FUNCTION FOR MULTICLASS COEFFICIENT FOR ALL 1VS1 CLASSIFIERS THE LAYOUT OF THE COEFFICIENTS IN THE MULTICLASS CASE IS SOMEWHAT NONTRIVIAL SEE THE SECTION ABOUT MULTICLASS CLASSIFICATION IN THE SVM SECTION OF THE USER GUIDE FOR DETAILS

COEF ARRAY SHAPE NCLASS NCLASS1 2 NFEATURES WEIGHTS ASSIGNED TO THE FEATURES COEFFICIENTS IN THE PRIMAL PROBLEM THIS IS ONLY AVAILABLE IN THE CASE OF A LINEAR KERNEL COEF IS A READONLY PROPERTY DERIVED FROM DUALCOEF ANDSUPPORTVECTORS

INTERCEPT ARRAY SHAPE NCLASS NCLASS1 2 CONSTANTS IN DECISION FUNCTION

FITSTATUS INT 0 IF CORRECTLY FITTED 1 OTHERWISE WILL RAISE WARNING

PROBA ARRAY SHAPE NCLASS NCLASS1 2

PROBB ARRAY SHAPE NCLASS NCLASS1 2 IF PROBABILITYTRUE THE PARAMETERS LEARNED IN PLATT SCALING TO PRODUCE PROBABILITY ESTIMATES FROM DECISION VALUES IF PROBABILITYFALSE AN EMPTY ARRAY PLATT SCALING USES THE LOGISTIC FUNCTION 1 1

EXPDECISIONVALUE PROBA PROBB WHEREPROBA ANDPROBB ARE LEARNED FROM THE DATASET R20C70293EF722 FOR MORE INFORMATION ON THE MULTICLASS CASE AND TRAINING PROCEDURE SEE SECTION 8 OF R20C70293EF721

SEE ALSO

SVR SUPPORT VECTOR MACHINE FOR REGRESSION IMPLEMENTED USING LIBSVM

LINEARSVC SCALABLE LINEAR SUPPORT VECTOR MACHINE FOR CLASSIFICATION IMPLEMENTED USING LIBLINEAR CHECK THE SEE ALSO SECTION OF LINEARSVC FOR MORE COMPARISON ELEMENT

REFERENCES

R20C70293EF721 R20C70293EF722

EXAMPLES

```
import numpy as np
X = np.array([1, 2, 1, 1, 1, 2, 1])
Y = np.array([1, 2, 2])
from sklearn.svm import SVC
clf = SVCGAMMAAUTO
clf.fit(X, Y)
svcc10 = cachesize200 classweightnone coef000
decisionfunctionshapeovr degree3 gammaauto kernelrbf
maxiter1 probabilityfalse randomstatenone shrinkingtrue
tol0001 verbosefalse
print(clf.predict(0.8))
```

1

2320 CHAPTER 6 API REFERENCE

METHODS

DECISIONFUNCTION SELF X EVALUATES THE DECISION FUNCTION FOR THE SAMPLES IN X  
FITSELF X Y SAMPLEWEIGHT FIT THE SVM MODEL ACCORDING TO THE GIVEN TRAINING DATA  
GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR  
PREDICT SELF X PERFORM CLASSIFICATION ON SAMPLES IN X  
SCORE SELF X Y SAMPLEWEIGHT RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR  
INIT SELF C10 KERNEL'RBF' DEGREE3 GAMMA'AUTODEPRECATED' COEF000 SHRINK  
INGTRUE PROBABILITYFALSE TOL0001 CACHESIZE200 CLASSWEIGHTNONE VER  
BOSEFALSE MAXITER1 DECISIONFUNCTIONSHAPE'OVR' RANDOMSTATENONE  
DECISIONFUNCTION SELF X  
EVALUATES THE DECISION FUNCTION FOR THE SAMPLES IN X  
PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES  
RETURNS  
XARRAYLIKE SHAPE NSAMPLES NCLASSES NCLASSES1 2 RETURNS THE DECISION FUNCTION OF THE SAMPLE FOR EACH CLASS IN THE MODEL IF DECISIONFUNCTIONSHAPE'OVR' THE SHAPE IS NSAMPLES NCLASSES

NOTES  
IF DECISIONFUNCTIONSHAPE'OVO' THE FUNCTION VALUES ARE PROPORTIONAL TO THE DISTANCE OF THE SAMPLES X TO THE SEPARATING HYPERPLANE IF THE EXACT DISTANCES ARE REQUIRED DIVIDE THE FUNCTION VALUES BY THE NORM OF THE WEIGHT VECTOR COEF SEE ALSO THIS QUESTION FOR FURTHER DETAILS IF DECISIONFUNCTIONSHAPE'OVR' THE DECISION FUNCTION IS A MONOTONIC TRANSFORMATION OF OVO DECISION FUNCTION

FITSELF X Y SAMPLEWEIGHT NONE  
FIT THE SVM MODEL ACCORDING TO THE GIVEN TRAINING DATA  
PARAMETERS  
XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES TRAINING VECTORS WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES FOR KERNEL"PRECOMPUTED" THE EXPECTED SHAPE OF X IS NSAMPLES NSAMPLES  
YARRAYLIKE SHAPE NSAMPLES TARGET VALUES CLASS LABELS IN CLASSIFICATION REAL NUMBERS IN REGRESSION  
SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES PERSAMPLE WEIGHTS RESCALE C PER SAMPLE HIGHER WEIGHTS FORCE THE CLASSIFIER TO PUT MORE EMPHASIS ON THESE POINTS

RETURNS  
SELF OBJECT  
NOTES  
IF X AND Y ARE NOT CORDERED AND CONTIGUOUS ARRAYS OF NPFLOAT64 AND X IS NOT A SCIPYSPARSECSRMATRIX X ANDOR Y MAY BE COPIED  
637SKLEARN SVM SUPPORT VECTOR MACHINES 2321

SCIKITLEARN USER GUIDE RELEASE 0213

IF X IS A DENSE ARRAY THEN THE OTHER METHODS WILL NOT SUPPORT SPARSE MATRICES AS INPUT

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELF

PERFORM CLASSIFICATION ON SAMPLES IN X

FOR AN ONECLASS MODEL 1 OR 1 IS RETURNED

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES FOR KERNEL"PRECOMPUTED"

THE EXPECTED SHAPE OF X IS NSAMPLESTEST NSAMPLESTRAIN

RETURNS

YPRED ARRAY SHAPE NSAMPLES CLASS LABELS FOR SAMPLES IN X

PREDICTLOGPROBA

COMPUTE LOG PROBABILITIES OF POSSIBLE OUTCOMES FOR SAMPLES IN X

THE MODEL NEED TO HAVE PROBABILITY INFORMATION COMPUTED AT TRAINING TIME FIT WITH ATTRIBUTE

PROBABILITY SET TO TRUE

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES FOR KERNEL"PRECOMPUTED" THE EXPECTED

SHAPE OF X IS NSAMPLESTEST NSAMPLESTRAIN

RETURNS

TARRAYLIKE SHAPE NSAMPLES NCLASSES RETURNS THE LOGPROBABILITIES OF THE SAMPLE FOR

EACH CLASS IN THE MODEL THE COLUMNS CORRESPOND TO THE CLASSES IN SORTED ORDER AS THEY

APPEAR IN THE ATTRIBUTE CLASSES

NOTES

THE PROBABILITY MODEL IS CREATED USING CROSS VALIDATION SO THE RESULTS CAN BE SLIGHTLY DIFFERENT THAN THOSE

OBTAINED BY PREDICT ALSO IT WILL PRODUCE MEANINGLESS RESULTS ON VERY SMALL DATASETS

PREDICTPROBA

COMPUTE PROBABILITIES OF POSSIBLE OUTCOMES FOR SAMPLES IN X

THE MODEL NEED TO HAVE PROBABILITY INFORMATION COMPUTED AT TRAINING TIME FIT WITH ATTRIBUTE

PROBABILITY SET TO TRUE

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES FOR KERNEL"PRECOMPUTED" THE EXPECTED

SHAPE OF X IS NSAMPLESTEST NSAMPLESTRAIN

RETURNS

2322 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

TARRAYLIKE SHAPE NSAMPLES NCLASSES RETURNS THE PROBABILITY OF THE SAMPLE FOR EACH CLASS IN THE MODEL THE COLUMNS CORRESPOND TO THE CLASSES IN SORTED ORDER AS THEY APPEAR IN THE ATTRIBUTE CLASSES

NOTES

THE PROBABILITY MODEL IS CREATED USING CROSS VALIDATION SO THE RESULTS CAN BE SLIGHTLY DIFFERENT THAN THOSE OBTAINED BY PREDICT ALSO IT WILL PRODUCE MEANINGLESS RESULTS ON VERY SMALL DATASETS

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS

IN MULTILABEL CLASSIFICATION THIS IS THE SUBSET ACCURACY WHICH IS A HARSH METRIC SINCE YOU REQUIRE FOR EACH SAMPLE THAT EACH LABEL SET BE CORRECTLY PREDICTED

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE LABELS FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT MEAN ACCURACY OF SELF PREDICT X WRT Y

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT

OF A NESTED OBJECT

RETURNS

SELF

EXAMPLES USING SKLEARN SVM SVC

- MULTILABEL CLASSIFICATION
- EXPLICIT FEATURE MAP APPROXIMATION FOR RBF KERNELS
- FACES RECOGNITION EXAMPLE USING EIGENFACES AND SVMs
- LIBSVM GUI
- RECOGNIZING HANDWRITTEN DIGITS
- PLOT CLASSIFICATION PROBABILITY
- CLASSIFIER COMPARISON
- CONCATENATING MULTIPLE FEATURE EXTRACTION METHODS
- PLOT THE DECISION BOUNDARIES OF A VOTING CLASSIFIER
- CROSSVALIDATION ON DIGITS DATASET EXERCISE
- SVM EXERCISE
- RECURSIVE FEATURE ELIMINATION

SCIKITLEARN USER GUIDE RELEASE 0213

- RECURSIVE FEATURE ELIMINATION WITH CROSSVALIDATION
- TEST WITH PERMUTATIONS THE SIGNIFICANCE OF A CLASSIFICATION SCORE
- UNIVARIATE FEATURE SELECTION
- PLOTING VALIDATION CURVES
- PARAMETER ESTIMATION USING GRID SEARCH WITH CROSSVALIDATION
- RECEIVER OPERATING CHARACTERISTIC ROC WITH CROSS VALIDATION
- NESTED VERSUS NONNESTED CROSSVALIDATION
- CONFUSION MATRIX
- RECEIVER OPERATING CHARACTERISTIC ROC
- PLOTING LEARNING CURVES
- FEATURE DISCRETIZATION
- DECISION BOUNDARY OF LABEL PROPAGATION VERSUS SVM ON THE IRIS DATASET
- SVM MAXIMUM MARGIN SEPARATING HYPERPLANE
- SVM WITH CUSTOM KERNEL
- SVM WEIGHTED SAMPLES
- SVM SEPARATING HYPERPLANE FOR UNBALANCED CLASSES
- SVMKERNELS
- SVMANOVA SVM WITH UNIVARIATE FEATURE SELECTION
- SVM MARGINS EXAMPLE
- PLOT DIFFERENT SVM CLASSIFIERS IN THE IRIS DATASET
- RBF SVM PARAMETERS

SKLEARN SVM SVR

CLASSSSKLEARN SVM SVRKERNEL' RBF' DEGREE3 GAMMA'AUTODEPRECATED' COEF000 TOL0001

C10 EPSILON01 SHRINKINGTRUE CACHESIZE200 VERBOSEFALSE

MAXITER1

EPSILONSUPPORT VECTOR REGRESSION

THE FREE PARAMETERS IN THE MODEL ARE C AND EPSILON

THE IMPLEMENTATION IS BASED ON LIBSVM THE FIT TIME COMPLEXITY IS MORE THAN QUADRATIC WITH THE NUMBER OF SAMPLES WHICH MAKES IT HARD TO SCALE TO DATASETS WITH MORE THAN A COUPLE OF 10000 SAMPLES FOR LARGE DATASETS CONSIDER USING SKLEARNLINEARMODELLINEARSVR ORSKLEARNLINEARMODELSGDREGRESSOR

INSTEAD POSSIBLY AFTER A SKLEARNKERNELAPPROXIMATIONNYSTROEM TRANSFORMER

READ MORE IN THE USER GUIDE

PARAMETERS

KERNEL STRING OPTIONAL DEFAULT' RBF' SPECIFIES THE KERNEL TYPE TO BE USED IN THE ALGORITHM

IT MUST BE ONE OF 'LINEAR' 'POLY' 'RBF' 'SIGMOID' 'PRECOMPUTED' OR A CALLABLE IF NONE IS

GIVEN 'RBF' WILL BE USED IF A CALLABLE IS GIVEN IT IS USED TO PRECOMPUTE THE KERNEL MATRIX

DEGREE INT OPTIONAL DEFAULT3 DEGREE OF THE POLYNOMIAL KERNEL FUNCTION 'POLY' IGNORED

BY ALL OTHER KERNELS

2324 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

GAMMA FLOAT OPTIONAL DEFAULT 'AUTO' KERNEL COEFFICIENT FOR 'RBF' 'POLY' AND 'SIGMOID' CURRENT DEFAULT IS 'AUTO' WHICH USES 1 NFEATURES IF GAMMASCALE IS PASSED THEN IT USES 1 NFEATURES XVAR AS VALUE OF GAMMA THE CURRENT DEFAULT OF GAMMA 'AUTO' WILL CHANGE TO 'SCALE' IN VERSION 022 'AUTODEPRECATED' A DEPRECATED VERSION OF 'AUTO' IS USED AS A DEFAULT INDICATING THAT NO EXPLICIT VALUE OF GAMMA WAS PASSED

COEF0 FLOAT OPTIONAL DEFAULT00 INDEPENDENT TERM IN KERNEL FUNCTION IT IS ONLY SIGNIFICANT IN 'POLY' AND 'SIGMOID'

TOLFLOAT OPTIONAL DEFAULT1E3 TOLERANCE FOR STOPPING CRITERION

CFLOAT OPTIONAL DEFAULT10 PENALTY PARAMETER C OF THE ERROR TERM

EPSILON FLOAT OPTIONAL DEFAULT01 EPSILON IN THE EPSILONSVR MODEL IT SPECIFIES THE EPSILONTUBE WITHIN WHICH NO PENALTY IS ASSOCIATED IN THE TRAINING LOSS FUNCTION WITH POINTS PREDICTED WITHIN A DISTANCE EPSILON FROM THE ACTUAL VALUE

SHRINKING BOOLEAN OPTIONAL DEFAULTTRUE WHETHER TO USE THE SHRINKING HEURISTIC

CACHESIZE FLOAT OPTIONAL SPECIFY THE SIZE OF THE KERNEL CACHE IN MB

VERBOSE BOOL DEFAULT FALSE ENABLE VERBOSE OUTPUT NOTE THAT THIS SETTING TAKES ADVANTAGE OF A PERPROCESS RUNTIME SETTING IN LIBSVM THAT IF ENABLED MAY NOT WORK PROPERLY IN A MULTITHREADED CONTEXT

MAXITER INT OPTIONAL DEFAULT1 HARD LIMIT ON ITERATIONS WITHIN SOLVER OR 1 FOR NO LIMIT

ATTRIBUTES

SUPPORT ARRAYLIKE SHAPE NSV INDICES OF SUPPORT VECTORS

SUPPORTVECTORS ARRAYLIKE SHAPE NSV NFEATURES SUPPORT VECTORS

DUALCOEF ARRAY SHAPE 1 NSV COEFFICIENTS OF THE SUPPORT VECTOR IN THE DECISION FUNCTION

COEF ARRAY SHAPE 1 NFEATURES WEIGHTS ASSIGNED TO THE FEATURES COEFFICIENTS IN THE PRIMAL PROBLEM THIS IS ONLY AVAILABLE IN THE CASE OF A LINEAR KERNEL

COEF IS READONLY PROPERTY DERIVED FROM DUALCOEF ANDSUPPORTVECTORS

INTERCEPT ARRAY SHAPE 1 CONSTANTS IN DECISION FUNCTION

SEE ALSO

NUSVR SUPPORT VECTOR MACHINE FOR REGRESSION IMPLEMENTED USING LIBSVM USING A PARAMETER TO CONTROL THE NUMBER OF SUPPORT VECTORS

LINEARSVR SCALABLE LINEAR SUPPORT VECTOR MACHINE FOR REGRESSION IMPLEMENTED USING LIBLINEAR

NOTES

REFERENCES LIBSVM A LIBRARY FOR SUPPORT VECTOR MACHINES

EXAMPLES

```
FROM SKLEARN SVM IMPORT SVR
IMPORT NUMPY AS NP
NSAMPLES NFEATURES 10 5
RNG NPRANDOMRANDOMSTATE0
637SKLEARN SVM SUPPORT VECTOR MACHINES 2325
```

SCIKITLEARN USER GUIDE RELEASE 0213

Y RNGRANDNNSAMPLES

X RNGRANDNNSAMPLES NFEATURES

CLF SVRGAMMASCALE C10 EPSILON02

CLFFITX Y

SVRC10 CACHESIZE200 COEF000 DEGREE3 EPSILON02 GAMMASCALE

KERNELRBF MAXITER1 SHRINKINGTRUE TOL0001 VERBOSEFALSE

METHODS

FITSELF X Y SAMPLEWEIGHT FIT THE SVM MODEL ACCORDING TO THE GIVEN TRAINING DATA

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PREDICT SELF X PERFORM REGRESSION ON SAMPLES IN X

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELFKERNEL'RBF' DEGREE3 GAMMA'AUTODEPRECATED' COEF000 TOL0001 C10 EPSILON01 SHRINKINGTRUE CACHESIZE200 VERBOSEFALSE MAXITER1

FITSELFXYSAMPLEWEIGHTNONE

FIT THE SVM MODEL ACCORDING TO THE GIVEN TRAINING DATA

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES TRAINING VECTORS WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES FOR KERNEL"PRECOMPUTED" THE EXPECTED SHAPE OF X IS NSAMPLES NSAMPLES

YARRAYLIKE SHAPE NSAMPLES TARGET VALUES CLASS LABELS IN CLASSIFICATION REAL NUMBERS IN REGRESSION

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES PERSAMPLE WEIGHTS RESCALE C PER SAMPLE HIGHER WEIGHTS FORCE THE CLASSIFIER TO PUT MORE EMPHASIS ON THESE POINTS

RETURNS

SELF OBJECT

NOTES

IF X AND Y ARE NOT CORDERED AND CONTIGUOUS ARRAYS OF NPFLOAT64 AND X IS NOT A SCIPYSPARSECSRMATRIX X ANDOR Y MAY BE COPIED

IF X IS A DENSE ARRAY THEN THE OTHER METHODS WILL NOT SUPPORT SPARSE MATRICES AS INPUT

GETPARAMS SELFDEEPTTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

2326 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

PREDICTSELF

PERFORM REGRESSION ON SAMPLES IN X

FOR AN ONECLASS MODEL 1 INLIER OR 1 OUTLIER IS RETURNED

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES FOR KERNEL"PRECOMPUTED"

THE EXPECTED SHAPE OF X IS NSAMPLESTEST NSAMPLESTRAIN

RETURNS

YPRED ARRAY SHAPE NSAMPLES

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{U}{V}$  WHERE U IS THE RESIDUAL SUM OF SQUARES YTRUE YPRED

2SUM AND V IS THE TOTAL SUM OF SQUARES YTRUE YTRUEMEAN 2SUM THE BEST POSSIBLE SCORE

IS 10 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS

PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 00

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES FOR SOME ESTIMATORS THIS MAY

BE A PRECOMPUTED KERNEL MATRIX INSTEAD SHAPE NSAMPLES NSAMPLESFITTED WHERE

NSAMPLESFITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE VALUES FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT R2 OF SELF-PREDICTX WRT Y

NOTES

THE R2 SCORE USED WHEN CALLING SCORE ON A REGRESSOR WILL USE MULTIOUTPUTUNIFORMAVERAGE

FROM VERSION 023 TO KEEP CONSISTENT WITH METRICSR2SCORE THIS WILL INFLUENCE THE SCORE

METHOD OF ALL THE MULTIOUTPUT REGRESSORS EXCEPT FOR MULTIOUTPUTMULTIOUTPUTREGRESSOR

TO SPECIFY THE DEFAULT VALUE MANUALLY AND AVOID THE WARNING PLEASE EITHER CALL METRICSR2SCORE

DIRECTLY OR MAKE A CUSTOM SCORER WITH METRICSMAKESCORER THE BUILTIN SCORER R2 USES

MULTIOUTPUTUNIFORMAVERAGE

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE

PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT

OF A NESTED OBJECT

RETURNS

SELF

EXAMPLES USING SKLEARN SVM SVR

- COMPARISON OF KERNEL RIDGE REGRESSION AND SVR
- PREDICTION LATENCY

637SKLEARN SVM SUPPORT VECTOR MACHINES 2327

SCIKITLEARN USER GUIDE RELEASE 0213

•SUPPORT VECTOR REGRESSION SVR USING LINEAR AND NONLINEAR KERNELS  
SVML1MINC X Y LOSS FITINTERCEPT RETURN THE LOWEST BOUND FOR C SUCH THAT FOR C IN L1MINC INFINITY THE MODEL IS GUARANTEED NOT TO BE EMPTY  
SKLEARN SVM L1MINC  
SKLEARN SVM L1MINC XYLOSS'SQUAREDHINGE' FITINTERCEPTTRUE INTERCEPTSCALING10  
RETURN THE LOWEST BOUND FOR C SUCH THAT FOR C IN L1MINC INFINITY THE MODEL IS GUARANTEED NOT TO BE EMPTY THIS APPLIES TO L1 PENALIZED CLASSIFIERS SUCH AS LINEARSVC WITH PENALTY'L1' AND LINEARLOGISTICREGRESSION WITH PENALTY'L1'  
THIS VALUE IS VALID IF CLASSWEIGHT PARAMETER IN FIT IS NOT SET  
PARAMETERS  
XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES TRAINING VECTOR WHERE NSAMPLES IN THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES  
YARRAY SHAPE NSAMPLES TARGET VECTOR RELATIVE TO X  
LOSS 'SQUAREDHINGE' 'LOG' DEFAULT 'SQUAREDHINGE' SPECIFIES THE LOSS FUNCTION WITH 'SQUAREDHINGE' IT IS THE SQUARED HINGE LOSS AKA L2 LOSS WITH 'LOG' IT IS THE LOSS OF LOGISTIC REGRESSION MODELS  
FITINTERCEPT BOOL DEFAULT TRUE SPECIFIES IF THE INTERCEPT SHOULD BE FITTED BY THE MODEL IT MUST MATCH THE FIT METHOD PARAMETER  
INTERCEPTSCALING FLOAT DEFAULT 1 WHEN FITINTERCEPT IS TRUE INSTANCE VECTOR X BECOMES X INTERCEPTSCALING IE A "SYNTHETIC" FEATURE WITH CONSTANT VALUE EQUALS TO INTERCEPTSCALING IS APPENDED TO THE INSTANCE VECTOR IT MUST MATCH THE FIT METHOD PARAMETER  
RETURNS  
L1MINC FLOAT MINIMUM VALUE FOR C  
EXAMPLES USING SKLEARN SVML1MINC  
•REGULARIZATION PATH OF L1 LOGISTIC REGRESSION  
6372 LOWLEVEL METHODS  
SVMLIBSVMCROSSVALIDATION BINDING OF THE CROSSVALIDATION ROUTINE LOWLEVEL ROUTINE  
SVMLIBSVMDECISIONFUNCTION PREDICT MARGIN LIBSVM NAME FOR THIS IS PREDICTVALUES  
SVMLIBSVMFIT TRAIN THE MODEL USING LIBSVM LOWLEVEL METHOD  
SVMLIBSVM PREDICT PREDICT TARGET VALUES OF X GIVEN A MODEL LOWLEVEL METHOD  
SVMLIBSVM PREDICTPROBA PREDICT PROBABILITIES  
SKLEARN SVMLIBSVM CROSSVALIDATION  
SKLEARN SVMLIBSVM CROSSVALIDATION  
BINDING OF THE CROSSVALIDATION ROUTINE LOWLEVEL ROUTINE  
PARAMETERS  
2328 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

XARRAYLIKE DTYPEFLOAT SIZENSAMPLES NFEATURES

YARRAY DTYPEFLOAT SIZENSAMPLES TARGET VECTOR

SVMTYPE 0 1 2 3 4 TYPE OF SVM C SVC NU SVC ONE CLASS EPSILON SVR NU SVR

KERNEL 'LINEAR' 'RBF' 'POLY' 'SIGMOID' 'PRECOMPUTED' KERNEL TO USE IN THE MODEL LINEAR

POLYNOMIAL RBF SIGMOID OR PRECOMPUTED

DEGREE INT DEGREE OF THE POLYNOMIAL KERNEL ONLY RELEVANT IF KERNEL IS SET TO POLYNOMIAL

GAMMA FLOAT GAMMA PARAMETER IN RBF POLY AND SIGMOID KERNELS IGNORED BY OTHER KERNELS

01 BY DEFAULT

COEF0 FLOAT INDEPENDENT PARAMETER IN POLYSIGMOID KERNEL

TOLFLOAT STOPPING CRITERIA

CFLOAT C PARAMETER IN CSUPPORT VECTOR CLASSIFICATION

NUFLOAT

CACHESIZE FLOAT

RANDOMSEED INT OPTIONAL SEED FOR THE RANDOM NUMBER GENERATOR USED FOR PROBABILITY ESTI

MATES 0 BY DEFAULT

RETURNS

TARGET ARRAY FLOAT

SKLEARNSVMLIBSVM DECISIONFUNCTION

SKLEARNSVMLIBSVM DECISIONFUNCTION

PREDICT MARGIN LIBSVM NAME FOR THIS IS PREDICTVALUES

WE HAVE TO RECONSTRUCT MODEL AND PARAMETERS TO MAKE SURE WE STAY IN SYNC WITH THE PYTHON OBJECT

SKLEARNSVMLIBSVM FIT

SKLEARNSVMLIBSVM FIT

TRAIN THE MODEL USING LIBSVM LOWLEVEL METHOD

PARAMETERS

XARRAYLIKE DTYPEFLOAT64 SIZENSAMPLES NFEATURES

YARRAY DTYPEFLOAT64 SIZENSAMPLES TARGET VECTOR

SVMTYPE 0 1 2 3 4 OPTIONAL TYPE OF SVM CSVC NUSVC ONECLASSSSVM EP

SILONSVR OR NUSVR RESPECTIVELY 0 BY DEFAULT

KERNEL 'LINEAR' 'RBF' 'POLY' 'SIGMOID' 'PRECOMPUTED' OPTIONAL KERNEL TO USE IN THE MODEL

LINEAR POLYNOMIAL RBF SIGMOID OR PRECOMPUTED 'RBF' BY DEFAULT

DEGREE INT32 OPTIONAL DEGREE OF THE POLYNOMIAL KERNEL ONLY RELEVANT IF KERNEL IS SET TO POLY

NOMIAL 3 BY DEFAULT

GAMMA FLOAT64 OPTIONAL GAMMA PARAMETER IN RBF POLY AND SIGMOID KERNELS IGNORED BY

OTHER KERNELS 01 BY DEFAULT

COEF0 FLOAT64 OPTIONAL INDEPENDENT PARAMETER IN POLYSIGMOID KERNEL 0 BY DEFAULT

TOLFLOAT64 OPTIONAL NUMERIC STOPPING CRITERION WRITEME 1E3 BY DEFAULT

637SKLEARN SVM SUPPORT VECTOR MACHINES 2329

SCIKITLEARN USER GUIDE RELEASE 0213

CFLOAT64 OPTIONAL C PARAMETER IN CSUPPORT VECTOR CLASSIFICATION 1 BY DEFAULT

NUFLOAT64 OPTIONAL 05 BY DEFAULT

EPSILON DOUBLE OPTIONAL 01 BY DEFAULT

CLASSWEIGHT ARRAY DTYPE FLOAT64 SHAPE NCLASSES OPTIONAL NEMPTY0 BY DEFAULT

SAMPLEWEIGHT ARRAY DTYPE FLOAT64 SHAPE NSAMPLES OPTIONAL NEMPTY0 BY DEFAULT

SHRINKING INT OPTIONAL 1 BY DEFAULT

PROBABILITY INT OPTIONAL 0 BY DEFAULT

CACHESIZE FLOAT64 OPTIONAL CACHE SIZE FOR GRAM MATRIX COLUMNS IN MEGABYTES 100 BY DEFAULT

MAXITER INT 1 FOR NO LIMIT OPTIONAL STOP SOLVER AFTER THIS MANY ITERATIONS REGARDLESS OF ACCURACY XXX CURRENTLY THERE IS NO API TO KNOW WHETHER THIS KICKED IN 1 BY DEFAULT

RANDOMSEED INT OPTIONAL SEED FOR THE RANDOM NUMBER GENERATOR USED FOR PROBABILITY ESTIMATES 0 BY DEFAULT

RETURNS

SUPPORT ARRAY SHAPENSUPPORT INDEX OF SUPPORT VECTORS

SUPPORTVECTORS ARRAY SHAPENSUPPORT NFEATURES SUPPORT VECTORS EQUIVALENT TO XSUPPORT WILL RETURN AN EMPTY ARRAY IN THE CASE OF PRECOMPUTED KERNEL

NCLASSSV ARRAY NUMBER OF SUPPORT VECTORS IN EACH CLASS

SVCOEF ARRAY COEFFICIENTS OF SUPPORT VECTORS IN DECISION FUNCTION

INTERCEPT ARRAY INTERCEPT IN DECISION FUNCTION

PROBA PROBB ARRAY PROBABILITY ESTIMATES EMPTY ARRAY FOR PROBABILITYFALSE

SKLEARNSVMLIBSVM PREDICT

SKLEARNSVMLIBSVM PREDICT

PREDICT TARGET VALUES OF X GIVEN A MODEL LOWLEVEL METHOD

PARAMETERS

XARRAYLIKE DTYPEFLOAT SIZENSAMPLES NFEATURES

SVMTYPE 0 1 2 3 4 TYPE OF SVM C SVC NU SVC ONE CLASS EPSILON SVR NU SVR

KERNEL 'LINEAR' 'RBF' 'POLY' 'SIGMOID' 'PRECOMPUTED' TYPE OF KERNEL

DEGREE INT DEGREE OF THE POLYNOMIAL KERNEL

GAMMA FLOAT GAMMA PARAMETER IN RBF POLY AND SIGMOID KERNELS IGNORED BY OTHER KERNELS 01 BY DEFAULT

COEF0 FLOAT INDEPENDENT PARAMETER IN POLYSIGMOID KERNEL

RETURNS

DECVALUES ARRAY PREDICTED VALUES

2330 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213  
SKLEARNVMLIBSVM PREDICTPROBA  
SKLEARNVMLIBSVM PREDICTPROBA  
PREDICT PROBABILITIES  
SVMMODEL STORES ALL PARAMETERS NEEDED TO PREDICT A GIVEN VALUE  
FOR SPEED ALL REAL WORK IS DONE AT THE C LEVEL IN FUNCTION COPYPREDICT LIBSVMHELPERC  
WE HAVE TO RECONSTRUCT MODEL AND PARAMETERS TO MAKE SURE WE STAY IN SYNC WITH THE PYTHON OBJECT  
SEE SKLEARNSVMPREDICT FOR A COMPLETE LIST OF PARAMETERS  
PARAMETERS  
XARRAYLIKE DTYPEFLOAT  
KERNEL 'LINEAR' 'RBF' 'POLY' 'SIGMOID' 'PRECOMPUTED'  
RETURNS  
DECVALUES ARRAY PREDICTED VALUES  
638SKLEARNTREE DECISION TREES  
THESKLEARNTREE MODULE INCLUDES DECISION TREEBASED MODELS FOR CLASSIFICATION AND REGRESSION  
USER GUIDE SEE THE DECISION TREES SECTION FOR FURTHER DETAILS  
TREEDECISIONTREECLASSIFIER CRITERION A DECISION TREE CLASSIFIER  
TREEDECISIONTREEREgressor CRITERION A DECISION TREE REGRESSOR  
TREEEXTRATREECLASSIFIER CRITERION AN EXTREMELY RANDOMIZED TREE CLASSIFIER  
TREEEXTRATREEREgressor CRITERION AN EXTREMELY RANDOMIZED TREE REGRESSOR  
6381SKLEARNTREE DECISIONTREECLASSIFIER  
CLASSSSKLEARNTREE DECISIONTREECLASSIFIER CRITERION'GINI' SPLITTER'BEST' MAXDEPTHNONE  
MINSAMPLESSPLIT2 MINSAMPLESLEAF1  
MINWEIGHTFRACTIONLEAF00  
MAXFEATURESNONE RAN  
DOMSTATENONE MAXLEAFNODESNONE  
MINIMPURITYDECREASE00  
MINIMPURITYSPLITNONE CLASSWEIGHTNONE  
PRESORTFALSE  
A DECISION TREE CLASSIFIER  
READ MORE IN THE USER GUIDE  
PARAMETERS  
CRITERION STRING OPTIONAL DEFAULT"GINI" THE FUNCTION TO MEASURE THE QUALITY OF A SPLIT SUP  
PORTED CRITERIA ARE "GINI" FOR THE GINI IMPURITY AND "ENTROPY" FOR THE INFORMATION GAIN  
SPLITTER STRING OPTIONAL DEFAULT"BEST" THE STRATEGY USED TO CHOOSE THE SPLIT AT EACH NODE  
SUPPORTED STRATEGIES ARE "BEST" TO CHOOSE THE BEST SPLIT AND "RANDOM" TO CHOOSE THE BEST  
RANDOM SPLIT  
MAXDEPTH INT OR NONE OPTIONAL DEFAULTNONE THE MAXIMUM DEPTH OF THE TREE IF NONE  
638SKLEARNTREE DECISION TREES 2331

SCIKITLEARN USER GUIDE RELEASE 0213

THEN NODES ARE EXPANDED UNTIL ALL LEAVES ARE PURE OR UNTIL ALL LEAVES CONTAIN LESS THAN MINSAMPLESSPLIT SAMPLES

MINSAMPLESSPLIT INT FLOAT OPTIONAL DEFAULT2 THE MINIMUM NUMBER OF SAMPLES REQUIRED TO SPLIT AN INTERNAL NODE

- IF INT THEN CONSIDER MINSAMPLESSPLIT AS THE MINIMUM NUMBER
- IF FLOAT THEN MINSAMPLESSPLIT IS A FRACTION AND CEILMINSAMPLESSPLIT

NSAMPLES ARE THE MINIMUM NUMBER OF SAMPLES FOR EACH SPLIT

CHANGED IN VERSION 018 ADDED FLOAT VALUES FOR FRACTIONS

MINSAMPLESLEAF INT FLOAT OPTIONAL DEFAULT1 THE MINIMUM NUMBER OF SAMPLES REQUIRED TO BE AT A LEAF NODE A SPLIT POINT AT ANY DEPTH WILL ONLY BE CONSIDERED IF IT LEAVES AT LEAST

MINSAMPLESLEAF TRAINING SAMPLES IN EACH OF THE LEFT AND RIGHT BRANCHES THIS MAY

HAVE THE EFFECT OF SMOOTHING THE MODEL ESPECIALLY IN REGRESSION

- IF INT THEN CONSIDER MINSAMPLESLEAF AS THE MINIMUM NUMBER
- IF FLOAT THEN MINSAMPLESLEAF IS A FRACTION AND CEILMINSAMPLESLEAF

NSAMPLES ARE THE MINIMUM NUMBER OF SAMPLES FOR EACH NODE

CHANGED IN VERSION 018 ADDED FLOAT VALUES FOR FRACTIONS

MINWEIGHTFRACTIONLEAF FLOAT OPTIONAL DEFAULT0 THE MINIMUM WEIGHTED FRACTION OF THE SUM TOTAL OF WEIGHTS OF ALL THE INPUT SAMPLES REQUIRED TO BE AT A LEAF NODE SAMPLES HAVE

EQUAL WEIGHT WHEN SAMPLEWEIGHT IS NOT PROVIDED

MAXFEATURES INT FLOAT STRING OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF FEATURES TO

CONSIDER WHEN LOOKING FOR THE BEST SPLIT

- IF INT THEN CONSIDER MAXFEATURES FEATURES AT EACH SPLIT
- IF FLOAT THEN MAXFEATURES IS A FRACTION AND INTMAXFEATURES

NFEATURES FEATURES ARE CONSIDERED AT EACH SPLIT

• IF "AUTO" THEN MAXFEATURESSQRTNFEATURES

• IF "SQRT" THEN MAXFEATURESSQRTNFEATURES

• IF "LOG2" THEN MAXFEATURESLOG2NFEATURES

• IF NONE THEN MAXFEATURESNFEATURES

NOTE THE SEARCH FOR A SPLIT DOES NOT STOP UNTIL AT LEAST ONE VALID PARTITION OF THE NODE SAMPLES

IS FOUND EVEN IF IT REQUIRES TO EFFECTIVELY INSPECT MORE THAN MAXFEATURES FEATURES

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RAN

DOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE

RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS

THE RANDOMSTATE INSTANCE USED BY NRANDOM

MAXLEAFNODES INT OR NONE OPTIONAL DEFAULTNONE GROW A TREE WITH MAXLEAFNODES

IN BESTFIRST FASHION BEST NODES ARE DEFINED AS RELATIVE REDUCTION IN IMPURITY IF NONE THEN

UNLIMITED NUMBER OF LEAF NODES

MINIMPURITYDECREASE FLOAT OPTIONAL DEFAULT0 A NODE WILL BE SPLIT IF THIS SPLIT INDUCES

A DECREASE OF THE IMPURITY GREATER THAN OR EQUAL TO THIS VALUE

THE WEIGHTED IMPURITY DECREASE EQUATION IS THE FOLLOWING

NT N IMPURITY NTR NT RIGHTIMPURITY

NTL NT LEFTIMPURITY

2332 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

WHERE N IS THE TOTAL NUMBER OF SAMPLES NTL IS THE NUMBER OF SAMPLES AT THE CURRENT NODE  
NTL IS THE NUMBER OF SAMPLES IN THE LEFT CHILD AND NTR IS THE NUMBER OF SAMPLES IN  
THE RIGHT CHILD

NNTNTR AND NTL ALL REFER TO THE WEIGHTED SUM IF SAMPLEWEIGHT IS PASSED  
NEW IN VERSION 019

MINIMPURITYSPLIT FLOAT DEFAULT 1E7 THRESHOLD FOR EARLY STOPPING IN TREE GROWTH A NODE  
WILL SPLIT IF ITS IMPURITY IS ABOVE THE THRESHOLD OTHERWISE IT IS A LEAF  
DEPRECATED SINCE VERSION 019 MINIMPURITYSPLIT HAS BEEN DEPRECATED IN FAVOR OF  
MINIMPURITYDECREASE IN 019 THE DEFAULT VALUE OF MINIMPURITYSPLIT  
WILL CHANGE FROM 1E7 TO 0 IN 023 AND IT WILL BE REMOVED IN 025 USE

MINIMPURITYDECREASE INSTEAD

CLASSWEIGHT DICT LIST OF DICTS "BALANCED" OR NONE DEFAULT NONE WEIGHTS ASSOCIATED WITH  
CLASSES IN THE FORM CLASSLABEL WEIGHT IF NOT GIVEN ALL CLASSES ARE SUPPOSED  
TO HAVE WEIGHT ONE FOR MULTIOUTPUT PROBLEMS A LIST OF DICTS CAN BE PROVIDED IN THE SAME  
ORDER AS THE COLUMNS OF Y

NOTE THAT FOR MULTIOUTPUT INCLUDING MULTILABEL WEIGHTS SHOULD BE DEFINED FOR EACH CLASS OF  
EVERY COLUMN IN ITS OWN DICT FOR EXAMPLE FOR FOURCLASS MULTILABEL CLASSIFICATION WEIGHTS  
SHOULD BE 0 1 1 1 0 1 1 5 0 1 1 1 0 1 1 1 INSTEAD OF 11 25

31 41

THE "BALANCED" MODE USES THE VALUES OF Y TO AUTOMATICALLY ADJUST WEIGHTS INVERSELY PRO  
PORTIONAL TO CLASS FREQUENCIES IN THE INPUT DATA AS NSAMPLES NCLASSES NP

BINCOUNTY  
FOR MULTIOUTPUT THE WEIGHTS OF EACH COLUMN OF Y WILL BE MULTIPLIED

NOTE THAT THESE WEIGHTS WILL BE MULTIPLIED WITH SAMPLEWEIGHT PASSED THROUGH THE FIT  
METHOD IF SAMPLEWEIGHT IS SPECIFIED

PRESORT BOOL OPTIONAL DEFAULT FALSE WHETHER TO PRESORT THE DATA TO SPEED UP THE FINDING OF  
BEST SPLITS IN FITTING FOR THE DEFAULT SETTINGS OF A DECISION TREE ON LARGE DATASETS SETTING  
THIS TO TRUE MAY SLOW DOWN THE TRAINING PROCESS WHEN USING EITHER A SMALLER DATASET OR A  
RESTRICTED DEPTH THIS MAY SPEED UP THE TRAINING

ATTRIBUTES

CLASSES ARRAY OF SHAPE NCLASSES OR A LIST OF SUCH ARRAYS THE CLASSES LABELS SINGLE OUTPUT  
PROBLEM OR A LIST OF ARRAYS OF CLASS LABELS MULTIOUTPUT PROBLEM

FEATUREIMPORTANCES ARRAY OF SHAPE NFEATURES RETURN THE FEATURE IMPORTANCES

MAXFEATURES INT THE INFERRED VALUE OF MAXFEATURES

NCLASSES INT OR LIST THE NUMBER OF CLASSES FOR SINGLE OUTPUT PROBLEMS OR A LIST CONTAINING  
THE NUMBER OF CLASSES FOR EACH OUTPUT FOR MULTIOUTPUT PROBLEMS

NFEATURES INT THE NUMBER OF FEATURES WHEN FIT IS PERFORMED

NOUTPUTS INT THE NUMBER OF OUTPUTS WHEN FIT IS PERFORMED

TREE TREE OBJECT THE UNDERLYING TREE OBJECT PLEASE REFER TO HELPSKLEARN TREE

TREETREE FOR ATTRIBUTES OF TREE OBJECT AND UNDERSTANDING THE DECISION TREE STRUCTURE  
FOR BASIC USAGE OF THESE ATTRIBUTES

SEE ALSO

DECISIONTREEREgressor

638SKLEARN TREE DECISION TREES 2333

SCIKITLEARN USER GUIDE RELEASE 0213

NOTES

THE DEFAULT VALUES FOR THE PARAMETERS CONTROLLING THE SIZE OF THE TREES EG MAXDEPTH MINSAMPLESLEAF ETC LEAD TO FULLY GROWN AND UNPRUNED TREES WHICH CAN POTENTIALLY BE VERY LARGE ON SOME DATA SETS TO REDUCE MEMORY CONSUMPTION THE COMPLEXITY AND SIZE OF THE TREES SHOULD BE CONTROLLED BY SETTING THOSE PARAMETER VALUES

THE FEATURES ARE ALWAYS RANDOMLY PERMUTED AT EACH SPLIT THEREFORE THE BEST FOUND SPLIT MAY VARY EVEN WITH THE SAME TRAINING DATA AND MAXFEATURESNFEATURES IF THE IMPROVEMENT OF THE CRITERION IS IDENTICAL FOR SEVERAL SPLITS ENUMERATED DURING THE SEARCH OF THE BEST SPLIT TO OBTAIN A DETERMINISTIC BEHAVIOUR DURING FITTING RANDOMSTATE HAS TO BE FIXED

REFERENCES

RB1EC977CD3071 RB1EC977CD3072 RB1EC977CD3073 RB1EC977CD3074

EXAMPLES

```
FROM SKLEARNDATASETS IMPORT LOADIRIS
FROM SKLEARNMODELSELECTION IMPORT CROSSVALSCORE
FROM SKLEARNTREE IMPORT DECISIONTREECLASSIFIER
CLF DECISIONTREECLASSIFIERRANDOMSTATE0
IRIS LOADIRIS
CROSSVALSCORECLF IRISDATA IRISTARGET CV10
```

ARRAY 1 093 086 093 093

093 093 1 093 1

METHODS

APPLY SELF X CHECKINPUT RETURNS THE INDEX OF THE LEAF THAT EACH SAMPLE IS PRE  
DICTED AS

DECISIONPATH SELF X CHECKINPUT RETURN THE DECISION PATH IN THE TREE

FITSELF X Y SAMPLEWEIGHT BUILD A DECISION TREE CLASSIFIER FROM THE TRAINING SET X  
Y

GETDEPTH SELF RETURNS THE DEPTH OF THE DECISION TREE

GETNLEAVES SELF RETURNS THE NUMBER OF LEAVES OF THE DECISION TREE

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PREDICT SELF X CHECKINPUT PREDICT CLASS OR REGRESSION VALUE FOR X

PREDICTLOGPROBA SELF X PREDICT CLASS LOGPROBABILITIES OF THE INPUT SAMPLES X

PREDICTPROBA SELF X CHECKINPUT PREDICT CLASS PROBABILITIES OF THE INPUT SAMPLES X

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND  
LABELS

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELF CRITERION'GINI' SPLITTER'BEST' MAXDEPTHNONE MINSAMPLESSPLIT2

MINSAMPLESLEAF1 MINWEIGHTFRACTIONLEAF00 MAXFEATURESNONE

RANDOMSTATENONE MAXLEAFNODESNONE MINIMPURITYDECREASE00

MINIMPURITYSPLITNONE CLASSWEIGHTNONE PRESORTFALSE

2334 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

APPLYSELFXCHECKINPUTTRUE

RETURNS THE INDEX OF THE LEAF THAT EACH SAMPLE IS PREDICTED AS  
NEW IN VERSION 017

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTER  
NALLY IT WILL BE CONVERTED TO DTYPENPFLOAT32 AND IF A SPARSE MATRIX IS PROVIDED TO  
A SPARSECSRMATRIX

CHECKINPUT BOOLEAN DEFAULTTRUE ALLOW TO BYPASS SEVERAL INPUT CHECKING DON'T USE  
THIS PARAMETER UNLESS YOU KNOW WHAT YOU DO

RETURNS

XLEAVES ARRAYLIKE SHAPE NSAMPLES FOR EACH DATAPOINT X IN X RETURN THE INDEX OF  
THE LEAF X ENDS UP IN LEAVES ARE NUMBERED WITHIN 0 SELFREENODECOUNT

POSSIBLY WITH GAPS IN THE NUMBERING

DECISIONPATH SELFXCHECKINPUTTRUE

RETURN THE DECISION PATH IN THE TREE

NEW IN VERSION 018

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTER  
NALLY IT WILL BE CONVERTED TO DTYPENPFLOAT32 AND IF A SPARSE MATRIX IS PROVIDED TO  
A SPARSECSRMATRIX

CHECKINPUT BOOLEAN DEFAULTTRUE ALLOW TO BYPASS SEVERAL INPUT CHECKING DON'T USE  
THIS PARAMETER UNLESS YOU KNOW WHAT YOU DO

RETURNS

INDICATOR SPARSE CSR ARRAY SHAPE NSAMPLES NNODES RETURN A NODE INDICATOR MATRIX  
WHERE NON ZERO ELEMENTS INDICATES THAT THE SAMPLES GOES THROUGH THE NODES

FEATUREIMPORTANCES

RETURN THE FEATURE IMPORTANCES

THE IMPORTANCE OF A FEATURE IS COMPUTED AS THE NORMALIZED TOTAL REDUCTION OF THE CRITERION BROUGHT BY THAT  
FEATURE IT IS ALSO KNOWN AS THE GINI IMPORTANCE

RETURNS

FEATUREIMPORTANCES ARRAY SHAPE NFEATURES

FITSELFXYSAMPLEWEIGHTNONE CHECKINPUTTRUE XIDXSORTEDNONE

BUILD A DECISION TREE CLASSIFIER FROM THE TRAINING SET X Y

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE TRAINING INPUT SAMPLES  
INTERNALLY IT WILL BE CONVERTED TO DTYPENPFLOAT32 AND IF A SPARSE MATRIX IS PROVIDED  
TO A SPARSESCMATRIX

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS THE TARGET VALUES CLASS LA  
BELS AS INTEGERS OR STRINGS

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OR NONE SAMPLE WEIGHTS IF NONE THEN

SAMPLES ARE EQUALLY WEIGHTED SPLITS THAT WOULD CREATE CHILD NODES WITH NET ZERO OR NEGA  
TIVE WEIGHT ARE IGNORED WHILE SEARCHING FOR A SPLIT IN EACH NODE SPLITS ARE ALSO IGNORED IF  
THEY WOULD RESULT IN ANY SINGLE CLASS CARRYING A NEGATIVE WEIGHT IN EITHER CHILD NODE

638SKLEARN TREE DECISION TREES 2335

SCIKITLEARN USER GUIDE RELEASE 0213

CHECKINPUT BOOLEAN DEFAULTTRUE ALLOW TO BYPASS SEVERAL INPUT CHECKING DON'T USE THIS PARAMETER UNLESS YOU KNOW WHAT YOU DO

XIDXSORTED ARRAYLIKE SHAPE NSAMPLES NFEATURES OPTIONAL THE INDEXES OF THE SORTED TRAINING INPUT SAMPLES IF MANY TREE ARE GROWN ON THE SAME DATASET THIS ALLOWS THE ORDERING TO BE CACHED BETWEEN TREES IF NONE THE DATA WILL BE SORTED HERE DON'T USE THIS PARAMETER UNLESS YOU KNOW WHAT TO DO

RETURNS  
SELF OBJECT

GETDEPTH SELF  
RETURNS THE DEPTH OF THE DECISION TREE  
THE DEPTH OF A TREE IS THE MAXIMUM DISTANCE BETWEEN THE ROOT AND ANY LEAF

GETNLEAVES SELF  
RETURNS THE NUMBER OF LEAVES OF THE DECISION TREE

GETPARAMS SELFDEEPTREE  
GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS  
DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS  
PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELFXCHECKINPUTTRUE  
PREDICT CLASS OR REGRESSION VALUE FOR X  
FOR A CLASSIFICATION MODEL THE PREDICTED CLASS FOR EACH SAMPLE IN X IS RETURNED FOR A REGRESSION MODEL THE PREDICTED VALUE BASED ON X IS RETURNED

PARAMETERS  
XARRAYLIKE OR SPARSE MATRIX OF SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTERNALLY IT WILL BE CONVERTED TO DTYPENPFLOAT32 AND IF A SPARSE MATRIX IS PROVIDED TO A SPARSECSRMATRIX

CHECKINPUT BOOLEAN DEFAULTTRUE ALLOW TO BYPASS SEVERAL INPUT CHECKING DON'T USE THIS PARAMETER UNLESS YOU KNOW WHAT YOU DO

RETURNS  
YARRAY OF SHAPE NSAMPLES OR NSAMPLES NOUTPUTS THE PREDICTED CLASSES OR THE PREDICT VALUES

PREDICTLOGPROBA SELF  
PREDICT CLASS LOGPROBABILITIES OF THE INPUT SAMPLES X

PARAMETERS  
XARRAYLIKE OR SPARSE MATRIX OF SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTERNALLY IT WILL BE CONVERTED TO DTYPENPFLOAT32 AND IF A SPARSE MATRIX IS PROVIDED TO A SPARSECSRMATRIX

RETURNS

2336 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

PARRAY OF SHAPE NSAMPLES NCLASSES OR A LIST OF NOUTPUTS SUCH ARRAYS IF NOUTPUTS  
1 THE CLASS LOGPROBABILITIES OF THE INPUT SAMPLES THE ORDER OF THE CLASSES CORRESPONDS  
TO THAT IN THE ATTRIBUTE CLASSES

PREDICTPROBA SELFXCHECKINPUTTRUE  
PREDICT CLASS PROBABILITIES OF THE INPUT SAMPLES X  
THE PREDICTED CLASS PROBABILITY IS THE FRACTION OF SAMPLES OF THE SAME CLASS IN A LEAF  
CHECKINPUT BOOLEAN DEFAULTTRUE ALLOW TO BYPASS SEVERAL INPUT CHECKING DON'T USE THIS PARAMETER  
UNLESS YOU KNOW WHAT YOU DO

PARAMETERS  
XARRAYLIKE OR SPARSE MATRIX OF SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTER  
NALLY IT WILL BE CONVERTED TO DTYPENPFLOAT32 AND IF A SPARSE MATRIX IS PROVIDED TO  
A SPARSECSRMATRIX  
CHECKINPUT BOOL RUN CHECKARRAY ON X

RETURNS  
PARRAY OF SHAPE NSAMPLES NCLASSES OR A LIST OF NOUTPUTS SUCH ARRAYS IF NOUTPUTS  
1 THE CLASS PROBABILITIES OF THE INPUT SAMPLES THE ORDER OF THE CLASSES CORRESPONDS TO  
THAT IN THE ATTRIBUTE CLASSES

SCORESELFXYSAMPLEWEIGHTNONE  
RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS  
IN MULTILABEL CLASSIFICATION THIS IS THE SUBSET ACCURACY WHICH IS A HARSH METRIC SINCE YOU REQUIRE FOR EACH  
SAMPLE THAT EACH LABEL SET BE CORRECTLY PREDICTED

PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES  
YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE LABELS FOR X  
SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS  
SCORE FLOAT MEAN ACCURACY OF SELFpredictX WRT Y

SETPARAMS SELFPARAMS  
SET THE PARAMETERS OF THIS ESTIMATOR  
THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE  
PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT  
OF A NESTED OBJECT

RETURNS  
SELF

EXAMPLES USING SKLEARNTREEDECISIONTREECLASSIFIER

- CLASSIFIER COMPARISON
- PLOT THE DECISION BOUNDARIES OF A VOTINGCLASSIFIER
- TWOCLASS ADABOOST
- MULTICLASS ADABOOSTED DECISION TREES

638SKLEARNTREE DECISION TREES 2337

SCIKITLEARN USER GUIDE RELEASE 0213

- DISCRETE VERSUS REAL ADABOOST
- PLOT THE DECISION SURFACES OF ENSEMBLES OF TREES ON THE IRIS DATASET
- DEMONSTRATION OF MULTIMETRIC EVALUATION ON CROSSVALSCORE AND GRIDSEARCHCV
- PLOT THE DECISION SURFACE OF A DECISION TREE ON THE IRIS DATASET
- UNDERSTANDING THE DECISION TREE STRUCTURE

6382SKLEARN TREE DECISION TREEREgressor  
CLASS SKLEARN TREE DECISION TREEREgressor CRITERION 'MSE' SPLITTER 'BEST' MAXDEPTH NONE

MINSAMPLESSPLIT 2 MINSAMPLESLEAF 1  
MINWEIGHTFRACTIONLEAF 0 MAXFEATURES NONE  
RANDOM STATE NONE MAXLEAFNODES NONE

MINIMPURITY DECREASE 0  
MINIMPURITYSPLIT NONE PRESORT FALSE

A DECISION TREE REGRESSOR  
READ MORE IN THE USER GUIDE  
PARAMETERS

CRITERION STRING OPTIONAL DEFAULT "MSE" THE FUNCTION TO MEASURE THE QUALITY OF A SPLIT SUPPORTED CRITERIA ARE "MSE" FOR THE MEAN SQUARED ERROR WHICH IS EQUAL TO VARIANCE REDUCTION AS FEATURE SELECTION CRITERION AND MINIMIZES THE L2 LOSS USING THE MEAN OF EACH TERMINAL NODE "FRIEDMAN MSE" WHICH USES MEAN SQUARED ERROR WITH FRIEDMAN'S IMPROVEMENT SCORE FOR POTENTIAL SPLITS AND "MAE" FOR THE MEAN ABSOLUTE ERROR WHICH MINIMIZES THE L1 LOSS USING THE MEDIAN OF EACH TERMINAL NODE

NEW IN VERSION 0.18 MEAN ABSOLUTE ERROR MAE CRITERION  
SPLITTER STRING OPTIONAL DEFAULT "BEST" THE STRATEGY USED TO CHOOSE THE SPLIT AT EACH NODE SUPPORTED STRATEGIES ARE "BEST" TO CHOOSE THE BEST SPLIT AND "RANDOM" TO CHOOSE THE BEST RANDOM SPLIT

MAXDEPTH INT OR NONE OPTIONAL DEFAULT NONE THE MAXIMUM DEPTH OF THE TREE IF NONE THEN NODES ARE EXPANDED UNTIL ALL LEAVES ARE PURE OR UNTIL ALL LEAVES CONTAIN LESS THAN MINSAMPLESSPLIT SAMPLES

MINSAMPLESSPLIT INT FLOAT OPTIONAL DEFAULT 2 THE MINIMUM NUMBER OF SAMPLES REQUIRED TO SPLIT AN INTERNAL NODE

- IF INT THEN CONSIDER MINSAMPLESSPLIT AS THE MINIMUM NUMBER
- IF FLOAT THEN MINSAMPLESSPLIT IS A FRACTION AND CEIL(MINSAMPLESSPLIT

NSAMPLES ARE THE MINIMUM NUMBER OF SAMPLES FOR EACH SPLIT

CHANGED IN VERSION 0.18 ADDED FLOAT VALUES FOR FRACTIONS

MINSAMPLESLEAF INT FLOAT OPTIONAL DEFAULT 1 THE MINIMUM NUMBER OF SAMPLES REQUIRED TO BE AT A LEAF NODE A SPLIT POINT AT ANY DEPTH WILL ONLY BE CONSIDERED IF IT LEAVES AT LEAST MINSAMPLESLEAF TRAINING SAMPLES IN EACH OF THE LEFT AND RIGHT BRANCHES THIS MAY HAVE THE EFFECT OF SMOOTHING THE MODEL ESPECIALLY IN REGRESSION

- IF INT THEN CONSIDER MINSAMPLESLEAF AS THE MINIMUM NUMBER
- IF FLOAT THEN MINSAMPLESLEAF IS A FRACTION AND CEIL(MINSAMPLESLEAF

NSAMPLES ARE THE MINIMUM NUMBER OF SAMPLES FOR EACH NODE

CHANGED IN VERSION 0.18 ADDED FLOAT VALUES FOR FRACTIONS

2338 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

MINWEIGHTFRACTIONLEAF FLOAT OPTIONAL DEFAULT0 THE MINIMUM WEIGHTED FRACTION OF THE SUM TOTAL OF WEIGHTS OF ALL THE INPUT SAMPLES REQUIRED TO BE AT A LEAF NODE SAMPLES HAVE EQUAL WEIGHT WHEN SAMPLEWEIGHT IS NOT PROVIDED

MAXFEATURES INT FLOAT STRING OR NONE OPTIONAL DEFAULTNONE THE NUMBER OF FEATURES TO CONSIDER WHEN LOOKING FOR THE BEST SPLIT

- IF INT THEN CONSIDER MAXFEATURES FEATURES AT EACH SPLIT
- IF FLOAT THEN MAXFEATURES IS A FRACTION AND INTMAXFEATURES NFEATURES FEATURES ARE CONSIDERED AT EACH SPLIT
- IF “AUTO” THEN MAXFEATURESNFEATURES
- IF “SQRT” THEN MAXFEATURESSQRTNFEATURES
- IF “LOG2” THEN MAXFEATURESLOG2NFEATURES
- IF NONE THEN MAXFEATURESNFEATURES

NOTE THE SEARCH FOR A SPLIT DOES NOT STOP UNTIL AT LEAST ONE VALID PARTITION OF THE NODE SAMPLES IS FOUND EVEN IF IT REQUIRES TO EFFECTIVELY INSPECT MORE THAN MAXFEATURES FEATURES

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RAN DOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NPRANDOM

MAXLEAFNODES INT OR NONE OPTIONAL DEFAULTNONE GROW A TREE WITH MAXLEAFNODES IN BESTFIRST FASHION BEST NODES ARE DEFINED AS RELATIVE REDUCTION IN IMPURITY IF NONE THEN UNLIMITED NUMBER OF LEAF NODES

MINIMPURITYDECREASE FLOAT OPTIONAL DEFAULT0 A NODE WILL BE SPLIT IF THIS SPLIT INDUCES A DECREASE OF THE IMPURITY GREATER THAN OR EQUAL TO THIS VALUE

THE WEIGHTED IMPURITY DECREASE EQUATION IS THE FOLLOWING

NT N IMPURITY NTR NT RIGHTIMPURITY

NTL NT LEFTIMPURITY

WHERE N IS THE TOTAL NUMBER OF SAMPLES NT IS THE NUMBER OF SAMPLES AT THE CURRENT NODE NTL IS THE NUMBER OF SAMPLES IN THE LEFT CHILD AND NTR IS THE NUMBER OF SAMPLES IN THE RIGHT CHILD

NNTNTR ANDNTL ALL REFER TO THE WEIGHTED SUM IF SAMPLEWEIGHT IS PASSED

NEW IN VERSION 019

MINIMPURITYSPLIT FLOAT DEFAULT1E7 THRESHOLD FOR EARLY STOPPING IN TREE GROWTH A NODE WILL SPLIT IF ITS IMPURITY IS ABOVE THE THRESHOLD OTHERWISE IT IS A LEAF

DEPRECATED SINCE VERSION 019 MINIMPURITYSPLIT HAS BEEN DEPRECATED IN FAVOR OF MINIMPURITYDECREASE IN 019 THE DEFAULT VALUE OF MINIMPURITYSPLIT WILL CHANGE FROM 1E7 TO 0 IN 023 AND IT WILL BE REMOVED IN 025 USE MINIMPURITYDECREASE INSTEAD

PRESORT BOOL OPTIONAL DEFAULTFALSE WHETHER TO PRESORT THE DATA TO SPEED UP THE FINDING OF BEST SPLITS IN FITTING FOR THE DEFAULT SETTINGS OF A DECISION TREE ON LARGE DATASETS SETTING THIS TO TRUE MAY SLOW DOWN THE TRAINING PROCESS WHEN USING EITHER A SMALLER DATASET OR A RESTRICTED DEPTH THIS MAY SPEED UP THE TRAINING

ATTRIBUTES

FEATUREIMPORTANCES ARRAY OF SHAPE NFEATURES RETURN THE FEATURE IMPORTANCES

638SKLEARN TREE DECISION TREES 2339

SCIKITLEARN USER GUIDE RELEASE 0213  
MAXFEATURES INT THE INFERRED VALUE OF MAXFEATURES  
NFEATURES INT THE NUMBER OF FEATURES WHEN FIT IS PERFORMED  
NOUTPUTS INT THE NUMBER OF OUTPUTS WHEN FIT IS PERFORMED  
TREE TREE OBJECT THE UNDERLYING TREE OBJECT PLEASE REFER TO HELPSKLEARNTREE  
TREETREE FOR ATTRIBUTES OF TREE OBJECT AND UNDERSTANDING THE DECISION TREE STRUCTURE  
FOR BASIC USAGE OF THESE ATTRIBUTES  
SEE ALSO  
DECISIONTREECLASSIFIER  
NOTES  
THE DEFAULT VALUES FOR THE PARAMETERS CONTROLLING THE SIZE OF THE TREES EG MAXDEPTH  
MINSAMPLESLEAF ETC LEAD TO FULLY GROWN AND UNPRUNED TREES WHICH CAN POTENTIALLY BE VERY LARGE ON  
SOME DATA SETS TO REDUCE MEMORY CONSUMPTION THE COMPLEXITY AND SIZE OF THE TREES SHOULD BE CONTROLLED BY  
SETTING THOSE PARAMETER VALUES  
THE FEATURES ARE ALWAYS RANDOMLY PERMUTED AT EACH SPLIT THEREFORE THE BEST FOUND SPLIT MAY VARY EVEN WITH  
THE SAME TRAINING DATA AND MAXFEATURESNFEATURES IF THE IMPROVEMENT OF THE CRITERION IS IDENTICAL FOR  
SEVERAL SPLITS ENUMERATED DURING THE SEARCH OF THE BEST SPLIT TO OBTAIN A DETERMINISTIC BEHAVIOUR DURING FITTING  
RANDOMSTATE HAS TO BE FIXED  
REFERENCES  
RA37B7E3ADB191 RA37B7E3ADB192 RA37B7E3ADB193 RA37B7E3ADB194  
EXAMPLES  
FROM SKLEARNDATASETS IMPORT LOADBOSTON  
FROM SKLEARNMODELSELECTION IMPORT CROSSVALSCORE  
FROM SKLEARNNTREE IMPORT DECISIONTREEREgressor  
BOSTON LOADBOSTON  
REGRESSOR DECISIONTREEREgressorRANDOMSTATE0  
CROSSVALSCOREREGRESSOR BOSTONDATA BOSTONTARGET CV10  
  
ARRAY 061 057 034 041 075  
007 029 033 142 177  
METHODS  
APPLY SELF X CHECKINPUT RETURNS THE INDEX OF THE LEAF THAT EACH SAMPLE IS PRE  
DICTED AS  
DECISIONPATH SELF X CHECKINPUT RETURN THE DECISION PATH IN THE TREE  
FITSELF X Y SAMPLEWEIGHT BUILD A DECISION TREE REGRESSOR FROM THE TRAINING SET X  
Y  
GETDEPTH SELF RETURNS THE DEPTH OF THE DECISION TREE  
CONTINUED ON NEXT PAGE  
2340 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

TABLE 6284 – CONTINUED FROM PREVIOUS PAGE

GETNLEAVES SELF RETURNS THE NUMBER OF LEAVES OF THE DECISION TREE

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PREDICT SELF X CHECKINPUT PREDICT CLASS OR REGRESSION VALUE FOR X

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELF CRITERION 'MSE' SPLITTER 'BEST' MAXDEPTH NONE MINSAMPLESSPLIT 2

MINSAMPLES LEAF 1 MINWEIGHT FRACTION LEAF 0 MAXFEATURES NONE

RANDOM STATE NONE MAXLEAF NODES NONE MINIMPURITY DECREASE 0

MINIMPURITY SPLIT NONE PRESORT FALSE

APPLY SELF X CHECK INPUT TRUE

RETURNS THE INDEX OF THE LEAF THAT EACH SAMPLE IS PREDICTED AS

NEW IN VERSION 017

PARAMETERS

X ARRAY LIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTERNALLY IT WILL BE CONVERTED TO DTYPENP FLOAT32 AND IF A SPARSE MATRIX IS PROVIDED TO A SPARSE CSR MATRIX

CHECK INPUT BOOLEAN DEFAULT TRUE ALLOW TO BYPASS SEVERAL INPUT CHECKING DON'T USE THIS PARAMETER UNLESS YOU KNOW WHAT YOU DO

RETURNS

X LEAVES ARRAY LIKE SHAPE NSAMPLES FOR EACH DATA POINT X IN X RETURN THE INDEX OF THE LEAF X ENDS UP IN LEAVES ARE NUMBERED WITHIN 0 SELF TREE NODE COUNT POSSIBLY WITH GAPS IN THE NUMBERING

DECISION PATH SELF X CHECK INPUT TRUE

RETURN THE DECISION PATH IN THE TREE

NEW IN VERSION 018

PARAMETERS

X ARRAY LIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTERNALLY IT WILL BE CONVERTED TO DTYPENP FLOAT32 AND IF A SPARSE MATRIX IS PROVIDED TO A SPARSE CSR MATRIX

CHECK INPUT BOOLEAN DEFAULT TRUE ALLOW TO BYPASS SEVERAL INPUT CHECKING DON'T USE THIS PARAMETER UNLESS YOU KNOW WHAT YOU DO

RETURNS

INDICATOR SPARSE CSR ARRAY SHAPE NSAMPLES NNODES RETURN A NODE INDICATOR MATRIX WHERE NON ZERO ELEMENTS INDICATES THAT THE SAMPLES GOES THROUGH THE NODES

FEATURE IMPORTANCES

RETURN THE FEATURE IMPORTANCES

THE IMPORTANCE OF A FEATURE IS COMPUTED AS THE NORMALIZED TOTAL REDUCTION OF THE CRITERION BROUGHT BY THAT FEATURE IT IS ALSO KNOWN AS THE GINI IMPORTANCE

RETURNS

FEATURE IMPORTANCES ARRAY SHAPE NFEATURES

638SKLEARN TREE DECISION TREES 2341

SCIKITLEARN USER GUIDE RELEASE 0213

FITSELFXYSAMPLEWEIGHTNONE CHECKINPUTTRUE XIDXSORTEDNONE

BUILD A DECISION TREE REGRESSOR FROM THE TRAINING SET X Y

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE TRAINING INPUT SAMPLES

INTERNALLY IT WILL BE CONVERTED TO DTYPENPFLOAT32 AND IF A SPARSE MATRIX IS PROVIDED

TO A SPARSECSMATRIX

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS THE TARGET VALUES REAL NUM

BERS USEDTPENPFLOAT64 ANDORDERC FOR MAXIMUM EFFICIENCY

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OR NONE SAMPLE WEIGHTS IF NONE THEN

SAMPLES ARE EQUALLY WEIGHTED SPLITS THAT WOULD CREATE CHILD NODES WITH NET ZERO OR NEGA

TIVE WEIGHT ARE IGNORED WHILE SEARCHING FOR A SPLIT IN EACH NODE

CHECKINPUT BOOLEAN DEFAULTTRUE ALLOW TO BYPASS SEVERAL INPUT CHECKING DON'T USE

THIS PARAMETER UNLESS YOU KNOW WHAT YOU DO

XIDXSORTED ARRAYLIKE SHAPE NSAMPLES NFEATURES OPTIONAL THE INDEXES OF THE

SORTED TRAINING INPUT SAMPLES IF MANY TREE ARE GROWN ON THE SAME DATASET THIS ALLOWS THE

ORDERING TO BE CACHED BETWEEN TREES IF NONE THE DATA WILL BE SORTED HERE DON'T USE THIS

PARAMETER UNLESS YOU KNOW WHAT TO DO

RETURNS

SELF OBJECT

GETDEPTH SELF

RETURNS THE DEPTH OF THE DECISION TREE

THE DEPTH OF A TREE IS THE MAXIMUM DISTANCE BETWEEN THE ROOT AND ANY LEAF

GETNLEAVES SELF

RETURNS THE NUMBER OF LEAVES OF THE DECISION TREE

GETPARAMS SELFDEEPTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELFXCHECKINPUTTRUE

PREDICT CLASS OR REGRESSION VALUE FOR X

FOR A CLASSIFICATION MODEL THE PREDICTED CLASS FOR EACH SAMPLE IN X IS RETURNED FOR A REGRESSION MODEL

THE PREDICTED VALUE BASED ON X IS RETURNED

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX OF SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTER

NALLY IT WILL BE CONVERTED TO DTYPENPFLOAT32 AND IF A SPARSE MATRIX IS PROVIDED TO

A SPARSECSRMATRIX

CHECKINPUT BOOLEAN DEFAULTTRUE ALLOW TO BYPASS SEVERAL INPUT CHECKING DON'T USE

THIS PARAMETER UNLESS YOU KNOW WHAT YOU DO

RETURNS

2342 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

`Y` ARRAY OF SHAPE `NSAMPLES` OR `NSAMPLES` NOUTPUTS THE PREDICTED CLASSES OR THE PREDICT VALUES

`SCORESELFXY` `SAMPLEWEIGHT` `NONE`

RETURNS THE COEFFICIENT OF DETERMINATION `R2` OF THE PREDICTION

THE COEFFICIENT `R2` IS DEFINED AS  $1 - \frac{U}{V}$  WHERE `U` IS THE RESIDUAL SUM OF SQUARES `YTRUE - YPRED` `2SUM` AND `V` IS THE TOTAL SUM OF SQUARES `YTRUE - YTRUEMEAN` `2SUM` THE BEST POSSIBLE SCORE IS 10 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS PREDICTS THE EXPECTED VALUE OF `Y` DISREGARDING THE INPUT FEATURES WOULD GET A `R2` SCORE OF 00

PARAMETERS

`X` ARRAYLIKE SHAPE `NSAMPLES` `NFEATURES` TEST SAMPLES FOR SOME ESTIMATORS THIS MAY BE A PRECOMPUTED KERNEL MATRIX INSTEAD SHAPE `NSAMPLES` `NSAMPLESFITTED` WHERE `NSAMPLESFITTED` IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR

`Y` ARRAYLIKE SHAPE `NSAMPLES` OR `NSAMPLES` NOUTPUTS TRUE VALUES FOR `X`

`SAMPLEWEIGHT` ARRAYLIKE SHAPE `NSAMPLES` OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT `R2` OF `SELF` `PREDICTX` WRT `Y`

NOTES

THE `R2` SCORE USED WHEN CALLING `SCORE` ON A REGRESSOR WILL USE `MULTIOUTPUTUNIFORMAVERAGE` FROM VERSION 023 TO KEEP CONSISTENT WITH `METRICSR2SCORE` THIS WILL INFLUENCE THE SCORE METHOD OF ALL THE `MULTIOUTPUT` REGRESSORS EXCEPT FOR `MULTIOUTPUTMULTIOUTPUTREGRESSOR` TO SPECIFY THE DEFAULT VALUE MANUALLY AND AVOID THE WARNING PLEASE EITHER CALL `METRICSR2SCORE` DIRECTLY OR MAKE A CUSTOM SCORER WITH `METRICSMAKESCORER` THE BUILTIN SCORER `R2` USES `MULTIOUTPUTUNIFORMAVERAGE`

`SETPARAMS` `SELF` `PARAMS`

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM `COMPONENTPARAMETER` SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

`SELF`

EXAMPLES USING `SKLEARNTREEDECISIONTREEREgressor`

- `DECISION TREE REGRESSION WITH ADABOOST`
- `SINGLE ESTIMATOR VERSUS BAGGING BIASVARIANCE DECOMPOSITION`
- `IMPUTING MISSING VALUES WITH VARIANTS OF ITERATIVEIMPUTER`
- `USING KBINSDISCRETIZER TO DISCRETIZE CONTINUOUS FEATURES`
- `DECISION TREE REGRESSION`
- `MULTIOUTPUT DECISION TREE REGRESSION`

638SKLEARNTREE DECISION TREES 2343

SCIKITLEARN USER GUIDE RELEASE 0213

6383SKLEARN TREE EXTRATREECLASSIFIER

CLASSSSKLEARN TREE EXTRATREECLASSIFIER CRITERION'GINI' SPLITTER'RANDOM'

MAXDEPTHNONE MINSAMPLESSPLIT2

MINSAMPLESLEAF1 MINWEIGHTFRACTIONLEAF00

MAXFEATURES'AUTO' RANDOMSTATENONE

MAXLEAFNODESNONE MINIMPURITYDECREASE00

MINIMPURITYSPLITNONE CLASSWEIGHTNONE

AN EXTREMELY RANDOMIZED TREE CLASSIFIER

EXTRATREES DIFFER FROM CLASSIC DECISION TREES IN THE WAY THEY ARE BUILT WHEN LOOKING FOR THE BEST SPLIT TO SEPARATE THE SAMPLES OF A NODE INTO TWO GROUPS RANDOM SPLITS ARE DRAWN FOR EACH OF THE MAXFEATURES RANDOMLY SELECTED FEATURES AND THE BEST SPLIT AMONG THOSE IS CHOSEN WHEN MAXFEATURES IS SET 1 THIS AMOUNTS TO BUILDING A TOTALLY RANDOM DECISION TREE

WARNING EXTRATREES SHOULD ONLY BE USED WITHIN ENSEMBLE METHODS

READ MORE IN THE USER GUIDE

PARAMETERS

CRITERION STRING OPTIONAL DEFAULT"GINI" THE FUNCTION TO MEASURE THE QUALITY OF A SPLIT SUPPORTED CRITERIA ARE "GINI" FOR THE GINI IMPURITY AND "ENTROPY" FOR THE INFORMATION GAIN

SPLITTER STRING OPTIONAL DEFAULT"RANDOM" THE STRATEGY USED TO CHOOSE THE SPLIT AT EACH NODE

SUPPORTED STRATEGIES ARE "BEST" TO CHOOSE THE BEST SPLIT AND "RANDOM" TO CHOOSE THE BEST RANDOM SPLIT

MAXDEPTH INT OR NONE OPTIONAL DEFAULTNONE THE MAXIMUM DEPTH OF THE TREE IF NONE THEN NODES ARE EXPANDED UNTIL ALL LEAVES ARE PURE OR UNTIL ALL LEAVES CONTAIN LESS THAN

MINSAMPLESSPLIT SAMPLES

MINSAMPLESSPLIT INT FLOAT OPTIONAL DEFAULT2 THE MINIMUM NUMBER OF SAMPLES REQUIRED TO SPLIT AN INTERNAL NODE

- IF INT THEN CONSIDER MINSAMPLESSPLIT AS THE MINIMUM NUMBER
- IF FLOAT THEN MINSAMPLESSPLIT IS A FRACTION AND CEILMINSAMPLESSPLIT

NSAMPLES ARE THE MINIMUM NUMBER OF SAMPLES FOR EACH SPLIT

CHANGED IN VERSION 018 ADDED FLOAT VALUES FOR FRACTIONS

MINSAMPLESLEAF INT FLOAT OPTIONAL DEFAULT1 THE MINIMUM NUMBER OF SAMPLES REQUIRED TO BE AT A LEAF NODE A SPLIT POINT AT ANY DEPTH WILL ONLY BE CONSIDERED IF IT LEAVES AT LEAST

MINSAMPLESLEAF TRAINING SAMPLES IN EACH OF THE LEFT AND RIGHT BRANCHES THIS MAY HAVE THE EFFECT OF SMOOTHING THE MODEL ESPECIALLY IN REGRESSION

- IF INT THEN CONSIDER MINSAMPLESLEAF AS THE MINIMUM NUMBER
- IF FLOAT THEN MINSAMPLESLEAF IS A FRACTION AND CEILMINSAMPLESLEAF

NSAMPLES ARE THE MINIMUM NUMBER OF SAMPLES FOR EACH NODE

CHANGED IN VERSION 018 ADDED FLOAT VALUES FOR FRACTIONS

MINWEIGHTFRACTIONLEAF FLOAT OPTIONAL DEFAULT0 THE MINIMUM WEIGHTED FRACTION OF THE SUM TOTAL OF WEIGHTS OF ALL THE INPUT SAMPLES REQUIRED TO BE AT A LEAF NODE SAMPLES HAVE EQUAL WEIGHT WHEN SAMPLEWEIGHT IS NOT PROVIDED

MAXFEATURES INT FLOAT STRING OR NONE OPTIONAL DEFAULT"AUTO" THE NUMBER OF FEATURES TO CONSIDER WHEN LOOKING FOR THE BEST SPLIT

- IF INT THEN CONSIDER MAXFEATURES FEATURES AT EACH SPLIT

2344 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

- IF FLOAT THEN MAXFEATURES IS A FRACTION AND INTMAXFEATURES NFEATURES FEATURES ARE CONSIDERED AT EACH SPLIT
- IF “AUTO” THEN MAXFEATURESSQRTNFEATURES
- IF “SQRT” THEN MAXFEATURESSQRTNFEATURES
- IF “LOG2” THEN MAXFEATURESLOG2NFEATURES
- IF NONE THEN MAXFEATURESNFEATURES

NOTE THE SEARCH FOR A SPLIT DOES NOT STOP UNTIL AT LEAST ONE VALID PARTITION OF THE NODE SAMPLES IS FOUND EVEN IF IT REQUIRES TO EFFECTIVELY INSPECT MORE THAN MAXFEATURES FEATURES

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RAN  
DOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE  
RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS  
THE RANDOMSTATE INSTANCE USED BY NPRANDOM

MAXLEAFNODES INT OR NONE OPTIONAL DEFAULTNONE GROW A TREE WITH MAXLEAFNODES  
IN BESTFIRST FASHION BEST NODES ARE DEFINED AS RELATIVE REDUCTION IN IMPURITY IF NONE THEN  
UNLIMITED NUMBER OF LEAF NODES

MINIMPURITYDECREASE FLOAT OPTIONAL DEFAULT0 A NODE WILL BE SPLIT IF THIS SPLIT INDUCES  
A DECREASE OF THE IMPURITY GREATER THAN OR EQUAL TO THIS VALUE  
THE WEIGHTED IMPURITY DECREASE EQUATION IS THE FOLLOWING

NT N IMPURITY NTR NT RIGHTIMPURITY  
NTL NT LEFTIMPURITY

WHEREIS THE TOTAL NUMBER OF SAMPLES NT IS THE NUMBER OF SAMPLES AT THE CURRENT NODE  
NTL IS THE NUMBER OF SAMPLES IN THE LEFT CHILD AND NTR IS THE NUMBER OF SAMPLES IN  
THE RIGHT CHILD

NNTNTR ANDNTL ALL REFER TO THE WEIGHTED SUM IF SAMPLEWEIGHT IS PASSED

NEW IN VERSION 019

MINIMPURITYSPLIT FLOAT DEFAULT1E7 THRESHOLD FOR EARLY STOPPING IN TREE GROWTH A NODE  
WILL SPLIT IF ITS IMPURITY IS ABOVE THE THRESHOLD OTHERWISE IT IS A LEAF

DEPRECATED SINCE VERSION 019 MINIMPURITYSPLIT HAS BEEN DEPRECATED IN FAVOR OF  
MINIMPURITYDECREASE IN 019 THE DEFAULT VALUE OF MINIMPURITYSPLIT

WILL CHANGE FROM 1E7 TO 0 IN 023 AND IT WILL BE REMOVED IN 025 USE

MINIMPURITYDECREASE INSTEAD

CLASSWEIGHT DICT LIST OF DICTS “BALANCED” OR NONE DEFAULTNONE WEIGHTS ASSOCIATED WITH  
CLASSES IN THE FORM CLASSLABEL WEIGHT IF NOT GIVEN ALL CLASSES ARE SUPPOSED  
TO HAVE WEIGHT ONE FOR MULTIOUTPUT PROBLEMS A LIST OF DICTS CAN BE PROVIDED IN THE SAME  
ORDER AS THE COLUMNS OF Y

NOTE THAT FOR MULTIOUTPUT INCLUDING MULTILABEL WEIGHTS SHOULD BE DEFINED FOR EACH CLASS OF  
EVERY COLUMN IN ITS OWN DICT FOR EXAMPLE FOR FOURCLASS MULTILABEL CLASSIFICATION WEIGHTS  
SHOULD BE 0 1 1 1 0 1 1 5 0 1 1 1 0 1 1 1 INSTEAD OF 11 25

31 41

THE “BALANCED” MODE USES THE VALUES OF Y TO AUTOMATICALLY ADJUST WEIGHTS INVERSELY PRO  
PORTIONAL TO CLASS FREQUENCIES IN THE INPUT DATA AS NSAMPLES NCLASSES NP

BINCOUNTY

FOR MULTIOUTPUT THE WEIGHTS OF EACH COLUMN OF Y WILL BE MULTIPLIED

638SKLEARNTREE DECISION TREES 2345

SCIKITLEARN USER GUIDE RELEASE 0213

NOTE THAT THESE WEIGHTS WILL BE MULTIPLIED WITH SAMPLEWEIGHT PASSED THROUGH THE FIT METHOD IF SAMPLEWEIGHT IS SPECIFIED

ATTRIBUTES

FEATUREIMPORTANCES RETURN THE FEATURE IMPORTANCES

SEE ALSO

EXTRATREEREgressor SKLEARNENSEMBLEEXTRATREESCLASSIFIER

SKLEARNENSEMBLEEXTRATREESREGRESSOR

NOTES

THE DEFAULT VALUES FOR THE PARAMETERS CONTROLLING THE SIZE OF THE TREES EG MAXDEPTH MINSAMPLESLEAF ETC LEAD TO FULLY GROWN AND UNPRUNED TREES WHICH CAN POTENTIALLY BE VERY LARGE ON SOME DATA SETS TO REDUCE MEMORY CONSUMPTION THE COMPLEXITY AND SIZE OF THE TREES SHOULD BE CONTROLLED BY SETTING THOSE PARAMETER VALUES

REFERENCES

RDD99A0224C6E1

METHODS

APPLY SELF X CHECKINPUT RETURNS THE INDEX OF THE LEAF THAT EACH SAMPLE IS PREDICTED AS

DECISIONPATH SELF X CHECKINPUT RETURN THE DECISION PATH IN THE TREE

FITSELF X Y SAMPLEWEIGHT BUILD A DECISION TREE CLASSIFIER FROM THE TRAINING SET X Y

GETDEPTH SELF RETURNS THE DEPTH OF THE DECISION TREE

GETNLEAVES SELF RETURNS THE NUMBER OF LEAVES OF THE DECISION TREE

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PREDICT SELF X CHECKINPUT PREDICT CLASS OR REGRESSION VALUE FOR X

PREDICTLOGPROBA SELF X PREDICT CLASS LOGPROBABILITIES OF THE INPUT SAMPLES X

PREDICTPROBA SELF X CHECKINPUT PREDICT CLASS PROBABILITIES OF THE INPUT SAMPLES X

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELF CRITERION'GINI' SPLITTER'RANDOM' MAXDEPTHNONE MINSAMPLESSPLIT2

MINSAMPLESLEAF1 MINWEIGHTFRACTIONLEAF00 MAXFEATURES'AUTO'

RANDOMSTATENONE MAXLEAFNODESNONE MINIMPURITYDECREASE00

MINIMPURITYSPLITNONE CLASSWEIGHTNONE

APPLYSELFXCHECKINPUTTRUE

RETURNS THE INDEX OF THE LEAF THAT EACH SAMPLE IS PREDICTED AS

NEW IN VERSION 017

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTER

2346 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

NALLY IT WILL BE CONVERTED TO DTYPENPFLOAT32 AND IF A SPARSE MATRIX IS PROVIDED TO A SPARSECSRMATRIX

CHECKINPUT BOOLEAN DEFAULTTRUE ALLOW TO BYPASS SEVERAL INPUT CHECKING DON'T USE THIS PARAMETER UNLESS YOU KNOW WHAT YOU DO

RETURNS

XLEAVES ARRAYLIKE SHAPE NSAMPLES FOR EACH DATAPOINT X IN X RETURN THE INDEX OF THE LEAF X ENDS UP IN LEAVES ARE NUMBERED WITHIN 0 SELFREENODECOUNT POSSIBLY WITH GAPS IN THE NUMBERING

DECISIONPATH SELFXCHECKINPUTTRUE

RETURN THE DECISION PATH IN THE TREE

NEW IN VERSION 018

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTER NALLY IT WILL BE CONVERTED TO DTYPENPFLOAT32 AND IF A SPARSE MATRIX IS PROVIDED TO A SPARSECSRMATRIX

CHECKINPUT BOOLEAN DEFAULTTRUE ALLOW TO BYPASS SEVERAL INPUT CHECKING DON'T USE THIS PARAMETER UNLESS YOU KNOW WHAT YOU DO

RETURNS

INDICATOR SPARSE CSR ARRAY SHAPE NSAMPLES NNODES RETURN A NODE INDICATOR MATRIX WHERE NON ZERO ELEMENTS INDICATES THAT THE SAMPLES GOES THROUGH THE NODES

FEATUREIMPORTANCES

RETURN THE FEATURE IMPORTANCES

THE IMPORTANCE OF A FEATURE IS COMPUTED AS THE NORMALIZED TOTAL REDUCTION OF THE CRITERION BROUGHT BY THAT FEATURE IT IS ALSO KNOWN AS THE GINI IMPORTANCE

RETURNS

FEATUREIMPORTANCES ARRAY SHAPE NFEATURES

FITSELFXYSAMPLEWEIGHTNONE CHECKINPUTTRUE XIDXSORTEDNONE

BUILD A DECISION TREE CLASSIFIER FROM THE TRAINING SET X Y

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE TRAINING INPUT SAMPLES INTERNALLY IT WILL BE CONVERTED TO DTYPENPFLOAT32 AND IF A SPARSE MATRIX IS PROVIDED TO A SPARSESCMATRIX

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS THE TARGET VALUES CLASS LABELS AS INTEGERS OR STRINGS

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OR NONE SAMPLE WEIGHTS IF NONE THEN SAMPLES ARE EQUALLY WEIGHTED SPLITS THAT WOULD CREATE CHILD NODES WITH NET ZERO OR NEGATIVE WEIGHT ARE IGNORED WHILE SEARCHING FOR A SPLIT IN EACH NODE SPLITS ARE ALSO IGNORED IF THEY WOULD RESULT IN ANY SINGLE CLASS CARRYING A NEGATIVE WEIGHT IN EITHER CHILD NODE

CHECKINPUT BOOLEAN DEFAULTTRUE ALLOW TO BYPASS SEVERAL INPUT CHECKING DON'T USE THIS PARAMETER UNLESS YOU KNOW WHAT YOU DO

XIDXSORTED ARRAYLIKE SHAPE NSAMPLES NFEATURES OPTIONAL THE INDEXES OF THE SORTED TRAINING INPUT SAMPLES IF MANY TREE ARE GROWN ON THE SAME DATASET THIS ALLOWS THE

638SKLEARNTREE DECISION TREES 2347

SCIKITLEARN USER GUIDE RELEASE 0213

ORDERING TO BE CACHED BETWEEN TREES IF NONE THE DATA WILL BE SORTED HERE DON'T USE THIS  
PARAMETER UNLESS YOU KNOW WHAT TO DO

RETURNS  
SELF OBJECT

GETDEPTH SELF  
RETURNS THE DEPTH OF THE DECISION TREE  
THE DEPTH OF A TREE IS THE MAXIMUM DISTANCE BETWEEN THE ROOT AND ANY LEAF

GETNLEAVES SELF  
RETURNS THE NUMBER OF LEAVES OF THE DECISION TREE

GETPARAMS SELFDEEPTRUE  
GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS  
DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED  
SUBOBJECTS THAT ARE ESTIMATORS

RETURNS  
PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELFXCHECKINPUTTRUE  
PREDICT CLASS OR REGRESSION VALUE FOR X  
FOR A CLASSIFICATION MODEL THE PREDICTED CLASS FOR EACH SAMPLE IN X IS RETURNED FOR A REGRESSION MODEL  
THE PREDICTED VALUE BASED ON X IS RETURNED

PARAMETERS  
XARRAYLIKE OR SPARSE MATRIX OF SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTER  
NALLY IT WILL BE CONVERTED TO DTYPENPFLOAT32 AND IF A SPARSE MATRIX IS PROVIDED TO  
A SPARSECSRMATRIX

CHECKINPUT BOOLEAN DEFAULTTRUE ALLOW TO BYPASS SEVERAL INPUT CHECKING DON'T USE  
THIS PARAMETER UNLESS YOU KNOW WHAT YOU DO

RETURNS  
YARRAY OF SHAPE NSAMPLES OR NSAMPLES NOUTPUTS THE PREDICTED CLASSES OR THE  
PREDICT VALUES

PREDICTLOGPROBA SELF  
PREDICT CLASS LOGPROBABILITIES OF THE INPUT SAMPLES X

PARAMETERS  
XARRAYLIKE OR SPARSE MATRIX OF SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTER  
NALLY IT WILL BE CONVERTED TO DTYPENPFLOAT32 AND IF A SPARSE MATRIX IS PROVIDED TO  
A SPARSECSRMATRIX

RETURNS  
PARRAY OF SHAPE NSAMPLES NCLASSES OR A LIST OF NOUTPUTS SUCH ARRAYS IF NOUTPUTS  
1 THE CLASS LOGPROBABILITIES OF THE INPUT SAMPLES THE ORDER OF THE CLASSES CORRESPONDS  
TO THAT IN THE ATTRIBUTE CLASSES

PREDICTPROBA SELFXCHECKINPUTTRUE  
PREDICT CLASS PROBABILITIES OF THE INPUT SAMPLES X  
THE PREDICTED CLASS PROBABILITY IS THE FRACTION OF SAMPLES OF THE SAME CLASS IN A LEAF

2348 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

CHECKINPUT BOOLEAN DEFAULTTRUE ALLOW TO BYPASS SEVERAL INPUT CHECKING DON'T USE THIS PARAMETER UNLESS YOU KNOW WHAT YOU DO

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX OF SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTER NALLY IT WILL BE CONVERTED TO DTYPENPFLOAT32 AND IF A SPARSE MATRIX IS PROVIDED TO A SPARSECSRMATRIX

CHECKINPUT BOOL RUN CHECKARRAY ON X

RETURNS

PARRAY OF SHAPE NSAMPLES NCLASSES OR A LIST OF NOUTPUTS SUCH ARRAYS IF NOUTPUTS 1 THE CLASS PROBABILITIES OF THE INPUT SAMPLES THE ORDER OF THE CLASSES CORRESPONDS TO THAT IN THE ATTRIBUTE CLASSES

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE MEAN ACCURACY ON THE GIVEN TEST DATA AND LABELS

IN MULTILABEL CLASSIFICATION THIS IS THE SUBSET ACCURACY WHICH IS A HARSH METRIC SINCE YOU REQUIRE FOR EACH SAMPLE THAT EACH LABEL SET BE CORRECTLY PREDICTED

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE LABELS FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS

SCORE FLOAT MEAN ACCURACY OF SELFpredictX WRT Y

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

6384SKLEARN TREE EXTRATREEREgressor

CLASSSSKLEARN TREE EXTRATREEREgressor CRITERION'MSE' SPLITTER'RANDOM' MAXDEPTHNONE

MINSAMPLESSPLIT2 MINSAMPLESLEAF1

MINWEIGHTFRACTIONLEAF00 MAXFEATURES'AUTO'

RANDOMSTATENONE MINIMPURITYDECREASE00

MINIMPURITYSPLITNONE MAXLEAFNODESNONE

AN EXTREMELY RANDOMIZED TREE REGRESSOR

EXTRATREES DIFFER FROM CLASSIC DECISION TREES IN THE WAY THEY ARE BUILT WHEN LOOKING FOR THE BEST SPLIT TO SEPARATE THE SAMPLES OF A NODE INTO TWO GROUPS RANDOM SPLITS ARE DRAWN FOR EACH OF THE MAXFEATURES RANDOMLY SELECTED FEATURES AND THE BEST SPLIT AMONG THOSE IS CHOSEN WHEN MAXFEATURES IS SET 1 THIS AMOUNTS TO BUILDING A TOTALLY RANDOM DECISION TREE

WARNING EXTRATREES SHOULD ONLY BE USED WITHIN ENSEMBLE METHODS

638SKLEARN TREE DECISION TREES 2349

SCIKITLEARN USER GUIDE RELEASE 0213

READ MORE IN THE USER GUIDE

PARAMETERS

CRITERION STRING OPTIONAL DEFAULT"mse" THE FUNCTION TO MEASURE THE QUALITY OF A SPLIT SUPPORTED CRITERIA ARE "mse" FOR THE MEAN SQUARED ERROR WHICH IS EQUAL TO VARIANCE REDUCTION AS FEATURE SELECTION CRITERION AND "mae" FOR THE MEAN ABSOLUTE ERROR

NEW IN VERSION 018 MEAN ABSOLUTE ERROR MAE CRITERION

SPLITTER STRING OPTIONAL DEFAULT"RANDOM" THE STRATEGY USED TO CHOOSE THE SPLIT AT EACH NODE SUPPORTED STRATEGIES ARE "BEST" TO CHOOSE THE BEST SPLIT AND "RANDOM" TO CHOOSE THE BEST

RANDOM SPLIT

MAXDEPTH INT OR NONE OPTIONAL DEFAULTNONE THE MAXIMUM DEPTH OF THE TREE IF NONE THEN NODES ARE EXPANDED UNTIL ALL LEAVES ARE PURE OR UNTIL ALL LEAVES CONTAIN LESS THAN MINSAMPLESSPLIT SAMPLES

MINSAMPLESSPLIT INT FLOAT OPTIONAL DEFAULT2 THE MINIMUM NUMBER OF SAMPLES REQUIRED TO SPLIT AN INTERNAL NODE

- IF INT THEN CONSIDER MINSAMPLESSPLIT AS THE MINIMUM NUMBER
- IF FLOAT THEN MINSAMPLESSPLIT IS A FRACTION AND CEILMINSAMPLESSPLIT

NSAMPLES ARE THE MINIMUM NUMBER OF SAMPLES FOR EACH SPLIT CHANGED IN VERSION 018 ADDED FLOAT VALUES FOR FRACTIONS

MINSAMPLESLEAF INT FLOAT OPTIONAL DEFAULT1 THE MINIMUM NUMBER OF SAMPLES REQUIRED TO BE AT A LEAF NODE A SPLIT POINT AT ANY DEPTH WILL ONLY BE CONSIDERED IF IT LEAVES AT LEAST MINSAMPLESLEAF TRAINING SAMPLES IN EACH OF THE LEFT AND RIGHT BRANCHES THIS MAY HAVE THE EFFECT OF SMOOTHING THE MODEL ESPECIALLY IN REGRESSION

- IF INT THEN CONSIDER MINSAMPLESLEAF AS THE MINIMUM NUMBER
- IF FLOAT THEN MINSAMPLESLEAF IS A FRACTION AND CEILMINSAMPLESLEAF

NSAMPLES ARE THE MINIMUM NUMBER OF SAMPLES FOR EACH NODE CHANGED IN VERSION 018 ADDED FLOAT VALUES FOR FRACTIONS

MINWEIGHTFRACTIONLEAF FLOAT OPTIONAL DEFAULT0 THE MINIMUM WEIGHTED FRACTION OF THE SUM TOTAL OF WEIGHTS OF ALL THE INPUT SAMPLES REQUIRED TO BE AT A LEAF NODE SAMPLES HAVE EQUAL WEIGHT WHEN SAMPLEWEIGHT IS NOT PROVIDED

MAXFEATURES INT FLOAT STRING OR NONE OPTIONAL DEFAULT"AUTO" THE NUMBER OF FEATURES TO CONSIDER WHEN LOOKING FOR THE BEST SPLIT

- IF INT THEN CONSIDER MAXFEATURES FEATURES AT EACH SPLIT
- IF FLOAT THEN MAXFEATURES IS A FRACTION AND INTMAXFEATURES

NFEATURES FEATURES ARE CONSIDERED AT EACH SPLIT

- IF "AUTO" THEN MAXFEATURESNFEATURES
- IF "SQRT" THEN MAXFEATURESSQRTNFEATURES
- IF "LOG2" THEN MAXFEATURESLOG2NFEATURES
- IF NONE THEN MAXFEATURESNFEATURES

NOTE THE SEARCH FOR A SPLIT DOES NOT STOP UNTIL AT LEAST ONE VALID PARTITION OF THE NODE SAMPLES IS FOUND EVEN IF IT REQUIRES TO EFFECTIVELY INSPECT MORE THAN MAXFEATURES FEATURES

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RAN

DOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE 2350 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NPRANDOM

MINIMPURITYDECREASE FLOAT OPTIONAL DEFAULT0 A NODE WILL BE SPLIT IF THIS SPLIT INDUCES A DECREASE OF THE IMPURITY GREATER THAN OR EQUAL TO THIS VALUE

THE WEIGHTED IMPURITY DECREASE EQUATION IS THE FOLLOWING

NT N IMPURITY NTR NT RIGHTIMPURITY

NTL NT LEFTIMPURITY

WHERE NIS THE TOTAL NUMBER OF SAMPLES NT IS THE NUMBER OF SAMPLES AT THE CURRENT NODE

NTL IS THE NUMBER OF SAMPLES IN THE LEFT CHILD AND NTR IS THE NUMBER OF SAMPLES IN THE RIGHT CHILD

NNTNTR ANDNTL ALL REFER TO THE WEIGHTED SUM IF SAMPLEWEIGHT IS PASSED

NEW IN VERSION 019

MINIMPURITYSPLIT FLOAT DEFAULT1E7 THRESHOLD FOR EARLY STOPPING IN TREE GROWTH A NODE WILL SPLIT IF ITS IMPURITY IS ABOVE THE THRESHOLD OTHERWISE IT IS A LEAF

DEPRECATED SINCE VERSION 019 MINIMPURITYSPLIT HAS BEEN DEPRECATED IN FAVOR OF

MINIMPURITYDECREASE IN 019 THE DEFAULT VALUE OF MINIMPURITYSPLIT

WILL CHANGE FROM 1E7 TO 0 IN 023 AND IT WILL BE REMOVED IN 025 USE

MINIMPURITYDECREASE INSTEAD

MAXLEAFNODES INT OR NONE OPTIONAL DEFAULTNONE GROW A TREE WITH MAXLEAFNODES

IN BESTFIRST FASHION BEST NODES ARE DEFINED AS RELATIVE REDUCTION IN IMPURITY IF NONE THEN

UNLIMITED NUMBER OF LEAF NODES

ATTRIBUTES

FEATUREIMPORTANCES RETURN THE FEATURE IMPORTANCES

SEE ALSO

EXTRATREECLASSIFIER SKLEARNENSEMBLEEXTRATREESCLASSIFIER

SKLEARNENSEMBLEEXTRATREESREGRESSOR

NOTES

THE DEFAULT VALUES FOR THE PARAMETERS CONTROLLING THE SIZE OF THE TREES EG MAXDEPTH

MINSAMPLESLEAF ETC LEAD TO FULLY GROWN AND UNPRUNED TREES WHICH CAN POTENTIALLY BE VERY LARGE ON

SOME DATA SETS TO REDUCE MEMORY CONSUMPTION THE COMPLEXITY AND SIZE OF THE TREES SHOULD BE CONTROLLED BY

SETTING THOSE PARAMETER VALUES

REFERENCES

R4939D63D5A491

METHODS

638SKLEARNTREE DECISION TREES 2351

SCIKITLEARN USER GUIDE RELEASE 0213

APPLY SELF X CHECKINPUT RETURNS THE INDEX OF THE LEAF THAT EACH SAMPLE IS PREDICTED AS

DECISIONPATH SELF X CHECKINPUT RETURN THE DECISION PATH IN THE TREE

FITSELF X Y SAMPLEWEIGHT BUILD A DECISION TREE REGRESSOR FROM THE TRAINING SET X Y

GETDEPTH SELF RETURNS THE DEPTH OF THE DECISION TREE

GETNLEAVES SELF RETURNS THE NUMBER OF LEAVES OF THE DECISION TREE

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

PREDICT SELF X CHECKINPUT PREDICT CLASS OR REGRESSION VALUE FOR X

SCORE SELF X Y SAMPLEWEIGHT RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT SELF CRITERION 'MSE' SPLITTER 'RANDOM' MAXDEPTH NONE MINSAMPLESSPLIT 2 MINSAMPLESLEAF 1 MINWEIGHTFRACTIONLEAF 0 MAXFEATURES 'AUTO' RANDOMSTATENONE MINIMPURITYDECREASE 0 MINIMPURITYSPLIT NONE MAXLEAFNODES NONE

APPLYSELF X CHECKINPUT TRUE

RETURNS THE INDEX OF THE LEAF THAT EACH SAMPLE IS PREDICTED AS

NEW IN VERSION 017

PARAMETERS

X ARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTERNALLY IT WILL BE CONVERTED TO DTYPE NPFLOAT32 AND IF A SPARSE MATRIX IS PROVIDED TO A SPARSECSR MATRIX

CHECKINPUT BOOLEAN DEFAULT TRUE ALLOW TO BYPASS SEVERAL INPUT CHECKING DON'T USE THIS PARAMETER UNLESS YOU KNOW WHAT YOU DO

RETURNS

X LEAVES ARRAYLIKE SHAPE NSAMPLES FOR EACH DATAPOINT X IN X RETURN THE INDEX OF THE LEAF X ENDS UP IN LEAVES ARE NUMBERED WITHIN 0 SELF TREENODECOUNT POSSIBLY WITH GAPS IN THE NUMBERING

DECISIONPATH SELF X CHECKINPUT TRUE

RETURN THE DECISION PATH IN THE TREE

NEW IN VERSION 018

PARAMETERS

X ARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTERNALLY IT WILL BE CONVERTED TO DTYPE NPFLOAT32 AND IF A SPARSE MATRIX IS PROVIDED TO A SPARSECSR MATRIX

CHECKINPUT BOOLEAN DEFAULT TRUE ALLOW TO BYPASS SEVERAL INPUT CHECKING DON'T USE THIS PARAMETER UNLESS YOU KNOW WHAT YOU DO

RETURNS

INDICATOR SPARSE CSR ARRAY SHAPE NSAMPLES NNODES RETURN A NODE INDICATOR MATRIX WHERE NON ZERO ELEMENTS INDICATES THAT THE SAMPLES GOES THROUGH THE NODES

FEATURE IMPORTANCES

RETURN THE FEATURE IMPORTANCES

2352 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

THE IMPORTANCE OF A FEATURE IS COMPUTED AS THE NORMALIZED TOTAL REDUCTION OF THE CRITERION BROUGHT BY THAT FEATURE IT IS ALSO KNOWN AS THE GINI IMPORTANCE

RETURNS

FEATUREIMPORTANCES ARRAY SHAPE NFEATURES

FITSELFXYSAMPLEWEIGHTNONE CHECKINPUTTRUE XIDXSORTEDNONE

BUILD A DECISION TREE REGRESSOR FROM THE TRAINING SET X Y

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE TRAINING INPUT SAMPLES

INTERNALLY IT WILL BE CONVERTED TO DTYPENPFLOAT32 AND IF A SPARSE MATRIX IS PROVIDED TO A SPARSESCMATRIX

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS THE TARGET VALUES REAL NUMBERS USEDTPENPFLOAT64 ANDORDERC FOR MAXIMUM EFFICIENCY

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OR NONE SAMPLE WEIGHTS IF NONE THEN SAMPLES ARE EQUALLY WEIGHTED SPLITS THAT WOULD CREATE CHILD NODES WITH NET ZERO OR NEGATIVE WEIGHT ARE IGNORED WHILE SEARCHING FOR A SPLIT IN EACH NODE

CHECKINPUT BOOLEAN DEFAULTTRUE ALLOW TO BYPASS SEVERAL INPUT CHECKING DON'T USE THIS PARAMETER UNLESS YOU KNOW WHAT YOU DO

XIDXSORTED ARRAYLIKE SHAPE NSAMPLES NFEATURES OPTIONAL THE INDEXES OF THE SORTED TRAINING INPUT SAMPLES IF MANY TREE ARE GROWN ON THE SAME DATASET THIS ALLOWS THE ORDERING TO BE CACHED BETWEEN TREES IF NONE THE DATA WILL BE SORTED HERE DON'T USE THIS PARAMETER UNLESS YOU KNOW WHAT TO DO

RETURNS

SELF OBJECT

GETDEPTH SELF

RETURNS THE DEPTH OF THE DECISION TREE

THE DEPTH OF A TREE IS THE MAXIMUM DISTANCE BETWEEN THE ROOT AND ANY LEAF

GETNLEAVES SELF

RETURNS THE NUMBER OF LEAVES OF THE DECISION TREE

GETPARAMS SELFDEEPTRUE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

PREDICTSELFXCHECKINPUTTRUE

PREDICT CLASS OR REGRESSION VALUE FOR X

FOR A CLASSIFICATION MODEL THE PREDICTED CLASS FOR EACH SAMPLE IN X IS RETURNED FOR A REGRESSION MODEL THE PREDICTED VALUE BASED ON X IS RETURNED

PARAMETERS

638SKLEARNTREE DECISION TREES 2353

SCIKITLEARN USER GUIDE RELEASE 0213

XARRAYLIKE OR SPARSE MATRIX OF SHAPE NSAMPLES NFEATURES THE INPUT SAMPLES INTER  
NALLY IT WILL BE CONVERTED TO DTYPENPFLOAT32 AND IF A SPARSE MATRIX IS PROVIDED TO  
A SPARSECSRMATRIX

CHECKINPUT BOOLEAN DEFAULTTRUE ALLOW TO BYPASS SEVERAL INPUT CHECKING DON'T USE  
THIS PARAMETER UNLESS YOU KNOW WHAT YOU DO

RETURNS  
YARRAY OF SHAPE NSAMPLES OR NSAMPLES NOUTPUTS THE PREDICTED CLASSES OR THE  
PREDICT VALUES

SCORESELFXYSAMPLEWEIGHTNONE

RETURNS THE COEFFICIENT OF DETERMINATION R2 OF THE PREDICTION

THE COEFFICIENT R2 IS DEFINED AS  $1 - \frac{U}{V}$  WHERE U IS THE RESIDUAL SUM OF SQUARES YTRUE YPRED  
2SUM AND V IS THE TOTAL SUM OF SQUARES YTRUE YTRUEMEAN 2SUM THE BEST POSSIBLE SCORE  
IS 10 AND IT CAN BE NEGATIVE BECAUSE THE MODEL CAN BE ARBITRARILY WORSE A CONSTANT MODEL THAT ALWAYS  
PREDICTS THE EXPECTED VALUE OF Y DISREGARDING THE INPUT FEATURES WOULD GET A R2 SCORE OF 00

PARAMETERS  
XARRAYLIKE SHAPE NSAMPLES NFEATURES TEST SAMPLES FOR SOME ESTIMATORS THIS MAY  
BE A PRECOMPUTED KERNEL MATRIX INSTEAD SHAPE NSAMPLES NSAMPLESFITTED WHERE  
NSAMPLESFITTED IS THE NUMBER OF SAMPLES USED IN THE FITTING FOR THE ESTIMATOR

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS TRUE VALUES FOR X

SAMPLEWEIGHT ARRAYLIKE SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS

RETURNS  
SCORE FLOAT R2 OF SELFpredictX WRT Y

NOTES  
THE R2 SCORE USED WHEN CALLING SCORE ON A REGRESSOR WILL USE MULTIOUTPUTUNIFORMAVERAGE  
FROM VERSION 023 TO KEEP CONSISTENT WITH METRICSR2SCORE THIS WILL INFLUENCE THE SCORE  
METHOD OF ALL THE MULTIOUTPUT REGRESSORS EXCEPT FOR MULTIOUTPUTMULTIOUTPUTREGRESSOR  
TO SPECIFY THE DEFAULT VALUE MANUALLY AND AVOID THE WARNING PLEASE EITHER CALL METRICSR2SCORE  
DIRECTLY OR MAKE A CUSTOM SCORER WITH METRICSMakesCORER THE BUILTIN SCORER R2 USES  
MULTIOUTPUTUNIFORMAVERAGE

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE  
PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT  
OF A NESTED OBJECT

RETURNS  
SELF

TREEEXPORTGRAPHVIZ DECISIONTREE EXPORT A DECISION TREE IN DOT FORMAT

TREEPLOTtree DECISIONTREE MAXDEPTH PLOT A DECISION TREE

TREEEXPORTTEXT DECISIONTREE BUILD A TEXT REPORT SHOWING THE RULES OF A DECISION TREE

2354 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

6385SKLEARNTREE EXPORTGRAPHVIZ

SKLEARNTREE EXPORTGRAPHVIZ DECISIONTREE OUTFILENONE MAXDEPTHNONE FEATURENAMESNONE CLASSNAMESNONE LABEL'ALL' FILLEDFALSE

LEAVESPARALLELFALSE IMPURITYTRUE NODEIDSFALSE

PROPORTIONFALSE ROTATEFALSE ROUNDEDFALSE SPECIALCHARACTERSFALSE PRECISION3

EXPORT A DECISION TREE IN DOT FORMAT

THIS FUNCTION GENERATES A GRAPHVIZ REPRESENTATION OF THE DECISION TREE WHICH IS THEN WRITTEN INTO OUTFILE

ONCE EXPORTED GRAPHICAL RENDERINGS CAN BE GENERATED USING FOR EXAMPLE

DOT TPS TREEDOT O TREEPS POSTSCRIPT FORMAT

DOT TPNG TREEDOT O TREEPNG PNG FORMAT

THE SAMPLE COUNTS THAT ARE SHOWN ARE WEIGHTED WITH ANY SAMPLEWEIGHTS THAT MIGHT BE PRESENT

READ MORE IN THE USER GUIDE

PARAMETERS

DECISIONTREE DECISION TREE CLASSIFIER THE DECISION TREE TO BE EXPORTED TO GRAPHVIZ

OUTFILE FILE OBJECT OR STRING OPTIONAL DEFAULTNONE HANDLE OR NAME OF THE OUTPUT FILE IF NONE THE RESULT IS RETURNED AS A STRING

CHANGED IN VERSION 020 DEFAULT OF OUTFILE CHANGED FROM "TREEDOT" TO NONE

MAXDEPTH INT OPTIONAL DEFAULTNONE THE MAXIMUM DEPTH OF THE REPRESENTATION IF NONE THE TREE IS FULLY GENERATED

FEATURENAMES LIST OF STRINGS OPTIONAL DEFAULTNONE NAMES OF EACH OF THE FEATURES

CLASSNAMES LIST OF STRINGS BOOL OR NONE OPTIONAL DEFAULTNONE NAMES OF EACH OF THE TARGET CLASSES IN ASCENDING NUMERICAL ORDER ONLY RELEVANT FOR CLASSIFICATION AND NOT SUPPORTED FOR MULTIOUTPUT IF TRUE SHOWS A SYMBOLIC REPRESENTATION OF THE CLASS NAME

LABEL 'ALL' 'ROOT' 'NONE' OPTIONAL DEFAULT'ALL' WHETHER TO SHOW INFORMATIVE LABELS FOR IMPURITY ETC OPTIONS INCLUDE 'ALL' TO SHOW AT EVERY NODE 'ROOT' TO SHOW ONLY AT THE TOP ROOT NODE OR 'NONE' TO NOT SHOW AT ANY NODE

FILLED BOOL OPTIONAL DEFAULTFALSE WHEN SET TO TRUE PAINT NODES TO INDICATE MAJORITY CLASS FOR CLASSIFICATION EXTREMITY OF VALUES FOR REGRESSION OR PURITY OF NODE FOR MULTIOUTPUT

LEAVESPARALLEL BOOL OPTIONAL DEFAULTFALSE WHEN SET TO TRUE DRAW ALL LEAF NODES AT THE BOTTOM OF THE TREE

IMPURITY BOOL OPTIONAL DEFAULTTRUE WHEN SET TO TRUE SHOW THE IMPURITY AT EACH NODE

NODEIDS BOOL OPTIONAL DEFAULTFALSE WHEN SET TO TRUE SHOW THE ID NUMBER ON EACH NODE

PROPORTION BOOL OPTIONAL DEFAULTFALSE WHEN SET TO TRUE CHANGE THE DISPLAY OF 'VALUES' ANDOR 'SAMPLES' TO BE PROPORTIONS AND PERCENTAGES RESPECTIVELY

ROTATE BOOL OPTIONAL DEFAULTFALSE WHEN SET TO TRUE ORIENT TREE LEFT TO RIGHT RATHER THAN TOPDOWN

ROUNDED BOOL OPTIONAL DEFAULTFALSE WHEN SET TO TRUE DRAW NODE BOXES WITH ROUNDED CORNERS AND USE HELVETICA FONTS INSTEAD OF TIMESROMAN

SPECIALCHARACTERS BOOL OPTIONAL DEFAULTFALSE WHEN SET TO FALSE IGNORE SPECIAL CHARACTERS FOR POSTSCRIPT COMPATIBILITY

6385SKLEARNTREE DECISION TREES 2355

SCIKITLEARN USER GUIDE RELEASE 0213

PRECISION INT OPTIONAL DEFAULT3 NUMBER OF DIGITS OF PRECISION FOR FLOATING POINT IN THE VALUES OF IMPURITY THRESHOLD AND VALUE ATTRIBUTES OF EACH NODE

RETURNS

DOTDATA STRING STRING REPRESENTATION OF THE INPUT TREE IN GRAPHVIZ DOT FORMAT ONLY RETURNED

IFOUTFILE IS NONE

NEW IN VERSION 018

EXAMPLES

FROM SKLEARN DATASETS IMPORT LOADIRIS

FROM SKLEARN IMPORT TREE

CLF TREEDECISIONTREECLASSIFIER

IRIS LOADIRIS

CLF CLFFITIRISDATA IRISTARGET

TREEEXPORTGRAPHVIZCLF

DIGRAPH TREE

6386SKLEARN TREE PLOTTREE

SKLEARN TREE PLOTTREE DECISIONTREE MAXDEPTHNONE FEATURENAMESNONE CLASSNAMESNONE

LABEL'ALL' FILLEDFALSE IMPURITYTRUE NODEIDFALSE PROPORTIONFALSE ROTATEFALSE ROUNDEDFALSE PRECISION3 AXNONE FONT

SIZE NONE

PLOT A DECISION TREE

THE SAMPLE COUNTS THAT ARE SHOWN ARE WEIGHTED WITH ANY SAMPLEWEIGHTS THAT MIGHT BE PRESENT THIS FUNCTION REQUIRES MATPLOTLIB AND WORKS BEST WITH MATPLOTLIB 15

THE VISUALIZATION IS FIT AUTOMATICALLY TO THE SIZE OF THE AXIS USE THE FIGSIZE OR DPI ARGUMENTS OF PLT FIGURE TO CONTROL THE SIZE OF THE RENDERING

READ MORE IN THE USER GUIDE

NEW IN VERSION 021

PARAMETERS

DECISIONTREE DECISION TREE REGRESSOR OR CLASSIFIER THE DECISION TREE TO BE EXPORTED TO GRAPHVIZ

MAXDEPTH INT OPTIONAL DEFAULTNONE THE MAXIMUM DEPTH OF THE REPRESENTATION IF NONE THE TREE IS FULLY GENERATED

FEATURENAMES LIST OF STRINGS OPTIONAL DEFAULTNONE NAMES OF EACH OF THE FEATURES

CLASSNAMES LIST OF STRINGS BOOL OR NONE OPTIONAL DEFAULTNONE NAMES OF EACH OF THE TARGET CLASSES IN ASCENDING NUMERICAL ORDER ONLY RELEVANT FOR CLASSIFICATION AND NOT SUPPORTED FOR MULTIOUTPUT IF TRUE SHOWS A SYMBOLIC REPRESENTATION OF THE CLASS NAME

LABEL 'ALL' 'ROOT' 'NONE' OPTIONAL DEFAULT'ALL' WHETHER TO SHOW INFORMATIVE LABELS FOR IMPURITY ETC OPTIONS INCLUDE 'ALL' TO SHOW AT EVERY NODE 'ROOT' TO SHOW ONLY AT THE TOP ROOT NODE OR 'NONE' TO NOT SHOW AT ANY NODE

2356 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

FILLED BOOL OPTIONAL DEFAULTFALSE WHEN SET TO TRUE PAINT NODES TO INDICATE MAJORITY CLASS FOR CLASSIFICATION EXTREMITY OF VALUES FOR REGRESSION OR PURITY OF NODE FOR MULTIOUTPUT IMPURITY BOOL OPTIONAL DEFAULTTRUE WHEN SET TO TRUE SHOW THE IMPURITY AT EACH NODE

NODEIDS BOOL OPTIONAL DEFAULTFALSE WHEN SET TO TRUE SHOW THE ID NUMBER ON EACH NODE

PROPORTION BOOL OPTIONAL DEFAULTFALSE WHEN SET TO TRUE CHANGE THE DISPLAY OF 'VALUES' ANDOR 'SAMPLES' TO BE PROPORTIONS AND PERCENTAGES RESPECTIVELY

ROTATE BOOL OPTIONAL DEFAULTFALSE WHEN SET TO TRUE ORIENT TREE LEFT TO RIGHT RATHER THAN TOPDOWN

ROUNDED BOOL OPTIONAL DEFAULTFALSE WHEN SET TO TRUE DRAW NODE BOXES WITH ROUNDED CORNERS AND USE HELVETICA FONTS INSTEAD OF TIMESROMAN

PRECISION INT OPTIONAL DEFAULT3 NUMBER OF DIGITS OF PRECISION FOR FLOATING POINT IN THE VALUES OF IMPURITY THRESHOLD AND VALUE ATTRIBUTES OF EACH NODE

AXMATPLOTLIB AXIS OPTIONAL DEFAULTNONE AXES TO PLOT TO IF NONE USE CURRENT AXIS ANY PREVIOUS CONTENT IS CLEARED

FONTSIZE INT OPTIONAL DEFAULTNONE SIZE OF TEXT FONT IF NONE DETERMINED AUTOMATICALLY TO FIT FIGURE

RETURNS

ANNOTATIONS LIST OF ARTISTS LIST CONTAINING THE ARTISTS FOR THE ANNOTATION BOXES MAKING UP THE TREE

EXAMPLES

FROM SKLEARNDATASETS IMPORT LOADIRIS

FROM SKLEARN IMPORT TREE

CLF TREEDECISIONTREECLASSIFIERRANDOMSTATEO

IRIS LOADIRIS

CLF CLFFITIRISDATA IRISTARGET

TREEPLOTTREECLF

TEXT2515345217X3 08

EXAMPLES USING SKLEARNTREEPLOTTREE

•PLOT THE DECISION SURFACE OF A DECISION TREE ON THE IRIS DATASET

6387SKLEARNTREE EXPORTTEXT

SKLEARNTREE EXPORTTEXT DECISIONTREE FEATURENAMESNONE MAXDEPTH10 SPACING3 DECIMALS2 SHOWWEIGHTSFALSE

BUILD A TEXT REPORT SHOWING THE RULES OF A DECISION TREE

NOTE THAT BACKWARDS COMPATIBILITY MAY NOT BE SUPPORTED

PARAMETERS

638SKLEARNTREE DECISION TREES 2357

SCIKITLEARN USER GUIDE RELEASE 0213

DECISIONTREE OBJECT THE DECISION TREE ESTIMATOR TO BE EXPORTED IT CAN BE AN INSTANCE OF  
DECISIONTREECLASSIFIER OR DECISIONTREEREgressor

FEATURENAMES LIST OPTIONAL DEFAULTNONE A LIST OF LENGTH NFEATURES CONTAINING THE FEATURE  
NAMES IF NONE GENERIC NAMES WILL BE USED "FEATURE0" "FEATURE1"

MAXDEPTH INT OPTIONAL DEFAULT10 ONLY THE FIRST MAXDEPTH LEVELS OF THE TREE ARE EXPORTED  
TRUNCATED BRANCHES WILL BE MARKED WITH " "

SPACING INT OPTIONAL DEFAULT3 NUMBER OF SPACES BETWEEN EDGES THE HIGHER IT IS THE WIDER  
THE RESULT

DECIMALS INT OPTIONAL DEFAULT2 NUMBER OF DECIMAL DIGITS TO DISPLAY

SHOWWEIGHTS BOOL OPTIONAL DEFAULTFALSE IF TRUE THE CLASSIFICATION WEIGHTS WILL BE EXPORTED  
ON EACH LEAF THE CLASSIFICATION WEIGHTS ARE THE NUMBER OF SAMPLES EACH CLASS

RETURNS

REPORT STRING TEXT SUMMARY OF ALL THE RULES IN THE DECISION TREE

EXAMPLES

```
FROM SKLEARNDATASETS IMPORT LOADIRIS
FROM SKLEARNTREE IMPORT DECISIONTREECLASSIFIER
FROM SKLEARNTREEEXPORT IMPORT EXPORTTEXT
IRIS LOADIRIS
X IRISDATA
Y IRISTARGET
DECISIONTREE DECISIONTREECLASSIFIERRANDOMSTATE0 MAXDEPTH2
DECISIONTREE DECISIONTREEFITX Y
R EXPORTTEXTDECISIONTREE FEATURENAMESIRISFEATURENAMES
PRINTR
PETAL WIDTH CM 080
CLASS 0
PETAL WIDTH CM 080
PETAL WIDTH CM 175
CLASS 1
PETAL WIDTH CM 175
CLASS 2
```

639SKLEARNUTILS UTILITIES

THESKLEARNUTILS MODULE INCLUDES VARIOUS UTILITIES

DEVELOPER GUIDE SEE THE UTILITIES FOR DEVELOPERS PAGE FOR FURTHER DETAILS

UTILSARRAYFUNCSCHOLESKYDELETE L

GOOUT

UTILSARRAYFUNCSMINPOS FIND THE MINIMUM VALUE OF AN ARRAY OVER POSITIVE VALUES

UTILSASFLOATARRAY X COPY FORCEALLFINITE CONVERTS AN ARRAYLIKE TO AN ARRAY OF FLOATS

UTILSASSERTALLFINITE X ALLOWNAN THROW A VALUEERROR IF X CONTAINS NAN OR INFINITY

UTILSCHECKXY X Y ACCEPTSPARSE INPUT VALIDATION FOR STANDARD ESTIMATORS

CONTINUED ON NEXT PAGE

2358 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

TABLE 6288 - CONTINUED FROM PREVIOUS PAGE

UTILSCHECKARRAY ARRAY ACCEPTSPARSE INPUT VALIDATION ON AN ARRAY LIST SPARSE MATRIX OR SIMILAR

UTILSCHECKSCALAR X NAME TARGETTYPE VALIDATE SCALAR PARAMETERS TYPE AND VALUE

UTILSCHECKCONSISTENTLENGTH ARRAYS CHECK THAT ALL ARRAYS HAVE CONSISTENT FIRST DIMENSIONS

UTILSCHECKRANDOMSTATE SEED TURN SEED INTO A NPRANDOMRANDOMSTATE INSTANCE

UTILSCCLASSWEIGHT

COMPUTECLASSWEIGHT ESTIMATE CLASS WEIGHTS FOR UNBALANCED DATASETS

UTILSCCLASSWEIGHT

COMPUTESAMPLEWEIGHT ESTIMATE SAMPLE WEIGHTS BY CLASS FOR UNBALANCED DATASETS

UTILSDEPRECATED EXTRA DECORATOR TO MARK A FUNCTION OR CLASS AS DEPRECATED

UTILSESTIMATORCHECKS

CHECKESTIMATOR ESTIMATORCHECK IF ESTIMATOR ADHERES TO SCIKITLEARN CONVENTIONS

UTILSEXTMATHSAFESPASEDOT A B DOT PRODUCT THAT HANDLE THE SPARSE MATRIX CASE CORRECTLY

UTILSEXTMATHRANDOMIZEDRANGEFINDER A

COMPUTES AN ORTHONORMAL MATRIX WHOSE RANGE APPROXI

MATES THE RANGE OF A

UTILSEXTMATHRANDOMIZEDSVD M

NCOMPONENTSCOMPUTES A TRUNCATED RANDOMIZED SVD

UTILSEXTMATHFASTLOGDET A COMPUTE LOGDETA FOR A SYMMETRIC

UTILSEXTMATHDENSITY W KWARGS COMPUTE DENSITY OF A SPARSE VECTOR

UTILSEXTMATHWEIGHTEDMODE A W AXIS RETURNS AN ARRAY OF THE WEIGHTED MODAL MOST COMMON

VALUE IN A

UTILSGENEVENSLICES N NPACKS NSAMPLES GENERATOR TO CREATE NPACKS SLICES GOING UP TO N

UTILSGRAPHSSINGLESOURCESHORTESTPATHLENGTH RETURN THE SHORTEST PATH LENGTH FROM SOURCE TO ALL REACHABLE

NODES

UTILSGRAPHSHORTESTPATH

GRAPHSHORTESTPATH PERFORM A SHORTESTPATH GRAPH SEARCH ON A POSITIVE DIRECTED

OR UNDIRECTED GRAPH

UTILSINDEXABLE ITERABLES MAKE ARRAYS INDEXABLE FOR CROSSVALIDATION

UTILSMETAESTIMATORS

IFDELEGATEHASMETHOD CREATE A DECORATOR FOR METHODS THAT ARE DELEGATED TO A SUB

ESTIMATOR

UTILSMULTICLASSTYPEOFTARGET Y DETERMINE THE TYPE OF DATA INDICATED BY THE TARGET

UTILSMULTICLASSISMULTILABEL Y CHECK IF YIS IN A MULTILABEL FORMAT

UTILSMULTICLASSUNIQUELABELS YS EXTRACT AN ORDERED ARRAY OF UNIQUE LABELS

UTILSMURMURHASH32 COMPUTE THE 32BIT MURMURHASH3 OF KEY AT SEED

UTILSRESAMPLE ARRAYS OPTIONS RESAMPLE ARRAYS OR SPARSE MATRICES IN A CONSISTENT WAY

UTILSSAFEINDEXING X INDICES RETURN ITEMS OR ROWS FROM X USING INDICES

UTILSSAFEMASK X MASK RETURN A MASK WHICH IS SAFE TO USE ON X

UTILSSAFESQR X COPY ELEMENT WISE SQUARING OF ARRAYLIKES AND SPARSE MATRICES

UTILSSHUFFLE ARRAYS OPTIONS SHUFFLE ARRAYS OR SPARSE MATRICES IN A CONSISTENT WAY

UTILSSPARSEFUNCS

INCRMEANVARIANCEAXIS X COMPUTE INCREMENTAL MEAN AND VARIANCE ALONG AN AXIX ON

A CSR OR CSC MATRIX

UTILSSPARSEFUNCS

INPLACECOLUMNSCALE X SCALEINPLACE COLUMN SCALING OF A CSCCSR MATRIX

UTILSSPARSEFUNCSINPLACEROWSCALE X

SCALEINPLACE ROW SCALING OF A CSR OR CSC MATRIX

UTILSSPARSEFUNCSINPLACESWAPROW X M

NSWAPS TWO ROWS OF A CSCCSR MATRIX INPLACE

UTILSSPARSEFUNCS

INPLACESWAPCOLUMN X M NSWAPS TWO COLUMNS OF A CSCCSR MATRIX INPLACE

UTILSSPARSEFUNCSMEANVARIANCEAXIS X

AXISCOMPUTE MEAN AND VARIANCE ALONG AN AXIX ON A CSR OR

CSC MATRIX

UTILSSPARSEFUNCS

INPLACECSROLUMNSCALE X INPLACE COLUMN SCALING OF A CSR MATRIX

CONTINUED ON NEXT PAGE

639SKLEARNUTILS UTILITIES 2359

SCIKITLEARN USER GUIDE RELEASE 0213

TABLE 6288 - CONTINUED FROM PREVIOUS PAGE

UTILSSPARSEFUNCSFAST

INPLACECSRROWNORMALIZEL1 INPLACE ROW NORMALIZE USING THE L1 NORM

UTILSSPARSEFUNCSFAST

INPLACECSRROWNORMALIZEL2 INPLACE ROW NORMALIZE USING THE L2 NORM

UTILSRANDOMSAMPLEWITHOUTREPLACEMENT SAMPLE INTEGERS WITHOUT REPLACEMENT

UTILSVALIDATIONCHECKISFITTED ESTIMATOR

PERFORM ISFITTED VALIDATION FOR ESTIMATOR

UTILSVALIDATIONCHECKMEMORY MEMORY CHECK THAT MEMORY IS JOBLIBMEMORYLIKE

UTILSVALIDATIONCHECKSYMMETRIC ARRAY

MAKE SURE THAT ARRAY IS 2D SQUARE AND SYMMETRIC

UTILSVALIDATIONCOLUMNOR1D Y WARN RAVEL COLUMN OR 1D NUMPY ARRAY ELSE RAISES AN ERROR

UTILSVALIDATIONHASFITPARAMETER CHECKS WHETHER THE ESTIMATOR'S FIT METHOD SUPPORTS THE GIVEN PARAMETER

UTILSTESTINGASSERTIN MEMBER CONTAINER

MSGJUST LIKE SELFASSERTTRUEA IN B BUT WITH A NICER DEFAULT MESSAGE

UTILSTESTINGASSERTNOTIN MEMBER CON

TAINERJUST LIKE SELFASSERTTRUEA NOT IN B BUT WITH A NICER DEFAULT MESSAGE

UTILSTESTINGASSERTRAISEMESSAGE HELPER FUNCTION TO TEST THE MESSAGE RAISED IN AN EXCEPTION

UTILSTESTINGALLESTIMATORS GET A LIST OF ALL ESTIMATORS FROM SKLEARN

6391SKLEARNUTILS ARRAYFUNCSCHOLESKYDELETE

6392SKLEARNUTILSARRAYFUNCS MINPOS

SKLEARNUTILSARRAYFUNCS MINPOS

FIND THE MINIMUM VALUE OF AN ARRAY OVER POSITIVE VALUES

RETURNS A HUGE VALUE IF NONE OF THE VALUES ARE POSITIVE

6393SKLEARNUTILS ASFLOATARRAY

SKLEARNUTILS ASFLOATARRAY XCOPYTRUE FORCEALLFINITETRUE

CONVERTS AN ARRAYLIKE TO AN ARRAY OF FLOATS

THE NEW DTYPE WILL BE NPFLOAT32 OR NPFLOAT64 DEPENDING ON THE ORIGINAL TYPE THE FUNCTION CAN CREATE A COPY OR MODIFY THE ARGUMENT DEPENDING ON THE ARGUMENT COPY

PARAMETERS

XARRAYLIKE SPARSE MATRIX

COPY BOOL OPTIONAL IF TRUE A COPY OF X WILL BE CREATED IF FALSE A COPY MAY STILL BE RETURNED

IF X'S DTYPE IS NOT A FLOATING POINT TYPE

FORCEALLFINITE BOOLEAN OR 'ALLOWNAN' DEFAULTTRUE WHETHER TO RAISE AN ERROR ON NPINF AND NPNAN IN X THE POSSIBILITIES ARE

- TRUE FORCE ALL VALUES OF X TO BE FINITE
- FALSE ACCEPT BOTH NPINF AND NPNAN IN X
- 'ALLOWNAN' ACCEPT ONLY NPNAN VALUES IN X VALUES CANNOT BE INFINITE

NEW IN VERSION 020 FORCEALLFINITE ACCEPTS THE STRING ALLOWNAN

RETURNS

2360 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

XT ARRAY SPARSE MATRIX AN ARRAY OF TYPE NPFLOAT

6394SKLEARNUTILS ASSERTALLFINITE

SKLEARNUTILS ASSERTALLFINITE XALLOWNANFALSE

THROW A VALUEERROR IF X CONTAINS NAN OR INFINITY

PARAMETERS

XARRAY OR SPARSE MATRIX

ALLOWNAN BOOL

6395SKLEARNUTILS CHECKXY

SKLEARNUTILS CHECKXY XYACCEPTSPARSEFALSE ACCEPTLARGESPARSETRUE DTYPE'NUMERIC'

ORDERNONE COPYFALSE FORCEALLFINITETRUE ENSURE2DTRUE

ALLOWNDFALSE MULTIOUTPUTFALSE ENSUREMINSAMPLES1 EN

SUREMINFEATURES1 YNUMERICFALSE WARNONDTYPENONE ESTIMA

TORNONE

INPUT VALIDATION FOR STANDARD ESTIMATORS

CHECKS X AND Y FOR CONSISTENT LENGTH ENFORCES X TO BE 2D AND Y 1D BY DEFAULT X IS CHECKED TO BE NONEMPTY

AND CONTAINING ONLY FINITE VALUES STANDARD INPUT CHECKS ARE ALSO APPLIED TO Y SUCH AS CHECKING THAT Y DOES NOT

HAVE NPNAN OR NPINF TARGETS FOR MULTILABEL Y SET MULTIOUTPUTTRUE TO ALLOW 2D AND SPARSE Y IF THE DTYPE OF

X IS OBJECT ATTEMPT CONVERTING TO FLOAT RAISING ON FAILURE

PARAMETERS

XNDARRAY LIST OR SPARSE MATRIX INPUT DATA

YNDARRAY LIST OR SPARSE MATRIX LABELS

ACCEPTSPARSE STRING BOOLEAN OR LIST OF STRING DEFAULTFALSE STRINGS REPRESENTING ALLOWED

SPARSE MATRIX FORMATS SUCH AS 'CSC' 'CSR' ETC IF THE INPUT IS SPARSE BUT NOT IN THE ALLOWED

FORMAT IT WILL BE CONVERTED TO THE FIRST LISTED FORMAT TRUE ALLOWS THE INPUT TO BE ANY FORMAT

FALSE MEANS THAT A SPARSE MATRIX INPUT WILL RAISE AN ERROR

ACCEPTLARGESPARSE BOOL DEFAULTTRUE IF A CSR CSC COO OR BSR SPARSE MATRIX IS SUP

PLIED AND ACCEPTED BY ACCEPTSPARSE ACCEPTLARGESPARSE WILL CAUSE IT TO BE ACCEPTED ONLY

IF ITS INDICES ARE STORED WITH A 32BIT DTYPE

NEW IN VERSION 020

DTYPE STRING TYPE LIST OF TYPES OR NONE DEFAULT"NUMERIC" DATA TYPE OF RESULT IF NONE THE

DTYPE OF THE INPUT IS PRESERVED IF "NUMERIC" DTYPE IS PRESERVED UNLESS ARRAYDTYPE IS OBJECT

IF DTYPE IS A LIST OF TYPES CONVERSION ON THE FIRST TYPE IS ONLY PERFORMED IF THE DTYPE OF THE

INPUT IS NOT IN THE LIST

ORDER 'F' 'C' OR NONE DEFAULTNONE WHETHER AN ARRAY WILL BE FORCED TO BE FORTRAN OR CSTYLE

COPY BOOLEAN DEFAULTFALSE WHETHER A FORCED COPY WILL BE TRIGGERED IF COPYFALSE A COPY

MIGHT BE TRIGGERED BY A CONVERSION

FORCEALLFINITE BOOLEAN OR 'ALLOWNAN' DEFAULTTRUE WHETHER TO RAISE AN ERROR ON NPINF

AND NPNAN IN X THIS PARAMETER DOES NOT INFLUENCE WHETHER Y CAN HAVE NPINF OR NPNAN

VALUES THE POSSIBILITIES ARE

- TRUE FORCE ALL VALUES OF X TO BE FINITE

6395SKLEARNUTILS UTILITIES 2361

SCIKITLEARN USER GUIDE RELEASE 0213

- FALSE ACCEPT BOTH NPINF AND NPNAN IN X
- ‘ALLOWNAN’ ACCEPT ONLY NPNAN VALUES IN X VALUES CANNOT BE INFINITE

NEW IN VERSION 020 FORCEALLFINITE ACCEPTS THE STRING ALLOWNAN

ENSURE2D BOOLEAN DEFAULTTRUE WHETHER TO RAISE A VALUE ERROR IF X IS NOT 2D

ALLOWND BOOLEAN DEFAULTFALSE WHETHER TO ALLOW XNDIM 2

MULTIOUTPUT BOOLEAN DEFAULTFALSE WHETHER TO ALLOW 2D Y ARRAY OR SPARSE MATRIX

IF FALSE Y WILL BE VALIDATED AS A VECTOR Y CANNOT HAVE NPNAN OR NPINF VALUES IF

MULTIOUTPUTTRUE

ENSUREMINSAMPLES INT DEFAULT1 MAKE SURE THAT X HAS A MINIMUM NUMBER OF SAMPLES IN

ITS FIRST AXIS ROWS FOR A 2D ARRAY

ENSUREMINFEATURES INT DEFAULT1 MAKE SURE THAT THE 2D ARRAY HAS SOME MINIMUM NUMBER

OF FEATURES COLUMNS THE DEFAULT VALUE OF 1 REJECTS EMPTY DATASETS THIS CHECK IS ONLY

ENFORCED WHEN X HAS EFFECTIVELY 2 DIMENSIONS OR IS ORIGINALLY 1D AND ENSURE2D IS TRUE

SETTING TO 0 DISABLES THIS CHECK

YNUMERIC BOOLEAN DEFAULTFALSE WHETHER TO ENSURE THAT Y HAS A NUMERIC TYPE IF DTYPE OF

Y IS OBJECT IT IS CONVERTED TO FLOAT64 SHOULD ONLY BE USED FOR REGRESSION ALGORITHMS

WARNONDDTYPE BOOLEAN OR NONE OPTIONAL DEFAULTNONE RAISE DATACONVERSIONWARNING IF

THE DTYPE OF THE INPUT DATA STRUCTURE DOES NOT MATCH THE REQUESTED DTYPE CAUSING A MEMORY

COPY

DEPRECATED SINCE VERSION 021 WARNONDDTYPE IS DEPRECATED IN VERSION 021 AND WILL BE

REMOVED IN 023

ESTIMATOR STR OR ESTIMATOR INSTANCE DEFAULTNONE IF PASSED INCLUDE THE NAME OF THE ESTIMATOR

IN WARNING MESSAGES

RETURNS

XCONVERTED OBJECT THE CONVERTED AND VALIDATED X

YCONVERTED OBJECT THE CONVERTED AND VALIDATED Y

6396SKLEARNUTILS CHECKARRAY

SKLEARNUTILS CHECKARRAY ARRAY ACCEPTSPARSEFALSE ACCEPTLARGESPARSETRUE

DTYPE’NUMERIC’ ORDERNONE COPYFALSE FORCEALLFINITETRUE

ENSURE2DTRUE ALLOWNDFALSE ENSUREMINSAMPLES1 EN

SUREMINFEATURES1 WARNONDDTYPENONE ESTIMATORNONE

INPUT VALIDATION ON AN ARRAY LIST SPARSE MATRIX OR SIMILAR

BY DEFAULT THE INPUT IS CHECKED TO BE A NONEMPTY 2D ARRAY CONTAINING ONLY FINITE VALUES IF THE DTYPE OF THE

ARRAY IS OBJECT ATTEMPT CONVERTING TO FLOAT RAISING ON FAILURE

PARAMETERS

ARRAY OBJECT INPUT OBJECT TO CHECK CONVERT

ACCEPTSPARSE STRING BOOLEAN OR LISTTUPLE OF STRINGS DEFAULTFALSE STRINGS REPRESENTING

ALLOWED SPARSE MATRIX FORMATS SUCH AS ‘CSC’ ‘CSR’ ETC IF THE INPUT IS SPARSE BUT NOT IN THE

ALLOWED FORMAT IT WILL BE CONVERTED TO THE FIRST LISTED FORMAT TRUE ALLOWS THE INPUT TO BE ANY

FORMAT FALSE MEANS THAT A SPARSE MATRIX INPUT WILL RAISE AN ERROR

2362 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

ACCEPTLARGESPARE BOOL DEFAULTTRUE IF A CSR CSC COO OR BSR SPARSE MATRIX IS SUP  
PLIED AND ACCEPTED BY ACCEPTSPARSE ACCEPTLARGESPAREFALSE WILL CAUSE IT TO BE ACCEPTED  
ONLY IF ITS INDICES ARE STORED WITH A 32BIT DTYPE

NEW IN VERSION 020

DTYPE STRING TYPE LIST OF TYPES OR NONE DEFAULT"NUMERIC" DATA TYPE OF RESULT IF NONE THE  
DTYPE OF THE INPUT IS PRESERVED IF "NUMERIC" DTYPE IS PRESERVED UNLESS ARRAYDTYPE IS OBJECT  
IF DTYPE IS A LIST OF TYPES CONVERSION ON THE FIRST TYPE IS ONLY PERFORMED IF THE DTYPE OF THE  
INPUT IS NOT IN THE LIST

ORDER 'F' 'C' OR NONE DEFAULTNONE WHETHER AN ARRAY WILL BE FORCED TO BE FORTRAN OR CSTYLE  
WHEN ORDER IS NONE DEFAULT THEN IF COPYFALSE NOTHING IS ENSURED ABOUT THE MEMORY  
LAYOUT OF THE OUTPUT ARRAY OTHERWISE COPYTRUE THE MEMORY LAYOUT OF THE RETURNED ARRAY  
IS KEPT AS CLOSE AS POSSIBLE TO THE ORIGINAL ARRAY

COPY BOOLEAN DEFAULTFALSE WHETHER A FORCED COPY WILL BE TRIGGERED IF COPYFALSE A COPY  
MIGHT BE TRIGGERED BY A CONVERSION

FORCEALLFINITE BOOLEAN OR 'ALLOWNAN' DEFAULTTRUE WHETHER TO RAISE AN ERROR ON NPINF  
AND NPNAN IN ARRAY THE POSSIBILITIES ARE

- TRUE FORCE ALL VALUES OF ARRAY TO BE FINITE
- FALSE ACCEPT BOTH NPINF AND NPNAN IN ARRAY
- 'ALLOWNAN' ACCEPT ONLY NPNAN VALUES IN ARRAY VALUES CANNOT BE INFINITE

FOR OBJECT DTYPED DATA ONLY NPNAN IS CHECKED AND NOT NPINF

NEW IN VERSION 020 FORCEALLFINITE ACCEPTS THE STRING ALLOWNAN

ENSURE2D BOOLEAN DEFAULTTRUE WHETHER TO RAISE A VALUE ERROR IF ARRAY IS NOT 2D

ALLOWND BOOLEAN DEFAULTFALSE WHETHER TO ALLOW ARRAYNDIM 2

ENSUREMINSAMPLES INT DEFAULT1 MAKE SURE THAT THE ARRAY HAS A MINIMUM NUMBER OF  
SAMPLES IN ITS FIRST AXIS ROWS FOR A 2D ARRAY SETTING TO 0 DISABLES THIS CHECK

ENSUREMINFEATURES INT DEFAULT1 MAKE SURE THAT THE 2D ARRAY HAS SOME MINIMUM NUMBER  
OF FEATURES COLUMNS THE DEFAULT VALUE OF 1 REJECTS EMPTY DATASETS THIS CHECK IS ONLY EN  
FORCED WHEN THE INPUT DATA HAS EFFECTIVELY 2 DIMENSIONS OR IS ORIGINALLY 1D AND ENSURE2D  
IS TRUE SETTING TO 0 DISABLES THIS CHECK

WARNONDDTYPE BOOLEAN OR NONE OPTIONAL DEFAULTNONE RAISE DATACONVERSIONWARNING IF  
THE DTYPE OF THE INPUT DATA STRUCTURE DOES NOT MATCH THE REQUESTED DTYPE CAUSING A MEMORY  
COPY

DEPRECATED SINCE VERSION 021 WARNONDDTYPE IS DEPRECATED IN VERSION 021 AND WILL BE  
REMOVED IN 023

ESTIMATOR STR OR ESTIMATOR INSTANCE DEFAULTNONE IF PASSED INCLUDE THE NAME OF THE ESTIMATOR  
IN WARNING MESSAGES

RETURNS

ARRAYCONVERTED OBJECT THE CONVERTED AND VALIDATED ARRAY

6397SKLEARNUTILS CHECKSCALAR

SKLEARNUTILS CHECKSCALAR XNAME TARGETTYPE MINVALNONE MAXVALNONE

VALIDATE SCALAR PARAMETERS TYPE AND VALUE

639SKLEARNUTILS UTILITIES 2363

SCIKITLEARN USER GUIDE RELEASE 0213

PARAMETERS

XOBJECT THE SCALAR PARAMETER TO VALIDATE

NAME STR THE NAME OF THE PARAMETER TO BE PRINTED IN ERROR MESSAGES

TARGETTYPE TYPE OR TUPLE ACCEPTABLE DATA TYPES FOR THE PARAMETER

MINVAL FLOAT OR INT OPTIONAL DEFAULTNONE THE MINIMUM VALID VALUE THE PARAMETER CAN

TAKE IF NONE DEFAULT IT IS IMPLIED THAT THE PARAMETER DOES NOT HAVE A LOWER BOUND

MAXVAL FLOAT OR INT OPTIONAL DEFAULTNONE THE MAXIMUM VALID VALUE THE PARAMETER CAN

TAKE IF NONE DEFAULT IT IS IMPLIED THAT THE PARAMETER DOES NOT HAVE AN UPPER BOUND

RAISES

TYPEERROR IF THE PARAMETER’S TYPE DOES NOT MATCH THE DESIRED TYPE

VALUEERROR IF THE PARAMETER’S VALUE VIOLATES THE GIVEN BOUNDS

6398SKLEARNUTILS CHECKCONSISTENTLENGTH

SKLEARNUTILS CHECKCONSISTENTLENGTH ARRAYS

CHECK THAT ALL ARRAYS HAVE CONSISTENT FIRST DIMENSIONS

CHECKS WHETHER ALL OBJECTS IN ARRAYS HAVE THE SAME SHAPE OR LENGTH

PARAMETERS

ARRAYS LIST OR TUPLE OF INPUT OBJECTS OBJECTS THAT WILL BE CHECKED FOR CONSISTENT LENGTH

6399SKLEARNUTILS CHECKRANDOMSTATE

SKLEARNUTILS CHECKRANDOMSTATE SEED

TURN SEED INTO A NPRANDOMRANDOMSTATE INSTANCE

PARAMETERS

SEED NONE INT INSTANCE OF RANDOMSTATE IF SEED IS NONE RETURN THE RANDOMSTATE SINGLETON

USED BY NPRANDOM IF SEED IS AN INT RETURN A NEW RANDOMSTATE INSTANCE SEEDED WITH SEED

IF SEED IS ALREADY A RANDOMSTATE INSTANCE RETURN IT OTHERWISE RAISE VALUEERROR

EXAMPLES USING SKLEARNUTILSCHECKRANDOMSTATE

- ISOTONIC REGRESSION
  - FACE COMPLETION WITH A MULTIOUTPUT ESTIMATORS
  - EMPIRICAL EVALUATION OF THE IMPACT OF KMEANS INITIALIZATION
  - MNIST CLASSFICATION USING MULTINOMIAL LOGISTIC L1
  - MANIFOLD LEARNING METHODS ON A SEVERED SPHERE
  - SCALING THE REGULARIZATION PARAMETER FOR SVCS
- 2364 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213  
63910SKLEARNUTILSCCLASSWEIGHT COMPUTECLASSWEIGHT  
SKLEARNUTILSCCLASSWEIGHT COMPUTECLASSWEIGHT CLASSSES Y  
ESTIMATE CLASS WEIGHTS FOR UNBALANCED DATASETS  
PARAMETERS  
CLASSWEIGHT DICT 'BALANCED' OR NONE IF 'BALANCED' CLASS WEIGHTS WILL BE GIVEN BY  
NSAMPLES NCLASSES NPBINCOUNTY IF A DICTIONARY IS GIVEN KEYS  
ARE CLASSES AND VALUES ARE CORRESPONDING CLASS WEIGHTS IF NONE IS GIVEN THE CLASS WEIGHTS  
WILL BE UNIFORM  
CLASSES NDARRAY ARRAY OF THE CLASSES OCCURRING IN THE DATA AS GIVEN BY NPUNIQUEYORG  
WITHYORG THE ORIGINAL CLASS LABELS  
YARRAYLIKE SHAPE NSAMPLES ARRAY OF ORIGINAL CLASS LABELS PER SAMPLE  
RETURNS  
CLASSWEIGHTVECT NDARRAY SHAPE NCLASSES ARRAY WITH CLASSWEIGHTVECTI THE WEIGHT FOR  
ITH CLASS  
REFERENCES  
THE "BALANCED" HEURISTIC IS INSPIRED BY LOGISTIC REGRESSION IN RARE EVENTS DATA KING ZEN 2001  
63911SKLEARNUTILSCCLASSWEIGHT COMPUTESAMPLEWEIGHT  
SKLEARNUTILSCCLASSWEIGHT COMPUTESAMPLEWEIGHT CLASSWEIGHT YINDICESNONE  
ESTIMATE SAMPLE WEIGHTS BY CLASS FOR UNBALANCED DATASETS  
PARAMETERS  
CLASSWEIGHT DICT LIST OF DICTS "BALANCED" OR NONE OPTIONAL WEIGHTS ASSOCIATED WITH CLASSES  
IN THE FORM CLASSLABEL WEIGHT IF NOT GIVEN ALL CLASSES ARE SUPPOSED TO HAVE  
WEIGHT ONE FOR MULTIOUTPUT PROBLEMS A LIST OF DICTS CAN BE PROVIDED IN THE SAME ORDER AS  
THE COLUMNS OF Y  
NOTE THAT FOR MULTIOUTPUT INCLUDING MULTILABEL WEIGHTS SHOULD BE DEFINED FOR EACH CLASS OF  
EVERY COLUMN IN ITS OWN DICT FOR EXAMPLE FOR FOURCLASS MULTILABEL CLASSIFICATION WEIGHTS  
SHOULD BE 0 1 1 1 0 1 1 5 0 1 1 1 0 1 1 1 INSTEAD OF 11 25  
31 41  
THE "BALANCED" MODE USES THE VALUES OF Y TO AUTOMATICALLY ADJUST WEIGHTS INVERSELY PRO  
PORTIONAL TO CLASS FREQUENCIES IN THE INPUT DATA NSAMPLES NCLASSES NP  
BINCOUNTY  
FOR MULTIOUTPUT THE WEIGHTS OF EACH COLUMN OF Y WILL BE MULTIPLIED  
YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NOUTPUTS ARRAY OF ORIGINAL CLASS LABELS PER  
SAMPLE  
INDICES ARRAYLIKE SHAPE NSUBSAMPLE OR NONE ARRAY OF INDICES TO BE USED IN A SUBSAMPLE  
CAN BE OF LENGTH LESS THAN NSAMPLES IN THE CASE OF A SUBSAMPLE OR EQUAL TO NSAMPLES IN  
THE CASE OF A BOOTSTRAP SUBSAMPLE WITH REPEATED INDICES IF NONE THE SAMPLE WEIGHT WILL  
BE CALCULATED OVER THE FULL SAMPLE ONLY "BALANCED" IS SUPPORTED FOR CLASSWEIGHT IF THIS IS  
PROVIDED  
RETURNS  
639SKLEARNUTILS UTILITIES 2365

SCIKITLEARN USER GUIDE RELEASE 0213

SAMPLEWEIGHTVECT NDARRAY SHAPE NSAMPLES ARRAY WITH SAMPLE WEIGHTS AS APPLIED TO THE ORIGINAL Y

63912SKLEARNUTILS DEPRECATED

SKLEARNUTILS DEPRECATED EXTRA''

DECORATOR TO MARK A FUNCTION OR CLASS AS DEPRECATED

ISSUE A WARNING WHEN THE FUNCTION IS CALLEDTHE CLASS IS INSTANTIATED AND ADDS A WARNING TO THE DOCSTRING

THE OPTIONAL EXTRA ARGUMENT WILL BE APPENDED TO THE DEPRECATION MESSAGE AND THE DOCSTRING NOTE TO USE THIS

WITH THE DEFAULT VALUE FOR EXTRA PUT IN AN EMPTY OF PARENTHESES

FROM SKLEARNUTILS IMPORT DEPRECATED

DEPRECATED

SKLEARNUTILSDEPRECATIONDEPRECATED OBJECT AT

DEPRECATED

DEF SOMEFUNCTION PASS

PARAMETERS

EXTRA STRING TO BE ADDED TO THE DEPRECATION MESSAGES

63913SKLEARNUTILSESTIMATORCHECKS CHECKESTIMATOR

SKLEARNUTILSESTIMATORCHECKS CHECKESTIMATOR ESTIMATOR

CHECK IF ESTIMATOR ADHERES TO SCIKITLEARN CONVENTIONS

THIS ESTIMATOR WILL RUN AN EXTENSIVE TESTSUITE FOR INPUT VALIDATION SHAPES ETC ADDITIONAL TESTS FOR CLASSIFIERS

REGRESSORS CLUSTERING OR TRANSFORMERS WILL BE RUN IF THE ESTIMATOR CLASS INHERITS FROM THE CORRESPONDING MIXIN

FROM SKLEARNBASE

THIS TEST CAN BE APPLIED TO CLASSES OR INSTANCES CLASSES CURRENTLY HAVE SOME ADDITIONAL TESTS THAT RELATED TO

CONSTRUCTION WHILE PASSING INSTANCES ALLOWS THE TESTING OF MULTIPLE OPTIONS

PARAMETERS

ESTIMATOR ESTIMATOR OBJECT OR CLASS ESTIMATOR TO CHECK ESTIMATOR IS A CLASS OBJECT OR INSTANCE

63914SKLEARNUTILSEXTMATH SAFESPASEDOT

SKLEARNUTILSEXTMATH SAFESPASEDOT ABDENSEOUTPUTFALSE

DOT PRODUCT THAT HANDLE THE SPARSE MATRIX CASE CORRECTLY

USES BLAS GEMM AS REPLACEMENT FOR NUMPYDOT WHERE POSSIBLE TO AVOID UNNECESSARY COPIES

PARAMETERS

AARRAY OR SPARSE MATRIX

BARRAY OR SPARSE MATRIX

DENSEOUTPUT BOOLEAN DEFAULT FALSE WHEN FALSE EITHER AORBBEING SPARSE WILL YIELD SPARSE

OUTPUT WHEN TRUE OUTPUT WILL ALWAYS BE AN ARRAY

RETURNS

2366 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

DOTPRODUCT ARRAY OR SPARSE MATRIX SPARSE IF AORBIS SPARSE AND DENSEOUTPUTFALSE

63915SKLEARNUTILSEXTMATH RANDOMIZEDRANGEFINDER

SKLEARNUTILSEXTMATH RANDOMIZEDRANGEFINDER A SIZE NITER

POWERITERATIONNORMALIZER'AUTO'

RANDOMSTATENONE

COMPUTES AN ORTHONORMAL MATRIX WHOSE RANGE APPROXIMATES THE RANGE OF A

PARAMETERS

A2D ARRAY THE INPUT DATA MATRIX

SIZE INTEGER SIZE OF THE RETURN ARRAY

NITER INTEGER NUMBER OF POWER ITERATIONS USED TO STABILIZE THE RESULT

POWERITERATIONNORMALIZER 'AUTO' DEFAULT 'QR' 'LU' 'NONE' WHETHER THE POWER ITERATIONS ARE NORMALIZED WITH STEPBYSTEP QR FACTORIZATION THE SLOWEST BUT MOST ACCURATE 'NONE' THE FASTEST BUT NUMERICALLY UNSTABLE WHEN NITER IS LARGE EG TYPICALLY 5 OR LARGER OR 'LU' FACTORIZATION NUMERICALLY STABLE BUT CAN LOSE SLIGHTLY IN ACCURACY THE 'AUTO' MODE APPLIES NO NORMALIZATION IF NITER > 2 AND SWITCHES TO LU OTHERWISE

NEW IN VERSION 018

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE THE SEED OF THE PSEUDO RANDOM NUMBER GENERATOR TO USE WHEN SHUFFLING THE DATA IF INT RANDOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NRANDOM

RETURNS

Q2D ARRAY A SIZE X SIZE PROJECTION MATRIX THE RANGE OF WHICH APPROXIMATES WELL THE RANGE OF THE INPUT MATRIX A

NOTES

FOLLOWS ALGORITHM 43 OF FINDING STRUCTURE WITH RANDOMNESS STOCHASTIC ALGORITHMS FOR CONSTRUCTING APPROXIMATE MATRIX DECOMPOSITIONS HALKO ET AL 2009 ARXIV909 [HTTPSARXIVORGPDF09094061PDF](https://arxiv.org/pdf/0909.4061.pdf)

AN IMPLEMENTATION OF A RANDOMIZED ALGORITHM FOR PRINCIPAL COMPONENT ANALYSIS A SZLAM ET AL 2014

63916SKLEARNUTILSEXTMATH RANDOMIZEDSVD

SKLEARNUTILSEXTMATH RANDOMIZEDSVD MNCOMPONENTS NOVERSAMPLES10 NITER'AUTO'

POWERITERATIONNORMALIZER'AUTO' TRANSPOSE'AUTO'

FLIPSIGNTRUE RANDOMSTATE0

COMPUTES A TRUNCATED RANDOMIZED SVD

PARAMETERS

MNDARRAY OR SPARSE MATRIX MATRIX TO DECOMPOSE

NCOMPONENTS INT NUMBER OF SINGULAR VALUES AND VECTORS TO EXTRACT

NOVERSAMPLES INT DEFAULT IS 10 ADDITIONAL NUMBER OF RANDOM VECTORS TO SAMPLE THE RANGE OF M SO AS TO ENSURE PROPER CONDITIONING THE TOTAL NUMBER OF RANDOM VECTORS USED TO FIND

6395SKLEARNUTILS UTILITIES 2367

SCIKITLEARN USER GUIDE RELEASE 0213

THE RANGE OF M IS NCOMPONENTS NOVERSAMPLES SMALLER NUMBER CAN IMPROVE SPEED BUT CAN NEGATIVELY IMPACT THE QUALITY OF APPROXIMATION OF SINGULAR VECTORS AND SINGULAR VALUES NITER INT OR 'AUTO' DEFAULT IS 'AUTO' NUMBER OF POWER ITERATIONS IT CAN BE USED TO DEAL WITH VERY NOISY PROBLEMS WHEN 'AUTO' IT IS SET TO 4 UNLESS NCOMPONENTS IS SMALL 1 MINXSHAPE NITER IN WHICH CASE IS SET TO 7 THIS IMPROVES PRECISION WITH FEW COMPONENTS

CHANGED IN VERSION 018

POWERITERATIONNORMALIZER 'AUTO' DEFAULT 'QR' 'LU' 'NONE' WHETHER THE POWER ITERATIONS ARE NORMALIZED WITH STEPBYSTEP QR FACTORIZATION THE SLOWEST BUT MOST ACCURATE 'NONE' THE FASTEST BUT NUMERICALLY UNSTABLE WHEN NITER IS LARGE EG TYPICALLY 5 OR LARGER OR 'LU' FACTORIZATION NUMERICALLY STABLE BUT CAN LOSE SLIGHTLY IN ACCURACY THE 'AUTO' MODE APPLIES NO NORMALIZATION IF NITER 2 AND SWITCHES TO LU OTHERWISE NEW IN VERSION 018

TRANSPOSE TRUE FALSE OR 'AUTO' DEFAULT WHETHER THE ALGORITHM SHOULD BE APPLIED TO MT INSTEAD OF M THE RESULT SHOULD APPROXIMATELY BE THE SAME THE 'AUTO' MODE WILL TRIGGER THE TRANSPOSITION IF MSHAPE1 MSHAPE0 SINCE THIS IMPLEMENTATION OF RANDOMIZED SVD TEND TO BE A LITTLE FASTER IN THAT CASE

CHANGED IN VERSION 018

FLIPSIGN BOOLEAN TRUE BY DEFAULT THE OUTPUT OF A SINGULAR VALUE DECOMPOSITION IS ONLY UNIQUE UP TO A PERMUTATION OF THE SIGNS OF THE SINGULAR VECTORS IF FLIPSIGN IS SET TO TRUE THE SIGN AMBIGUITY IS RESOLVED BY MAKING THE LARGEST LOADINGS FOR EACH COMPONENT IN THE LEFT SINGULAR VECTORS POSITIVE

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE THE SEED OF THE PSEUDO RANDOM NUMBER GENERATOR TO USE WHEN SHUFFLING THE DATA IF INT RANDOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NPRANDOM

NOTES

THIS ALGORITHM FINDS A USUALLY VERY GOOD APPROXIMATE TRUNCATED SINGULAR VALUE DECOMPOSITION USING RANDOMIZATION TO SPEED UP THE COMPUTATIONS IT IS PARTICULARLY FAST ON LARGE MATRICES ON WHICH YOU WISH TO EXTRACT ONLY A SMALL NUMBER OF COMPONENTS IN ORDER TO OBTAIN FURTHER SPEED UP NITER CAN BE SET 2 AT THE COST OF LOSS OF PRECISION

REFERENCES

- FINDING STRUCTURE WITH RANDOMNESS STOCHASTIC ALGORITHMS FOR CONSTRUCTING APPROXIMATE MATRIX DECOMPOSITIONS HALKO ET AL 2009 [HTTPSARXIVORGABS09094061](https://arxiv.org/abs/0909.4061)
  - A RANDOMIZED ALGORITHM FOR THE DECOMPOSITION OF MATRICES PERGUNNAR MARTINSSON VLADIMIR ROKHLIN AND MARK TYGERT
  - AN IMPLEMENTATION OF A RANDOMIZED ALGORITHM FOR PRINCIPAL COMPONENT ANALYSIS A SZLAM ET AL 2014
- 2368 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213  
63917SKLEARNUTILSEXTMATH FASTLOGDET  
SKLEARNUTILSEXTMATH FASTLOGDET A  
COMPUTE LOGDETA FOR A SYMMETRIC  
EQUIVALENT TO NPLOGNLDETA BUT MORE ROBUST IT RETURNS INF IF DETA IS NON POSITIVE OR IS NOT DEFINED  
PARAMETERS  
AARRAYLIKE THE MATRIX  
63918SKLEARNUTILSEXTMATH DENSITY  
SKLEARNUTILSEXTMATH DENSITYWKWARGS  
COMPUTE DENSITY OF A SPARSE VECTOR  
PARAMETERS  
WARRAYLIKE THE SPARSE VECTOR  
RETURNS  
FLOAT THE DENSITY OF W BETWEEN 0 AND 1  
EXAMPLES USING SKLEARNUTILSEXTMATHDENSITY  
•CLASSIFICATION OF TEXT DOCUMENTS USING SPARSE FEATURES  
63919SKLEARNUTILSEXTMATH WEIGHTEDMODE  
SKLEARNUTILSEXTMATH WEIGHTEDMODE AWAXIS0  
RETURNS AN ARRAY OF THE WEIGHTED MODAL MOST COMMON VALUE IN A  
IF THERE IS MORE THAN ONE SUCH VALUE ONLY THE FIRST IS RETURNED THE BINCOUNT FOR THE MODAL BINS IS ALSO RETURNED  
THIS IS AN EXTENSION OF THE ALGORITHM IN SCIPYSTATSMODE  
PARAMETERS  
AARRAYLIKE NDIMENSIONAL ARRAY OF WHICH TO FIND MODES  
WARRAYLIKE NDIMENSIONAL ARRAY OF WEIGHTS FOR EACH VALUE  
AXIS INT OPTIONAL AXIS ALONG WHICH TO OPERATE DEFAULT IS 0 IE THE FIRST AXIS  
RETURNS  
VALS NDARRAY ARRAY OF MODAL VALUES  
SCORE NDARRAY ARRAY OF WEIGHTED COUNTS FOR EACH MODE  
SEE ALSO  
SCIPYSTATSMODE  
EXAMPLES  
639SKLEARNUTILS UTILITIES 2369

SCIKITLEARN USER GUIDE RELEASE 0213  
 FROM SKLEARNUTILSEXTMATH IMPORT WEIGHTEDMODE  
 X 4 1 4 2 4 2  
 WEIGHTS 1 1 1 1 1 1  
 WEIGHTEDMODEX WEIGHTS  
 ARRAY4 ARRAY3  
 THE VALUE 4 APPEARS THREE TIMES WITH UNIFORM WEIGHTS THE RESULT IS SIMPLY THE MODE OF THE DISTRIBUTION  
 WEIGHTS 1 3 05 15 1 2 DEWEIGHT THE 4S  
 WEIGHTEDMODEX WEIGHTS  
 ARRAY2 ARRAY35  
 THE VALUE 2 HAS THE HIGHEST SCORE IT APPEARS TWICE WITH WEIGHTS OF 15 AND 2 THE SUM OF THESE IS 35  
 63920SKLEARNUTILS GENEVENSLICES  
 SKLEARNUTILS GENEVENSLICES NNPACKS NSAMPLESNONE  
 GENERATOR TO CREATE NPACKS SLICES GOING UP TO N  
 PARAMETERS  
 NINT  
 NPACKS INT NUMBER OF SLICES TO GENERATE  
 NSAMPLES INT OR NONE DEFAULT NONE NUMBER OF SAMPLES PASS NSAMPLES WHEN THE SLICES  
 ARE TO BE USED FOR SPARSE MATRIX INDEXING SLICING OFFTHEEND RAISES AN EXCEPTION WHILE IT  
 WORKS FOR NUMPY ARRAYS  
 YIELDS  
 SLICE  
 EXAMPLES  
 FROM SKLEARNUTILS IMPORT GENEVENSLICES  
 LISTGENEVENSLICES10 1  
 SLICE0 10 NONE  
 LISTGENEVENSLICES10 10  
 SLICE0 1 NONE SLICE1 2 NONE SLICE9 10 NONE  
 LISTGENEVENSLICES10 5  
 SLICE0 2 NONE SLICE2 4 NONE SLICE8 10 NONE  
 LISTGENEVENSLICES10 3  
 SLICE0 4 NONE SLICE4 7 NONE SLICE7 10 NONE  
 63921SKLEARNUTILSGRAPH SINGLESOURCESHORTESTPATHLENGTH  
 SKLEARNUTILSGRAPH SINGLESOURCESHORTESTPATHLENGTH GRAPH SOURCE CUT  
 OFFNONE  
 RETURN THE SHORTEST PATH LENGTH FROM SOURCE TO ALL REACHABLE NODES  
 RETURNS A DICTIONARY OF SHORTEST PATH LENGTHS KEYED BY TARGET  
 PARAMETERS  
 GRAPH SPARSE MATRIX OR 2D ARRAY PREFERABLY LIL MATRIX ADJACENCY MATRIX OF THE GRAPH  
 2370 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213  
SOURCE INTEGER STARTING NODE FOR PATH  
CUTOFF INTEGER OPTIONAL DEPTH TO STOP THE SEARCH ONLY PATHS OF LENGTH CUTOFF ARE RETURNED  
EXAMPLES  
FROM SKLEARNUTILSGRAPH IMPORT SINGLESOURCESHORTESTPATHLENGTH  
IMPORT NUMPY AS NP  
GRAPH NPARRAY 0 1 0 0  
1 0 1 0  
0 1 0 1  
0 0 1 0  
LISTSORTEDSINGLESOURCESHORTESTPATHLENGTHGRAPH 0ITEMS  
0 0 1 1 2 2 3 3  
GRAPH NPONES6 6  
LISTSORTEDSINGLESOURCESHORTESTPATHLENGTHGRAPH 2ITEMS  
0 1 1 1 2 0 3 1 4 1 5 1  
63922SKLEARNUTILSGRAPHSHORTESTPATH GRAPHSHORTESTPATH  
SKLEARNUTILSGRAPHSHORTESTPATH GRAPHSHORTESTPATH  
PERFORM A SHORTESTPATH GRAPH SEARCH ON A POSITIVE DIRECTED OR UNDIRECTED GRAPH  
PARAMETERS  
DISTMATRIX ARRAYLIKE OR SPARSE MATRIX SHAPE NN ARRAY OF POSITIVE DISTANCES IF VERTEX  
I IS CONNECTED TO VERTEX J THEN DISTMATRIXIJ GIVES THE DISTANCE BETWEEN THE VERTICES IF  
VERTEX I IS NOT CONNECTED TO VERTEX J THEN DISTMATRIXIJ 0  
DIRECTED BOOLEAN IF TRUE THEN FIND THE SHORTEST PATH ON A DIRECTED GRAPH ONLY PROGRESS FROM  
A POINT TO ITS NEIGHBORS NOT THE OTHER WAY AROUND IF FALSE THEN FIND THE SHORTEST PATH ON AN  
UNDIRECTED GRAPH THE ALGORITHM CAN PROGRESS FROM A POINT TO ITS NEIGHBORS AND VICE VERSA  
METHOD STRING 'AUTO' 'FW' 'D' METHOD TO USE OPTIONS ARE 'AUTO' ATTEMPT TO CHOOSE THE BEST  
METHOD FOR THE CURRENT PROBLEM 'FW' FLOYDWARSHALL ALGORITHM ON3 'D' DIJKSTRA'S  
ALGORITHM WITH FIBONACCI STACKS OKLOGNN2  
RETURNS  
GNPNDARRAY FLOAT SHAPE NN GIJ GIVES THE SHORTEST DISTANCE FROM POINT I TO POINT J  
ALONG THE GRAPH  
NOTES  
AS CURRENTLY IMPLEMENTED DIJKSTRA'S ALGORITHM DOES NOT WORK FOR GRAPHS WITH DIRECTIONDEPENDENT DISTANCES  
WHEN DIRECTED FALSE IE IF DISTMATRIXIJ AND DISTMATRIXJI ARE NOT EQUAL AND BOTH ARE NONZERO  
METHOD'D' WILL NOT NECESSARILY YIELD THE CORRECT RESULT  
ALSO THESE ROUTINES HAVE NOT BEEN TESTED FOR GRAPHS WITH NEGATIVE DISTANCES NEGATIVE DISTANCES CAN LEAD TO  
INFINITE CYCLES THAT MUST BE HANDLED BY SPECIALIZED ALGORITHMS  
63923SKLEARNUTILS INDEXABLE  
SKLEARNUTILS INDEXABLE ITERABLES  
MAKE ARRAYS INDEXABLE FOR CROSSVALIDATION  
639SKLEARNUTILS UTILITIES 2371

SCIKITLEARN USER GUIDE RELEASE 0213

CHECKS CONSISTENT LENGTH PASSES THROUGH NONE AND ENSURES THAT EVERYTHING CAN BE INDEXED BY CONVERTING SPARSE MATRICES TO CSR AND CONVERTING NONINTERABLE OBJECTS TO ARRAYS

PARAMETERS

ITERABLES LISTS DATAFRAMES ARRAYS SPARSE MATRICES LIST OF OBJECTS TO ENSURE SLICEABILITY

63924SKLEARNUTILSMETAESTIMATORS IFDELEGATEHASMETHOD

SKLEARNUTILSMETAESTIMATORS IFDELEGATEHASMETHOD DELEGATE

CREATE A DECORATOR FOR METHODS THAT ARE DELEGATED TO A SUBESTIMATOR

THIS ENABLES DUCKTYPING BY HASATTR RETURNING TRUE ACCORDING TO THE SUBESTIMATOR

PARAMETERS

DELEGATE STRING LIST OF STRINGS OR TUPLE OF STRINGS NAME OF THE SUBESTIMATOR THAT CAN BE ACCESED AS AN ATTRIBUTE OF THE BASE OBJECT IF A LIST OR A TUPLE OF NAMES ARE PROVIDED THE FIRST SUBESTIMATOR THAT IS AN ATTRIBUTE OF THE BASE OBJECT WILL BE USED

EXAMPLES USING SKLEARNUTILSMETAESTIMATORSIFDELEGATEHASMETHOD

- INDUCTIVE CLUSTERING

63925SKLEARNUTILSMULTICLASS TYPEOFTARGET

SKLEARNUTILSMULTICLASS TYPEOFTARGET Y

DETERMINE THE TYPE OF DATA INDICATED BY THE TARGET

NOTE THAT THIS TYPE IS THE MOST SPECIFIC TYPE THAT CAN BE INFERRED FOR EXAMPLE

- BINARY IS MORE SPECIFIC BUT COMPATIBLE WITH MULTICLASS
- MULTICLASS OF INTEGERS IS MORE SPECIFIC BUT COMPATIBLE WITH CONTINUOUS
- MULTILABELINDICATOR IS MORE SPECIFIC BUT COMPATIBLE WITH MULTICLASSMULTIOUTPUT

PARAMETERS

YARRAYLIKE

RETURNS

TARGETTYPE STRING ONE OF

- ‘CONTINUOUS’ YIS AN ARRAYLIKE OF FLOATS THAT ARE NOT ALL INTEGERS AND IS 1D OR A COLUMN VECTOR
- ‘CONTINUOUSMULTIOUTPUT’ YIS A 2D ARRAY OF FLOATS THAT ARE NOT ALL INTEGERS AND BOTH DIMENSIONS ARE OF SIZE 1
- ‘BINARY’ YCONTAINS 2 DISCRETE VALUES AND IS 1D OR A COLUMN VECTOR
- ‘MULTICLASS’ YCONTAINS MORE THAN TWO DISCRETE VALUES IS NOT A SEQUENCE OF SEQUENCES AND IS 1D OR A COLUMN VECTOR
- ‘MULTICLASSMULTIOUTPUT’ YIS A 2D ARRAY THAT CONTAINS MORE THAN TWO DISCRETE VALUES IS NOT A SEQUENCE OF SEQUENCES AND BOTH DIMENSIONS ARE OF SIZE 1
- ‘MULTILABELINDICATOR’ YIS A LABEL INDICATOR MATRIX AN ARRAY OF TWO DIMENSIONS WITH AT LEAST TWO COLUMNS AND AT MOST 2 UNIQUE VALUES

2372 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

- ‘UNKNOWN’ YIS ARRAYLIKE BUT NONE OF THE ABOVE SUCH AS A 3D ARRAY SEQUENCE OF SEQUENCES OR AN ARRAY OF NONSEQUENCE OBJECTS

EXAMPLES

```
import numpy as np
TYPEOFTARGET01 06
CONTINUOUS
TYPEOFTARGET1 1 1 1
BINARY
TYPEOFTARGETA B A
BINARY
TYPEOFTARGET10 20
BINARY
TYPEOFTARGET1 0 2
MULTICLASS
TYPEOFTARGET10 00 30
MULTICLASS
TYPEOFTARGETA B C
MULTICLASS
TYPEOFTARGETNPARRAY1 2 3 1
MULTICLASSMULTIOUTPUT
TYPEOFTARGET1 2
MULTICLASSMULTIOUTPUT
TYPEOFTARGETNPARRAY15 20 30 16
CONTINUOUSMULTIOUTPUT
TYPEOFTARGETNPARRAY0 1 1 1
MULTILABELINDICATOR
63926SKLEARNUTILSMULTICLASS ISMULTILABEL
SKLEARNUTILSMULTICLASS ISMULTILABEL Y
CHECK IFYIS IN A MULTILABEL FORMAT
PARAMETERS
YNUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES
RETURNS
OUT BOOL RETURN TRUE IFYIS IN A MULTILABEL FORMAT ELSE FALSE
EXAMPLES
import numpy as np
from SKLEARNUTILSMULTICLASS import ISMULTILABEL
ISMULTILABEL0 1 0 1
FALSE
ISMULTILABEL1 0 2
FALSE
ISMULTILABELNPARRAY1 0 0 0
TRUE
ISMULTILABELNPARRAY1 0 0
FALSE
639SKLEARNUTILS UTILITIES 2373
```

SCIKITLEARN USER GUIDE RELEASE 0213

ISMULTILABELNPARRAY1 0 0

TRUE

63927SKLEARNUTILSMULTICLASS UNIQUELABELS

SKLEARNUTILSMULTICLASS UNIQUELABELS YS

EXTRACT AN ORDERED ARRAY OF UNIQUE LABELS

WE DON'T ALLOW

- MIX OF MULTILABEL AND MULTICLASS SINGLE LABEL TARGETS
- MIX OF LABEL INDICATOR MATRIX AND ANYTHING ELSE BECAUSE THERE ARE NO EXPLICIT LABELS
- MIX OF LABEL INDICATOR MATRICES OF DIFFERENT SIZES
- MIX OF STRING AND INTEGER LABELS

AT THE MOMENT WE ALSO DON'T ALLOW "MULTICLASSMULTIOUTPUT" INPUT TYPE

PARAMETERS

YS ARRAYLIKES

RETURNS

OUT NUMPY ARRAY OF SHAPE NUNIQUELABELS AN ORDERED ARRAY OF UNIQUE LABELS

EXAMPLES

FROM SKLEARNUTILSMULTICLASS IMPORT UNIQUELABELS

UNIQUELABELS3 5 5 5 7 7

ARRAY3 5 7

UNIQUELABELS1 2 3 4 2 2 3 4

ARRAY1 2 3 4

UNIQUELABELS1 2 10 5 11

ARRAY 1 2 5 10 11

EXAMPLES USING SKLEARNUTILSMULTICLASSUNIQUELABELS

- CONFUSION MATRIX

63928SKLEARNUTILS MURMURHASH332

SKLEARNUTILS MURMURHASH332

COMPUTE THE 32BIT MURMURHASH3 OF KEY AT SEED

THE UNDERLYING IMPLEMENTATION IS MURMURHASH3X8632 GENERATING LOW LATENCY 32BITS HASH SUITABLE FOR IM

PLEMENTING LOOKUP TABLES BLOOM FILTERS COUNT MIN SKETCH OR FEATURE HASHING

PARAMETERS

KEY INT32 BYTES UNICODE OR NDARRAY WITH DTYPE INT32 THE PHYSICAL OBJECT TO HASH

SEED INT OPTIONAL DEFAULT IS 0 INTEGER SEED FOR THE HASHING ALGORITHM

POSITIVE BOOLEAN OPTIONAL DEFAULT IS FALSE

2374 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213  
 TRUE THE RESULTS IS CASTED TO AN UNSIGNED INT FROM 0 TO 2<sup>32</sup> - 1  
 FALSE THE RESULTS IS CASTED TO A SIGNED INT FROM 2<sup>31</sup> TO 2<sup>31</sup> - 1  
 63929SKLEARNUTILS RESAMPLE  
 SKLEARNUTILS RESAMPLE ARRAYS OPTIONS  
 RESAMPLE ARRAYS OR SPARSE MATRICES IN A CONSISTENT WAY  
 THE DEFAULT STRATEGY IMPLEMENTS ONE STEP OF THE BOOTSTRAPPING PROCEDURE  
 PARAMETERS  
 ARRAYS SEQUENCE OF INDEXABLE DATASTRUCTURES INDEXABLE DATASTRUCTURES CAN BE ARRAYS LISTS  
 DATAFRAMES OR SCIPY SPARSE MATRICES WITH CONSISTENT FIRST DIMENSION  
 RETURNS  
 RESAMPLEDARRAYS SEQUENCE OF INDEXABLE DATASTRUCTURES SEQUENCE OF RESAMPLED COPIES OF  
 THE COLLECTIONS THE ORIGINAL ARRAYS ARE NOT IMPACTED  
 OTHER PARAMETERS  
 REPLACE BOOLEAN TRUE BY DEFAULT IMPLEMENTS RESAMPLING WITH REPLACEMENT IF FALSE THIS WILL  
 IMPLEMENT SLICED RANDOM PERMUTATIONS  
 NSAMPLES INT NONE BY DEFAULT NUMBER OF SAMPLES TO GENERATE IF LEFT TO NONE THIS IS AUTO  
 MATICALLY SET TO THE FIRST DIMENSION OF THE ARRAYS IF REPLACE IS FALSE IT SHOULD NOT BE LARGER  
 THAN THE LENGTH OF ARRAYS  
 RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE THE SEED OF THE  
 PSEUDO RANDOM NUMBER GENERATOR TO USE WHEN SHUFFLING THE DATA IF INT RANDOMSTATE IS  
 THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS  
 THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE  
 INSTANCE USED BY NPRANDOM  
 STRATIFY ARRAYLIKE OR NONE DEFAULTNONE IF NOT NONE DATA IS SPLIT IN A STRATIFIED FASHION  
 USING THIS AS THE CLASS LABELS  
 SEE ALSO  
 SKLEARNUTILSSHUFFLE  
 EXAMPLES  
 IT IS POSSIBLE TO MIX SPARSE AND DENSE ARRAYS IN THE SAME RUN  
 X NPARRAY1 0 2 1 0 0  
 Y NPARRAY0 1 2  
 FROM SCIPYSPARSE IMPORT COOMATRIX  
 XSPARSE COOMATRIX  
 FROM SKLEARNUTILS IMPORT RESAMPLE  
 X XSPARSE Y RESAMPLEX XSPARSE Y RANDOMSTATE0  
 X  
 ARRAY1 0  
 2 1  
 1 0  
 639SKLEARNUTILS UTILITIES 2375

SCIKITLEARN USER GUIDE RELEASE 0213

XSPARSE

3X2 SPARSE MATRIX OF TYPE NUMPYFLOAT64

WITH 4 STORED ELEMENTS IN COMPRESSED SPARSE ROW FORMAT

XSPARSETOARRAY

ARRAY1 0

2 1

1 0

Y

ARRAY0 1 0

RESAMPLEY NSAMPLES2 RANDOMSTATE0

ARRAY0 1

EXAMPLE USING STRATIFICATION

Y 0 0 1 1 1 1 1 1 1

RESAMPLEY NSAMPLES5 REPLACE FALSE STRATIFY

RANDOMSTATE0

1 1 1 0 1

63930SKLEARNUTILS SAFEINDEXING

SKLEARNUTILS SAFEINDEXING XINDICES

RETURN ITEMS OR ROWS FROM X USING INDICES

ALLOWS SIMPLE INDEXING OF LISTS OR ARRAYS

PARAMETERS

XARRAYLIKE SPARSEMATRIX LIST PANDASDATAFRAME PANDASSERIES DATA FROM WHICH TO SAMPLE ROWS OR ITEMS

INDICES ARRAYLIKE OF INT INDICES ACCORDING TO WHICH X WILL BE SUBSAMPLED

RETURNS

SUBSET SUBSET OF X ON FIRST AXIS

NOTES

CSR CSC AND LIL SPARSE MATRICES ARE SUPPORTED COO SPARSE MATRICES ARE NOT SUPPORTED

63931SKLEARNUTILS SAFEMASK

SKLEARNUTILS SAFEMASK XMASK

RETURN A MASK WHICH IS SAFE TO USE ON X

PARAMETERS

XARRAYLIKE SPARSE MATRIX DATA ON WHICH TO APPLY MASK

MASK ARRAY MASK TO BE USED ON X

RETURNS

2376 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

MASK

63932SKLEARNUTILS SAFESQR

SKLEARNUTILS SAFESQR XCOPYTRUE

ELEMENT WISE SQUARING OF ARRAYLIKES AND SPARSE MATRICES

PARAMETERS

XARRAY LIKE MATRIX SPARSE MATRIX

COPY BOOLEAN OPTIONAL DEFAULT TRUE WHETHER TO CREATE A COPY OF X AND OPERATE ON IT OR TO

PERFORM INPLACE COMPUTATION DEFAULT BEHAVIOUR

RETURNS

X 2 ELEMENT WISE SQUARE

63933SKLEARNUTILS SHUFFLE

SKLEARNUTILS SHUFFLEARRAYS OPTIONS

SHUFFLE ARRAYS OR SPARSE MATRICES IN A CONSISTENT WAY

THIS IS A CONVENIENCE ALIAS TO RESAMPLE ARRAYS REPLACEFALSE TO DO RANDOM PERMUTATIONS OF THE

COLLECTIONS

PARAMETERS

ARRAYS SEQUENCE OF INDEXABLE DATASTRUCTURES INDEXABLE DATASTRUCTURES CAN BE ARRAYS LISTS

DATAFRAMES OR SCIPY SPARSE MATRICES WITH CONSISTENT FIRST DIMENSION

RETURNS

SHUFFLEDARRAYS SEQUENCE OF INDEXABLE DATASTRUCTURES SEQUENCE OF SHUFFLED COPIES OF THE COL

LECTIONS THE ORIGINAL ARRAYS ARE NOT IMPACTED

OTHER PARAMETERS

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE THE SEED OF THE

PSEUDO RANDOM NUMBER GENERATOR TO USE WHEN SHUFFLING THE DATA IF INT RANDOMSTATE IS

THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS

THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE

INSTANCE USED BY NPRANDOM

NSAMPLES INT NONE BY DEFAULT NUMBER OF SAMPLES TO GENERATE IF LEFT TO NONE THIS IS AUTO

MATICALLY SET TO THE FIRST DIMENSION OF THE ARRAYS

SEE ALSO

SKLEARNUTILSRESAMPLE

EXAMPLES

IT IS POSSIBLE TO MIX SPARSE AND DENSE ARRAYS IN THE SAME RUN

6395SKLEARNUTILS UTILITIES 2377

SCIKITLEARN USER GUIDE RELEASE 0213

X NPARRAY1 0 2 1 0 0

Y NPARRAY0 1 2

FROM SCIPYSPARSE IMPORT COOMATRIX

XSPARSE COOMATRIX

FROM SKLEARNUTILS IMPORT SHUFFLE

X XSPARSE Y SHUFFLEX XSPARSE Y RANDOMSTATE0

X

ARRAY0 0

2 1

1 0

XSPARSE

3X2 SPARSE MATRIX OF TYPE NUMPYFLOAT64

WITH 3 STORED ELEMENTS IN COMPRESSED SPARSE ROW FORMAT

XSPARSETOARRAY

ARRAY0 0

2 1

1 0

Y

ARRAY2 1 0

SHUFFLEY NSAMPLES2 RANDOMSTATE0

ARRAY0 1

EXAMPLES USING SKLEARNUTILSSHUFFLE

- MODEL COMPLEXITY INFLUENCE
- PREDICTION LATENCY
- COLOR QUANTIZATION USING KMEANS
- EMPIRICAL EVALUATION OF THE IMPACT OF KMEANS INITIALIZATION
- GRADIENT BOOSTING REGRESSION
- EARLY STOPPING OF STOCHASTIC GRADIENT DESCENT

63934SKLEARNUTILSSPARSEFUNCS INCRMEANVARIANCEAXIS

SKLEARNUTILSSPARSEFUNCS INCRMEANVARIANCEAXIS XAXISLASTMEAN LASTVAR LASTN

COMPUTE INCREMENTAL MEAN AND VARIANCE ALONG AN AXIX ON A CSR OR CSC MATRIX

LASTMEAN LASTVAR ARE THE STATISTICS COMPUTED AT THE LAST STEP BY THIS FUNCTION BOTH MUST BE INITIALIZED TO 0

ARRAYS OF THE PROPER SIZE IE THE NUMBER OF FEATURES IN X LASTN IS THE NUMBER OF SAMPLES ENCOUNTERED UNTIL

NOW

PARAMETERS

XCSR OR CSC SPARSE MATRIX SHAPE NSAMPLES NFEATURES INPUT DATA

AXIS INT EITHER 0 OR 1 AXIS ALONG WHICH THE AXIS SHOULD BE COMPUTED

LASTMEAN FLOAT ARRAY WITH SHAPE NFEATURES ARRAY OF FEATUREWISE MEANS TO UPDATE WITH THE

NEW DATA X

2378 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

LASTVAR FLOAT ARRAY WITH SHAPE NFEATURES ARRAY OF FEATUREWISE VAR TO UPDATE WITH THE NEW DATA X

LASTN INT WITH SHAPE NFEATURES NUMBER OF SAMPLES SEEN SO FAR EXCLUDED X

RETURNS

MEANS FLOAT ARRAY WITH SHAPE NFEATURES UPDATED FEATUREWISE MEANS

VARIANCES FLOAT ARRAY WITH SHAPE NFEATURES UPDATED FEATUREWISE VARIANCES

NINT WITH SHAPE NFEATURES UPDATED NUMBER OF SEEN SAMPLES

NOTES

NANS ARE IGNORED IN THE ALGORITHM

63935SKLEARNUTILSSPARSEFUNCS INPLACECOLUMNSCALE

SKLEARNUTILSSPARSEFUNCS INPLACECOLUMNSCALE XSCALE

INPLACE COLUMN SCALING OF A CSCCSR MATRIX

SCALE EACH FEATURE OF THE DATA MATRIX BY MULTIPLYING WITH SPECIFIC SCALE PROVIDED BY THE CALLER ASSUMING A NSAMPLES NFEATURES SHAPE

PARAMETERS

XCSR OR CSR MATRIX WITH SHAPE NSAMPLES NFEATURES MATRIX TO NORMALIZE USING THE VARIANCE OF THE FEATURES

SCALE FLOAT ARRAY WITH SHAPE NFEATURES ARRAY OF PRECOMPUTED FEATUREWISE VALUES TO USE FOR SCALING

63936SKLEARNUTILSSPARSEFUNCS INPLACEROWSCALE

SKLEARNUTILSSPARSEFUNCS INPLACEROWSCALE XSCALE

INPLACE ROW SCALING OF A CSR OR CSC MATRIX

SCALE EACH ROW OF THE DATA MATRIX BY MULTIPLYING WITH SPECIFIC SCALE PROVIDED BY THE CALLER ASSUMING A NSAMPLES NFEATURES SHAPE

PARAMETERS

XCSR OR CSC SPARSE MATRIX SHAPE NSAMPLES NFEATURES MATRIX TO BE SCALED

SCALE FLOAT ARRAY WITH SHAPE NFEATURES ARRAY OF PRECOMPUTED SAMPLEWISE VALUES TO USE FOR SCALING

63937SKLEARNUTILSSPARSEFUNCS INPLACESWAPROW

SKLEARNUTILSSPARSEFUNCS INPLACESWAPROW XMN

SWAPS TWO ROWS OF A CSCCSR MATRIX INPLACE

PARAMETERS

XCSR OR CSC SPARSE MATRIX SHAPENSAMPLES NFEATURES MATRIX WHOSE TWO ROWS ARE TO BE SWAPPED

639SKLEARNUTILS UTILITIES 2379

SCIKITLEARN USER GUIDE RELEASE 0213

MINT INDEX OF THE ROW OF X TO BE SWAPPED

NINT INDEX OF THE ROW OF X TO BE SWAPPED

63938SKLEARNUTILSSPARSEFUNCS INPLACESWAPCOLUMN

SKLEARNUTILSSPARSEFUNCS INPLACESWAPCOLUMN XMN

SWAPS TWO COLUMNS OF A CSCCSR MATRIX INPLACE

PARAMETERS

XCSR OR CSC SPARSE MATRIX SHAPENSAMPLES NFEATURES MATRIX WHOSE TWO COLUMNS ARE TO BE SWAPPED

MINT INDEX OF THE COLUMN OF X TO BE SWAPPED

NINT INDEX OF THE COLUMN OF X TO BE SWAPPED

63939SKLEARNUTILSSPARSEFUNCS MEANVARIANCEAXIS

SKLEARNUTILSSPARSEFUNCS MEANVARIANCEAXIS XAXIS

COMPUTE MEAN AND VARIANCE ALONG AN AXIX ON A CSR OR CSC MATRIX

PARAMETERS

XCSR OR CSC SPARSE MATRIX SHAPE NSAMPLES NFEATURES INPUT DATA

AXIS INT EITHER 0 OR 1 AXIS ALONG WHICH THE AXIS SHOULD BE COMPUTED

RETURNS

MEANS FLOAT ARRAY WITH SHAPE NFEATURES FEATUREWISE MEANS

VARIANCES FLOAT ARRAY WITH SHAPE NFEATURES FEATUREWISE VARIANCES

63940SKLEARNUTILSSPARSEFUNCS INPLACECSRCOLUMNSCALE

SKLEARNUTILSSPARSEFUNCS INPLACECSRCOLUMNSCALE XSCALE

INPLACE COLUMN SCALING OF A CSR MATRIX

SCALE EACH FEATURE OF THE DATA MATRIX BY MULTIPLYING WITH SPECIFIC SCALE PROVIDED BY THE CALLER ASSUMING A NSAMPLES NFEATURES SHAPE

PARAMETERS

XCSR MATRIX WITH SHAPE NSAMPLES NFEATURES MATRIX TO NORMALIZE USING THE VARIANCE OF THE FEATURES

SCALE FLOAT ARRAY WITH SHAPE NFEATURES ARRAY OF PRECOMPUTED FEATUREWISE VALUES TO USE FOR SCALING

63941SKLEARNUTILSSPARSEFUNCSFAST INPLACECSRROWNORMALIZEL1

SKLEARNUTILSSPARSEFUNCSFAST INPLACECSRROWNORMALIZEL1

INPLACE ROW NORMALIZE USING THE L1 NORM

2380 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

63942SKLEARNUTILSSPARSEFUNCSFAST INPLACECSRROWNORMALIZEL2  
SKLEARNUTILSSPARSEFUNCSFAST INPLACECSRROWNORMALIZEL2  
INPLACE ROW NORMALIZE USING THE L2 NORM

63943SKLEARNUTILSRANDOM SAMPLEWITHOUTREPLACEMENT  
SKLEARNUTILSRANDOM SAMPLEWITHOUTREPLACEMENT  
SAMPLE INTEGERS WITHOUT REPLACEMENT  
SELECT NSAMPLES INTEGERS FROM THE SET 0 NPOPULATION WITHOUT REPLACEMENT

PARAMETERS

NPOPULATION INT THE SIZE OF THE SET TO SAMPLE FROM  
NSAMPLES INT THE NUMBER OF INTEGER TO SAMPLE  
RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULTNONE IF INT RAN  
DOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE  
RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS  
THE RANDOMSTATE INSTANCE USED BY NPRANDOM  
METHOD “AUTO” “TRACKINGSELECTION” “RESERVOIRSAMPLING” OR “POOL” IF METHOD “AUTO”  
THE RATIO OF NSAMPLES NPOPULATION IS USED TO DETERMINE WHICH ALGORITHM TO USE IF RA  
TIO IS BETWEEN 0 AND 001 TRACKING SELECTION IS USED IF RATIO IS BETWEEN 001 AND 099  
NUMPYRANDOMPERMUTATION IS USED IF RATIO IS GREATER THAN 099 RESERVOIR SAMPLING IS USED  
THE ORDER OF THE SELECTED INTEGERS IS UNDEFINED IF A RANDOM ORDER IS DESIRED THE SELECTED  
SUBSET SHOULD BE SHUFFLED  
IF METHOD “TRACKINGSELECTION” A SET BASED IMPLEMENTATION IS USED WHICH IS SUITABLE FOR  
NSAMPLES NPOPULATION  
IF METHOD “RESERVOIRSAMPLING” A RESERVOIR SAMPLING ALGORITHM IS USED WHICH IS SUITABLE  
FOR HIGH MEMORY CONSTRAINT OR WHEN O NSAMPLES ONPOPULATION THE ORDER OF THE  
SELECTED INTEGERS IS UNDEFINED IF A RANDOM ORDER IS DESIRED THE SELECTED SUBSET SHOULD BE  
SHUFFLED  
IF METHOD “POOL” A POOL BASED ALGORITHM IS PARTICULARLY FAST EVEN FASTER THAN THE TRACKING  
SELECTION METHOD HOVEWER A VECTOR CONTAINING THE ENTIRE POPULATION HAS TO BE INITIALIZED  
IF NSAMPLES NPOPULATION THE RESERVOIR SAMPLING METHOD IS FASTER

RETURNS

OUT ARRAY OF SIZE NSAMPLES THE SAMPLED SUBSETS OF INTEGER THE SUBSET OF SELECTED INTEGER  
MIGHT NOT BE RANDOMIZED SEE THE METHOD ARGUMENT

63944SKLEARNUTILSVALIDATION CHECKISFITTED  
SKLEARNUTILSVALIDATION CHECKISFITTED ESTIMATOR ATTRIBUTES MSGNONE  
ALLORANYBUILTIN FUNCTION ALL  
PERFORM ISFITTED VALIDATION FOR ESTIMATOR  
CHECKS IF THE ESTIMATOR IS FITTED BY VERIFYING THE PRESENCE OF “ALLORANY” OF THE PASSED ATTRIBUTES AND RAISES A  
NOTFITTEDERROR WITH THE GIVEN MESSAGE

PARAMETERS

ESTIMATOR ESTIMATOR INSTANCE ESTIMATOR INSTANCE FOR WHICH THE CHECK IS PERFORMED

6395SKLEARNUTILS UTILITIES 2381

SCIKITLEARN USER GUIDE RELEASE 0213

ATTRIBUTES ATTRIBUTE NAMES GIVEN AS STRING OR A LISTTUPLE OF STRINGS

EGCOEF ESTIMATOR COEF

MSG STRING THE DEFAULT ERROR MESSAGE IS “THIS NAMES INSTANCE IS NOT FITTED YET CALL ‘FIT’ WITH APPROPRIATE ARGUMENTS BEFORE USING THIS METHOD”

FOR CUSTOM MESSAGES IF “NAMES” IS PRESENT IN THE MESSAGE STRING IT IS SUBSTITUTED FOR THE ESTIMATOR NAME

EG “ESTIMATOR NAMES MUST BE FITTED BEFORE SPARSIFYING”

ALLORANY CALLABLE ALL ANY DEFAULT ALL SPECIFY WHETHER ALL OR ANY OF THE GIVEN ATTRIBUTES MUST EXIST

RETURNS

NONE

RAISES

NOTFITTEDERROR IF THE ATTRIBUTES ARE NOT FOUND

63945SKLEARNUTILSVALIDATION CHECKMEMORY

SKLEARNUTILSVALIDATION CHECKMEMORY MEMORY

CHECK THAT MEMORY IS JOBLIBMEMORYLIKE

JOBLIBMEMORYLIKE MEANS THAT MEMORY CAN BE CONVERTED INTO A JOBLIBMEMORY INSTANCE TYPICALLY A STR DENOTING THELOCATION OR HAS THE SAME INTERFACE HAS A CACHE METHOD

PARAMETERS

MEMORY NONE STR OR OBJECT WITH THE JOBLIBMEMORY INTERFACE

RETURNS

MEMORY OBJECT WITH THE JOBLIBMEMORY INTERFACE

RAISES

VALUEERROR IFMEMORY IS NOT JOBLIBMEMORYLIKE

63946SKLEARNUTILSVALIDATION CHECKSYMMETRIC

SKLEARNUTILSVALIDATION CHECKSYMMETRIC ARRAY TOL1E10 RAISEWARNINGTRUE

RAISEEXCEPTIONFALSE

MAKE SURE THAT ARRAY IS 2D SQUARE AND SYMMETRIC

IF THE ARRAY IS NOT SYMMETRIC THEN A SYMMETRIZED VERSION IS RETURNED OPTIONALLY A WARNING OR EXCEPTION IS RAISED IF THE MATRIX IS NOT SYMMETRIC

PARAMETERS

ARRAY NDARRAY OR SPARSE MATRIX INPUT OBJECT TO CHECK CONVERT MUST BE TWODIMENSIONAL AND SQUARE OTHERWISE A VALUEERROR WILL BE RAISED

TOLFLOAT ABSOLUTE TOLERANCE FOR EQUIVALENCE OF ARRAYS DEFAULT 1E10

RAISEWARNING BOOLEAN DEFAULTTRUE IF TRUE THEN RAISE A WARNING IF CONVERSION IS REQUIRED

RAISEEXCEPTION BOOLEAN DEFAULTFALSE IF TRUE THEN RAISE AN EXCEPTION IF ARRAY IS NOT SYM

METRIC

2382 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

RETURNS  
ARRAYSYM NDARRAY OR SPARSE MATRIX SYMMETRIZED VERSION OF THE INPUT ARRAY IE THE AVERAGE  
OF ARRAY AND ARRAYTRANSPPOSE IF SPARSE THEN DUPLICATE ENTRIES ARE FIRST SUMMED AND ZEROS  
ARE ELIMINATED

63947SKLEARNUTILSVALIDATION COLUMNOR1D  
SKLEARNUTILSVALIDATION COLUMNOR1D YWARNFALSE  
RAVEL COLUMN OR 1D NUMPY ARRAY ELSE RAISES AN ERROR

PARAMETERS  
YARRAYLIKE  
WARN BOOLEAN DEFAULT FALSE TO CONTROL DISPLAY OF WARNINGS

RETURNS  
YARRAY  
63948SKLEARNUTILSVALIDATION HASFITPARAMETER  
SKLEARNUTILSVALIDATION HASFITPARAMETER ESTIMATOR PARAMETER  
CHECKS WHETHER THE ESTIMATOR'S FIT METHOD SUPPORTS THE GIVEN PARAMETER

PARAMETERS  
ESTIMATOR OBJECT AN ESTIMATOR TO INSPECT  
PARAMETER STR THE SEARCHED PARAMETER  
RETURNS  
ISPARAMETER BOOL WHETHER THE PARAMETER WAS FOUND TO BE A NAMED PARAMETER OF THE ESTIMA  
TOR'S FIT METHOD

EXAMPLES  
FROM SKLEARN SVM IMPORT SVC  
HASFITPARAMETERSVC SAMPLEWEIGHT  
TRUE  
63949SKLEARNUTILSTESTING ASSERTIN  
SKLEARNUTILSTESTING ASSERTIN MEMBER CONTAINER MSGNONE  
JUST LIKE SELFASSERTTRUEA IN B BUT WITH A NICER DEFAULT MESSAGE  
63950SKLEARNUTILSTESTING ASSERTNOTIN  
SKLEARNUTILSTESTING ASSERTNOTIN MEMBER CONTAINER MSGNONE  
JUST LIKE SELFASSERTTRUEA NOT IN B BUT WITH A NICER DEFAULT MESSAGE  
639SKLEARNUTILS UTILITIES 2383

SCIKITLEARN USER GUIDE RELEASE 0213

63951SKLEARNUTILSTESTING ASSERTRAISEMESSAGE  
SKLEARNUTILSTESTING ASSERTRAISEMESSAGE EXCEPTIONS MESSAGE FUNCTION ARGS  
KWARGS  
HELPER FUNCTION TO TEST THE MESSAGE RAISED IN AN EXCEPTION  
GIVEN AN EXCEPTION A CALLABLE TO RAISE THE EXCEPTION AND A MESSAGE STRING TESTS THAT THE CORRECT EXCEPTION IS  
RAISED AND THAT THE MESSAGE IS A SUBSTRING OF THE ERROR THROWN USED TO TEST THAT THE SPECIFIC MESSAGE THROWN  
DURING AN EXCEPTION IS CORRECT  
PARAMETERS  
EXCEPTIONS EXCEPTION OR TUPLE OF EXCEPTION AN EXCEPTION OBJECT  
MESSAGE STR THE ERROR MESSAGE OR A SUBSTRING OF THE ERROR MESSAGE  
FUNCTION CALLABLE CALLABLE OBJECT TO RAISE ERROR  
ARGS THE POSITIONAL ARGUMENTS TO FUNCTION  
KWARGS THE KEYWORD ARGUMENTS TO FUNCTION

63952SKLEARNUTILSTESTING ALLESTIMATORS  
SKLEARNUTILSTESTING ALLESTIMATORS INCLUDEMETAESTIMATORSNONE INCLUDEOTHERNONE  
TYPEFILTERNONE INCLUDEDONTTESTNONE  
GET A LIST OF ALL ESTIMATORS FROM SKLEARN  
THIS FUNCTION CRAWLS THE MODULE AND GETS ALL CLASSES THAT INHERIT FROM BASEESTIMATOR CLASSES THAT ARE DEFINED IN  
TESTMODULES ARE NOT INCLUDED BY DEFAULT METAESTIMATORS SUCH AS GRIDSEARCHCV ARE ALSO NOT INCLUDED  
PARAMETERS  
INCLUDEMETAESTIMATORS BOOLEAN DEFAULTFALSE DEPRECATED IGNORED DEPRECATED 021  
INCLUDEMETAESTIMATORS HAS BEEN DEPRECATED AND HAS NO EFFECT IN 021 AND  
WILL BE REMOVED IN 023  
INCLUDEOTHER BOOLEAN DEFAULTFALSE DEPRECATED IGNORED DEPRECATED 021  
INCLUDEOTHER HAS BEEN DEPRECATED AND HAS NOT EFFECT IN 021 AND WILL BE REMOVED  
IN 023  
TYPEFILTER STRING LIST OF STRING OR NONE DEFAULTNONE WHICH KIND OF ESTIMATORS SHOULD BE  
RETURNED IF NONE NO FILTER IS APPLIED AND ALL ESTIMATORS ARE RETURNED POSSIBLE VALUES ARE  
'CLASSIFIER' 'REGRESSOR' 'CLUSTER' AND 'TRANSFORMER' TO GET ESTIMATORS ONLY OF THESE SPECIFIC  
TYPES OR A LIST OF THESE TO GET THE ESTIMATORS THAT FIT AT LEAST ONE OF THE TYPES  
INCLUDEDONTTEST BOOLEAN DEFAULTFALSE DEPRECATED IGNORED DEPRECATED 021  
INCLUDEDONTTEST HAS BEEN DEPRECATED AND HAS NO EFFECT IN 021 AND WILL BE  
REMOVED IN 023  
RETURNS  
ESTIMATORS LIST OF TUPLES LIST OF NAME CLASS WHERE NAME IS THE CLASS NAME AS STRING AND  
CLASS IS THE ACTUALL TYPE OF THE CLASS  
UTILITIES FROM JOBLIB

2384 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

UTILSPARALLELBACKEND BACKEND NJOBS CHANGE THE DEFAULT BACKEND USED BY PARALLEL INSIDE A WITH BLOCK

UTILSREGISTERPARALLELBACKEND NAME FAC

TORYREGISTER A NEW PARALLEL BACKEND FACTORY

63953SKLEARNUTILS PARALLELBACKEND

SKLEARNUTILS PARALLELBACKEND BACKEND NJOBS1 BACKENDPARAMS

CHANGE THE DEFAULT BACKEND USED BY PARALLEL INSIDE A WITH BLOCK

IFBACKEND IS A STRING IT MUST MATCH A PREVIOUSLY REGISTERED IMPLEMENTATION USING THE REGISTERPARALLELBACKEND FUNCTION

BY DEFAULT THE FOLLOWING BACKENDS ARE AVAILABLE

- ‘LOKY’ SINGLEHOST PROCESSBASED PARALLELISM USED BY DEFAULT
- ‘THREADING’ SINGLEHOST THREADBASED PARALLELISM
- ‘MULTIPROCESSING’ LEGACY SINGLEHOST PROCESSBASED PARALLELISM

‘LOKY’ IS RECOMMENDED TO RUN FUNCTIONS THAT MANIPULATE PYTHON OBJECTS ‘THREADING’ IS A LOWOVERHEAD ALTERNATIVE THAT IS MOST EFFICIENT FOR FUNCTIONS THAT RELEASE THE GLOBAL INTERPRETER LOCK EG IOBOUND CODE OR CPUBOUND CODE IN A FEW CALLS TO NATIVE CODE THAT EXPLICITLY RELEASES THE GIL

IN ADDITION IF THE DASK ANDDISTRIBUTED PYTHON PACKAGES ARE INSTALLED IT IS POSSIBLE TO USE THE ‘DASK’ BACKEND FOR BETTER SCHEDULING OF NESTED PARALLEL CALLS WITHOUT OVERSUBSCRIPTION AND POTENTIALLY DISTRIBUTE PARALLEL CALLS OVER A NETWORKED CLUSTER OF SEVERAL HOSTS

ALTERNATIVELY THE BACKEND CAN BE PASSED DIRECTLY AS AN INSTANCE

BY DEFAULT ALL AVAILABLE WORKERS WILL BE USED NJOBS1 UNLESS THE CALLER PASSES AN EXPLICIT VALUE FOR THE NJOBS PARAMETER

THIS IS AN ALTERNATIVE TO PASSING A BACKENDBACKENDNAME ARGUMENT TO THE PARALLEL CLASS CONSTRUC TOR IT IS PARTICULARLY USEFUL WHEN CALLING INTO LIBRARY CODE THAT USES JOBLIB INTERNALLY BUT DOES NOT EXPOSE THE BACKEND ARGUMENT IN ITS OWN API

```
FROM OPERATOR IMPORT NEG
WITH PARALLELBACKENDTHREADING
PRINTPARALLELDELAYEDNEGI 1 FORIINRANGE5
```

1 2 3 4 5

WARNING THIS FUNCTION IS EXPERIMENTAL AND SUBJECT TO CHANGE IN A FUTURE VERSION OF JOBLIB

NEW IN VERSION 010

63954SKLEARNUTILS REGISTERPARALLELBACKEND

SKLEARNUTILS REGISTERPARALLELBACKEND NAME FACTORY MAKEDEFAULTFALSE

REGISTER A NEW PARALLEL BACKEND FACTORY

THE NEW BACKEND CAN THEN BE SELECTED BY PASSING ITS NAME AS THE BACKEND ARGUMENT TO THE PARALLEL CLASS

MOREOVER THE DEFAULT BACKEND CAN BE OVERWRITTEN GLOBALLY BY SETTING MAKEDEFAULTTRUE

THE FACTORY CAN BE ANY CALLABLE THAT TAKES NO ARGUMENT AND RETURN AN INSTANCE OF PARALLELBACKENDBASE

WARNING THIS FUNCTION IS EXPERIMENTAL AND SUBJECT TO CHANGE IN A FUTURE VERSION OF JOBLIB

639SKLEARNUTILS UTILITIES 2385

SCIKITLEARN USER GUIDE RELEASE 0213  
NEW IN VERSION 010  
640 RECENTLY DEPRECATED  
6401 TO BE REMOVED IN 023  
UTILSMEMORY ARGS KWARGS  
ATTRIBUTES  
UTILSPARALLEL ARGS KWARGS  
METHODS  
SKLEARNUTILS MEMORY  
WARNING DEPRECATED  
CLASSSSKLEARNUTILS MEMORYARGS KWARGS  
ATTRIBUTES  
CACHEDIR  
METHODS  
CACHE SELF FUNC IGNORE VERBOSE MMAPMODE DECORATES THE GIVEN FUNCTION FUNC TO ONLY COMPUTE ITS  
RETURN VALUE FOR INPUT ARGUMENTS NOT CACHED ON DISK  
CLEAR SELF WARN ERASE THE COMPLETE CACHE DIRECTORY  
EVAL SELF FUNC ARGS KWARGS EVAL FUNCTION FUNC WITH ARGUMENTS ARGS AND  
KWARGS IN THE CONTEXT OF THE MEMORY  
FORMAT SELF OBJ INDENT RETURN THE FORMATTED REPRESENTATION OF THE OBJECT  
REDUCESIZE SELF REMOVE CACHE ELEMENTS TO MAKE CACHE SIZE FIT IN  
BYTESLIMIT  
DEBUG  
WARN  
INIT ARGS KWARGS  
DEPRECATED DEPRECATED IN VERSION 0201 TO BE REMOVED IN VERSION 023 PLEASE IMPORT THIS FUNCTION  
ALITY DIRECTLY FROM JOBLIB WHICH CAN BE INSTALLED WITH PIP INSTALL JOBLIB  
CACHESELFUNCFUNCNONE IGNORENONE VERBOSENONE MMAPMODEFALSE  
DECORATES THE GIVEN FUNCTION FUNC TO ONLY COMPUTE ITS RETURN VALUE FOR INPUT ARGUMENTS NOT CACHED ON DISK  
PARAMETERS  
2386 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

FUNC CALLABLE OPTIONAL THE FUNCTION TO BE DECORATED

IGNORE LIST OF STRINGS A LIST OF ARGUMENTS NAME TO IGNORE IN THE HASHING

VERBOSE INTEGER OPTIONAL THE VERBOSITY MODE OF THE FUNCTION BY DEFAULT THAT OF THE MEMORY OBJECT IS USED

MMAPMODE NONE 'R' 'R' 'W' 'C' OPTIONAL THE MEMMAPPING MODE USED WHEN

LOADING FROM CACHE NUMPY ARRAYS SEE NUMPYLOAD FOR THE MEANING OF THE ARGUMENTS

BY DEFAULT THAT OF THE MEMORY OBJECT IS USED

RETURNS

DECORATEDFUNC MEMORIZEDFUNC OBJECT THE RETURNED OBJECT IS A MEMORIZEDFUNC OBJECT THAT IS CALLABLE BEHAVES LIKE A FUNCTION BUT OFFERS EXTRA METHODS FOR CACHE LOOKUP AND MANAGEMENT SEE THE DOCUMENTATION FOR JOBLIBMEMORYMEMORIZEDFUNC

CLEARSELFWARNTRUE

ERASE THE COMPLETE CACHE DIRECTORY

EVALSELFFUNCARGS KWARGS

EVAL FUNCTION FUNC WITH ARGUMENTS ARGS ANDKWARGS IN THE CONTEXT OF THE MEMORY

THIS METHOD WORKS SIMILARLY TO THE BUILTIN APPLY EXCEPT THAT THE FUNCTION IS CALLED ONLY IF THE CACHE IS NOT UP TO DATE

FORMATSELFBJINDENT0

RETURN THE FORMATTED REPRESENTATION OF THE OBJECT

REDUCESIZE SELF

REMOVE CACHE ELEMENTS TO MAKE CACHE SIZE FIT IN BYTESLIMIT

SKLEARNUTILS PARALLEL

WARNING DEPRECATED

CLASSSSKLEARNUTILS PARALLEL ARGS KWARGS

METHODS

CALL SELF ITERABLE

DISPATCHNEXT SELF DISPATCH MORE DATA FOR PARALLEL PROCESSING

DISPATCHONEBATCH SELF ITERATOR PREFETCH THE TASKS FOR THE NEXT BATCH AND DISPATCH THEM

FORMAT SELF OBJ INDENT RETURN THE FORMATTED REPRESENTATION OF THE OBJECT

PRINTPROGRESS SELF DISPLAY THE PROCESS OF THE PARALLEL EXECUTION ONLY A FRACTION OF TIME CONTROLLED BY SELFVERBOSE

DEBUG

RETRIEVE

WARN

INIT ARGS KWARGS

DEPRECATED DEPRECATED IN VERSION 0201 TO BE REMOVED IN VERSION 023 PLEASE IMPORT THIS FUNCTION

640 RECENTLY DEPRECATED 2387

SCIKITLEARN USER GUIDE RELEASE 0213

ALITY DIRECTLY FROM JOBLIB WHICH CAN BE INSTALLED WITH PIP INSTALL JOBLIB

DISPATCHNEXT SELF

DISPATCH MORE DATA FOR PARALLEL PROCESSING

THIS METHOD IS MEANT TO BE CALLED CONCURRENTLY BY THE MULTIPROCESSING CALLBACK WE RELY ON THE THREAD SAFETY OF DISPATCHONEBATCH TO PROTECT AGAINST CONCURRENT CONSUMPTION OF THE UNPROTECTED ITERATOR

DISPATCHONEBATCH SELFITERATOR

PREFETCH THE TASKS FOR THE NEXT BATCH AND DISPATCH THEM

THE EFFECTIVE SIZE OF THE BATCH IS COMPUTED HERE IF THERE ARE NO MORE JOBS TO DISPATCH RETURN FALSE ELSE RETURN TRUE

THE ITERATOR CONSUMPTION AND DISPATCHING IS PROTECTED BY THE SAME LOCK SO CALLING THIS FUNCTION SHOULD BE THREAD SAFE

FORMATSELF0BJINDENT0

RETURN THE FORMATTED REPRESENTATION OF THE OBJECT

PRINTPROGRESS SELF

DISPLAY THE PROCESS OF THE PARALLEL EXECUTION ONLY A FRACTION OF TIME CONTROLLED BY SELFVERBOSE

UTILSCPUCOUNT DEPRECATED DEPRECATED IN VERSION 0201 TO BE RE

MOVED IN VERSION 023

UTILSDELAYED FUNCTION CHECKPICKLE DEPRECATED DEPRECATED IN VERSION 0201 TO BE RE

MOVED IN VERSION 023

METRICSCALINSKI HARABAZSCORE X LABELS DEPRECATED FUNCTION 'CALINSKI HARABAZSCORE' HAS BEEN RENAMED TO 'CALINSKI HARABASZSCORE' AND WILL BE RE

MOVED IN VERSION 023

METRICSJACCARDSIMILARITYSCORE YTRUE

YPREDJACCARD SIMILARITY COEFFICIENT SCORE

LINEARMODELLOGISTICREGRESSIONPATH X

YDEPRECATED LOGISTICREGRESSIONPATH WAS DEPRECATED

IN VERSION 021 AND WILL BE REMOVED IN VERSION 0230

SKLEARNUTILS CPUCOUNT

WARNING DEPRECATED

SKLEARNUTILS CPUCOUNT

DEPRECATED DEPRECATED IN VERSION 0201 TO BE REMOVED IN VERSION 023 PLEASE IMPORT THIS FUNCTIONALITY DIRECTLY FROM JOBLIB WHICH CAN BE INSTALLED WITH PIP INSTALL JOBLIB

RETURN THE NUMBER OF CPUS

SKLEARNUTILS DELAYED

WARNING DEPRECATED

SKLEARNUTILS DELAYEDFUNCTION CHECKPICKLENONE

DEPRECATED DEPRECATED IN VERSION 0201 TO BE REMOVED IN VERSION 023 PLEASE IMPORT THIS FUNCTIONALITY DIRECTLY FROM JOBLIB WHICH CAN BE INSTALLED WITH PIP INSTALL JOBLIB

2388 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213  
DECORATOR USED TO CAPTURE THE ARGUMENTS OF A FUNCTION  
SKLEARNMETRICS CALINSKI HARABAZ SCORE  
WARNING DEPRECATED  
SKLEARNMETRICS CALINSKI HARABAZ SCORE XLABELS  
DEPRECATED FUNCTION 'CALINSKI HARABAZ SCORE' HAS BEEN RENAMED TO 'CALINSKI HARABASZ SCORE' AND WILL BE  
REMOVED IN VERSION 023  
SKLEARNMETRICS JACCARD SIMILARITY SCORE  
WARNING DEPRECATED  
SKLEARNMETRICS JACCARD SIMILARITY SCORE YTRUE YPRED NORMALIZE TRUE SAMPLEWEIGHT NONE  
JACCARD SIMILARITY COEFFICIENT SCORE  
DEPRECATED SINCE VERSION 021 THIS IS DEPRECATED TO BE REMOVED IN 023 SINCE ITS HANDLING OF BINARY AND  
MULTICLASS INPUTS WAS BROKEN JACCARD SCORE HAS AN API THAT IS CONSISTENT WITH PRECISION SCORE F SCORE  
ETC  
READ MORE IN THE USER GUIDE  
PARAMETERS  
YTRUE 1D ARRAYLIKE OR LABEL INDICATOR ARRAY SPARSE MATRIX GROUND TRUTH CORRECT LABELS  
YPRED 1D ARRAYLIKE OR LABEL INDICATOR ARRAY SPARSE MATRIX PREDICTED LABELS AS RETURNED BY  
A CLASSIFIER  
NORMALIZE BOOL OPTIONAL DEFAULT TRUE IF FALSE RETURN THE SUM OF THE JACCARD SIMILARITY  
COEFFICIENT OVER THE SAMPLE SET OTHERWISE RETURN THE AVERAGE OF JACCARD SIMILARITY COEFFI  
CIENT  
SAMPLEWEIGHT ARRAYLIKE OF SHAPE NSAMPLES OPTIONAL SAMPLE WEIGHTS  
RETURNS  
SCORE FLOAT IF NORMALIZE TRUE RETURN THE AVERAGE JACCARD SIMILARITY COEFFICIENT ELSE  
IT RETURNS THE SUM OF THE JACCARD SIMILARITY COEFFICIENT OVER THE SAMPLE SET  
THE BEST PERFORMANCE IS 1 WITH NORMALIZE TRUE AND THE NUMBER OF SAMPLES WITH  
NORMALIZE FALSE  
SEE ALSO  
ACCURACY SCORE HAMMING LOSS ZERO ONE LOSS  
NOTES  
IN BINARY AND MULTICLASS CLASSIFICATION THIS FUNCTION IS EQUIVALENT TO THE ACCURACY SCORE IT DIFFERS IN THE  
MULTILABEL CLASSIFICATION PROBLEM  
640 RECENTLY DEPRECATED 2389

SCIKITLEARN USER GUIDE RELEASE 0213

REFERENCES

1

SKLEARNLINEARMODEL LOGISTICREGRESSIONPATH

WARNING DEPRECATED

SKLEARNLINEARMODEL LOGISTICREGRESSIONPATH XY POSCLASSNONE CS10

FITINTERCEPTTRUE MAXITER100

TOL00001 VERBOSE0 SOLVER'LBFGS'

COEFNONE CLASSWEIGHTNONE

DUALFALSE PENALTY'L2' INTER

CEPTSCALING10 MULTICLASS'WARN'

RANDOMSTATENONE CHECKINPUTTRUE

MAXSQUAREDSSUMNONE SAM

PLEWEIGHTNONE L1RATIONONE

DEPRECATED LOGISTICREGRESSIONPATH WAS DEPRECATED IN VERSION 021 AND WILL BE REMOVED IN VERSION 0230

COMPUTE A LOGISTIC REGRESSION MODEL FOR A LIST OF REGULARIZATION PARAMETERS

THIS IS AN IMPLEMENTATION THAT USES THE RESULT OF THE PREVIOUS MODEL TO SPEED UP COMPUTATIONS ALONG THE SET OF SOLUTIONS MAKING IT FASTER THAN SEQUENTIALLY CALLING LOGISTICREGRESSION FOR THE DIFFERENT PARAMETERS

NOTE THAT THERE WILL BE NO SPEEDUP WITH LIBLINEAR SOLVER SINCE IT DOES NOT HANDLE WARMSTARTING

DEPRECATED SINCE VERSION 021 LOGISTICREGRESSIONPATH WAS DEPRECATED IN VERSION 021 AND

WILL BE REMOVED IN 023

READ MORE IN THE USER GUIDE

PARAMETERS

XARRAYLIKE OR SPARSE MATRIX SHAPE NSAMPLES NFEATURES

INPUT DATA

YARRAYLIKE SHAPE NSAMPLES OR NSAMPLES NTARGETS INPUT DATA TARGET VALUES

POSCLASS INT NONE THE CLASS WITH RESPECT TO WHICH WE PERFORM A ONEVSALL FIT IF NONE

THEN IT IS ASSUMED THAT THE GIVEN PROBLEM IS BINARY

CSINT ARRAYLIKE SHAPE NCS LIST OF VALUES FOR THE REGULARIZATION PARAMETER OR INTEGER

SPECIFYING THE NUMBER OF REGULARIZATION PARAMETERS THAT SHOULD BE USED IN THIS CASE THE

PARAMETERS WILL BE CHOSEN IN A LOGARITHMIC SCALE BETWEEN 1E4 AND 1E4

FITINTERCEPT BOOL WHETHER TO FIT AN INTERCEPT FOR THE MODEL IN THIS CASE THE SHAPE OF THE

RETURNED ARRAY IS NCS NFEATURES 1

MAXITER INT MAXIMUM NUMBER OF ITERATIONS FOR THE SOLVER

TOLFLOAT STOPPING CRITERION FOR THE NEWTONCG AND LBFGS SOLVERS THE ITERATION WILL STOP

WHENMAXGI I 1 N TOL WHEREGI IS THE ITH COMPONENT

OF THE GRADIENT

VERBOSE INT FOR THE LIBLINEAR AND LBFGS SOLVERS SET VERBOSE TO ANY POSITIVE NUMBER FOR

VERBOSITY

SOLVER 'LBFGS' 'NEWTONCG' 'LIBLINEAR' 'SAG' 'SAGA' NUMERICAL SOLVER TO USE

2390 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

COEF ARRAYLIKE SHAPE NFEATURES DEFAULT NONE INITIALIZATION VALUE FOR COEFFICIENTS OF LOGISTIC REGRESSION USELESS FOR LIBLINEAR SOLVER

CLASSWEIGHT DICT OR 'BALANCED' OPTIONAL WEIGHTS ASSOCIATED WITH CLASSES IN THE FORM CLASSLABEL WEIGHT IF NOT GIVEN ALL CLASSES ARE SUPPOSED TO HAVE WEIGHT ONE

THE "BALANCED" MODE USES THE VALUES OF Y TO AUTOMATICALLY ADJUST WEIGHTS INVERSELY PROPORTIONAL TO CLASS FREQUENCIES IN THE INPUT DATA AS NSAMPLES NCLASSES

NPBINCOUNTY

NOTE THAT THESE WEIGHTS WILL BE MULTIPLIED WITH SAMPLEWEIGHT PASSED THROUGH THE FIT METHOD IF SAMPLEWEIGHT IS SPECIFIED

DUAL BOOL DUAL OR PRIMAL FORMULATION DUAL FORMULATION IS ONLY IMPLEMENTED FOR L2 PENALTY WITH LIBLINEAR SOLVER PREFER DUALFALSE WHEN NSAMPLES NFEATURES

PENALTY STR 'L1' 'L2' OR 'ELASTICNET' USED TO SPECIFY THE NORM USED IN THE PENALIZATION THE 'NEWTONCG' 'SAG' AND 'LBFGS' SOLVERS SUPPORT ONLY L2 PENALTIES 'ELASTICNET' IS ONLY SUPPORTED BY THE 'SAGA' SOLVER

INTERCEPTSCALING FLOAT DEFAULT 1 USEFUL ONLY WHEN THE SOLVER 'LIBLINEAR' IS USED AND SELFFITINTERCEPT IS SET TO TRUE IN THIS CASE X BECOMES X SELFINTERCEPTSCALING IE A "SYNTHETIC" FEATURE WITH CONSTANT VALUE EQUAL TO INTERCEPTSCALING IS APPENDED TO THE INSTANCE VECTOR THE INTERCEPT BECOMES INTERCEPTSCALING

SYNTHETICFEATUREWEIGHT

NOTE THE SYNTHETIC FEATURE WEIGHT IS SUBJECT TO L1L2 REGULARIZATION AS ALL OTHER FEATURES TO LESSEN THE EFFECT OF REGULARIZATION ON SYNTHETIC FEATURE WEIGHT AND THEREFORE ON THE INTERCEPT INTERCEPTSCALING HAS TO BE INCREASED

MULTICLASS STR 'OVR' 'MULTINOMIAL' 'AUTO' DEFAULT 'OVR' IF THE OPTION CHOSEN IS 'OVR' THEN A BINARY PROBLEM IS FIT FOR EACH LABEL FOR 'MULTINOMIAL' THE LOSS MINIMISED IS THE MULTINOMIAL LOSS FIT ACROSS THE ENTIRE PROBABILITY DISTRIBUTION EVEN WHEN THE DATA IS BINARY 'MULTINOMIAL' IS UNAVAILABLE WHEN SOLVER'LIBLINEAR' 'AUTO' SELECTS 'OVR' IF THE DATA IS BINARY OR IF SOLVER'LIBLINEAR' AND OTHERWISE SELECTS 'MULTINOMIAL'

NEW IN VERSION 018 STOCHASTIC AVERAGE GRADIENT DESCENT SOLVER FOR 'MULTINOMIAL' CASE CHANGED IN VERSION 020 DEFAULT WILL CHANGE FROM 'OVR' TO 'AUTO' IN 022

RANDOMSTATE INT RANDOMSTATE INSTANCE OR NONE OPTIONAL DEFAULT NONE THE SEED OF THE PSEUDO RANDOM NUMBER GENERATOR TO USE WHEN SHUFFLING THE DATA IF INT RANDOMSTATE IS THE SEED USED BY THE RANDOM NUMBER GENERATOR IF RANDOMSTATE INSTANCE RANDOMSTATE IS THE RANDOM NUMBER GENERATOR IF NONE THE RANDOM NUMBER GENERATOR IS THE RANDOMSTATE INSTANCE USED BY NPRANDOM USED WHEN SOLVER 'SAG' OR 'LIBLINEAR'

CHECKINPUT BOOL DEFAULT TRUE IF FALSE THE INPUT ARRAYS X AND Y WILL NOT BE CHECKED

MAXSQUARED SUM FLOAT DEFAULT NONE MAXIMUM SQUARED SUM OF X OVER SAMPLES USED ONLY IN SAG SOLVER IF NONE IT WILL BE COMPUTED GOING THROUGH ALL THE SAMPLES THE VALUE SHOULD BE PRECOMPUTED TO SPEED UP CROSS VALIDATION

SAMPLEWEIGHT ARRAYLIKE SHAPENSAMPLES OPTIONAL ARRAY OF WEIGHTS THAT ARE ASSIGNED TO INDIVIDUAL SAMPLES IF NOT PROVIDED THEN EACH SAMPLE IS GIVEN UNIT WEIGHT

L1RATIO FLOAT OR NONE OPTIONAL DEFAULTNONE THE ELASTICNET MIXING PARAMETER WITH0 L1RATIO 1 ONLY USED IF PENALTYELASTICNET SETTING

L1RATIO0 IS EQUIVALENT TO USING PENALTYL2 WHILE SETTING L1RATIO1 IS EQUIVALENT TO USING PENALTYL1 FOR0 L1RATIO 1 THE PENALTY IS A COMBINATION OF L1 AND L2

640 RECENTLY DEPRECATED 2391

SCIKITLEARN USER GUIDE RELEASE 0213

RETURNS

COEFS NDARRAY SHAPE NCS NFEATURES OR NCS NFEATURES 1  
LIST OF COEFFICIENTS FOR THE LOGISTIC REGRESSION MODEL IF FITINTERCEPT IS SET TO TRUE  
THEN THE SECOND DIMENSION WILL BE NFEATURES 1 WHERE THE LAST ITEM REPRESENTS  
THE INTERCEPT FOR MULTICLASSMULTINOMIAL THE SHAPE IS NCLASSES NCS  
NFEATURES OR NCLASSES NCS NFEATURES 1  
CSNDARRAY GRID OF CS USED FOR CROSSVALIDATION  
NITER ARRAY SHAPE NCS ACTUAL NUMBER OF ITERATION FOR EACH CS

NOTES

YOU MIGHT GET SLIGHTLY DIFFERENT RESULTS WITH THE SOLVER LIBLINEAR THAN WITH THE OTHERS SINCE THIS USES LIBLINEAR  
WHICH PENALIZES THE INTERCEPT  
CHANGED IN VERSION 019 THE “COPY” PARAMETER WAS REMOVED  
ENSEMBLEPARTIALDEPENDENCE  
PARTIALDEPENDENCE DEPRECATED THE FUNCTION ENSEM  
BLEPARTIALDEPENDENCE HAS BEEN DEPRECATED IN FAVOUR  
OF INSPECTIONPARTIALDEPENDENCE IN 021 AND WILL BE  
REMOVED IN 023  
ENSEMBLEPARTIALDEPENDENCE  
PLOTPARTIALDEPENDENCE DEPRECATED THE FUNCTION ENSEM  
BLEPLOTPARTIALDEPENDENCE HAS BEEN DEPRECATED IN  
FAVOUR OF SKLEARNINSPECTIONPLOTPARTIALDEPENDENCE IN  
021 AND WILL BE REMOVED IN 023  
SKLEARNENSEMBLEPARTIALDEPENDENCE PARTIALDEPENDENCE  
SKLEARNENSEMBLEPARTIALDEPENDENCE PARTIALDEPENDENCE GBRT TARGETVARIABLES  
GRIDNONE XNONE  
PERCENTILES005 095  
GRIDRESOLUTION100  
DEPRECATED THE FUNCTION ENSEMBLEPARTIALDEPENDENCE HAS BEEN DEPRECATED IN FAVOUR OF INSPEC  
TIONPARTIALDEPENDENCE IN 021 AND WILL BE REMOVED IN 023  
PARTIAL DEPENDENCE OF TARGETVARIABLES  
PARTIAL DEPENDENCE PLOTS SHOW THE DEPENDENCE BETWEEN THE JOINT VALUES OF THE  
TARGETVARIABLES AND THE FUNCTION REPRESENTED BY THE GBRT  
READ MORE IN THE USER GUIDE  
DEPRECATED SINCE VERSION 021 THIS FUNCTION WAS DEPRECATED IN VERSION 021 IN FAVOR OF SKLEARN  
INSPECTIONPARTIALDEPENDENCE AND WILL BE REMOVED IN 023  
PARAMETERS  
GBRT BASEGRADIENTBOOSTING  
A FITTED GRADIENT BOOSTING MODEL  
TARGETVARIABLES ARRAYLIKE DTYPEINT THE TARGET FEATURES FOR WHICH THE PARTIAL DEPENDEN  
CY SHOULD BE COMPUTED SIZE SHOULD BE SMALLER THAN 3 FOR VISUAL RENDERINGS  
2392 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

GRID ARRAYLIKE SHAPENPOINTS LENTARGETVARIABLES THE GRID OF  
TARGETVARIABLES VALUES FOR WHICH THE PARTIAL DEPENDENCY SHOULD BE EVALU  
ATED EITHER GRID ORXMUST BE SPECIFIED

XARRAYLIKE SHAPENSAMPLES NFEATURES THE DATA ON WHICH GBRT WAS TRAINED IT  
IS USED TO GENERATE A GRID FOR THETARGETVARIABLES THEGRID COMPRISES  
GRIDRESOLUTION EQUALLY SPACED POINTS BETWEEN THE TWO PERCENTILES  
PERCENTILES LOW HIGH DEFAULT005 095 THE LOWER AND UPPER PERCENTILE USED CREATE  
THE EXTREME VALUES FOR THE GRID ONLY IFXIS NOT NONE

GRIDRESOLUTION INT DEFAULT100 THE NUMBER OF EQUALLY SPACED POINTS ON THE GRID  
RETURNS

PDP ARRAY SHAPENCLASSES NPOINTS  
THE PARTIAL DEPENDENCE FUNCTION EVALUATED ON THE GRID FOR REGRESSION AND BINARY  
CLASSIFICATION NCLASSES1  
AXES SEQ OF NDARRAY OR NONE THE AXES WITH WHICH THE GRID HAS BEEN CREATED OR NONE IF THE  
GRID HAS BEEN GIVEN

EXAMPLES

SAMPLES 0 0 2 1 0 0  
LABELS 0 1

FROM SKLEARNENSEMBLE IMPORT GRADIENTBOOSTINGCLASSIFIER  
GB GRADIENTBOOSTINGCLASSIFIERRANDOMSTATE0FITSAMPLES LABELS  
KWARGS DICTXSAMPLES PERCENTILES0 1 GRIDRESOLUTION2  
PARTIALDEPENDENCEGB 0 KWARGS  
ARRAY452 452 ARRAY 0 1  
SKLEARNENSEMBLEPARTIALDEPENDENCE PLOTPARTIALDEPENDENCE  
SKLEARNENSEMBLEPARTIALDEPENDENCE PLOTPARTIALDEPENDENCE GBRTXFEATURES FEA  
TURENAMESNONE  
LABELNONE  
NCOLS3  
GRIDRESOLUTION100  
PERCENTILES005  
095 NJOBSNONE  
VERBOSE0  
AXNONE  
LINEKWNONE  
CONTOURKWNONE  
FIGKW

DEPRECATED THE FUNCTION ENSEMBLEPLOTPARTIALDEPENDENCE HAS BEEN DEPRECATED IN FAVOUR OF  
SKLEARNINSPECTIONPLOTPARTIALDEPENDENCE IN 021 AND WILL BE REMOVED IN 023

PARTIAL DEPENDENCE PLOTS FOR FEATURES

THELENFEATURES PLOTS ARE ARRANGED IN A GRID WITH NCOLS COLUMNS TWOWAY PARTIAL DEPENDENCE PLOTS ARE PLOTTED AS CONTOUR PLOTS

READ MORE IN THE USER GUIDE

640 RECENTLY DEPRECATED 2393

SCIKITLEARN USER GUIDE RELEASE 0213  
DEPRECATED SINCE VERSION 021 THIS FUNCTION WAS DEPRECATED IN VERSION 021 IN FAVOR OF SKLEARN  
INSPECTIONPLOT  
PARTIALDEPENDENCE AND WILL BE REMOVED IN 023

PARAMETERS

GBRT BASEGRADIENTBOOSTING  
A FITTED GRADIENT BOOSTING MODEL

XARRAYLIKE SHAPENSAMPLES NFEATURES THE DATA ON WHICH GBRT WAS TRAINED  
FEATURES SEQ OF INTS STRINGS OR TUPLES OF INTS OR STRINGS IF SEQI IS AN INT OR A TUPLE  
WITH ONE INT VALUE A ONEWAY PDP IS CREATED IF SEQI IS A TUPLE OF TWO INTS A TWO  
WAY PDP IS CREATED IF FEATURENAMES IS SPECIFIED AND SEQI IS AN INT SEQI MUST BE  
LENFEATURENAMES IF SEQI IS A STRING FEATURENAMES MUST BE SPECIFIED AND SEQI MUST  
BE IN FEATURENAMES  
FEATURENAMES SEQ OF STR NAME OF EACH FEATURE FEATURENAMESI HOLDS THE NAME OF THE  
FEATURE WITH INDEX I  
LABEL OBJECT THE CLASS LABEL FOR WHICH THE PDPS SHOULD BE COMPUTED ONLY IF GBRT IS A  
MULTICLASS MODEL MUST BE IN GBRTCLASSES  
NCOLS INT THE NUMBER OF COLUMNS IN THE GRID PLOT DEFAULT 3  
GRIDRESOLUTION INT DEFAULT100 THE NUMBER OF EQUALLY SPACED POINTS ON THE AXES  
PERCENTILES LOW HIGH DEFAULT005 095 THE LOWER AND UPPER PERCENTILE USED TO CREATE  
THE EXTREME VALUES FOR THE PDP AXES  
NJOBS INT OR NONE OPTIONAL DEFAULTNONE NONE MEANS 1 UNLESS IN A JOBLIB  
PARALLELBACKEND CONTEXT1MEANS USING ALL PROCESSORS SEE GLOSSARY FOR MORE  
DETAILS  
VERBOSE INT VERBOSE OUTPUT DURING PD COMPUTATIONS DEFAULTS TO 0  
AXMATPLOTLIB AXIS OBJECT DEFAULT NONE AN AXIS OBJECT ONTO WHICH THE PLOTS WILL BE DRAWN  
LINEKW DICT DICT WITH KEYWORDS PASSED TO THE MATPLOTLIBPYPLOTPLOT CALL FOR  
ONEWAY PARTIAL DEPENDENCE PLOTS  
CONTOURKW DICT DICT WITH KEYWORDS PASSED TO THE MATPLOTLIBPYPLOTPLOT CALL  
FOR TWOWAY PARTIAL DEPENDENCE PLOTS  
FIGKW DICT DICT WITH KEYWORDS PASSED TO THE FIGURE CALL NOTE THAT ALL KEYWORDS NOT  
RECOGNIZED ABOVE WILL BE AUTOMATICALLY INCLUDED HERE

RETURNS

FIGFIGURE  
THE MATPLOTLIB FIGURE OBJECT  
AXS SEQ OF AXIS OBJECTS A SEQ OF AXIS OBJECTS ONE FOR EACH SUBPLOT

EXAMPLES

```
FROM SKLEARNDATASETS IMPORT MAKEFRIEDMAN1  
FROM SKLEARNENSEMBLE IMPORT GRADIENTBOOSTINGREGRESSOR  
X Y MAKEFRIEDMAN1  
CLF GRADIENTBOOSTINGREGRESSORNESTIMATORS10FITX Y
```

2394 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213  
FIG AXS PLOT PARTIAL DEPENDENCE CLF X 0 0 1

6402 TO BE REMOVED IN 022  
COVARIANCE GRAPH LASSO ARGS KWARGS SPARSE INVERSE COVARIANCE ESTIMATION WITH AN L1 PENALIZED ESTIMATOR  
COVARIANCE GRAPH LASSO CV ARGS KWARGS SPARSE INVERSE COVARIANCE W CROSS VALIDATED CHOICE OF THE L1 PENALTY  
PREPROCESSING IMPUTER ARGS KWARGS IMPUTATION TRANSFORMER FOR COMPLETING MISSING VALUES  
UTIL TESTING MOCK ML DATA URLOPEN ARGS  
OBJECT THAT MOCKS THE URLOPEN FUNCTION TO FAKE REQUESTS TO  
ML DATA  
SKLEARN COVARIANCE GRAPH LASSO  
WARNING DEPRECATED  
CLASS SKLEARN COVARIANCE GRAPH LASSO ARGS KWARGS  
SPARSE INVERSE COVARIANCE ESTIMATION WITH AN L1 PENALIZED ESTIMATOR  
THIS CLASS IMPLEMENTS THE GRAPHICAL LASSO ALGORITHM  
READ MORE IN THE USER GUIDE  
PARAMETERS  
ALPHA POSITIVE FLOAT DEFAULT 0.01 THE REGULARIZATION PARAMETER THE HIGHER ALPHA THE MORE REGULARIZATION THE SPARSER THE INVERSE COVARIANCE  
MODE 'CD' 'LARS' DEFAULT 'CD' THE LASSO SOLVER TO USE COORDINATE DESCENT OR LARS USE LARS FOR VERY SPARSE UNDERLYING GRAPHS WHERE  $p \gg n$  ELSEWHERE PREFER CD WHICH IS MORE NUMERICALLY STABLE  
TOL POSITIVE FLOAT DEFAULT 1e4 THE TOLERANCE TO DECLARE CONVERGENCE IF THE DUAL GAP GOES BELOW THIS VALUE ITERATIONS ARE STOPPED  
ENET\_TOL POSITIVE FLOAT OPTIONAL THE TOLERANCE FOR THE ELASTIC NET SOLVER USED TO CALCULATE THE DESCENT DIRECTION THIS PARAMETER CONTROLS THE ACCURACY OF THE SEARCH DIRECTION FOR A GIVEN COLUMN UPDATE NOT OF THE OVERALL PARAMETER ESTIMATE ONLY USED FOR MODE 'CD'  
MAX\_ITER INTEGER DEFAULT 100 THE MAXIMUM NUMBER OF ITERATIONS  
VERBOSE BOOLEAN DEFAULT FALSE IF VERBOSE IS TRUE THE OBJECTIVE FUNCTION AND DUAL GAP ARE PLOTTED AT EACH ITERATION  
ASSUME\_CENTERED BOOLEAN DEFAULT FALSE IF TRUE DATA ARE NOT CENTERED BEFORE COMPUTATION USEFUL WHEN WORKING WITH DATA WHOSE MEAN IS ALMOST BUT NOT EXACTLY ZERO IF FALSE DATA ARE CENTERED BEFORE COMPUTATION  
ATTRIBUTES  
COVARIANCE ARRAYLIKE SHAPE (n\_features, n\_features) ESTIMATED COVARIANCE MATRIX  
PRECISION ARRAYLIKE SHAPE (n\_features, n\_features) ESTIMATED PSEUDO INVERSE MATRIX  
N\_ITER INT NUMBER OF ITERATIONS RUN  
640 RECENTLY DEPRECATED 2395

SCIKITLEARN USER GUIDE RELEASE 0213

SEE ALSO

GRAPHLASSO GRAPHLASSOCV

METHODS

ERRORNORM SELF COMPCOV NORM SCALING COMPUTES THE MEAN SQUARED ERROR BETWEEN TWO COVARIANCE ESTIMATORS

FITSELF X Y FITS THE GRAPHICALASSO MODEL TO X

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

GETPRECISION SELF GETTER FOR THE PRECISION MATRIX

MAHALANOBIS SELF X COMPUTES THE SQUARED MAHALANOBIS DISTANCES OF GIVEN OBSERVATIONS

SCORE SELF XTEST Y COMPUTES THE LOGLIKELIHOOD OF A GAUSSIAN DATA SET WITH SELFCOVARIANCE AS AN ESTIMATOR OF ITS COVARIANCE MATRIX

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT ARGS KWARGS

DEPRECATED THE 'GRAPHLASSO' WAS RENAMED TO 'GRAPHICALASSO' IN VERSION 020 AND WILL BE REMOVED IN 022

ERRORNORM SELFCOMPCOV NORM'FROBENIUS' SCALINGTRUE SQUAREDTRUE COMPUTES THE MEAN SQUARED ERROR BETWEEN TWO COVARIANCE ESTIMATORS IN THE SENSE OF THE FROBENIUS NORM

PARAMETERS

COMPCOV ARRAYLIKE SHAPE NFEATURES NFEATURES THE COVARIANCE TO COMPARE WITH

NORM STR THE TYPE OF NORM USED TO COMPUTE THE ERROR AVAILABLE ERROR TYPES 'FROBENIUS'

DEFAULT SQRTTRATA 'SPECTRAL' SQRTMAXEIGENVALUESATA WHERE A IS THE ERROR

COMPCOV SELFCOVARIANCE

SCALING BOOL IF TRUE DEFAULT THE SQUARED ERROR NORM IS DIVIDED BY NFEATURES IF FALSE THE SQUARED ERROR NORM IS NOT RESCALED

SQUARED BOOL WHETHER TO COMPUTE THE SQUARED ERROR NORM OR THE ERROR NORM IF TRUE DEFAULT THE SQUARED ERROR NORM IS RETURNED IF FALSE THE ERROR NORM IS RETURNED

RETURNS

THE MEAN SQUARED ERROR IN THE SENSE OF THE FROBENIUS NORM BETWEEN SELF ANDCOMPCOV COVARIANCE ESTIMATORS

FITSELFXYNONE

FITS THE GRAPHICALASSO MODEL TO X

PARAMETERS

XNDARRAY SHAPE NSAMPLES NFEATURES DATA FROM WHICH TO COMPUTE THE COVARIANCE ESTIMATE

YIGNORED

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

2396 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

GETPRECISION SELF

GETTER FOR THE PRECISION MATRIX

RETURNS

PRECISION ARRAYLIKE THE PRECISION MATRIX ASSOCIATED TO THE CURRENT COVARIANCE OBJECT

MAHALANOBIS SELF

COMPUTES THE SQUARED MAHALANOBIS DISTANCES OF GIVEN OBSERVATIONS

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES THE OBSERVATIONS THE MAHALANOBIS DISTANCES OF THE WHICH WE COMPUTE OBSERVATIONS ARE ASSUMED TO BE DRAWN FROM THE SAME DISTRIBUTION THAN THE DATA USED IN FIT

RETURNS

DIST ARRAY SHAPE NSAMPLES SQUARED MAHALANOBIS DISTANCES OF THE OBSERVATIONS

SCORESELFXTTEST YNONE

COMPUTES THE LOGLIKELIHOOD OF A GAUSSIAN DATA SET WITH SELFCOVARIANCE AS AN ESTIMATOR OF ITS COVARIANCE MATRIX

PARAMETERS

XTTEST ARRAYLIKE SHAPE NSAMPLES NFEATURES TEST DATA OF WHICH WE COMPUTE THE LIKELIHOOD WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES XTTEST IS ASSUMED TO BE DRAWN FROM THE SAME DISTRIBUTION THAN THE DATA USED IN FIT INCLUDING CENTERING

YNOT USED PRESENT FOR API CONSISTENCE PURPOSE

RETURNS

RESFLOAT THE LIKELIHOOD OF THE DATA SET WITH SELFCOVARIANCE AS AN ESTIMATOR OF ITS COVARIANCE MATRIX

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT'S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

SKLEARNCOVARIANCE GRAPHLASSOCV

WARNING DEPRECATED

640 RECENTLY DEPRECATED 2397

SCIKITLEARN USER GUIDE RELEASE 0213

CLASSSSKLEARNCOVARIANCE GRAPHLASSOCV ARGS KWARGS

SPARSE INVERSE COVARIANCE W CROSSVALIDATED CHOICE OF THE L1 PENALTY

SEE GLOSSARY ENTRY FOR CROSSVALIDATION ESTIMATOR

THIS CLASS IMPLEMENTS THE GRAPHICAL LASSO ALGORITHM

READ MORE IN THE USER GUIDE

PARAMETERS

ALPHAS INTEGER OR LIST POSITIVE FLOAT OPTIONAL IF AN INTEGER IS GIVEN IT FIXES THE NUMBER OF POINTS ON THE GRIDS OF ALPHA TO BE USED IF A LIST IS GIVEN IT GIVES THE GRID TO BE USED SEE THE NOTES IN THE CLASS DOCSTRING FOR MORE DETAILS

NREFINEMENTS STRICTLY POSITIVE INTEGER THE NUMBER OF TIMES THE GRID IS REFINED NOT USED IF EXPLICIT VALUES OF ALPHAS ARE PASSED

CVINT CROSSVALIDATION GENERATOR OR AN ITERABLE OPTIONAL DETERMINES THE CROSSVALIDATION SPLITTING STRATEGY POSSIBLE INPUTS FOR CV ARE

- NONE TO USE THE DEFAULT 3FOLD CROSSVALIDATION
- INTEGER TO SPECIFY THE NUMBER OF FOLDS
- CV SPLITTER
- AN ITERABLE YIELDING TRAIN TEST SPLITS AS ARRAYS OF INDICES

FOR INTEGERNONE INPUTS KFOLD IS USED

REFER USER GUIDE FOR THE VARIOUS CROSSVALIDATION STRATEGIES THAT CAN BE USED HERE

CHANGED IN VERSION 020 CVDEFAULT VALUE IF NONE WILL CHANGE FROM 3FOLD TO 5FOLD IN V022

TOLPOSITIVE FLOAT OPTIONAL THE TOLERANCE TO DECLARE CONVERGENCE IF THE DUAL GAP GOES BELOW THIS VALUE ITERATIONS ARE STOPPED

ENETTOL POSITIVE FLOAT OPTIONAL THE TOLERANCE FOR THE ELASTIC NET SOLVER USED TO CALCULATE THE DESCENT DIRECTION THIS PARAMETER CONTROLS THE ACCURACY OF THE SEARCH DIRECTION FOR A GIVEN COLUMN UPDATE NOT OF THE OVERALL PARAMETER ESTIMATE ONLY USED FOR MODE'CD'

MAXITER INTEGER OPTIONAL MAXIMUM NUMBER OF ITERATIONS

MODE 'CD' 'LARS' THE LASSO SOLVER TO USE COORDINATE DESCENT OR LARS USE LARS FOR VERY SPARSE UNDERLYING GRAPHS WHERE NUMBER OF FEATURES IS GREATER THAN NUMBER OF SAMPLES ELSEWHERE PREFER CD WHICH IS MORE NUMERICALLY STABLE

NJOBS INT OR NONE OPTIONAL DEFAULTNONE NUMBER OF JOBS TO RUN IN PARALLEL NONE MEANS 1 UNLESS IN A JOBLIBPARALLELBKEND CONTEXT1MEANS USING ALL PROCESSORS SEE GLOSSARY FOR MORE DETAILS

VERBOSE BOOLEAN OPTIONAL IF VERBOSE IS TRUE THE OBJECTIVE FUNCTION AND DUALITY GAP ARE PRINTED AT EACH ITERATION

ASSUMECENTERED BOOLEAN IF TRUE DATA ARE NOT CENTERED BEFORE COMPUTATION USEFUL WHEN WORKING WITH DATA WHOSE MEAN IS ALMOST BUT NOT EXACTLY ZERO IF FALSE DATA ARE CENTERED BEFORE COMPUTATION

ATTRIBUTES

COVARIANCE NUMPYNDARRAY SHAPE NFEATURES NFEATURES ESTIMATED COVARIANCE MATRIX

PRECISION NUMPYNDARRAY SHAPE NFEATURES NFEATURES ESTIMATED PRECISION MATRIX INVERSE COVARIANCE

2398 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

ALPHA FLOAT PENALIZATION PARAMETER SELECTED

CVALPHAS LIST OF FLOAT ALL PENALIZATION PARAMETERS EXPLORED

GRIDScores 2D NUMPYNDARRAY NALPHAS NFOLDS LOGLIKELIHOOD SCORE ON LEFTOUT DATA

ACROSS FOLDS

NITER INT NUMBER OF ITERATIONS RUN FOR THE OPTIMAL ALPHA

SEE ALSO

GRAPHLASSO GRAPHLASSO

NOTES

THE SEARCH FOR THE OPTIMAL PENALIZATION PARAMETER ALPHA IS DONE ON AN ITERATIVELY REFINED GRID FIRST THE CROSS

VALIDATED SCORES ON A GRID ARE COMPUTED THEN A NEW REFINED GRID IS CENTERED AROUND THE MAXIMUM AND SO ON

ONE OF THE CHALLENGES WHICH IS FACED HERE IS THAT THE SOLVERS CAN FAIL TO CONVERGE TO A WELLCONDITIONED ESTIMATE

THE CORRESPONDING VALUES OF ALPHA THEN COME OUT AS MISSING VALUES BUT THE OPTIMUM MAY BE CLOSE TO THESE

MISSING VALUES

METHODS

ERRORNORM SELF COMPCOV NORM SCALING COMPUTES THE MEAN SQUARED ERROR BETWEEN TWO COVARIANCE ESTIMATORS

FITSELF X Y FITS THE GRAPHICALASSO COVARIANCE MODEL TO X

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

GETPRECISION SELF GETTER FOR THE PRECISION MATRIX

MAHALANOBIS SELF X COMPUTES THE SQUARED MAHALANOBIS DISTANCES OF GIVEN OBSERVATIONS

SCORE SELF XTEST Y COMPUTES THE LOGLIKELIHOOD OF A GAUSSIAN DATA SET WITH SELF COVARIANCE AS AN ESTIMATOR OF ITS COVARIANCE MATRIX

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

INIT ARGS KWARGS

DEPRECATED THE 'GRAPHLASSOCV' WAS RENAMED TO 'GRAPHICALASSOCV' IN VERSION 020 AND WILL BE REMOVED IN 022

ERRORNORM SELF COMPCOV NORM 'FROBENIUS' SCALING TRUE SQUARED TRUE

COMPUTES THE MEAN SQUARED ERROR BETWEEN TWO COVARIANCE ESTIMATORS IN THE SENSE OF THE FROBENIUS NORM

PARAMETERS

COMPCOV ARRAYLIKE SHAPE NFEATURES NFEATURES THE COVARIANCE TO COMPARE WITH

NORM STR THE TYPE OF NORM USED TO COMPUTE THE ERROR AVAILABLE ERROR TYPES 'FROBENIUS'

DEFAULT SQRTTRATA 'SPECTRAL' SQRTMAXEIGENVALUESATA WHERE A IS THE ERROR

COMPCOV SELF COVARIANCE

SCALING BOOL IF TRUE DEFAULT THE SQUARED ERROR NORM IS DIVIDED BY NFEATURES IF FALSE THE SQUARED ERROR NORM IS NOT RESCALED

SQUARED BOOL WHETHER TO COMPUTE THE SQUARED ERROR NORM OR THE ERROR NORM IF TRUE

640 RECENTLY DEPRECATED 2399

SCIKITLEARN USER GUIDE RELEASE 0213

DEFAULT THE SQUARED ERROR NORM IS RETURNED IF FALSE THE ERROR NORM IS RETURNED

RETURNS

THE MEAN SQUARED ERROR IN THE SENSE OF THE FROBENIUS NORM BETWEEN

SELF ANDCOMPCOV COVARIANCE ESTIMATORS

FITSELFXYNONE

FITS THE GRAPHICALASSO COVARIANCE MODEL TO X

PARAMETERS

XNDARRAY SHAPE NSAMPLES NFEATURES DATA FROM WHICH TO COMPUTE THE COVARIANCE ES

TIMATE

YIGNORED

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

GETPRECISION SELF

GETTER FOR THE PRECISION MATRIX

RETURNS

PRECISION ARRAYLIKE THE PRECISION MATRIX ASSOCIATED TO THE CURRENT COVARIANCE OBJECT

MAHALANOBIS SELF

COMPUTES THE SQUARED MAHALANOBIS DISTANCES OF GIVEN OBSERVATIONS

PARAMETERS

XARRAYLIKE SHAPE NSAMPLES NFEATURES THE OBSERVATIONS THE MAHALANOBIS DISTANCES

OF THE WHICH WE COMPUTE OBSERVATIONS ARE ASSUMED TO BE DRAWN FROM THE SAME DISTRIBU

TION THAN THE DATA USED IN FIT

RETURNS

DIST ARRAY SHAPE NSAMPLES SQUARED MAHALANOBIS DISTANCES OF THE OBSERVATIONS

SCORESELFXTTEST YNONE

COMPUTES THE LOGLIKELIHOOD OF A GAUSSIAN DATA SET WITH SELFCOVARIANCE AS AN ESTIMATOR OF ITS

COVARIANCE MATRIX

PARAMETERS

XTTEST ARRAYLIKE SHAPE NSAMPLES NFEATURES TEST DATA OF WHICH WE COMPUTE THE

LIKELIHOOD WHERE NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF

FEATURES XTTEST IS ASSUMED TO BE DRAWN FROM THE SAME DISTRIBUTION THAN THE DATA USED IN

FIT INCLUDING CENTERING

YNOT USED PRESENT FOR API CONSISTENCE PURPOSE

RETURNS

RESFLOAT THE LIKELIHOOD OF THE DATA SET WITH SELFCOVARIANCE AS AN ESTIMATOR OF ITS

COVARIANCE MATRIX

2400 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS

SELF

SKLEARNPREPROCESSING IMPUTER

WARNING DEPRECATED

CLASSSSKLEARNPREPROCESSING IMPUTERARGS KWARGS

IMPUTATION TRANSFORMER FOR COMPLETING MISSING VALUES

READ MORE IN THE USER GUIDE

PARAMETERS

MISSINGVALUES INTEGER OR “NAN” OPTIONAL DEFAULT”NAN” THE PLACEHOLDER FOR THE MISSING VALUES ALL OCCURRENCES OF MISSINGVALUES WILL BE IMPUTED FOR MISSING VALUES ENCODED AS NPNAN USE THE STRING VALUE “NAN”

STRATEGY STRING OPTIONAL DEFAULT”MEAN” THE IMPUTATION STRATEGY

- IF “MEAN” THEN REPLACE MISSING VALUES USING THE MEAN ALONG THE AXIS
- IF “MEDIAN” THEN REPLACE MISSING VALUES USING THE MEDIAN ALONG THE AXIS
- IF “MOSTFREQUENT” THEN REPLACE MISSING USING THE MOST FREQUENT VALUE ALONG THE AXIS

AXIS INTEGER OPTIONAL DEFAULT0 THE AXIS ALONG WHICH TO IMPUTE

- IFAXIS0 THEN IMPUTE ALONG COLUMNS
- IFAXIS1 THEN IMPUTE ALONG ROWS

VERBOSE INTEGER OPTIONAL DEFAULT0 CONTROLS THE VERBOSITY OF THE IMPUTER

COPY BOOLEAN OPTIONAL DEFAULTTRUE IF TRUE A COPY OF X WILL BE CREATED IF FALSE IMPUTATION WILL BE DONE INPLACE WHENEVER POSSIBLE NOTE THAT IN THE FOLLOWING CASES A NEW COPY WILL ALWAYS BE MADE EVEN IF COPYFALSE

- IF X IS NOT AN ARRAY OF FLOATING VALUES
- IF X IS SPARSE AND MISSINGVALUES0
- IFAXIS0 AND X IS ENCODED AS A CSR MATRIX
- IFAXIS1 AND X IS ENCODED AS A CSC MATRIX

ATTRIBUTES

STATISTICS ARRAY OF SHAPE NFEATURES THE IMPUTATION FILL VALUE FOR EACH FEATURE IF AXIS 0

640 RECENTLY DEPRECATED 2401

SCIKITLEARN USER GUIDE RELEASE 0213

NOTES

- WHENAXISO COLUMNS WHICH ONLY CONTAINED MISSING VALUES AT FIT ARE DISCARDED UPON TRANSFORM
- WHENAXIS1 AN EXCEPTION IS RAISED IF THERE ARE ROWS FOR WHICH IT IS NOT POSSIBLE TO FILL IN THE MISSING VALUES EG BECAUSE THEY ONLY CONTAIN MISSING VALUES

METHODS

FITSELF X Y FIT THE IMPUTER ON X

FITTRANSFORM SELF X Y FIT TO DATA THEN TRANSFORM IT

GETPARAMS SELF DEEP GET PARAMETERS FOR THIS ESTIMATOR

SETPARAMS SELF PARAMS SET THE PARAMETERS OF THIS ESTIMATOR

TRANSFORM SELF X IMPUTE ALL MISSING VALUES IN X

INIT ARGS KWARGS

DEPRECATED IMPUTER WAS DEPRECATED IN VERSION 020 AND WILL BE REMOVED IN 022 IMPORT IM

PUTESIMPLEIMPUTER FROM SKLEARN INSTEAD

FITSELFXYNONE

FIT THE IMPUTER ON X

PARAMETERS

XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES INPUT DATA WHERE

NSAMPLES IS THE NUMBER OF SAMPLES AND NFEATURES IS THE NUMBER OF FEATURES

RETURNS

SELF IMPUTER

FITTRANSFORM SELFXYNONE FITPARAMS

FIT TO DATA THEN TRANSFORM IT

FITS TRANSFORMER TO X AND Y WITH OPTIONAL PARAMETERS FITPARAMS AND RETURNS A TRANSFORMED VERSION OF X

PARAMETERS

XNUMPY ARRAY OF SHAPE NSAMPLES NFEATURES TRAINING SET

YNUMPY ARRAY OF SHAPE NSAMPLES TARGET VALUES

RETURNS

XNEW NUMPY ARRAY OF SHAPE NSAMPLES NFEATURESNEW TRANSFORMED ARRAY

GETPARAMS SELFDEEPTURE

GET PARAMETERS FOR THIS ESTIMATOR

PARAMETERS

DEEP BOOLEAN OPTIONAL IF TRUE WILL RETURN THE PARAMETERS FOR THIS ESTIMATOR AND CONTAINED

SUBOBJECTS THAT ARE ESTIMATORS

RETURNS

PARAMS MAPPING OF STRING TO ANY PARAMETER NAMES MAPPED TO THEIR VALUES

SETPARAMS SELFPARAMS

SET THE PARAMETERS OF THIS ESTIMATOR

2402 CHAPTER 6 API REFERENCE



SCIKITLEARN USER GUIDE RELEASE 0213

THE METHOD WORKS ON SIMPLE ESTIMATORS AS WELL AS ON NESTED OBJECTS SUCH AS PIPELINES THE LATTER HAVE PARAMETERS OF THE FORM COMPONENTPARAMETER SO THAT IT’S POSSIBLE TO UPDATE EACH COMPONENT OF A NESTED OBJECT

RETURNS  
SELF

TRANSFORM SELF  
IMPUTE ALL MISSING VALUES IN X

PARAMETERS  
XARRAYLIKE SPARSE MATRIX SHAPE NSAMPLES NFEATURES THE INPUT DATA TO COMPLETE

SKLEARNUTILSTESTING MOCKMLDATAURLOPEN

WARNING DEPRECATED

CLASSSKLEARNUTILSTESTING MOCKMLDATAURLOPEN ARGS KWARGS

OBJECT THAT MOCKS THE URLOPEN FUNCTION TO FAKE REQUESTS TO MLDATA

WHEN REQUESTING A DATASET WITH A NAME THAT IS IN MOCKDATASETS THIS OBJECT CREATES A FAKE DATASET IN A STRINGIO OBJECT AND RETURNS IT OTHERWISE IT RAISES AN HTTPERROR

DEPRECATED SINCE VERSION 020 WILL BE REMOVED IN VERSION 022

PARAMETERS  
MOCKDATASETS DICT A DICTIONARY OF DATASETNAME DATADICT OR DATASETNAME DATADICT

ORDERING DATADICT ITSELF IS A DICTIONARY OF COLUMNNAME DATAARRAY AND

ORDERING IS A LIST OF COLUMNNAMES TO DETERMINE THE ORDERING IN THE DATA SET SEE

FAKEMLDATA FOR DETAILS

METHODS  
CALL SELF URLNAME

PARAMETERS  
INIT ARGS KWARGS

DEPRECATED DEPRECATED IN VERSION 020 TO BE REMOVED IN VERSION 022

COVARIANCEGRAPHLASSO EMPCOV ALPHA DEPRECATED THE ‘GRAPHLASSO’ WAS RENAMED TO ‘GRAPHICALLASSO’ IN VERSION 020 AND WILL BE REMOVED IN 022

DATASETSFETCHMLDATA DATANAME DEPRECATED FETCHMLDATA WAS DEPRECATED IN VERSION 020 AND WILL BE REMOVED IN VERSION 022

DATASETSMMLDATAFILENAME DATANAME DEPRECATED MLDATAFILENAME WAS DEPRECATED IN VERSION 020 AND WILL BE REMOVED IN VERSION 022

640 RECENTLY DEPRECATED 2403

SCIKITLEARN USER GUIDE RELEASE 0213  
SKLEARNCOVARIANCE GRAPHLASSO  
WARNING DEPRECATED  
SKLEARNCOVARIANCE GRAPHLASSO EMPCOV ALPHA COVINITNONE MODE'CD' TOL00001  
ENETTOL00001 MAXITER100 VERBOSEFALSE RE  
TURNCOSTSFALSE EPS2220446049250313E16 RE  
TURNNITERFALSE  
DEPRECATED THE 'GRAPHLASSO' WAS RENAMED TO 'GRAPHICALASSO' IN VERSION 020 AND WILL BE REMOVED IN  
022  
L1PENALIZED COVARIANCE ESTIMATOR  
READ MORE IN THE USER GUIDE  
PARAMETERS  
EMPCOV 2D NDARRAY SHAPE NFEATURES NFEATURES  
EMPIRICAL COVARIANCE FROM WHICH TO COMPUTE THE COVARIANCE ESTIMATE  
ALPHA POSITIVE FLOAT THE REGULARIZATION PARAMETER THE HIGHER ALPHA THE MORE REGULARIZA  
TION THE SPARSER THE INVERSE COVARIANCE  
COVINIT 2D ARRAY NFEATURES NFEATURES OPTIONAL THE INITIAL GUESS FOR THE COVARIANCE  
MODE 'CD' 'LARS' THE LASSO SOLVER TO USE COORDINATE DESCENT OR LARS USE LARS  
FOR VERY SPARSE UNDERLYING GRAPHS WHERE P N ELSEWHERE PREFER CD WHICH IS MORE  
NUMERICALLY STABLE  
TOLPOSITIVE FLOAT OPTIONAL THE TOLERANCE TO DECLARE CONVERGENCE IF THE DUAL GAP GOES  
BELOW THIS VALUE ITERATIONS ARE STOPPED  
ENETTOL POSITIVE FLOAT OPTIONAL THE TOLERANCE FOR THE ELASTIC NET SOLVER USED TO CALCULATE  
THE DESCENT DIRECTION THIS PARAMETER CONTROLS THE ACCURACY OF THE SEARCH DIRECTION FOR A  
GIVEN COLUMN UPDATE NOT OF THE OVERALL PARAMETER ESTIMATE ONLY USED FOR MODE'CD'  
MAXITER INTEGER OPTIONAL THE MAXIMUM NUMBER OF ITERATIONS  
VERBOSE BOOLEAN OPTIONAL IF VERBOSE IS TRUE THE OBJECTIVE FUNCTION AND DUAL GAP ARE  
PRINTED AT EACH ITERATION  
RETURNCOSTS BOOLEAN OPTIONAL IF RETURNCOSTS IS TRUE THE OBJECTIVE FUNCTION AND DUAL  
GAP AT EACH ITERATION ARE RETURNED  
EPS FLOAT OPTIONAL THE MACHINEPRECISION REGULARIZATION IN THE COMPUTATION OF THE  
CHOLESKY DIAGONAL FACTORS INCREASE THIS FOR VERY ILLCONDITIONED SYSTEMS  
RETURNNITER BOOL OPTIONAL WHETHER OR NOT TO RETURN THE NUMBER OF ITERATIONS  
RETURNS  
COVARIANCE 2D NDARRAY SHAPE NFEATURES NFEATURES  
THE ESTIMATED COVARIANCE MATRIX  
PRECISION 2D NDARRAY SHAPE NFEATURES NFEATURES THE ESTIMATED SPARSE PRECISION  
MATRIX  
2404 CHAPTER 6 API REFERENCE

SCIKITLEARN USER GUIDE RELEASE 0213

COSTS LIST OF OBJECTIVE DUALGAP PAIRS THE LIST OF VALUES OF THE OBJECTIVE FUNCTION AND THE DUAL GAP AT EACH ITERATION RETURNED ONLY IF RETURNCOSTS IS TRUE

NITER INT NUMBER OF ITERATIONS RETURNED ONLY IF RETURNNITER IS SET TO TRUE

NOTES

THE ALGORITHM EMPLOYED TO SOLVE THIS PROBLEM IS THE GLASSO ALGORITHM FROM THE FRIEDMAN 2008 BIOSTATISTICS PAPER IT IS THE SAME ALGORITHM AS IN THE R GLASSO PACKAGE

ONE POSSIBLE DIFFERENCE WITH THE GLASSO R PACKAGE IS THAT THE DIAGONAL COEFFICIENTS ARE NOT PENALIZED

SKLEARNDATASETS FETCHMLDATA

WARNING DEPRECATED

SKLEARNDATASETS FETCHMLDATA DATANAME TARGETNAME'LABEL' DATANAME'DATA' TRANS

POSEDATATRUE DATAHOMENONE

DEPRECATED FETCHMLDATA WAS DEPRECATED IN VERSION 020 AND WILL BE REMOVED IN VERSION 022 PLEASE USE FETCHOPENML

FETCH AN MLDATAORG DATA SET

MLDATAORG IS NO LONGER OPERATIONAL

IF THE FILE DOES NOT EXIST YET IT IS DOWNLOADED FROM MLDATAORG

MLDATAORG DOES NOT HAVE AN ENFORCED CONVENTION FOR STORING DATA OR NAMING THE COLUMNS IN A DATA SET THE DEFAULT BEHAVIOR OF THIS FUNCTION WORKS WELL WITH THE MOST COMMON CASES

1 DATA VALUES ARE STORED IN THE COLUMN 'DATA' AND TARGET VALUES IN THE COLUMN 'LABEL'

2 ALTERNATIVELY THE FIRST COLUMN STORES TARGET VALUES AND THE SECOND DATA VALUES

3 THE DATA ARRAY IS STORED AS NFEATURES X NSAMPLES AND THUS NEEDS TO BE TRANSPOSED TO MATCH THESKLEARN STANDARD

KEYWORD ARGUMENTS ALLOW TO ADAPT THESE DEFAULTS TO SPECIFIC DATA SETS SEE PARAMETERS

TARGETNAME DATANAME TRANSPOSEDATA AND THE EXAMPLES BELOW

MLDATAORG DATA SETS MAY HAVE MULTIPLE COLUMNS WHICH ARE STORED IN THE BUNCH OBJECT WITH THEIR ORIGINAL NAME

DEPRECATED SINCE VERSION 020 WILL BE REMOVED IN VERSION 022

PARAMETERS

DATANAME STR

NAME OF THE DATA SET ON MLDATAORG EG "LEUKEMIA" "WHISTLER DAILY SNOWFALL" ETC

THE RAW NAME IS AUTOMATICALLY CONVERTED TO A MLDATAORG URL

TARGETNAME OPTIONAL DEFAULT 'LABEL' NAME OR INDEX OF THE COLUMN CONTAINING THE TARGET VALUES

DATANAME OPTIONAL DEFAULT 'DATA' NAME OR INDEX OF THE COLUMN CONTAINING THE DATA

TRANSPOSEDATA OPTIONAL DEFAULT TRUE IF TRUE TRANSPOSE THE DOWNLOADED DATA ARRAY

640 RECENTLY DEPRECATED 2405

SCIKITLEARN USER GUIDE RELEASE 0213  
DATAHOME OPTIONAL DEFAULT NONE SPECIFY ANOTHER DOWNLOAD AND CACHE FOLDER FOR THE  
DATA SETS BY DEFAULT ALL SCIKITLEARN DATA IS STORED IN 'SCIKITLEARNDATA' SUBFOLDERS  
RETURNS  
DATA BUNCH DICTIONARYLIKE OBJECT THE INTERESTING ATTRIBUTES ARE 'DATA' THE DATA TO LEARN  
'TARGET' THE CLASSIFICATION LABELS 'DESCR' THE FULL DESCRIPTION OF THE DATASET AND  
'COLNAMES' THE ORIGINAL NAMES OF THE DATASET COLUMNS  
SKLEARNDATASETS MLDATAFILENAME  
WARNING DEPRECATED  
SKLEARNDATASETS MLDATAFILENAME DATANAME  
DEPRECATED MLDATAFILENAME WAS DEPRECATED IN VERSION 020 AND WILL BE REMOVED IN VERSION 022 PLEASE  
USE FETCHOPENML  
CONVERT A RAW NAME FOR A DATA SET IN A MLDATAORG FILENAME  
DEPRECATED SINCE VERSION 020 WILL BE REMOVED IN VERSION 022  
PARAMETERS  
DATANAME STR NAME OF DATASET  
RETURNS  
FNAME STR THE CONVERTED DATANAME  
2406 CHAPTER 6 API REFERENCE

CHAPTER  
SEVEN  
DEVELOPER’S GUIDE  
71 CONTRIBUTING  
THIS PROJECT IS A COMMUNITY EFFORT AND EVERYONE IS WELCOME TO CONTRIBUTE  
THE PROJECT IS HOSTED ON [HTTPSGITHUBCOMSCIKITLEARNSCIKITLEARN](https://github.com/scikitlearn/scikitlearn)  
THE DECISION MAKING PROCESS AND GOVERNANCE STRUCTURE OF SCIKITLEARN IS LAID OUT IN THE GOVERNANCE DOCUMENT SCIKIT  
LEARN GOVERNANCE AND DECISIONMAKING  
SCIKITLEARN IS SOMEWHAT SELECTIVE WHEN IT COMES TO ADDING NEW ALGORITHMS AND THE BEST WAY TO CONTRIBUTE AND TO HELP  
THE PROJECT IS TO START WORKING ON KNOWN ISSUES SEE ISSUES FOR NEW CONTRIBUTORS TO GET STARTED  
OUR COMMUNITY OUR VALUES  
WE ARE A COMMUNITY BASED ON OPENNESS AND FRIENDLY DIDACTIC DISCUSSIONS  
WE ASPIRE TO TREAT EVERYBODY EQUALLY AND VALUE THEIR CONTRIBUTIONS  
DECISIONS ARE MADE BASED ON TECHNICAL MERIT AND CONSENSUS  
CODE IS NOT THE ONLY WAY TO HELP THE PROJECT REVIEWING PULL REQUESTS ANSWERING QUESTIONS TO HELP OTHERS ON  
MAILING LISTS OR ISSUES ORGANIZING AND TEACHING TUTORIALS WORKING ON THE WEBSITE IMPROVING THE DOCUMENTATION ARE  
ALL PRICELESS CONTRIBUTIONS  
WE ABIDE BY THE PRINCIPLES OF OPENNESS RESPECT AND CONSIDERATION OF OTHERS OF THE PYTHON SOFTWARE FOUNDATION  
[HTTPSWWWPYTHONORGPSFCODEOFCONDUCT](https://www.python.org/psf/codeofconduct)  
IN CASE YOU EXPERIENCE ISSUES USING THIS PACKAGE DO NOT HESITATE TO SUBMIT A TICKET TO THE GITHUB ISSUE TRACKER YOU  
ARE ALSO WELCOME TO POST FEATURE REQUESTS OR PULL REQUESTS  
711 WAYS TO CONTRIBUTE  
THERE ARE MANY WAYS TO CONTRIBUTE TO SCIKITLEARN WITH THE MOST COMMON ONES BEING CONTRIBUTION OF CODE OR DOCU  
MENTATION TO THE PROJECT IMPROVING THE DOCUMENTATION IS NO LESS IMPORTANT THAN IMPROVING THE LIBRARY ITSELF IF YOU  
FIND A TYPO IN THE DOCUMENTATION OR HAVE MADE IMPROVEMENTS DO NOT HESITATE TO SEND AN EMAIL TO THE MAILING LIST OR  
PREFERABLY SUBMIT A GITHUB PULL REQUEST FULL DOCUMENTATION CAN BE FOUND UNDER THE DOC DIRECTORY  
BUT THERE ARE MANY OTHER WAYS TO HELP IN PARTICULAR ANSWERING QUERIES ON THE ISSUE TRACKER INVESTIGATING BUGS  
ANDREVIEWING OTHER DEVELOPERS’ PULL REQUESTS ARE VERY VALUABLE CONTRIBUTIONS THAT DECREASE THE BURDEN ON THE PROJECT  
MAINTAINERS  
2407

SCIKITLEARN USER GUIDE RELEASE 0213

ANOTHER WAY TO CONTRIBUTE IS TO REPORT ISSUES YOU'RE FACING AND GIVE A "THUMBS UP" ON ISSUES THAT OTHERS REPORTED AND THAT ARE RELEVANT TO YOU IT ALSO HELPS US IF YOU SPREAD THE WORD REFERENCE THE PROJECT FROM YOUR BLOG AND ARTICLES LINK TO IT FROM YOUR WEBSITE OR SIMPLY STAR TO SAY "I USE IT"

IN CASE A CONTRIBUTIONISSUE INVOLVES CHANGES TO THE API PRINCIPLES OR CHANGES TO DEPENDENCIES OR SUPPORTED VERSIONS IT MUST BE BACKED BY A ENHANCEMENT PROPOSALS SLEPS WHERE A SLEP MUST BE SUBMITTED AS A PULLREQUEST TO ENHANCE MENT PROPOSALS USING THE SLEP TEMPLATE AND FOLLOWS THE DECISIONMAKING PROCESS OUTLINED IN SCIKITLEARN GOVERNANCE AND DECISIONMAKING

CONTRIBUTING TO RELATED PROJECTS

SCIKITLEARN THRIVES IN AN ECOSYSTEM OF SEVERAL RELATED PROJECTS WHICH ALSO MAY HAVE RELEVANT ISSUES TO WORK ON INCLUDING SMALLER PROJECTS SUCH AS

- SCIKITLEARNCONTRIB
- JOBLIB
- SPHINXGALLERY
- NUMPYDOC
- LIACARFF

AND LARGER PROJECTS

- NUMPY
- SCIPY
- MATPLOTLIB
- AND SO ON

LOOK FOR ISSUES MARKED "HELP WANTED" OR SIMILAR HELPING THESE PROJECTS MAY HELP SCIKITLEARN TOO SEE ALSO RELATED PROJECTS

712 SUBMITTING A BUG REPORT OR A FEATURE REQUEST

WE USE GITHUB ISSUES TO TRACK ALL BUGS AND FEATURE REQUESTS FEEL FREE TO OPEN AN ISSUE IF YOU HAVE FOUND A BUG OR WISH TO SEE A FEATURE IMPLEMENTED

IN CASE YOU EXPERIENCE ISSUES USING THIS PACKAGE DO NOT HESITATE TO SUBMIT A TICKET TO THE BUG TRACKER YOU ARE ALSO WELCOME TO POST FEATURE REQUESTS OR PULL REQUESTS

IT IS RECOMMENDED TO CHECK THAT YOUR ISSUE COMPLIES WITH THE FOLLOWING RULES BEFORE SUBMITTING

- VERIFY THAT YOUR ISSUE IS NOT BEING CURRENTLY ADDRESSED BY OTHER ISSUES OR PULL REQUESTS
- IF YOU ARE SUBMITTING AN ALGORITHM OR FEATURE REQUEST PLEASE VERIFY THAT THE ALGORITHM FULFILLS OUR NEW ALGORITHM REQUIREMENTS

• IF YOU ARE SUBMITTING A BUG REPORT WE STRONGLY ENCOURAGE YOU TO FOLLOW THE GUIDELINES IN HOW TO MAKE A GOOD BUG REPORT

HOW TO MAKE A GOOD BUG REPORT

WHEN YOU SUBMIT AN ISSUE TO GITHUB PLEASE DO YOUR BEST TO FOLLOW THESE GUIDELINES THIS WILL MAKE IT A LOT EASIER TO PROVIDE YOU WITH GOOD FEEDBACK

SCIKITLEARN USER GUIDE RELEASE 0213

- THE IDEAL BUG REPORT CONTAINS A SHORT REPRODUCIBLE CODE SNIPPET THIS WAY ANYONE CAN TRY TO REPRODUCE THE BUG EASILY SEE THIS FOR MORE DETAILS IF YOUR SNIPPET IS LONGER THAN AROUND 50 LINES PLEASE LINK TO A GIST OR A GITHUB REPO
- IF NOT FEASIBLE TO INCLUDE A REPRODUCIBLE SNIPPET PLEASE BE SPECIFIC ABOUT WHAT ESTIMATORS ANDOR FUNCTIONS ARE INVOLVED AND THE SHAPE OF THE DATA
- IF AN EXCEPTION IS RAISED PLEASE PROVIDE THE FULL TRACEBACK
- PLEASE INCLUDE YOUR OPERATING SYSTEM TYPE AND VERSION NUMBER AS WELL AS YOUR PYTHON SCIKITLEARN NUMPY AND SCIPY VERSIONS THIS INFORMATION CAN BE FOUND BY RUNNING THE FOLLOWING CODE SNIPPET

```
import sklearn
sklearn.show_versions()
```

NOTE THIS UTILITY FUNCTION IS ONLY AVAILABLE IN SCIKITLEARN V020 FOR PREVIOUS VERSIONS ONE HAS TO EXPLICITLY RUN

```
import platform; print(platform.platform())
import sys; print(sys.version)
import numpy; print(numpy.__version__)
import scipy; print(scipy.__version__)
import sklearn; print(sklearn.__version__)
```

- PLEASE ENSURE ALL CODE SNIPPETS AND ERROR MESSAGES ARE FORMATTED IN APPROPRIATE CODE BLOCKS SEE CREATING AND HIGHLIGHTING CODE BLOCKS FOR MORE DETAILS

713 CONTRIBUTING CODE

NOTE TO AVOID DUPLICATING WORK IT IS HIGHLY ADVISED THAT YOU SEARCH THROUGH THE ISSUE TRACKER AND THE PR LIST IF IN DOUBT ABOUT DUPLICATED WORK OR IF YOU WANT TO WORK ON A NONTRIVIAL FEATURE IT’S RECOMMENDED TO FIRST OPEN AN ISSUE IN THE ISSUE TRACKER TO GET SOME FEEDBACKS FROM CORE DEVELOPERS

HOW TO CONTRIBUTE

THE PREFERRED WAY TO CONTRIBUTE TO SCIKITLEARN IS TO FORK THE MAIN REPOSITORY ON GITHUB THEN SUBMIT A “PULL REQUEST” PR

- 1 CREATE AN ACCOUNT ON GITHUB IF YOU DO NOT ALREADY HAVE ONE
  - 2 FORK THE PROJECT REPOSITORY CLICK ON THE ‘FORK’ BUTTON NEAR THE TOP OF THE PAGE THIS CREATES A COPY OF THE CODE UNDER YOUR ACCOUNT ON THE GITHUB USER ACCOUNT FOR MORE DETAILS ON HOW TO FORK A REPOSITORY SEE THIS GUIDE
  - 3 CLONE YOUR FORK OF THE SCIKITLEARN REPO FROM YOUR GITHUB ACCOUNT TO YOUR LOCAL DISK  
GIT CLONE GITGITHUBCOMYOURLOGINSCIKITLEARNGIT  
CD SCIKITLEARN
  - 4 INSTALL THE DEVELOPMENT DEPENDENCIES  
PIP INSTALL CYTHON PYTEST FLAKE8
  - 5 INSTALL SCIKITLEARN IN EDITABLE MODE
- 71 CONTRIBUTING 2409

SCIKITLEARN USER GUIDE RELEASE 0213

PIP INSTALL EDITABLE

FOR MORE DETAILS ABOUT ADVANCED INSTALLATION SEE THE BUILDING FROM SOURCE SECTION

6 ADD THE UPSTREAM REMOTE THIS SAVES A REFERENCE TO THE MAIN SCIKITLEARN REPOSITORY WHICH YOU CAN USE TO KEEP YOUR REPOSITORY SYNCHRONIZED WITH THE LATEST CHANGES

GIT REMOTE ADD UPSTREAM HTTPSGITHUBCOMSCIKITLEARNSCIKITLEARNGIT

7 FETCH THE UPSTREAM AND THEN CREATE A BRANCH TO HOLD YOUR DEVELOPMENT CHANGES

GIT FETCH UPSTREAM

GIT CHECKOUT B MYFEATURE UPSTREAMMASTER

AND START MAKING CHANGES ALWAYS USE A FEATURE BRANCH IT'S GOOD PRACTICE TO NEVER WORK ON THE MASTER BRANCH

8 DEVELOP THE FEATURE ON YOUR FEATURE BRANCH ON YOUR COMPUTER USING GIT TO DO THE VERSION CONTROL WHEN YOU'RE DONE EDITING ADD CHANGED FILES USING GIT ADD AND THEN GIT COMMIT FILES

GIT ADD MODIFIEDFILES

GIT COMMIT

TO RECORD YOUR CHANGES IN GIT THEN PUSH THE CHANGES TO YOUR GITHUB ACCOUNT WITH

GIT PUSH U ORIGIN MYFEATURE

9 FOLLOW THESE INSTRUCTIONS TO CREATE A PULL REQUEST FROM YOUR FORK THIS WILL SEND AN EMAIL TO THE COMMITTERS YOU MAY WANT TO CONSIDER SENDING AN EMAIL TO THE MAILING LIST FOR MORE VISIBILITY

NOTE IF YOU ARE MODIFYING A CYTHON MODULE YOU HAVE TO RERUN STEP 5 AFTER MODIFICATIONS AND BEFORE TESTING THEM

IT IS OFTEN HELPFUL TO KEEP YOUR LOCAL BRANCH SYNCHRONIZED WITH THE LATEST CHANGES OF THE MAIN SCIKITLEARN REPOSITORY

GIT FETCH UPSTREAM

GIT MERGE UPSTREAMMASTER

SUBSEQUENTLY YOU MIGHT NEED TO SOLVE THE CONFLICTS YOU CAN REFER TO THE GIT DOCUMENTATION RELATED TO RESOLVING MERGE CONFLICT USING THE COMMAND LINE

LEARNING GIT

THE GIT DOCUMENTATION AND HTTPTRYGITHUBIO ARE EXCELLENT RESOURCES TO GET STARTED WITH GIT AND UNDERSTANDING ALL OF THE COMMANDS SHOWN HERE

PULL REQUEST CHECKLIST

BEFORE A PR CAN BE MERGED IT NEEDS TO BE APPROVED BY TWO CORE DEVELOPERS PLEASE PREFIX THE TITLE OF YOUR PULL REQUEST WITHMRG IF THE CONTRIBUTION IS COMPLETE AND SHOULD BE SUBJECTED TO A DETAILED REVIEW AN INCOMPLETE CONTRIBUTION - WHERE YOU EXPECT TO DO MORE WORK BEFORE RECEIVING A FULL REVIEW - SHOULD BE PREFIXED WIP TO INDICATE A WORK IN PROGRESS AND CHANGED TO MRG WHEN IT MATURES WIPS MAY BE USEFUL TO INDICATE YOU ARE WORKING ON SOMETHING TO AVOID DUPLICATED WORK REQUEST BROAD REVIEW OF FUNCTIONALITY OR API OR SEEK COLLABORATORS WIPS OFTEN BENEFIT FROM THE INCLUSION OF A TASK LIST IN THE PR DESCRIPTION

2410 CHAPTER 7 DEVELOPER'S GUIDE



SCIKITLEARN USER GUIDE RELEASE 0213

IN ORDER TO EASE THE REVIEWING PROCESS WE RECOMMEND THAT YOUR CONTRIBUTION COMPLIES WITH THE FOLLOWING RULES BEFORE MARKING A PR AS MRG THE BOLDDED ONES ARE ESPECIALLY IMPORTANT

1GIVE YOUR PULL REQUEST A HELPFUL TITLE THAT SUMMARISES WHAT YOUR CONTRIBUTION DOES THIS TITLE WILL OFTEN BECOME THE COMMIT MESSAGE ONCE MERGED SO IT SHOULD SUMMARISE YOUR CONTRIBUTION FOR POSTERITY IN SOME CASES “FIX ISSUE TITLE” IS ENOUGH “FIX ISSUE NUMBER” IS NEVER A GOOD TITLE

2MAKE SURE YOUR CODE PASSES THE TESTS THE WHOLE TEST SUITE CAN BE RUN WITH PYTEST BUT IT IS USUALLY NOT RECOMMENDED SINCE IT TAKES A LONG TIME IT IS OFTEN ENOUGH TO ONLY RUN THE TEST RELATED TO YOUR CHANGES FOR EXAMPLE IF YOU CHANGED SOMETHING IN SKLEARNLINEARMODELLOGISTICPY RUNNING THE FOLLOWING COMMANDS WILL USUALLY BE ENOUGH

- PYTEST SKLEARNLINEARMODELLOGISTICPY TO MAKE SURE THE DOCTEST EXAMPLES ARE CORRECT
- PYTEST SKLEARNLINEARMODELTESTSTESTLOGISTICPY TO RUN THE TESTS SPECIFIC TO THE FILE
- PYTEST SKLEARNLINEARMODEL TO TEST THE WHOLE GENERALIZED LINEAR MODELS MODULE
- PYTEST SKLEARNDOCLINEARMODELRST TO MAKE SURE THE USER GUIDE EXAMPLES ARE CORRECT
- PYTEST SKLEARNTESTSTESTCOMMONPY K LOGISTICREGRESSION TO RUN ALL OUR ESTI

MATOR CHECKS SPECIFICALLY FOR LOGISTICREGRESSION IF THAT’S THE ESTIMATOR YOU CHANGED THERE MAY BE OTHER FAILING TESTS BUT THEY WILL BE CAUGHT BY THE CI SO YOU DON’T NEED TO RUN THE WHOLE TEST SUITE LOCALLY YOU CAN READ MORE IN TESTING AND IMPROVING TEST COVERAGE

3MAKE SURE YOUR CODE IS PROPERLY COMMENTED AND DOCUMENTED AND MAKE SURE THE DOCUMENTATION RENDERS PROPERLY TO BUILD THE DOCUMENTATION PLEASE REFER TO OUR DOCUMENTATION GUIDELINES THE CI WILL ALSO BUILD THE DOCS PLEASE REFER TO GENERATED DOCUMENTATION ON CIRCLECI

4TESTS ARE NECESSARY FOR ENHANCEMENTS TO BE ACCEPTED BUGFIXES OR NEW FEATURES SHOULD BE PROVIDED WITH NON REGRESSION TESTS THESE TESTS VERIFY THE CORRECT BEHAVIOR OF THE FIX OR FEATURE IN THIS MANNER FURTHER MODIFICATIONS ON THE CODE BASE ARE GRANTED TO BE CONSISTENT WITH THE DESIRED BEHAVIOR IN THE CASE OF BUG FIXES AT THE TIME OF THE PR THE NONREGRESSION TESTS SHOULD FAIL FOR THE CODE BASE IN THE MASTER BRANCH AND PASS FOR THE PR CODE

5MAKE SURE THAT YOUR PR DOES NOT ADD PEP8 VIOLATIONS ON A UNIXLIKE SYSTEM YOU CAN RUN MAKE FLAKE8DIFF FLAKE8 PATHTOFILE WOULD WORK FOR ANY SYSTEM BUT PLEASE AVOID REFORMATTING PARTS OF THE FILE THAT YOUR PULL REQUEST DOESN’T CHANGE AS IT DISTRACTS FROM CODE REVIEW

6 FOLLOW THE CODINGGUIDELINES SEE BELOW

7 WHEN APPLICABLE USE THE VALIDATION TOOLS AND SCRIPTS IN THE SKLEARNUTILS SUBMODULE A LIST OF UTILITY ROUTINES AVAILABLE FOR DEVELOPERS CAN BE FOUND IN THE UTILITIES FOR DEVELOPERS PAGE

8 OFTEN PULL REQUESTS RESOLVE ONE OR MORE OTHER ISSUES OR PULL REQUESTS IF MERGING YOUR PULL REQUEST MEANS THAT SOME OTHER ISSUESPRS SHOULD BE CLOSED YOU SHOULD USE KEYWORDS TO CREATE LINK TO THEM EG FIXES 1234

MULTIPLE ISSUESPRS ARE ALLOWED AS LONG AS EACH ONE IS PRECEDED BY A KEYWORD UPON MERGING THOSE ISSUESPRS WILL AUTOMATICALLY BE CLOSED BY GITHUB IF YOUR PULL REQUEST IS SIMPLY RELATED TO SOME OTHER ISSUESPRS CREATE A LINK TO THEM WITHOUT USING THE KEYWORDS EG SEE ALSO 1234

9 PRS SHOULD OFTEN SUBSTANTIATE THE CHANGE THROUGH BENCHMARKS OF PERFORMANCE AND EFFICIENCY OR THROUGH EXAM PLES OF USAGE EXAMPLES ALSO ILLUSTRATE THE FEATURES AND INTRICACIES OF THE LIBRARY TO USERS HAVE A LOOK AT OTHER EXAMPLES IN THE EXAMPLES DIRECTORY FOR REFERENCE EXAMPLES SHOULD DEMONSTRATE WHY THE NEW FUNCTIONALITY IS USEFUL IN PRACTICE AND IF POSSIBLE COMPARE IT TO OTHER METHODS AVAILABLE IN SCIKITLEARN

10 NEW FEATURES OFTEN NEED TO BE ILLUSTRATED WITH NARRATIVE DOCUMENTATION IN THE USER GUIDE WITH SMALL CODE SNIPETS IF RELEVANT PLEASE ALSO ADD REFERENCES IN THE LITERATURE WITH PDF LINKS WHEN POSSIBLE

11 THE USER GUIDE SHOULD ALSO INCLUDE EXPECTED TIME AND SPACE COMPLEXITY OF THE ALGORITHM AND SCALABILITY EG “THIS ALGORITHM CAN SCALE TO A LARGE NUMBER OF SAMPLES 100000 BUT DOES NOT SCALE IN DIMENSIONALITY NFEATURES IS EXPECTED TO BE LOWER THAN 100”

71 CONTRIBUTING 2411

SCIKITLEARN USER GUIDE RELEASE 0213

YOU CAN ALSO CHECK OUR CODE REVIEW GUIDELINES TO GET AN IDEA OF WHAT REVIEWERS WILL EXPECT

YOU CAN CHECK FOR COMMON PROGRAMMING ERRORS WITH THE FOLLOWING TOOLS

- CODE WITH A GOOD UNITTEST COVERAGE AT LEAST 80 BETTER 100 CHECK WITH  
PIP INSTALL PYTEST PYTESTCOV  
PYTEST COV SKLEARN PATHTOTESTSFORPACKAGE

SEE ALSO TESTING AND IMPROVING TEST COVERAGE

BONUS POINTS FOR CONTRIBUTIONS THAT INCLUDE A PERFORMANCE ANALYSIS WITH A BENCHMARK SCRIPT AND PROFILING OUTPUT PLEASE REPORT ON THE MAILING LIST OR ON THE GITHUB ISSUE

ALSO CHECK OUT THE HOW TO OPTIMIZE FOR SPEED GUIDE FOR MORE DETAILS ON PROFILING AND CYTHON OPTIMIZATIONS

NOTE THE CURRENT STATE OF THE SCIKITLEARN CODE BASE IS NOT COMPLIANT WITH ALL OF THOSE GUIDELINES BUT WE EXPECT THAT ENFORCING THOSE CONSTRAINTS ON ALL NEW CONTRIBUTIONS WILL GET THE OVERALL CODE BASE QUALITY IN THE RIGHT DIRECTION

NOTE FOR TWO VERY WELL DOCUMENTED AND MORE DETAILED GUIDES ON DEVELOPMENT WORKFLOW PLEASE PAY A VISIT TO THE SCIPY DEVELOPMENT WORKFLOW AND THE ASTROPY WORKFLOW FOR DEVELOPERS SECTIONS

CONTINUOUS INTEGRATION CI

- AZURE PIPELINES ARE USED FOR TESTING SCIKITLEARN ON LINUX MAC AND WINDOWS WITH DIFFERENT DEPENDENCIES AND SETTINGS
- CIRCLECI IS USED TO BUILD THE DOCS FOR VIEWING FOR LINTING WITH FLAKE8 AND FOR TESTING WITH PYPY ON LINUX

PLEASE NOTE THAT IF ONE OF THE FOLLOWING MARKERS APPEAR IN THE LATEST COMMIT MESSAGE THE FOLLOWING ACTIONS ARE TAKEN

COMMIT MESSAGE

MARKERACTION TAKEN BY CI

SCIPYDEV ADD A TRAVIS BUILD WITH OUR DEPENDENCIES NUMPY SCIPY ETC    DEVELOP

MENT BUILDS

CI SKIP CI IS SKIPPED COMPLETELY

DOC SKIP DOCS ARE NOT BUILT

DOC QUICK DOCS BUILT BUT EXCLUDES EXAMPLE GALLERY PLOTS

DOC BUILD DOCS BUILT INCLUDING EXAMPLE GALLERY PLOTS

STALLED PULL REQUESTS

AS CONTRIBUTING A FEATURE CAN BE A LENGTHY PROCESS SOME PULL REQUESTS APPEAR INACTIVE BUT UNFINISHED IN SUCH A CASE TAKING THEM OVER IS A GREAT SERVICE FOR THE PROJECT

A GOOD ETIQUETTE TO TAKE OVER IS

- DETERMINE IF A PR IS STALLED

-A PULL REQUEST MAY HAVE THE LABEL “STALLED” OR “HELP WANTED” IF WE HAVE ALREADY IDENTIFIED IT AS A CANDIDATE FOR OTHER CONTRIBUTORS

2412 CHAPTER 7 DEVELOPER’S GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

-TO DECIDE WHETHER AN INACTIVE PR IS STALLED ASK THE CONTRIBUTOR IF SHEHE PLANS TO CONTINUE WORKING ON THE PR IN THE NEAR FUTURE FAILURE TO RESPOND WITHIN 2 WEEKS WITH AN ACTIVITY THAT MOVES THE PR FORWARD SUGGESTS THAT THE PR IS STALLED AND WILL RESULT IN TAGGING THAT PR WITH “HELP WANTED”  
NOTE THAT IF A PR HAS RECEIVED EARLIER COMMENTS ON THE CONTRIBUTION THAT HAVE HAD NO REPLY IN A MONTH IT IS SAFE TO ASSUME THAT THE PR IS STALLED AND TO SHORTEN THE WAIT TIME TO ONE DAY

AFTER A SPRINT FOLLOWUP FOR UNMERGED PRS OPENED DURING SPRINT WILL BE COMMUNICATED TO PARTICIPANTS AT THE SPRINT AND THOSE PRS WILL BE TAGGED “SPRINT” PRS TAGGED WITH “SPRINT” CAN BE REASSIGNED OR DECLARED STALLED BY SPRINT LEADERS

•TAKING OVER A STALLED PR TO TAKE OVER A PR IT IS IMPORTANT TO COMMENT ON THE STALLED PR THAT YOU ARE TAKING OVER AND TO LINK FROM THE NEW PR TO THE OLD ONE THE NEW PR SHOULD BE CREATED BY PULLING FROM THE OLD ONE  
ISSUES FOR NEW CONTRIBUTORS

NEW CONTRIBUTORS SHOULD LOOK FOR THE FOLLOWING TAGS WHEN LOOKING FOR ISSUES WE STRONGLY RECOMMEND THAT NEW CONTRIBUTORS TACKLE “EASY” ISSUES FIRST THIS HELPS THE CONTRIBUTOR BECOME FAMILIAR WITH THE CONTRIBUTION WORKFLOW AND FOR THE CORE DEVS TO BECOME ACQUAINTED WITH THE CONTRIBUTOR BESIDES WHICH WE FREQUENTLY UNDERESTIMATE HOW EASY AN ISSUE IS TO SOLVE

GOOD FIRST ISSUE TAG

A GREAT WAY TO START CONTRIBUTING TO SCIKITLEARN IS TO PICK AN ITEM FROM THE LIST OF GOOD FIRST ISSUES IN THE ISSUE TRACKER RESOLVING THESE ISSUES ALLOW YOU TO START CONTRIBUTING TO THE PROJECT WITHOUT MUCH PRIOR KNOWLEDGE IF YOU HAVE ALREADY CONTRIBUTED TO SCIKITLEARN YOU SHOULD LOOK AT EASY ISSUES INSTEAD

EASY TAG

IF YOU HAVE ALREADY CONTRIBUTED TO SCIKITLEARN ANOTHER GREAT WAY TO CONTRIBUTE TO SCIKITLEARN IS TO PICK AN ITEM FROM THE LIST OF EASY ISSUES IN THE ISSUE TRACKER YOUR ASSISTANCE IN THIS AREA WILL BE GREATLY APPRECIATED BY THE MORE EXPERIENCED DEVELOPERS AS IT HELPS FREE UP THEIR TIME TO CONCENTRATE ON OTHER ISSUES

HELP WANTED TAG

WE OFTEN USE THE HELP WANTED TAG TO MARK ISSUES REGARDLESS OF DIFFICULTY ADDITIONALLY WE USE THE HELP WANTED TAG TO MARK PULL REQUESTS WHICH HAVE BEEN ABANDONED BY THEIR ORIGINAL CONTRIBUTOR AND ARE AVAILABLE FOR SOMEONE TO PICK UP WHERE THE ORIGINAL CONTRIBUTOR LEFT OFF THE LIST OF ISSUES WITH THE HELP WANTED TAG CAN BE FOUND HERE

NOTE THAT NOT ALL ISSUES WHICH NEED CONTRIBUTORS WILL HAVE THIS TAG

714 DOCUMENTATION

WE ARE GLAD TO ACCEPT ANY SORT OF DOCUMENTATION FUNCTION DOCSTRINGS RESTRUCTUREDTEXT DOCUMENTS LIKE THIS ONE TUTORIALS ETC RESTRUCTUREDTEXT DOCUMENTS LIVE IN THE SOURCE CODE REPOSITORY UNDER THE DOC DIRECTORY YOU CAN EDIT THE DOCUMENTATION USING ANY TEXT EDITOR AND THEN GENERATE THE HTML OUTPUT BY TYPING MAKE FROM THEDOC DIRECTORY ALTERNATIVELY MAKE HTML MAY BE USED TO GENERATE THE DOCUMENTATION WITH THE EXAMPLE GALLERY WHICH TAKES QUITE SOME TIME THE RESULTING HTML FILES WILL BE PLACED IN BUILDHTMLSTABLE AND ARE VIEWABLE IN A WEB BROWSER

71 CONTRIBUTING 2413

SCIKITLEARN USER GUIDE RELEASE 0213

BUILDING THE DOCUMENTATION

FIRST MAKE SURE YOU HAVE PROPERLY INSTALLED THE DEVELOPMENT VERSION

BUILDING THE DOCUMENTATION REQUIRES INSTALLING SOME ADDITIONAL PACKAGES

PIP INSTALL SPHINX SPHINXGALLERY NUMPYDOC MATPLOTLIB PILLOW PANDAS SCIKITIMAGE

TO BUILD THE DOCUMENTATION YOU NEED TO BE IN THE DOC FOLDER

CD DOC

IN THE VAST MAJORITY OF CASES YOU ONLY NEED TO GENERATE THE FULL WEB SITE WITHOUT THE EXAMPLE GALLERY

MAKE

THE DOCUMENTATION WILL BE GENERATED IN THE BUILDHTMLSTABLE DIRECTORY TO ALSO GENERATE THE EXAMPLE GALLERY

YOU CAN USE

MAKE HTML

THIS WILL RUN ALL THE EXAMPLES WHICH TAKES A WHILE IF YOU ONLY WANT TO GENERATE A FEW EXAMPLES YOU CAN USE

EXAMPLESPATTERNYOURREGEXGOESHERE MAKE HTML

THIS IS PARTICULARLY USEFUL IF YOU ARE MODIFYING A FEW EXAMPLES

SET THE ENVIRONMENT VARIABLE NOMATHJAX1 IF YOU INTEND TO VIEW THE DOCUMENTATION IN AN OFFLINE SETTING

TO BUILD THE PDF MANUAL RUN

MAKE LATEXPDF

WARNING SPHINX VERSION

WHILE WE DO OUR BEST TO HAVE THE DOCUMENTATION BUILD UNDER AS MANY VERSIONS OF SPHINX AS POSSIBLE THE DIFFERENT

VERSIONS TEND TO BEHAVE SLIGHTLY DIFFERENTLY TO GET THE BEST RESULTS YOU SHOULD USE THE SAME VERSION AS THE ONE WE

USED ON CIRCLECI LOOK AT THIS GITHUB SEARCH TO KNOW THE EXACT VERSION

GUIDELINES FOR WRITING DOCUMENTATION

IT IS IMPORTANT TO KEEP A GOOD COMPROMISE BETWEEN MATHEMATICAL AND ALGORITHMIC DETAILS AND GIVE INTUITION TO THE

READER ON WHAT THE ALGORITHM DOES

BASICALLY TO ELABORATE ON THE ABOVE IT IS BEST TO ALWAYS START WITH A SMALL PARAGRAPH WITH A HANDWAVING EXPLANATION OF

WHAT THE METHOD DOES TO THE DATA THEN IT IS VERY HELPFUL TO POINT OUT WHY THE FEATURE IS USEFUL AND WHEN IT SHOULD BE

USED THE LATTER ALSO INCLUDING “BIG O”  $\mathcal{O}$  COMPLEXITIES OF THE ALGORITHM AS OPPOSED TO JUST RULES OF THUMB AS

THE LATTER CAN BE VERY MACHINEDEPENDENT IF THOSE COMPLEXITIES ARE NOT AVAILABLE THEN RULES OF THUMB MAY BE PROVIDED

INSTEAD

SECONDLY A GENERATED FIGURE FROM AN EXAMPLE AS MENTIONED IN THE PREVIOUS PARAGRAPH SHOULD THEN BE INCLUDED TO

FURTHER PROVIDE SOME INTUITION

NEXT ONE OR TWO SMALL CODE EXAMPLES TO SHOW ITS USE CAN BE ADDED

2414 CHAPTER 7 DEVELOPER’S GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

NEXT ANY MATH AND EQUATIONS FOLLOWED BY REFERENCES CAN BE ADDED TO FURTHER THE DOCUMENTATION NOT STARTING THE DOCUMENTATION WITH THE MATHS MAKES IT MORE FRIENDLY TOWARDS USERS THAT ARE JUST INTERESTED IN WHAT THE FEATURE WILL DO AS OPPOSED TO HOW IT WORKS “UNDER THE HOOD”

FINALLY FOLLOW THE FORMATTING RULES BELOW TO MAKE IT CONSISTENTLY GOOD

- ADD “SEE ALSO” IN DOCSTRINGS FOR RELATED CLASSESFUNCTIONS
- “SEE ALSO” IN DOCSTRINGS SHOULD BE ONE LINE PER REFERENCE WITH A COLON AND AN EXPLANATION FOR EXAMPLE SEE ALSO

SELECTKBEST SELECT FEATURES BASED ON THE K HIGHEST SCORES

SELECTFPR SELECT FEATURES BASED ON A FALSE POSITIVE RATE TEST

- FOR UNWRITTEN FORMATTING RULES TRY TO FOLLOW EXISTING GOOD WORKS
- FOR “REFERENCES” IN DOCSTRINGS SEE THE SILHOUETTE COEFFICIENT SKLEARNMETRICS

SILHOUETTESCORE

- WHEN EDITING RESTRUCTUREDTEXT RST FILES TRY TO KEEP LINE LENGTH UNDER 80 CHARACTERS WHEN POSSIBLE EXCEPTIONS INCLUDE LINKS AND TABLES

GENERATED DOCUMENTATION ON CIRCLECI

WHEN YOU CHANGE THE DOCUMENTATION IN A PULL REQUEST CIRCLECI AUTOMATICALLY BUILDS IT TO VIEW THE DOCUMENTATION GENERATED BY CIRCLECI

- NAVIGATE TO THE BOTTOM OF YOUR PULL REQUEST PAGE TO SEE THE CI STATUSES YOU MAY NEED TO CLICK ON “SHOW ALL CHECKS” TO SEE ALL THE CI STATUSES

- CLICK ON THE CIRCLECI STATUS WITH “DOC” IN THE TITLE

- ADDARTIFACTS AT THE END OF THE URL NOTE YOU NEED TO WAIT FOR THE CIRCLECI BUILD TO FINISH BEFORE BEING ABLE TO LOOK AT THE ARTIFACTS

- ONCE THE ARTIFACTS ARE VISIBLE NAVIGATE TO DOCCHANGEDHTML TO SEE A LIST OF DOCUMENTATION PAGES THAT ARE LIKELY TO BE AFFECTED BY YOUR PULL REQUEST NAVIGATE TO DOCINDEXHTML TO SEE THE FULL GENERATED HTML DOCUMENTATION

IF YOU OFTEN NEED TO LOOK AT THE DOCUMENTATION GENERATED BY CIRCLECI EG WHEN REVIEWING PULL REQUESTS YOU MAY FIND THIS TIP VERY HANDY

715 TESTING AND IMPROVING TEST COVERAGE

HIGHQUALITY UNIT TESTING IS A CORNERSTONE OF THE SCIKITLEARN DEVELOPMENT PROCESS FOR THIS PURPOSE WE USE THE PYTEST PACKAGE THE TESTS ARE FUNCTIONS APPROPRIATELY NAMED LOCATED IN TESTS SUBDIRECTORIES THAT CHECK THE VALIDITY OF THE ALGORITHMS AND THE DIFFERENT OPTIONS OF THE CODE

THE FULL SCIKITLEARN TESTS CAN BE RUN USING ‘MAKE’ IN THE ROOT FOLDER ALTERNATIVELY RUNNING ‘PYTEST’ IN A FOLDER WILL RUN ALL THE TESTS OF THE CORRESPONDING SUBPACKAGES

WE EXPECT CODE COVERAGE OF NEW FEATURES TO BE AT LEAST AROUND 90

FOR GUIDELINES ON HOW TO USE PYTEST EFFICIENTLY SEE THE USEFUL PYTEST ALIASES AND FLAGS

71 CONTRIBUTING 2415

SCIKITLEARN USER GUIDE RELEASE 0213

WRITING MATPLOTLIB RELATED TESTS

TEST FIXTURES ENSURE THAT A SET OF TESTS WILL BE EXECUTING WITH THE APPROPRIATE INITIALIZATION AND CLEANUP THE SCIKITLEARN TEST SUITE IMPLEMENTS A FIXTURE WHICH CAN BE USED WITH MATPLOTLIB

PYPLOT THEPYPLOT FIXTURE SHOULD BE USED WHEN A TEST FUNCTION IS DEALING WITH MATPLOTLIB MATPLOTLIB IS A SOFT DEPENDENCY AND IS NOT REQUIRED THIS FIXTURE IS IN CHARGE OF SKIPPING THE TESTS IF MATPLOTLIB IS NOT INSTALLED IN ADDITION FIGURES CREATED DURING THE TESTS WILL BE AUTOMATICALLY CLOSED ONCE THE TEST FUNCTION HAS BEEN EXECUTED

TO USE THIS FIXTURE IN A TEST FUNCTION ONE NEEDS TO PASS IT AS AN ARGUMENT

DEFTESTREQUIRINGMPLFIXTUREPYPLOT

YOU CAN NOW SAFELY USE MATPLOTLIB

WORKFLOW TO IMPROVE TEST COVERAGE

TO TEST CODE COVERAGE YOU NEED TO INSTALL THE COVERAGE PACKAGE IN ADDITION TO PYTEST

1RUN 'MAKE TESTCOVERAGE' THE OUTPUT LISTS FOR EACH FILE THE LINE NUMBERS THAT ARE NOT TESTED

2FIND A LOW HANGING FRUIT LOOKING AT WHICH LINES ARE NOT TESTED WRITE OR ADAPT A TEST SPECIFICALLY FOR THESE LINES

3 LOOP

ISSUE TRACKER TAGS

ALL ISSUES AND PULL REQUESTS ON THE GITHUB ISSUE TRACKER SHOULD HAVE AT LEAST ONE OF THE FOLLOWING TAGS

BUG CRASH SOMETHING IS HAPPENING THAT CLEARLY SHOULDN'T HAPPEN WRONG RESULTS AS WELL AS UNEXPECTED ERRORS FROM ESTIMATORS GO HERE

CLEANUP ENHANCEMENT IMPROVING PERFORMANCE USABILITY CONSISTENCY

DOCUMENTATION MISSING INCORRECT OR SUBSTANDARD DOCUMENTATIONS AND EXAMPLES

NEW FEATURE FEATURE REQUESTS AND PULL REQUESTS IMPLEMENTING A NEW FEATURE

THERE ARE FOUR OTHER TAGS TO HELP NEW CONTRIBUTORS

GOOD FIRST ISSUE THIS ISSUE IS IDEAL FOR A FIRST CONTRIBUTION TO SCIKITLEARN ASK FOR HELP IF THE FORMULATION IS UNCLEAR IF YOU HAVE ALREADY CONTRIBUTED TO SCIKITLEARN LOOK AT EASY ISSUES INSTEAD

EASY THIS ISSUE CAN BE TACKLED WITHOUT MUCH PRIOR EXPERIENCE

MODERATE MIGHT NEED SOME KNOWLEDGE OF MACHINE LEARNING OR THE PACKAGE BUT IS STILL APPROACHABLE FOR SOMEONE NEW TO THE PROJECT

HELP WANTED THIS TAG MARKS AN ISSUE WHICH CURRENTLY LACKS A CONTRIBUTOR OR A PR THAT NEEDS ANOTHER CONTRIBUTOR TO TAKE OVER THE WORK THESE ISSUES CAN RANGE IN DIFFICULTY AND MAY NOT BE APPROACHABLE FOR NEW CONTRIBUTORS NOTE THAT NOT ALL ISSUES WHICH NEED CONTRIBUTORS WILL HAVE THIS TAG

716 CODING GUIDELINES

THE FOLLOWING ARE SOME GUIDELINES ON HOW NEW CODE SHOULD BE WRITTEN OF COURSE THERE ARE SPECIAL CASES AND THERE WILL BE EXCEPTIONS TO THESE RULES HOWEVER FOLLOWING THESE RULES WHEN SUBMITTING NEW CODE MAKES THE REVIEW EASIER SO NEW CODE CAN BE INTEGRATED IN LESS TIME

2416 CHAPTER 7 DEVELOPER'S GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

UNIFORMLY FORMATTED CODE MAKES IT EASIER TO SHARE CODE OWNERSHIP THE SCIKITLEARN PROJECT TRIES TO CLOSELY FOLLOW THE OFFICIAL PYTHON GUIDELINES DETAILED IN PEP8 THAT DETAIL HOW CODE SHOULD BE FORMATTED AND INDENTED PLEASE READ IT AND FOLLOW IT

IN ADDITION WE ADD THE FOLLOWING GUIDELINES

- USE UNDERSCORES TO SEPARATE WORDS IN NON CLASS NAMES NSAMPLES RATHER THAN NSAMPLES
- AVOID MULTIPLE STATEMENTS ON ONE LINE PREFER A LINE RETURN AFTER A CONTROL FLOW STATEMENT IFFOR
- USE RELATIVE IMPORTS FOR REFERENCES INSIDE SCIKITLEARN
- UNIT TESTS ARE AN EXCEPTION TO THE PREVIOUS RULE THEY SHOULD USE ABSOLUTE IMPORTS EXACTLY AS CLIENT CODE WOULD

A COROLLARY IS THAT IF SKLEARNFOO EXPORTS A CLASS OR FUNCTION THAT IS IMPLEMENTED IN SKLEARNFOOBAR BAZ THE TEST SHOULD IMPORT IT FROM SKLEARNFOO

- PLEASE DON'T USE IMPORTIN ANY CASE IT IS CONSIDERED HARMFUL BY THE OFFICIAL PYTHON RECOMMENDATIONS IT MAKES THE CODE HARDER TO READ AS THE ORIGIN OF SYMBOLS IS NO LONGER EXPLICITLY REFERENCED BUT MOST IMPORTANT IT PREVENTS USING A STATIC ANALYSIS TOOL LIKE PYFLAKES TO AUTOMATICALLY FIND BUGS IN SCIKITLEARN
- USE THE NUMPY DOCSTRING STANDARD IN ALL YOUR DOCSTRINGS

A GOOD EXAMPLE OF CODE THAT WE LIKE CAN BE FOUND HERE

INPUT VALIDATION

THE MODULE SKLEARNUTILS CONTAINS VARIOUS FUNCTIONS FOR DOING INPUT VALIDATION AND CONVERSION SOMETIMES NPASARRAY SUFFICES FOR VALIDATION DO NOTUSENPASANYARRAY ORNPATLEAST2D SINCE THOSE LET NUMPY'S NPMATRIX THROUGH WHICH HAS A DIFFERENT API EG MEANS DOT PRODUCT ON NPMATRIX BUT HADAMARD PRODUCT ONNPNDDARRAY

IN OTHER CASES BE SURE TO CALL CHECKARRAY ON ANY ARRAYLIKE ARGUMENT PASSED TO A SCIKITLEARN API FUNCTION THE EXACT PARAMETERS TO USE DEPENDS MAINLY ON WHETHER AND WHICH SCIPYSPARSE MATRICES MUST BE ACCEPTED FOR MORE INFORMATION REFER TO THE UTILITIES FOR DEVELOPERS PAGE

RANDOM NUMBERS

IF YOUR CODE DEPENDS ON A RANDOM NUMBER GENERATOR DO NOT USE NUMPYRANDOMRANDOM OR SIMILAR ROUTINES TO ENSURE REPEATABILITY IN ERROR CHECKING THE ROUTINE SHOULD ACCEPT A KEYWORD RANDOMSTATE AND USE THIS TO CONSTRUCT ANUMPYRANDOMRANDOMSTATE OBJECT SEE SKLEARNUTILSCHECKRANDOMSTATE INUTILITIES FOR DEVELOPERS

HERE'S A SIMPLE EXAMPLE OF CODE USING SOME OF THE ABOVE GUIDELINES

```
FROM SKLEARNUTILS IMPORT CHECKARRAY CHECKRANDOMSTATE
DEFCHOOSERANDOMSAMPLE(X, RANDOMSTATE)

    CHOOSE A RANDOM POINT FROM X
    PARAMETERS

    X: ARRAYLIKE SHAPE (NSAMPLES, NFEATURES)
    ARRAY REPRESENTING THE DATA
    RANDOMSTATE: RANDOMSTATE OR AN INT SEED 0 BY DEFAULT
    A RANDOM NUMBER GENERATOR INSTANCE TO DEFINE THE STATE OF THE
    RANDOM PERMUTATIONS GENERATOR
```

71 CONTRIBUTING 2417

SCIKITLEARN USER GUIDE RELEASE 0213  
RETURNS

X NUMPY ARRAY SHAPE NFEATURES  
A RANDOM POINT SELECTED FROM X

X CHECKARRAYX  
RANDOMSTATE CHECKRANDOMSTATERANDOMSTATE  
I RANDOMSTATERANDINTXSHAPE0  
RETURNXI

IF YOU USE RANDOMNESS IN AN ESTIMATOR INSTEAD OF A FREESTANDING FUNCTION SOME ADDITIONAL GUIDELINES APPLY  
FIRST OFF THE ESTIMATOR SHOULD TAKE A RANDOMSTATE ARGUMENT TO ITS INIT WITH A DEFAULT VALUE OF  
NONE IT SHOULD STORE THAT ARGUMENT’S VALUE UNMODIFIED IN AN ATTRIBUTE RANDOMSTATE FIT CAN CALL  
CHECKRANDOMSTATE ON THAT ATTRIBUTE TO GET AN ACTUAL RANDOM NUMBER GENERATOR IF FOR SOME REASON RAN  
DOMNESS IS NEEDED AFTER FIT THE RNG SHOULD BE STORED IN AN ATTRIBUTE RANDOMSTATE THE FOLLOWING EXAMPLE  
SHOULD MAKE THIS CLEAR  
CLASS GAUSSIANNNOISE BASEESTIMATOR TRANSFORMERMIXIN  
THIS ESTIMATOR IGNORES ITS INPUT AND RETURNS RANDOM GAUSSIAN NOISE  
IT ALSO DOES NOT ADHERE TO ALL SCIKITLEARN CONVENTIONS  
BUT SHOWCASES HOW TO HANDLE RANDOMNESS

DEFINITSELF NCOMPONENTS100 RANDOMSTATE NONE  
SELFRANDOMSTATE RANDOMSTATE  
THE ARGUMENTS ARE IGNORED ANYWAY SO WE MAKE THEM OPTIONAL  
DEFFITSELF X NONE YNONE  
SELFRANDOMSTATE CHECKRANDOMSTATESELFRANDOMSTATE  
DEFTRANSFORMSELF X  
NSAMPLES XSHAPE0  
RETURNSELFRANDOMSTATERANDNNSAMPLES NCOMPONENTS  
THE REASON FOR THIS SETUP IS REPRODUCIBILITY WHEN AN ESTIMATOR IS FIT TWICE TO THE SAME DATA IT SHOULD PRODUCE AN  
IDENTICAL MODEL BOTH TIMES HENCE THE VALIDATION IN FIT NOTINIT  
DEPRECATION

IF ANY PUBLICLY ACCESSIBLE METHOD FUNCTION ATTRIBUTE OR PARAMETER IS RENAMED WE STILL SUPPORT THE OLD ONE FOR TWO  
RELEASES AND ISSUE A DEPRECATION WARNING WHEN IT IS CALLEDPASSEDACCESSED EG IF THE FUNCTION ZEROONE IS RE  
NAMED TOZEROONELOSS WE ADD THE DECORATOR DEPRECATED FROMSKLEARNUTILS TOZEROONE AND CALL  
ZEROONELOSS FROM THAT FUNCTION  
FROM UTILS IMPORT DEPRECATED  
DEFZEROONELOSSYTRUE YPRED NORMALIZE TRUE  
ACTUAL IMPLEMENTATION  
PASS  
DEPRECATED FUNCTION ZEROONE WAS RENAMED TO ZEROONELOSS  
IN VERSION 013 AND WILL BE REMOVED IN RELEASE 015  
DEFAULT BEHAVIOR IS CHANGED FROM NORMALIZEFALSE TO  
NORMALIZETRUE  
2418 CHAPTER 7 DEVELOPER’S GUIDE



SCIKITLEARN USER GUIDE RELEASE 0213

DEFZEROONEYTRUE YPRED NORMALIZE FALSE

RETURNZEROONELOSSYTRUE YPRED NORMALIZE

IF AN ATTRIBUTE IS TO BE DEPRECATED USE THE DECORATOR DEPRECATED ON A PROPERTY PLEASE NOTE THAT THE PROPERTY DECORATOR SHOULD BE PLACED BEFORE THE DEPRECATED DECORATOR FOR THE DOCSTRINGS TO BE RENDERED PROPERLY EG RENAMING AN ATTRIBUTE LABELS TOCLASSES CAN BE DONE AS

DEPRECATED ATTRIBUTE LABELS WAS DEPRECATED IN VERSION 013 AND

WILL BE REMOVED IN 015 USE CLASSES INSTEAD

PROPERTY

DEFLABELSSELF

RETURNSELFCLASSES

IF A PARAMETER HAS TO BE DEPRECATED USE DEPRECATIONWARNING APPROPRIATELY IN THE FOLLOWING EXAMPLE K IS DEPRECATED AND RENAMED TO NCLUSTERS

IMPORT WARNINGS

DEFEXAMPLEFUNCTIONNCLUSTERS8 KNOTUSED

IFK NOTUSED

WARNINGSWARNK WAS RENAMED TO NCLUSTERS IN VERSION 013 AND

WILL BE REMOVED IN 015 DEPRECATIONWARNING

NCLUSTERS K

WHEN THE CHANGE IS IN A CLASS WE VALIDATE AND RAISE WARNING IN FIT

IMPORT WARNINGS

CLASS EXAMPLEESTIMATOR BASEESTIMATOR

DEFINITSELF NCLUSTERS8 KNOTUSED

SELFNCLUSTERS NCLUSTERS

SELFK K

DEFFITSELF X Y

IFSELFK NOTUSED

WARNINGSWARNK WAS RENAMED TO NCLUSTERS IN VERSION 013 AND

WILL BE REMOVED IN 015 DEPRECATIONWARNING

SELFNCLUSTERS SELFK

ELSE

SELFNCLUSTERS SELFNCLUSTERS

AS IN THESE EXAMPLES THE WARNING MESSAGE SHOULD ALWAYS GIVE BOTH THE VERSION IN WHICH THE DEPRECATION HAPPENED AND THE VERSION IN WHICH THE OLD BEHAVIOR WILL BE REMOVED IF THE DEPRECATION HAPPENED IN VERSION 0XDEV THE MESSAGE SHOULD SAY DEPRECATION OCCURRED IN VERSION 0X AND THE REMOVAL WILL BE IN 0X2 SO THAT USERS WILL HAVE ENOUGH TIME TO ADAPT THEIR CODE TO THE NEW BEHAVIOUR FOR EXAMPLE IF THE DEPRECATION HAPPENED IN VERSION 018DEV THE MESSAGE SHOULD SAY IT HAPPENED IN VERSION 018 AND THE OLD BEHAVIOR WILL BE REMOVED IN VERSION 020

IN ADDITION A DEPRECATION NOTE SHOULD BE ADDED IN THE DOCSTRING RECALLING THE SAME INFORMATION AS THE DEPRECATION WARNING AS EXPLAINED ABOVE USE THE DEPRECATED DIRECTIVE

DEPRECATED 013

K WAS RENAMED TO NCLUSTERS IN VERSION 013 AND WILL BE REMOVED

IN 015

WHAT’S MORE A DEPRECATION REQUIRES A TEST WHICH ENSURES THAT THE WARNING IS RAISED IN RELEVANT CASES BUT NOT IN OTHER CASES THE WARNING SHOULD BE CAUGHT IN ALL OTHER TESTS USING EG PYTESTMARKFILTERWARNINGS AND THERE SHOULD BE NO WARNING IN THE EXAMPLES

71 CONTRIBUTING 2419

SCIKITLEARN USER GUIDE RELEASE 0213

CHANGE THE DEFAULT VALUE OF A PARAMETER

IF THE DEFAULT VALUE OF A PARAMETER NEEDS TO BE CHANGED PLEASE REPLACE THE DEFAULT VALUE WITH A SPECIFIC VALUE EG

WARN AND RAISEFUTUREWARNING WHEN USERS ARE USING THE DEFAULT VALUE IN THE FOLLOWING EXAMPLE WE CHANGE THE

DEFAULT VALUE OF NCLUSTERS FROM 5 TO 10 CURRENT VERSION IS 020

IMPORT WARNINGS

DEFEXAMPLEFUNCTIONNCLUSTERSWARN

IFNCLUSTERS WARN

WARNINGSWARNTHE DEFAULT VALUE OF NCLUSTERS WILL CHANGE FROM

5 TO 10 IN 022 FUTUREWARNING

NCLUSTERS 5

WHEN THE CHANGE IS IN A CLASS WE VALIDATE AND RAISE WARNING IN FIT

IMPORT WARNINGS

CLASS EXAMPLEESTIMATOR

DEFINITSELF NCLUSTERSWARN

SELFNCLUSTERS NCLUSTERS

DEFFITSELF X Y

IFSELFNCLUSTERS WARN

WARNINGSWARNTHE DEFAULT VALUE OF NCLUSTERS WILL CHANGE FROM

5 TO 10 IN 022 FUTUREWARNING

SELFNCLUSTERS 5

SIMILAR TO DEPRECATIONS THE WARNING MESSAGE SHOULD ALWAYS GIVE BOTH THE VERSION IN WHICH THE CHANGE HAPPENED AND

THE VERSION IN WHICH THE OLD BEHAVIOR WILL BE REMOVED THE DOCSTRING NEEDS TO BE UPDATED ACCORDINGLY WE NEED A TEST

WHICH ENSURES THAT THE WARNING IS RAISED IN RELEVANT CASES BUT NOT IN OTHER CASES THE WARNING SHOULD BE CAUGHT IN ALL

OTHER TESTS USING EG PYTESTMARKFILTERWARNINGS AND THERE SHOULD BE NO WARNING IN THE EXAMPLES

PYTHON VERSIONS SUPPORTED

SINCE SCIKITLEARN 021 ONLY PYTHON 35 AND NEWER IS SUPPORTED

717 CODE REVIEW GUIDELINES

REVIEWING CODE CONTRIBUTED TO THE PROJECT AS PRS IS A CRUCIAL COMPONENT OF SCIKITLEARN DEVELOPMENT WE ENCOURAGE

ANYONE TO START REVIEWING CODE OF OTHER DEVELOPERS THE CODE REVIEW PROCESS IS OFTEN HIGHLY EDUCATIONAL FOR EVERYBODY

INVOLVED THIS IS PARTICULARLY APPROPRIATE IF IT IS A FEATURE YOU WOULD LIKE TO USE AND SO CAN RESPOND CRITICALLY ABOUT

WHETHER THE PR MEETS YOUR NEEDS WHILE EACH PULL REQUEST NEEDS TO BE SIGNED OFF BY TWO CORE DEVELOPERS YOU CAN

SPEED UP THIS PROCESS BY PROVIDING YOUR FEEDBACK

HERE ARE A FEW IMPORTANT ASPECTS THAT NEED TO BE COVERED IN ANY CODE REVIEW FROM HIGHLEVEL QUESTIONS TO A MORE

DETAILED CHECKLIST

- DO WE WANT THIS IN THE LIBRARY IS IT LIKELY TO BE USED DO YOU AS A SCIKITLEARN USER LIKE THE CHANGE AND INTEND TO USE IT IS IT IN THE SCOPE OF SCIKITLEARN WILL THE COST OF MAINTAINING A NEW FEATURE BE WORTH ITS BENEFITS
- IS THE CODE CONSISTENT WITH THE API OF SCIKITLEARN ARE PUBLIC FUNCTIONSCLASSESPARAMETERS WELL NAMED AND INTUITIVELY DESIGNED
- ARE ALL PUBLIC FUNCTIONSCLASSES AND THEIR PARAMETERS RETURN TYPES AND STORED ATTRIBUTES NAMED ACCORDING TO SCIKITLEARN CONVENTIONS AND DOCUMENTED CLEARLY

2420 CHAPTER 7 DEVELOPER’S GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

- IS ANY NEW FUNCTIONALITY DESCRIBED IN THE USERGUIDE AND ILLUSTRATED WITH EXAMPLES
- IS EVERY PUBLIC FUNCTIONCLASS TESTED ARE A REASONABLE SET OF PARAMETERS THEIR VALUES VALUE TYPES AND COMBINATIONS TESTED DO THE TESTS VALIDATE THAT THE CODE IS CORRECT IE DOING WHAT THE DOCUMENTATION SAYS IT DOES IF THE CHANGE IS A BUGFIX IS A NONREGRESSION TEST INCLUDED LOOK AT THIS TO GET STARTED WITH TESTING IN PYTHON
- DO THE TESTS PASS IN THE CONTINUOUS INTEGRATION BUILD IF APPROPRIATE HELP THE CONTRIBUTOR UNDERSTAND WHY TESTS FAILED
- DO THE TESTS COVER EVERY LINE OF CODE SEE THE COVERAGE REPORT IN THE BUILD LOG IF NOT ARE THE LINES MISSING COVERAGE GOOD EXCEPTIONS
- IS THE CODE EASY TO READ AND LOW ON REDUNDANCY SHOULD VARIABLE NAMES BE IMPROVED FOR CLARITY OR CONSISTENCY SHOULD COMMENTS BE ADDED SHOULD COMMENTS BE REMOVED AS UNHELPFUL OR EXTRANEIOUS
- COULD THE CODE EASILY BE REWRITTEN TO RUN MUCH MORE EFFICIENTLY FOR RELEVANT SETTINGS
- IS THE CODE BACKWARDS COMPATIBLE WITH PREVIOUS VERSIONS OR IS A DEPRECATION CYCLE NECESSARY
- WILL THE NEW CODE ADD ANY DEPENDENCIES ON OTHER LIBRARIES THIS IS UNLIKELY TO BE ACCEPTED
- DOES THE DOCUMENTATION RENDER PROPERLY SEE THE DOCUMENTATION SECTION FOR MORE DETAILS AND ARE THE PLOTS INSTRUCTIVE

STANDARD REPLIES FOR REVIEWING INCLUDES SOME FREQUENT COMMENTS THAT REVIEWERS MAY MAKE

718 APIS OF SCIKITLEARN OBJECTS

TO HAVE A UNIFORM API WE TRY TO HAVE A COMMON BASIC API FOR ALL THE OBJECTS IN ADDITION TO AVOID THE PROLIFERATION OF FRAMEWORK CODE WE TRY TO ADOPT SIMPLE CONVENTIONS AND LIMIT TO A MINIMUM THE NUMBER OF METHODS AN OBJECT MUST IMPLEMENT

ELEMENTS OF THE SCIKITLEARN API ARE DESCRIBED MORE DEFINITELY IN THE GLOSSARY OF COMMON TERMS AND API ELEMENTS DIFFERENT OBJECTS

THE MAIN OBJECTS IN SCIKITLEARN ARE ONE CLASS CAN IMPLEMENT MULTIPLE INTERFACES

ESTIMATOR THE BASE OBJECT IMPLEMENTS A FIT METHOD TO LEARN FROM DATA EITHER

ESTIMATOR ESTIMATORFITDATA TARGETS

OR

ESTIMATOR ESTIMATORFITDATA

PREDICTOR FOR SUPERVISED LEARNING OR SOME UNSUPERVISED PROBLEMS IMPLEMENTS

PREDICTION PREDICTORPREDICTDATA

CLASSIFICATION ALGORITHMS USUALLY ALSO OFFER A WAY TO QUANTIFY CERTAINTY OF A PREDICTION EITHER USING

DECISIONFUNCTION ORPREDICTPROBA

PROBABILITY PREDICTORPREDICTPROBADATA

TRANSFORMER FOR FILTERING OR MODIFYING THE DATA IN A SUPERVISED OR UNSUPERVISED WAY IMPLEMENTS

NEWDATA TRANSFORMERTRANSFORMDATA

71 CONTRIBUTING 2421

SCIKITLEARN USER GUIDE RELEASE 0213

WHEN FITTING AND TRANSFORMING CAN BE PERFORMED MUCH MORE EFFICIENTLY TOGETHER THAN SEPARATELY IMPLEMENTS

NEWDATA TRANSFORMERFITTRANSFORMDATA

MODEL A MODEL THAT CAN GIVE A GOODNESS OF FIT MEASURE OR A LIKELIHOOD OF UNSEEN DATA IMPLEMENTS HIGHER IS BETTER

SCORE MODELSCOREDATA

ESTIMATORS

THE API HAS ONE PREDOMINANT OBJECT THE ESTIMATOR A ESTIMATOR IS AN OBJECT THAT FITS A MODEL BASED ON SOME TRAINING DATA AND IS CAPABLE OF INFERRING SOME PROPERTIES ON NEW DATA IT CAN BE FOR INSTANCE A CLASSIFIER OR A REGRESSOR ALL ESTIMATORS IMPLEMENT THE FIT METHOD

ESTIMATORFITX Y

ALL BUILTIN ESTIMATORS ALSO HAVE A SETPARAMS METHOD WHICH SETS DATAINDEPENDENT PARAMETERS OVERRIDING PREVIOUS PARAMETER VALUES PASSED TO INIT

ALL ESTIMATORS IN THE MAIN SCIKITLEARN CODEBASE SHOULD INHERIT FROM SKLEARNBASEBASEESTIMATOR

INSTANTIATION

THIS CONCERNS THE CREATION OF AN OBJECT THE OBJECT’S INIT METHOD MIGHT ACCEPT CONSTANTS AS ARGUMENTS THAT DETERMINE THE ESTIMATOR’S BEHAVIOR LIKE THE C CONSTANT IN SVMs IT SHOULD NOT HOWEVER TAKE THE ACTUAL TRAINING DATA AS AN ARGUMENT AS THIS IS LEFT TO THE FIT METHOD

CLF2 SVCC23

CLF3 SVC1 2 2 3 1 1 WRONG

THE ARGUMENTS ACCEPTED BY INIT SHOULD ALL BE KEYWORD ARGUMENTS WITH A DEFAULT VALUE IN OTHER WORDS A USER SHOULD BE ABLE TO INSTANTIATE AN ESTIMATOR WITHOUT PASSING ANY ARGUMENTS TO IT THE ARGUMENTS SHOULD ALL CORRESPOND TO HYPERPARAMETERS DESCRIBING THE MODEL OR THE OPTIMISATION PROBLEM THE ESTIMATOR TRIES TO SOLVE THESE INITIAL ARGUMENTS OR PARAMETERS ARE ALWAYS REMEMBERED BY THE ESTIMATOR ALSO NOTE THAT THEY SHOULD NOT BE DOCUMENTED UNDER THE “ATTRIBUTES” SECTION BUT RATHER UNDER THE “PARAMETERS” SECTION FOR THAT ESTIMATOR

IN ADDITION EVERY KEYWORD ARGUMENT ACCEPTED BY INIT SHOULD CORRESPOND TO AN ATTRIBUTE ON THE INSTANCE SCIKITLEARN RELIES ON THIS TO FIND THE RELEVANT ATTRIBUTES TO SET ON AN ESTIMATOR WHEN DOING MODEL SELECTION

TO SUMMARIZE AN INIT SHOULD LOOK LIKE

```
DEFINITSELF PARAM11 PARAM22
SELFPARAM1 PARAM1
SELFPARAM2 PARAM2
```

THERE SHOULD BE NO LOGIC NOT EVEN INPUT VALIDATION AND THE PARAMETERS SHOULD NOT BE CHANGED THE CORRESPONDING LOGIC SHOULD BE PUT WHERE THE PARAMETERS ARE USED TYPICALLY IN FIT THE FOLLOWING IS WRONG

```
DEFINITSELF PARAM11 PARAM22 PARAM33
WRONG PARAMETERS SHOULD NOT BE MODIFIED
IFPARAM1 1
PARAM2 1
SELFPARAM1 PARAM1
```

2422 CHAPTER 7 DEVELOPER’S GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

WRONG THE OBJECTS ATTRIBUTES SHOULD HAVE EXACTLY THE NAME OF  
THE ARGUMENT IN THE CONSTRUCTOR

SELFPARAM3 PARAM2

THE REASON FOR POSTPONING THE VALIDATION IS THAT THE SAME VALIDATION WOULD HAVE TO BE PERFORMED IN SETPARAMS  
WHICH IS USED IN ALGORITHMS LIKE GRIDSEARCHCV

FITTING

THE NEXT THING YOU WILL PROBABLY WANT TO DO IS TO ESTIMATE SOME PARAMETERS IN THE MODEL THIS IS IMPLEMENTED IN THE  
FIT METHOD

THEFIT METHOD TAKES THE TRAINING DATA AS ARGUMENTS WHICH CAN BE ONE ARRAY IN THE CASE OF UNSUPERVISED LEARNING  
OR TWO ARRAYS IN THE CASE OF SUPERVISED LEARNING

NOTE THAT THE MODEL IS FITTED USING X AND Y BUT THE OBJECT HOLDS NO REFERENCE TO X AND Y THERE ARE HOWEVER SOME  
EXCEPTIONS TO THIS AS IN THE CASE OF PRECOMPUTED KERNELS WHERE THIS DATA MUST BE STORED FOR USE BY THE PREDICT METHOD

PARAMETERS

X ARRAYLIKE SHAPE NSAMPLES NFEATURES

Y ARRAY SHAPE NSAMPLES

KWARGS OPTIONAL DATADEPENDENT PARAMETERS

XSHAPE0 SHOULD BE THE SAME AS YSHAPE0 IF THIS REQUISITE IS NOT MET AN EXCEPTION OF TYPE VALUEERROR  
SHOULD BE RAISED

YMIGHT BE IGNORED IN THE CASE OF UNSUPERVISED LEARNING HOWEVER TO MAKE IT POSSIBLE TO USE THE ESTIMATOR AS PART OF  
A PIPELINE THAT CAN MIX BOTH SUPERVISED AND UNSUPERVISED TRANSFORMERS EVEN UNSUPERVISED ESTIMATORS NEED TO ACCEPT  
AYNONE KEYWORD ARGUMENT IN THE SECOND POSITION THAT IS JUST IGNORED BY THE ESTIMATOR FOR THE SAME REASON

FITPREDICT FITTRANSFORM SCORE ANDPARTIALFIT METHODS NEED TO ACCEPT A YARGUMENT IN THE SECOND  
PLACE IF THEY ARE IMPLEMENTED

THE METHOD SHOULD RETURN THE OBJECT SELF THIS PATTERN IS USEFUL TO BE ABLE TO IMPLEMENT QUICK ONE LINERS IN AN  
IPYTHON SESSION SUCH AS

YPREDICTED SVCC100FITXTRAIN YTRAINPREDICTXTEST

DEPENDING ON THE NATURE OF THE ALGORITHM FIT CAN SOMETIMES ALSO ACCEPT ADDITIONAL KEYWORDS ARGUMENTS HOWEVER  
ANY PARAMETER THAT CAN HAVE A VALUE ASSIGNED PRIOR TO HAVING ACCESS TO THE DATA SHOULD BE AN INIT KEYWORD  
ARGUMENT FIT PARAMETERS SHOULD BE RESTRICTED TO DIRECTLY DATA DEPENDENT VARIABLES FOR INSTANCE A GRAM MATRIX  
OR AN AFFINITY MATRIX WHICH ARE PRECOMPUTED FROM THE DATA MATRIX XARE DATA DEPENDENT A TOLERANCE STOPPING CRITERION  
TOL IS NOT DIRECTLY DATA DEPENDENT ALTHOUGH THE OPTIMAL VALUE ACCORDING TO SOME SCORING FUNCTION PROBABLY IS  
WHENFIT IS CALLED ANY PREVIOUS CALL TO FIT SHOULD BE IGNORED IN GENERAL CALLING ESTIMATORFITX1 AND  
THENESTIMATORFITX2 SHOULD BE THE SAME AS ONLY CALLING ESTIMATORFITX2 HOWEVER THIS MAY NOT  
BE TRUE IN PRACTICE WHEN FIT DEPENDS ON SOME RANDOM PROCESS SEE RANDOMSTATE ANOTHER EXCEPTION TO THIS RULE IS  
WHEN THE HYPERPARAMETER WARMSTART IS SET TOTRUE FOR ESTIMATORS THAT SUPPORT IT WARMSTARTTRUE MEANS  
THAT THE PREVIOUS STATE OF THE TRAINABLE PARAMETERS OF THE ESTIMATOR ARE REUSED INSTEAD OF USING THE DEFAULT INITIALIZATION  
STRATEGY

ESTIMATED ATTRIBUTES

ATTRIBUTES THAT HAVE BEEN ESTIMATED FROM THE DATA MUST ALWAYS HAVE A NAME ENDING WITH TRAILING UNDERSCORE FOR  
EXAMPLE THE COEFFICIENTS OF SOME REGRESSION ESTIMATOR WOULD BE STORED IN A COEF ATTRIBUTE AFTER FIT HAS BEEN CALLED

71 CONTRIBUTING 2423

SCIKITLEARN USER GUIDE RELEASE 0213

THE ESTIMATED ATTRIBUTES ARE EXPECTED TO BE OVERRIDDEN WHEN YOU CALL FIT A SECOND TIME  
OPTIONAL ARGUMENTS

IN ITERATIVE ALGORITHMS THE NUMBER OF ITERATIONS SHOULD BE SPECIFIED BY AN INTEGER CALLED NITER  
PAIRWISE ATTRIBUTES

AN ESTIMATOR THAT ACCEPT XOF SHAPENSAMPLES NSAMPLES AND DEFINES A PAIRWISE PROPERTY EQUAL TO TRUE  
ALLOWS FOR CROSSVALIDATION OF THE DATASET EG WHEN XIS A PRECOMPUTED KERNEL MATRIX SPECIFICALLY THE PAIRWISE  
PROPERTY IS USED BY UTILSMETAESTIMATORSSAFESPLIT TO SLICE ROWS AND COLUMNS

719 ROLLING YOUR OWN ESTIMATOR

IF YOU WANT TO IMPLEMENT A NEW ESTIMATOR THAT IS SCIKITLEARNCOMPATIBLE WHETHER IT IS JUST FOR YOU OR FOR CONTRIBUTING IT  
TO SCIKITLEARN THERE ARE SEVERAL INTERNALS OF SCIKITLEARN THAT YOU SHOULD BE AWARE OF IN ADDITION TO THE SCIKITLEARN API  
OUTLINED ABOVE YOU CAN CHECK WHETHER YOUR ESTIMATOR ADHERES TO THE SCIKITLEARN INTERFACE AND STANDARDS BY RUNNING  
UTILSESTIMATORCHECKSCHECKESTIMATOR ON THE CLASS

FROM SKLEARNUTILSESTIMATORCHECKS IMPORT CHECKESTIMATOR

FROM SKLEARN SVM IMPORT LINEARSVC

CHECKESTIMATORLINEARSVC PASSES

THE MAIN MOTIVATION TO MAKE A CLASS COMPATIBLE TO THE SCIKITLEARN ESTIMATOR INTERFACE MIGHT BE THAT YOU WANT  
TO USE IT TOGETHER WITH MODEL EVALUATION AND SELECTION TOOLS SUCH AS MODELSELECTIONGRIDSEARCHCV AND  
PIPELINEPIPELINE

BEFORE DETAILING THE REQUIRED INTERFACE BELOW WE DESCRIBE TWO WAYS TO ACHIEVE THE CORRECT INTERFACE MORE EASILY  
PROJECT TEMPLATE

WE PROVIDE A PROJECT TEMPLATE WHICH HELPS IN THE CREATION OF PYTHON PACKAGES CONTAINING SCIKITLEARN COMPATIBLE  
ESTIMATORS IT PROVIDES

- AN INITIAL GIT REPOSITORY WITH PYTHON PACKAGE DIRECTORY STRUCTURE
- A TEMPLATE OF A SCIKITLEARN ESTIMATOR
- AN INITIAL TEST SUITE INCLUDING USE OF CHECKESTIMATOR
- DIRECTORY STRUCTURES AND SCRIPTS TO COMPILE DOCUMENTATION AND EXAMPLE GALLERIES
- SCRIPTS TO MANAGE CONTINUOUS INTEGRATION TESTING ON LINUX AND WINDOWS
- INSTRUCTIONS FROM GETTING STARTED TO PUBLISHING ON PYPI

BASEESTIMATOR AND MIXINS

WE TEND TO USE “DUCK TYPING” SO BUILDING AN ESTIMATOR WHICH FOLLOWS THE API SUFFICES FOR COMPATIBILITY WITHOUT  
NEEDING TO INHERIT FROM OR EVEN IMPORT ANY SCIKITLEARN CLASSES

```
SCIKITLEARN USER GUIDE RELEASE 0213
HOWEVER IF A DEPENDENCY ON SCIKITLEARN IS ACCEPTABLE IN YOUR CODE YOU CAN PREVENT A LOT OF BOILERPLATE CODE BY
DERIVING A CLASS FROM BASEESTIMATOR AND OPTIONALLY THE MIXIN CLASSES IN SKLEARNBASE FOR EXAMPLE BELOW
IS A CUSTOM CLASSIFIER WITH MORE EXAMPLES INCLUDED IN THE SCIKITLEARNCONTRIB PROJECT TEMPLATE
IMPORT NUMPY AS NP
FROM SKLEARNBASE IMPORT BASEESTIMATOR CLASSIFIERMIXIN
FROM SKLEARNUTILSVALIDATION IMPORT CHECKXY CHECKARRAY CHECKISFITTED
FROM SKLEARNUTILSMULTICLASS IMPORT UNIQUELABELS
FROM SKLEARNMETRICS IMPORT EUCLIDEANDISTANCES
CLASS TEMPLATECLASSIFIER BASEESTIMATOR CLASSIFIERMIXIN

    DEF INITSELF DEMOPARAMDEMO
    SELFDEMOPARAM DEMOPARAM

    DEF FITSELF X Y

        CHECK THAT X AND Y HAVE CORRECT SHAPE
        X Y CHECKXYX Y
        STORE THE CLASSES SEEN DURING FIT
        SELFCLASSES UNIQUELABELSY

        SELFX X
        SELF Y
        RETURN THE CLASSIFIER
        RETURN SELF

    DEF PREDICTSELF X

        CHECK IS FIT HAD BEEN CALLED
        CHECKISFITTEDSELF X Y

        INPUT VALIDATION
        X CHECKARRAYX

    CLOSEST NPARGMINEUCLIDEANDISTANCESX SELFX AXIS1
    RETURN SELFYCLOSEST
    GETPARAMS AND SETPARAMS
ALL SCIKITLEARN ESTIMATORS HAVE GETPARAMS ANDSETPARAMS FUNCTIONS THE GETPARAMS FUNCTION TAKES NO
ARGUMENTS AND RETURNS A DICT OF THE INIT PARAMETERS OF THE ESTIMATOR TOGETHER WITH THEIR VALUES IT MUST TAKE
ONE KEYWORD ARGUMENT DEEP WHICH RECEIVES A BOOLEAN VALUE THAT DETERMINES WHETHER THE METHOD SHOULD RETURN THE
PARAMETERS OF SUBESTIMATORS FOR MOST ESTIMATORS THIS CAN BE IGNORED THE DEFAULT VALUE FOR DEEP SHOULD BE TRUE
THESETPARAMS ON THE OTHER HAND TAKES AS INPUT A DICT OF THE FORM PARAMETER VALUE AND SETS THE
PARAMETER OF THE ESTIMATOR USING THIS DICT RETURN VALUE MUST BE ESTIMATOR ITSELF
WHILE THEGETPARAMS MECHANISM IS NOT ESSENTIAL SEE CLONING BELOW THE SETPARAMS FUNCTION IS NECESSARY AS
IT IS USED TO SET PARAMETERS DURING GRID SEARCHES
THE EASIEST WAY TO IMPLEMENT THESE FUNCTIONS AND TO GET A SENSIBLE REPR METHOD IS TO INHERIT FROM SKLEARN
BASEBASEESTIMATOR IF YOU DO NOT WANT TO MAKE YOUR CODE DEPENDENT ON SCIKITLEARN THE EASIEST WAY TO
IMPLEMENT THE INTERFACE IS
DEFGETPARAMSSSELF DEEP TRUE
    SUPPOSE THIS ESTIMATOR HAS PARAMETERS ALPHA AND RECURSIVE
    RETURNALPHA SELFALPHA RECURSIVE SELFRECURSIVE
71 CONTRIBUTING 2425
```

SCIKITLEARN USER GUIDE RELEASE 0213

DEFSETPARAMS  
SELF PARAMETERS  
FORPARAMETER VALUE INPARAMETERSITEMS  
SETATTRSELF PARAMETER VALUE  
RETURNSELF

PARAMETERS AND INIT

ASMODELSELECTIONGRIDSEARCHCV USESSETPARAMS TO APPLY PARAMETER SETTING TO ESTIMATORS IT IS ESSENTIAL THAT CALLING SETPARAMS HAS THE SAME EFFECT AS SETTING PARAMETERS USING THE INIT METHOD THE EASIEST AND RECOMMENDED WAY TO ACCOMPLISH THIS IS TO NOT DO ANY PARAMETER VALIDATION IN INIT ALL LOGIC BEHIND ESTIMATOR PARAMETERS LIKE TRANSLATING STRING ARGUMENTS INTO FUNCTIONS SHOULD BE DONE IN FIT ALSO IT IS EXPECTED THAT PARAMETERS WITH TRAILING ARE NOT TO BE SET INSIDE THE INIT METHOD ALL AND ONLY THE PUBLIC ATTRIBUTES SET BY FIT HAVE A TRAILING AS A RESULT THE EXISTENCE OF PARAMETERS WITH TRAILING IS USED TO CHECK IF THE ESTIMATOR HAS BEEN FITTED

CLONING

FOR USE WITH THE MODELSELECTION MODULE AN ESTIMATOR MUST SUPPORT THE BASECLONE FUNCTION TO REPLICATE AN ESTIMATOR THIS CAN BE DONE BY PROVIDING A GETPARAMS METHOD IF GETPARAMS IS PRESENT THEN CLONEESTIMATOR WILL BE AN INSTANCE OF TYPEESTIMATOR ON WHICHSETPARAMS HAS BEEN CALLED WITH CLONES OF THE RESULT OF ESTIMATORGETPARAMS OBJECTS THAT DO NOT PROVIDE THIS METHOD WILL BE DEEPCOPIED USING THE PYTHON STANDARD FUNCTION COPYDEEPCOPY IFSAFEFALSE IS PASSED TO CLONE

PIPELINE COMPATIBILITY

FOR AN ESTIMATOR TO BE USABLE TOGETHER WITH PIPELINEPIPELINE IN ANY BUT THE LAST STEP IT NEEDS TO PROVIDE A FIT ORFITTRANSFORM FUNCTION TO BE ABLE TO EVALUATE THE PIPELINE ON ANY DATA BUT THE TRAINING SET IT ALSO NEEDS TO PROVIDE ATRANSFORM FUNCTION THERE ARE NO SPECIAL REQUIREMENTS FOR THE LAST STEP IN A PIPELINE EXCEPT THAT IT HAS A FIT FUNCTION ALL FIT ANDFITTRANSFORM FUNCTIONS MUST TAKE ARGUMENTS X Y EVEN IF Y IS NOT USED SIMILARLY FORSCORE TO BE USABLE THE LAST STEP OF THE PIPELINE NEEDS TO HAVE A SCORE FUNCTION THAT ACCEPTS AN OPTIONAL Y

ESTIMATOR TYPES

SOME COMMON FUNCTIONALITY DEPENDS ON THE KIND OF ESTIMATOR PASSED FOR EXAMPLE CROSSVALIDATION IN MODELSELECTIONGRIDSEARCHCV ANDMODELSELECTIONCROSSVALSCORE DEFAULTS TO BEING STRATIFIED WHEN USED ON A CLASSIFIER BUT NOT OTHERWISE SIMILARLY SCORERS FOR AVERAGE PRECISION THAT TAKE A CONTINUOUS PREDICTION NEED TO CALL DECISIONFUNCTION FOR CLASSIFIERS BUT PREDICT FOR REGRESSORS THIS DISTINCTION BETWEEN CLASSIFIERS AND REGRESSORS IS IMPLEMENTED USING THE ESTIMATOR TYPE ATTRIBUTE WHICH TAKES A STRING VALUE IT SHOULD BE CLASSIFIER FOR CLASSIFIERS AND REGRESSOR FOR REGRESSORS AND CLUSTERER FOR CLUSTERING METHODS TO WORK AS EXPECTED INHERITING FROM CLASSIFIERMIXIN REGRESSORMIXIN ORCLUSTERMIXIN WILL SET THE ATTRIBUTE AUTOMATICALLY WHEN A METAESTIMATOR NEEDS TO DISTINGUISH AMONG ESTIMATOR TYPES INSTEAD OF CHECKING ESTIMATOR TYPE DIRECTLY HELPERS LIKE BASEISCLASSIFIER SHOULD BE USED

SPECIFIC MODELS

CLASSIFIERS SHOULD ACCEPT YTARGET ARGUMENTS TO FIT THAT ARE SEQUENCES LISTS ARRAYS OF EITHER STRINGS OR INTEGERS THEY SHOULD NOT ASSUME THAT THE CLASS LABELS ARE A CONTIGUOUS RANGE OF INTEGERS INSTEAD THEY SHOULD STORE A LIST OF

2426 CHAPTER 7 DEVELOPER'S GUIDE



SCIKITLEARN USER GUIDE RELEASE 0213

CLASSES IN A CLASSES ATTRIBUTE OR PROPERTY THE ORDER OF CLASS LABELS IN THIS ATTRIBUTE SHOULD MATCH THE ORDER IN WHICH PREDICTPROBA PREDICTLOGPROBA ANDDECISIONFUNCTION RETURN THEIR VALUES THE EASIEST WAY TO ACHIEVE THIS IS TO PUT

SELFCLASSES Y NPUNIQUEY RETURNINVERSE TRUE

INFIT THIS RETURNS A NEW YTHAT CONTAINS CLASS INDEXES RATHER THAN LABELS IN THE RANGE 0 NCLASSES

A CLASSIFIER’S PREDICT METHOD SHOULD RETURN ARRAYS CONTAINING CLASS LABELS FROM CLASSES IN A CLASSIFIER THAT IMPLEMENTS DECISIONFUNCTION THIS CAN BE ACHIEVED WITH

DEFPREDICTSELF X

D SELFDECISIONFUNCTIONX

RETURNSELFCLASSESNPARGMAXD AXIS1

IN LINEAR MODELS COEFFICIENTS ARE STORED IN AN ARRAY CALLED COEF AND THE INDEPENDENT TERM IS STORED IN INTERCEPT

SKLEARNLINEARMODELBASE CONTAINS A FEW BASE CLASSES AND MIXINS THAT IMPLEMENT COMMON LINEAR MODEL PATTERNS

THESKLEARNUTILSMULTICLASS MODULE CONTAINS USEFUL FUNCTIONS FOR WORKING WITH MULTICLASS AND MULTILABEL PROBLEMS

ESTIMATOR TAGS

WARNING THE ESTIMATOR TAGS ARE EXPERIMENTAL AND THE API IS SUBJECT TO CHANGE

SCIKITLEARN INTRODUCED ESTIMATOR TAGS IN VERSION 021 THESE ARE ANNOTATIONS OF ESTIMATORS THAT ALLOW PROGRAMMATIC INSPECTION OF THEIR CAPABILITIES SUCH AS SPARSE MATRIX SUPPORT SUPPORTED OUTPUT TYPES AND SUPPORTED METHODS THE ESTIMATOR TAGS ARE A DICTIONARY RETURNED BY THE METHOD GETTAGS THESE TAGS ARE USED BY THE COMMON TESTS AND THESKLEARNUTILSESTIMATORCHECKSCHECKESTIMATOR FUNCTION TO DECIDE WHAT TESTS TO RUN AND WHAT INPUT DATA IS APPROPRIATE TAGS CAN DEPEND ON ESTIMATOR PARAMETERS OR EVEN SYSTEM ARCHITECTURE AND CAN IN GENERAL ONLY BE DETERMINED AT RUNTIME

THE DEFAULT VALUE OF ALL TAGS EXCEPT FOR XTYPES ISFALSE THESE ARE DEFINED IN THE BASEESTIMATOR CLASS

THE CURRENT SET OF ESTIMATOR TAGS ARE

NONDETERMINISTIC WHETHER THE ESTIMATOR IS NOT DETERMINISTIC GIVEN A FIXED RANDOMSTATE

REQUIRESPOSITIVEDATA UNUSED FOR NOW WHETHER THE ESTIMATOR REQUIRES POSITIVE X

NOVALIDATION WHETHER THE ESTIMATOR SKIPS INPUTVALIDATION THIS IS ONLY MEANT FOR STATELESS AND DUMMY TRANSFORMERS

MULTIOUTPUT UNUSED FOR NOW WHETHER A REGRESSOR SUPPORTS MULTITARGET OUTPUTS OR A CLASSIFIER SUPPORTS MULTICLASS MULTIOUTPUT

MULTILABEL WHETHER THE ESTIMATOR SUPPORTS MULTILABEL OUTPUT

STATELESS WHETHER THE ESTIMATOR NEEDS ACCESS TO DATA FOR FITTING EVEN THOUGH AN ESTIMATOR IS STATELESS IT MIGHT STILL NEED A CALL TO FIT FOR INITIALIZATION

ALLOWNAN WHETHER THE ESTIMATOR SUPPORTS DATA WITH MISSING VALUES ENCODED AS NPNAN

POORSORE WHETHER THE ESTIMATOR FAILS TO PROVIDE A “REASONABLE” TESTSET SCORE WHICH CURRENTLY FOR REGRESSION IS AN R2 OF 05 ON A SUBSET OF THE BOSTON HOUSING DATASET AND FOR CLASSIFICATION AN ACCURACY OF 083 ON MAKEBLOBSNSAMPLES300 RANDOMSTATE0 THESE DATASETS AND VALUES ARE BASED ON CURRENT ESTIMATORS IN SKLEARN AND MIGHT BE REPLACED BY SOMETHING MORE SYSTEMATIC

MULTIOUTPUTONLY WHETHER ESTIMATOR SUPPORTS ONLY MULTIOUTPUT CLASSIFICATION OR REGRESSION

71 CONTRIBUTING 2427

SCIKITLEARN USER GUIDE RELEASE 0213

SKIPTTEST WHETHER TO SKIP COMMON TESTS ENTIRELY DON'T USE THIS UNLESS YOU HAVE A VERY GOOD REASON

XTYPES SUPPORTED INPUT TYPES FOR X AS LIST OF STRINGS TESTS ARE CURRENTLY ONLY RUN IF '2DARRAY' IS CONTAINED IN THE LIST SIGNIFYING THAT THE ESTIMATOR TAKES CONTINUOUS 2D NUMPY ARRAYS AS INPUT THE DEFAULT VALUE IS '2DARRAY' OTHER POSSIBLE TYPES ARE STRING SPARSE CATEGORICAL DICT 1DLABELS AND2DLABELS

THE GOAL IS THAT IN THE FUTURE THE SUPPORTED INPUT TYPE WILL DETERMINE THE DATA USED DURING TESTING IN PARTICULAR FORSTRING SPARSE ANDCATEGORICAL DATA FOR NOW THE TEST FOR SPARSE DATA DO NOT MAKE USE OF THEPARSE TAG

TO OVERRIDE THE TAGS OF A CHILD CLASS ONE MUST DEFINE THE MORETAGS METHOD AND RETURN A DICT WITH THE DESIRED TAGS EG

```
CLASS MYMULTIOUTPUTESTIMATOR BASEESTIMATOR
DEFMORETAGSSELF
RETURNMULTIOUTPUTONLY TRUE
NONDETERMINISTIC TRUE
```

IN ADDITION TO THE TAGS ESTIMATORS ALSO NEED TO DECLARE ANY NONOPTIONAL PARAMETERS TO INIT IN THE REQUIREDPARAMETERS CLASS ATTRIBUTE WHICH IS A LIST OR TUPLE IF REQUIREDPARAMETERS IS ONLY ESTIMATOR ORBASEESTIMATOR THEN THE ESTIMATOR WILL BE INSTANTIATED WITH AN INSTANCE OF LINEARDISCRIMINANTANALYSIS ORRIDGEREGRESSION IF THE ESTIMATOR IS A REGRESSOR IN THE TESTS THE CHOICE OF THESE TWO MODELS IS SOMEWHAT IDIOSYNCRATIC BUT BOTH SHOULD PROVIDE ROBUST CLOSEDFORM SOLUTIONS

7110 READING THE EXISTING CODE BASE

READING AND DIGESTING AN EXISTING CODE BASE IS ALWAYS A DIFFICULT EXERCISE THAT TAKES TIME AND EXPERIENCE TO MASTER EVEN THOUGH WE TRY TO WRITE SIMPLE CODE IN GENERAL UNDERSTANDING THE CODE CAN SEEM OVERWHELMING AT FIRST GIVEN THE SHEER SIZE OF THE PROJECT HERE IS A LIST OF TIPS THAT MAY HELP MAKE THIS TASK EASIER AND FASTER IN NO PARTICULAR ORDER

- GET ACQUAINTED WITH THE APIS OF SCIKITLEARN OBJECTS UNDERSTAND WHAT FITPREDICT TRANSFORM ETC ARE USED FOR
- BEFORE DIVING INTO READING THE CODE OF A FUNCTION CLASS GO THROUGH THE DOCSTRINGS FIRST AND TRY TO GET AN IDEA OF WHAT EACH PARAMETER ATTRIBUTE IS DOING IT MAY ALSO HELP TO STOP A MINUTE AND THINK HOW WOULD I DO THIS MYSELF IF I HAD TO
- THE TRICKIEST THING IS OFTEN TO IDENTIFY WHICH PORTIONS OF THE CODE ARE RELEVANT AND WHICH ARE NOT IN SCIKIT LEARN A LOT OF INPUT CHECKING IS PERFORMED ESPECIALLY AT THE BEGINNING OF THE FITMETHODS SOMETIMES ONLY A VERY SMALL PORTION OF THE CODE IS DOING THE ACTUAL JOB FOR EXAMPLE LOOKING AT THE FIT METHOD OFSKLEARNLINEARMODELLINEARREGRESSION WHAT YOU'RE LOOKING FOR MIGHT JUST BE THE CALL THE SCIPYLINALGLSTSQ BUT IT IS BURIED INTO MULTIPLE LINES OF INPUT CHECKING AND THE HANDLING OF DIFFERENT KINDS OF PARAMETERS
- DUE TO THE USE OF INHERITANCE SOME METHODS MAY BE IMPLEMENTED IN PARENT CLASSES ALL ESTIMATORS INHERIT AT LEAST FROM BASEESTIMATOR AND FROM A MIXIN CLASS EG CLASSIFIERMIXIN THAT ENABLES DEFAULT BEHAVIOUR DEPENDING ON THE NATURE OF THE ESTIMATOR CLASSIFIER REGRESSOR TRANSFORMER ETC
- SOMETIMES READING THE TESTS FOR A GIVEN FUNCTION WILL GIVE YOU AN IDEA OF WHAT ITS INTENDED PURPOSE IS YOU CAN USEGIT GREP SEE BELOW TO FIND ALL THE TESTS WRITTEN FOR A FUNCTION MOST TESTS FOR A SPECIFIC FUNCTIONCLASS ARE PLACED UNDER THE TESTS FOLDER OF THE MODULE
- YOU'LL OFTEN SEE CODE LOOKING LIKE THIS OUT PARALLEL

```
DELAYEDSOMEFUNCTIONPARAM FOR PARAM IN SOMEITERABLE THIS RUNS
SOMEFUNCTION IN PARALLEL USING JOBLIB OUT IS THEN AN ITERABLE CONTAINING THE VALUES RETURNED BY
SOMEFUNCTION FOR EACH CALL
```

- WE USE CYTHON TO WRITE FAST CODE CYTHON CODE IS LOCATED IN PYX ANDPXD FILES CYTHON CODE HAS A MORE CLIKE FLAVOR WE USE POINTERS PERFORM MANUAL MEMORY ALLOCATION ETC HAVING SOME MINIMAL EXPERIENCE IN C

2428 CHAPTER 7 DEVELOPER'S GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

C IS PRETTY MUCH MANDATORY HERE

• MASTER YOUR TOOLS

-WITH SUCH A BIG PROJECT BEING EFFICIENT WITH YOUR FAVORITE EDITOR OR IDE GOES A LONG WAY TOWARDS DIGESTING THE CODE BASE BEING ABLE TO QUICKLY JUMP OR PEEK TO A FUNCTIONCLASSATTRIBUTE DEFINITION HELPS A LOT SO DOES BEING ABLE TO QUICKLY SEE WHERE A GIVEN NAME IS USED IN A FILE

-GIT ALSO HAS SOME BUILTIN KILLER FEATURES IT IS OFTEN USEFUL TO UNDERSTAND HOW A FILE CHANGED OVER TIME USING EGGIT BLAME MANUAL THIS CAN ALSO BE DONE DIRECTLY ON GITHUB GIT GREP EXAMPLES IS ALSO EXTREMELY USEFUL TO SEE EVERY OCCURRENCE OF A PATTERN EG A FUNCTION CALL OR A VARIABLE IN THE CODE BASE

72 DEVELOPERS’ TIPS AND TRICKS

721 PRODUCTIVITY AND SANITYPRESERVING TIPS

IN THIS SECTION WE GATHER SOME USEFUL ADVICE AND TOOLS THAT MAY INCREASE YOUR QUALITYOFLIFE WHEN REVIEWING PULL REQUESTS RUNNING UNIT TESTS AND SO FORTH SOME OF THESE TRICKS CONSIST OF USERSCRIPTS THAT REQUIRE A BROWSER EXTENSION SUCH AS TAMPERMONKEY OR GREASEMONKEY TO SET UP USERSCRIPTS YOU MUST HAVE ONE OF THESE EXTENSIONS INSTALLED ENABLED AND RUNNING WE PROVIDE USERSCRIPTS AS GITHUB GISTS TO INSTALL THEM CLICK ON THE “RAW” BUTTON ON THE GIST PAGE

VIEWING THE RENDERED HTML DOCUMENTATION FOR A PULL REQUEST

WE USE CIRCLECI TO BUILD THE HTML DOCUMENTATION FOR EVERY PULL REQUEST TO ACCESS THAT DOCUMENTATION INSTRUCTIONS ARE PROVIDED IN THE DOCUMENTATION SECTION OF THE CONTRIBUTOR GUIDE TO SAVE YOU A FEW CLICKS WE PROVIDE A USERSCRIPT THAT ADDS A BUTTON TO EVERY PR AFTER INSTALLING THE USERSCRIPT NAVIGATE TO ANY GITHUB PR A NEW BUTTON LABELED “SEE CIRCLECI DOC FOR THIS PR” SHOULD APPEAR IN THE TOPRIGHT AREA

FOLDING AND UNFOLDING OUTDATED DIFFS ON PULL REQUESTS

GITHUB HIDES DISCUSSIONS ON PRS WHEN THE CORRESPONDING LINES OF CODE HAVE BEEN CHANGED IN THE MEAN WHILE THIS USERSCRIPT PROVIDES A SHORTCUT CONTROLALTP AT THE TIME OF WRITING BUT LOOK AT THE CODE TO BE SURE TO UNFOLD ALL SUCH HIDDEN DISCUSSIONS AT ONCE SO YOU CAN CATCH UP

CHECKING OUT PULL REQUESTS AS REMOTETRACKING BRANCHES

IN YOUR LOCAL FORK ADD TO YOUR GITCONFIG UNDER THEREMOTE UPSTREAM HEADING THE LINE

FETCH REFSPULL HEADREFSREMOTESUPSTREAMPR

YOU MAY THEN USE GIT CHECKOUT PRPRNUMBER TO NAVIGATE TO THE CODE OF THE PULLREQUEST WITH THE GIVEN NUMBER READ MORE IN THIS GIST

DISPLAY CODE COVERAGE IN PULL REQUESTS

TO OVERLAY THE CODE COVERAGE REPORTS GENERATED BY THE CODECOV CONTINUOUS INTEGRATION CONSIDER THIS BROWSER EXTEN SION THE COVERAGE OF EACH LINE WILL BE DISPLAYED AS A COLOR BACKGROUND BEHIND THE LINE NUMBER

72 DEVELOPERS’ TIPS AND TRICKS 2429

SCIKITLEARN USER GUIDE RELEASE 0213

USEFUL PYTEST ALIASES AND FLAGS

THE FULL TEST SUITE TAKES FAIRLY LONG TO RUN FOR FASTER ITERATIONS IT IS POSSIBLY TO SELECT A SUBSET OF TESTS USING PYTEST SELECTORS IN PARTICULAR ONE CAN RUN A SINGLE TEST BASED ON ITS NODE ID

PYTEST V SKLEARNLINEARMODELTESTSTESTLOGISTICPYTESTSPARSIFY

OR USE THE K PYTEST PARAMETER TO SELECT TESTS BASED ON THEIR NAME FOR INSTANCE

PYTEST SKLEARNTTESTSTESTCOMMONPY V K LOGISTICREGRESSION

WILL RUN ALL COMMON TESTS FOR THELOGISTICREGRESSION ESTIMATOR

WHEN A UNIT TEST FAILS THE FOLLOWING TRICKS CAN MAKE DEBUGGING EASIER

1 THE COMMAND LINE ARGUMENT PYTEST L INSTRUCTS PYTEST TO PRINT THE LOCAL VARIABLES WHEN A FAILURE OCCURS

2 THE ARGUMENT PYTEST PDB DROPS INTO THE PYTHON DEBUGGER ON FAILURE TO INSTEAD DROP INTO THE RICH IPYTHON DEBUGGERIPDB YOU MAY SET UP A SHELL ALIAS TO

PYTEST PDBCLSIPTYTHONTERMINALDEBUGGERTERMINALPDB CAPTURE NO

OTHERPYTEST OPTIONS THAT MAY BECOME USEFUL INCLUDE

- XWHICH EXITS ON THE FIRST FAILED TEST
- LF TO RERUN THE TESTS THAT FAILED ON THE PREVIOUS RUN
- FF TO RERUN ALL PREVIOUS TESTS RUNNING THE ONES THAT FAILED FIRST
- SSO THAT PYTEST DOES NOT CAPTURE THE OUTPUT OF PRINT STATEMENTS
- TBSHORT ORTBLINE TO CONTROL THE LENGTH OF THE LOGS

SINCE OUR CONTINUOUS INTEGRATION TESTS WILL ERROR IF DEPRECATIONWARNING ORFUTUREWARNING AREN'T PROPERLY CAUGHT IT IS ALSO RECOMMENDED TO RUN PYTEST ALONG WITH THE WERRORDEPRECATIONWARNING AND WERRORFUTUREWARNING FLAGS

STANDARD REPLIES FOR REVIEWING

IT MAY BE HELPFUL TO STORE SOME OF THESE IN GITHUB'S SAVED REPLIES FOR REVIEWING

ISSUE USAGE QUESTIONS

YOU'RE ASKING A USAGE QUESTION THE ISSUE TRACKER IS MAINLY FOR BUGS AND NEW

→FEATURES FOR USAGE QUESTIONS IT IS RECOMMENDED TO TRY STACK OVERFLOWHTTPS

→STACKOVERFLOWCOMQUESTIONSTAGGEDSCIKITLEARN OR THE MAILING LISTHTTPS

→MAILPYTHONORGMAILMANLISTINFOSCIKITLEARN

ISSUE YOU'RE WELCOME TO UPDATE THE DOCS

PLEASE FEEL FREE TO OFFER A PULL REQUEST UPDATING THE DOCUMENTATION IF YOU FEEL

→IT COULD BE IMPROVED

ISSUE SELFCONTAINED EXAMPLE FOR BUG

PLEASE PROVIDE SELFCONTAINED EXAMPLE CODEHTTPSSTACKOVERFLOWCOMHELPMCVE

→INCLUDING IMPORTS AND DATA IF POSSIBLE SO THAT OTHER CONTRIBUTORS CAN JUST

→RUN IT AND REPRODUCE YOUR ISSUE IDEALLY YOUR EXAMPLE CODE SHOULD BE MINIMAL

2430 CHAPTER 7 DEVELOPER'S GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213  
ISSUE SOFTWARE VERSIONS  
TO HELP DIAGNOSE YOUR ISSUE PLEASE PASTE THE OUTPUT OF  
PY  
IMPORT SKLEARN SKLEARNSHOWVERSIONS

THANKS  
ISSUE CODE BLOCKS  
READABILITY CAN BE GREATLY IMPROVED IF YOU FORMATHTTPSHELPGITHUBCOM  
(→ARTICLESCREATINGANDHIGHLIGHTINGCODEBLOCKS YOUR CODE SNIPPETS AND  
(→COMPLETE ERROR MESSAGES APPROPRIATELY FOR EXAMPLE  
PYTHON  
PRINTSOMETHING

GENERATES  
PYTHON  
PRINTSOMETHING

AND  
PYTB  
TRACEBACK MOST RECENT CALL LAST  
FILE STDIN LINE 1 IN MODULE  
IMPORTERROR NO MODULE NAMED HELLO

GENERATES  
PYTB  
TRACEBACK MOST RECENT CALL LAST  
FILE STDIN LINE 1 IN MODULE  
IMPORTERROR NO MODULE NAMED HELLO

YOU CAN EDIT YOUR ISSUE DESCRIPTIONS AND COMMENTS AT ANY TIME TO IMPROVE  
(→READABILITY THIS HELPS MAINTAINERS A LOT THANKS  
ISSUECOMMENT LINKING TO CODE  
FRIENDLY ADVICE FOR CLARITY'S SAKE YOU CAN LINK TO CODE LIKE THISHTTPS  
(→HELPGITHUBCOMARTICLESCREATINGAPERMANENTLINKTOACODESNIPPET  
ISSUECOMMENT LINKING TO COMMENTS  
PLEASE USE LINKS TO COMMENTS WHICH MAKE IT A LOT EASIER TO SEE WHAT YOU ARE  
(→REFERRING TO RATHER THAN JUST LINKING TO THE ISSUE SEE THISHTTPS  
(→STACKOVERFLOWCOMQUESTIONS25163598HOWDOIREFERENCEASPECIFICISSUE  
(→COMMENTONGITHUB FOR MORE DETAILS  
PRNEW BETTER DESCRIPTION  
THANKS FOR THE PULL REQUEST PLEASE MAKE THE TITLE OF THE PR DESCRIPTIVE SO THAT  
(→WE CAN EASILY RECALL THE ISSUE IT IS RESOLVING YOU SHOULD STATE WHAT ISSUE OR  
(→PR IT FIXESRESOLVES IN THE DESCRIPTION SEE HEREHTTPSCIKITLEARNORG  
(→DEVDEVELOPERSCONTRIBUTINGHTMLCONTRIBUTINGPULLREQUESTS  
PRNEW FIX  
72 DEVELOPERS' TIPS AND TRICKS 2431

SCIKITLEARN USER GUIDE RELEASE 0213

PLEASE USE FIX ISSUENUMBER IN YOUR PR DESCRIPTION AND YOU CAN DO IT MORE THAN

→ONCE THIS WAY THE ASSOCIATED ISSUE GETS CLOSED AUTOMATICALLY WHEN THE PR IS

→MERGED FOR MORE DETAILS LOOK AT THISHTTPSGITHUBCOMBLOG1506CLOSING

→ISSUESVIAPULLREQUESTS

PRNEW OR ISSUE MAINTENANCE COST

EVERY FEATURE WE INCLUDE HAS A MAINTENANCE COSTHTTPSCIKITLEARNORGDEVFAQ

→HTMLWHYAREYOUSOSELECTIVEONWHATALGORITHMSYOUINCLUDEINSCIKITLEARN

→OUR MAINTAINERS ARE MOSTLY VOLUNTEERS FOR A NEW FEATURE TO BE INCLUDED WE

→NEED EVIDENCE THAT IT IS OFTEN USEFUL AND IDEALLY WELLESTABLISHEDHTTP

→SCIKITLEARNORGDEVFAQHTMLWHATARETHEINCLUSIONCRITERIAFORNEW

→ALGORITHMS IN THE LITERATURE OR IN PRACTICE THAT DOESNT STOP YOU

→IMPLEMENTING IT FOR YOURSELF AND PUBLISHING IT IN A SEPARATE REPOSITORY OR

→EVEN SCIKITLEARNCONTRIBHTTPSSCIKITLEARNCONTRIBGITHUBIO

PRWIP WHAT’S NEEDED BEFORE MERGE

PLEASE CLARIFY PERHAPS AS A TODO LIST IN THE PR DESCRIPTION WHAT WORK YOU

→BELIEVE STILL NEEDS TO BE DONE BEFORE IT CAN BE REVIEWED FOR MERGE WHEN IT IS

→READY PLEASE PREFIX THE PR TITLE WITH MRG

PRWIP REGRESSION TEST NEEDED

PLEASE ADD A NONREGRESSION TESTHTTPSENWIKIPEDIAORGWIKINONREGRESSION

→TESTING THAT WOULD FAIL AT MASTER BUT PASS IN THIS PR

PRWIP PEP8

YOU HAVE SOME PEP8HTTPSWWWPYTHONORGDEVPEPSPEP0008 VIOLATIONS WHOSE

→DETAILS YOU CAN SEE IN THE CIRCLE CI LINT JOB IT MIGHT BE WORTH CONFIGURING

→YOUR CODE EDITOR TO CHECK FOR SUCH ERRORS ON THE FLY SO YOU CAN CATCH THEM

→BEFORE COMMITTING

PRMRG PATIENCE

BEFORE MERGING WE GENERALLY REQUIRE TWO CORE DEVELOPERS TO AGREE THAT YOUR PULL

→REQUEST IS DESIRABLE AND READY PLEASE BE PATIENTHTTPSCIKITLEARNORGDEV

→FAQHTMLWHYISMYPULLREQUESTNOTGETTINGANYATTENTION AS WE MOSTLY RELY

→ON VOLUNTEERED TIME FROM BUSY CORE DEVELOPERS YOU ARE ALSO WELCOME TO HELP US

→OUT WITH REVIEWING OTHER PRSHTTPSCIKITLEARNORGDEVDEVELOPERS

→CONTRIBUTINGHTMLCODEREVIEWGUIDELINES

PRMRG ADD TO WHAT’S NEW

PLEASE ADD AN ENTRY TO THE CHANGE LOG AT DOCWHATSNEWV RST LIKE THE OTHER

→ENTRIES THERE PLEASE REFERENCE THIS PULL REQUEST WITH PR AND CREDIT

→YOURSELF AND OTHER CONTRIBUTORS IF APPLICABLE WITH USER

PR DON’T CHANGE UNRELATED

PLEASE DO NOT CHANGE UNRELATED LINES IT MAKES YOUR CONTRIBUTION HARDER TO REVIEW

→AND MAY INTRODUCE MERGE CONFLICTS TO OTHER PULL REQUESTS

2432 CHAPTER 7 DEVELOPER’S GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

722 DEBUGGING MEMORY ERRORS IN CYTHON WITH VALGRIND

WHILE PYTHONNUMPY'S BUILTIN MEMORY MANAGEMENT IS RELATIVELY ROBUST IT CAN LEAD TO PERFORMANCE PENALTIES FOR SOME ROUTINES FOR THIS REASON MUCH OF THE HIGHPERFORMANCE CODE IN SCIKITLEARN IS WRITTEN IN CYTHON THIS PERFORMANCE GAIN COMES WITH A TRADEOFF HOWEVER IT IS VERY EASY FOR MEMORY BUGS TO CROP UP IN CYTHON CODE ESPECIALLY IN SITUATIONS WHERE THAT CODE RELIES HEAVILY ON POINTER ARITHMETIC MEMORY ERRORS CAN MANIFEST THEMSELVES A NUMBER OF WAYS THE EASIEST ONES TO DEBUG ARE OFTEN SEGMENTATION FAULTS AND RELATED GLIBC ERRORS UNINITIALIZED VARIABLES CAN LEAD TO UNEXPECTED BEHAVIOR THAT IS DIFFICULT TO TRACK DOWN A VERY USEFUL TOOL WHEN DEBUGGING THESE SORTS OF ERRORS IS VALGRIND

VALGRIND IS A COMMANDLINE TOOL THAT CAN TRACE MEMORY ERRORS IN A VARIETY OF CODE FOLLOW THESE STEPS

1 INSTALL VALGRIND ON YOUR SYSTEM

2 DOWNLOAD THE PYTHON VALGRIND SUPPRESSION FILE VALGRINDPYTHONSUPP

3 FOLLOW THE DIRECTIONS IN THE READMEVALGRIND FILE TO CUSTOMIZE YOUR PYTHON SUPPRESSIONS IF YOU DON'T YOU WILL HAVE SPURIOUS OUTPUT COMING RELATED TO THE PYTHON INTERPRETER INSTEAD OF YOUR OWN CODE

4 RUN VALGRIND AS FOLLOWS

VALGRIND V SUPPRESSIONSVALGRINDPYTHONSUPP PYTHON MYTESTSCRIPTPY

THE RESULT WILL BE A LIST OF ALL THE MEMORYRELATED ERRORS WHICH REFERENCE LINES IN THE CCODE GENERATED BY CYTHON FROM YOUR PYX FILE IF YOU EXAMINE THE REFERENCED LINES IN THE C FILE YOU WILL SEE COMMENTS WHICH INDICATE THE CORRESPONDING LOCATION IN YOUR PYX SOURCE FILE HOPEFULLY THE OUTPUT WILL GIVE YOU CLUES AS TO THE SOURCE OF YOUR MEMORY ERROR

FOR MORE INFORMATION ON VALGRIND AND THE ARRAY OF OPTIONS IT HAS SEE THE TUTORIALS AND DOCUMENTATION ON THE VALGRIND WEB SITE

73 UTILITIES FOR DEVELOPERS

SCIKITLEARN CONTAINS A NUMBER OF UTILITIES TO HELP WITH DEVELOPMENT THESE ARE LOCATED IN SKLEARNUTILS AND INCLUDE TOOLS IN A NUMBER OF CATEGORIES ALL THE FOLLOWING FUNCTIONS AND CLASSES ARE IN THE MODULE SKLEARNUTILS WARNING THESE UTILITIES ARE MEANT TO BE USED INTERNALLY WITHIN THE SCIKITLEARN PACKAGE THEY ARE NOT GUARANTEED TO BE STABLE BETWEEN VERSIONS OF SCIKITLEARN BACKPORTS IN PARTICULAR WILL BE REMOVED AS THE SCIKITLEARN DEPENDENCIES EVOLVE

731 VALIDATION TOOLS

THESE ARE TOOLS USED TO CHECK AND VALIDATE INPUT WHEN YOU WRITE A FUNCTION WHICH ACCEPTS ARRAYS MATRICES OR SPARSE MATRICES AS ARGUMENTS THE FOLLOWING SHOULD BE USED WHEN APPLICABLE

- ASSERTALLFINITE THROW AN ERROR IF ARRAY CONTAINS NANS OR INFS
- ASFLOATARRAY CONVERT INPUT TO AN ARRAY OF FLOATS IF A SPARSE MATRIX IS PASSED A SPARSE MATRIX WILL BE RETURNED
- CHECKARRAY CHECK THAT INPUT IS A 2D ARRAY RAISE ERROR ON SPARSE MATRICES ALLOWED SPARSE MATRIX FORMATS CAN BE GIVEN OPTIONALLY AS WELL AS ALLOWING 1D OR NDIMENSIONAL ARRAYS CALLS ASSERTALLFINITE BY DEFAULT

73 UTILITIES FOR DEVELOPERS 2433

SCIKITLEARN USER GUIDE RELEASE 0213

- CHECKXY CHECK THAT X AND Y HAVE CONSISTENT LENGTH CALLS CHECKARRAY ON X AND COLUMNOR1D ON Y FOR MULTILABEL CLASSIFICATION OR MULTITARGET REGRESSION SPECIFY MULTIOUTPUTTRUE IN WHICH CASE CHECKARRAY WILL BE CALLED ON Y
- INDEXABLE CHECK THAT ALL INPUT ARRAYS HAVE CONSISTENT LENGTH AND CAN BE SLICED OR INDEXED USING SAFEINDEX THIS IS USED TO VALIDATE INPUT FOR CROSSVALIDATION
- VALIDATIONCHECKMEMORY CHECKS THAT INPUT IS JOBLIBMEMORY LIKE WHICH MEANS THAT IT CAN BE CONVERTED INTO A SKLEARNUTILSMEMORY INSTANCE TYPICALLY A STR DENOTING THE CACHEDIR OR HAS THE SAME INTERFACE

IF YOUR CODE RELIES ON A RANDOM NUMBER GENERATOR IT SHOULD NEVER USE FUNCTIONS LIKE NUMPYRANDOMRANDOM OR NUMPYRANDOMNORMAL THIS APPROACH CAN LEAD TO REPEATABILITY ISSUES IN UNIT TESTS INSTEAD A NUMPY RANDOMRANDOMSTATE OBJECT SHOULD BE USED WHICH IS BUILT FROM A RANDOMSTATE ARGUMENT PASSED TO THE CLASS OR FUNCTION THE FUNCTION CHECKRANDOMSTATE BELOW CAN THEN BE USED TO CREATE A RANDOM NUMBER GENERATOR OBJECT

- CHECKRANDOMSTATE CREATE ANRANDOMRANDOMSTATE OBJECT FROM A PARAMETER RANDOMSTATE
- IFRANDOMSTATE ISNONE ORNPRANDOM THEN A RANDOMLYINITIALIZED RANDOMSTATE OBJECT IS RETURNED

-IFRANDOMSTATE IS AN INTEGER THEN IT IS USED TO SEED A NEW RANDOMSTATE OBJECT

-IFRANDOMSTATE IS ARANDOMSTATE OBJECT THEN IT IS PASSED THROUGH

FOR EXAMPLE

```
FROM SKLEARNUTILS IMPORT CHECKRANDOMSTATE
RANDOMSTATE 0
RANDOMSTATE CHECKRANDOMSTATERANDOMSTATE
RANDOMSTATERAND4
```

ARRAY05488135 071518937 060276338 054488318

WHEN DEVELOPING YOUR OWN SCIKITLEARN COMPATIBLE ESTIMATOR THE FOLLOWING HELPERS ARE AVAILABLE

- VALIDATIONCHECKISFITTED CHECK THAT THE ESTIMATOR HAS BEEN FITTED BEFORE CALLING TRANSFORM PREDICT OR SIMILAR METHODS THIS HELPER ALLOWS TO RAISE A STANDARDIZED ERROR MESSAGE ACROSS ESTIMATOR
- VALIDATIONHASFITPARAMETER CHECK THAT A GIVEN PARAMETER IS SUPPORTED IN THE FIT METHOD OF A GIVEN ESTIMATOR

732 EFFICIENT LINEAR ALGEBRA ARRAY OPERATIONS

- EXTMATHRANDOMIZEDRANGEFINDER CONSTRUCT AN ORTHONORMAL MATRIX WHOSE RANGE APPROXIMATES THE RANGE OF THE INPUT THIS IS USED IN EXTMATHRANDOMIZEDSVD BELOW
- EXTMATHRANDOMIZEDSVD COMPUTE THE KTRUNCATED RANDOMIZED SVD THIS ALGORITHM FINDS THE EXACT TRUNCATED SINGULAR VALUES DECOMPOSITION USING RANDOMIZATION TO SPEED UP THE COMPUTATIONS IT IS PARTICULARLY FAST ON LARGE MATRICES ON WHICH YOU WISH TO EXTRACT ONLY A SMALL NUMBER OF COMPONENTS
- ARRAYFUNCSCHOLESKYDELETE USED INSKLEARNLINEARMODELLARSPATH REMOVE AN ITEM FROM A CHOLESKY FACTORIZATION
- ARRAYFUNCSMINPOS USED INSKLEARNLINEARMODELLEASTANGLE FIND THE MINIMUM OF THE POSITIVE VALUES WITHIN AN ARRAY
- EXTMATHFASTLOGDET EFFICIENTLY COMPUTE THE LOG OF THE DETERMINANT OF A MATRIX
- EXTMATHDENSITY EFFICIENTLY COMPUTE THE DENSITY OF A SPARSE VECTOR



SCIKITLEARN USER GUIDE RELEASE 0213

- EXTMATHSAFESParsedOT DOT PRODUCT WHICH WILL CORRECTLY HANDLE SCIPYSPARSE INPUTS IF THE INPUTS ARE DENSE IT IS EQUIVALENT TO NUMPYDOT
  - EXTMATHWEIGHTEDMODE AN EXTENSION OF SCIPYSTATSMODE WHICH ALLOWS EACH ITEM TO HAVE A REAL VALUED WEIGHT
  - RESAMPLE RESAMPLE ARRAYS OR SPARSE MATRICES IN A CONSISTENT WAY USED IN SHUFFLE BELOW
  - SHUFFLE SHUFFLE ARRAYS OR SPARSE MATRICES IN A CONSISTENT WAY USED IN SKLEARNCLUSTERKMEANS
- 733 EFFICIENT RANDOM SAMPLING
- RANDOMSAMPLEWITHOUTREPLACEMENT IMPLEMENTS EFFICIENT ALGORITHMS FOR SAMPLING NSAMPLES INTEGERS FROM A POPULATION OF SIZE NPOPULATION WITHOUT REPLACEMENT
- 734 EFFICIENT ROUTINES FOR SPARSE MATRICES
- THE SKLEARNUTILSSPARSEFUNCS CYTHON MODULE HOSTS COMPILED EXTENSIONS TO EFFICIENTLY PROCESS SCIPY SPARSE DATA
- SPARSEFUNCSMEANVARIANCEAXIS COMPUTE THE MEANS AND VARIANCES ALONG A SPECIFIED AXIS OF A CSR MATRIX USED FOR NORMALIZING THE TOLERANCE STOPPING CRITERION IN SKLEARNCLUSTERKMEANS
  - SPARSEFUNCSFASTINPLACECSRROWNORMALIZEL1 ANDSPARSEFUNCSFASTINPLACECSRROWNORMALIZEL2 CAN BE USED TO NORMALIZE INDIVIDUAL SPARSE SAMPLES TO UNIT L1 OR L2 NORM AS DONE IN SKLEARNPREPROCESSINGNORMALIZER
  - SPARSEFUNCSINPLACECSRCOLUMNSCALE CAN BE USED TO MULTIPLY THE COLUMNS OF A CSR MATRIX BY A CONSTANT SCALE ONE SCALE PER COLUMN USED FOR SCALING FEATURES TO UNIT STANDARD DEVIATION IN SKLEARN PREPROCESSINGSTANDARDSCALER
- 735 GRAPH ROUTINES
- GRAPHSINGLESOURCESHORTESTPATHLENGTH NOT CURRENTLY USED IN SCIKITLEARN RETURN THE SHORTEST PATH FROM A SINGLE SOURCE TO ALL CONNECTED NODES ON A GRAPH CODE IS ADAPTED FROM NETWORKX IF THIS IS EVER NEEDED AGAIN IT WOULD BE FAR FASTER TO USE A SINGLE ITERATION OF DIJKSTRA'S ALGORITHM FROM GRAPHSHORTESTPATH
  - GRAPHSHORTESTPATHGRAPHSHORTESTPATH USED IN SKLEARNMANIFOLDISOMAP RETURN THE SHORTEST PATH BETWEEN ALL PAIRS OF CONNECTED POINTS ON A DIRECTED OR UNDIRECTED GRAPH BOTH THE FLOYD WARSHALL ALGORITHM AND DIJKSTRA'S ALGORITHM ARE AVAILABLE THE ALGORITHM IS MOST EFFICIENT WHEN THE CONNECTIVITY MATRIX IS ASCIPYSPARSECSRMATRIX
- 736 TESTING FUNCTIONS
- TESTINGASSERTIN TESTINGASSERTNOTIN ASSERTIONS FOR CONTAINER MEMBERSHIP DESIGNED FOR FORWARD COMPATIBILITY WITH NOSE 10
  - TESTINGASSERTRAISEMESSAGE ASSERTIONS FOR CHECKING THE ERROR RAISE MESSAGE
  - TESTINGMOCKMLDATAURLOPEN MOCKS THE URLOPEN FUNCTION TO FAKE REQUESTS TO MLDATAORG USED IN TESTS OFSKLEARNDATASETS
  - TESTINGALLESTIMATORS RETURNS A LIST OF ALL ESTIMATORS IN SCIKITLEARN TO TEST FOR CONSISTENT BEHAVIOR AND INTERFACES
- 73 UTILITIES FOR DEVELOPERS 2435

SCIKITLEARN USER GUIDE RELEASE 0213

737 MULTICLASS AND MULTILABEL UTILITY FUNCTION

- MULTICLASSISMULTILABEL HELPER FUNCTION TO CHECK IF THE TASK IS A MULTILABEL CLASSIFICATION ONE
- MULTICLASSUNIQUELABELS HELPER FUNCTION TO EXTRACT AN ORDERED ARRAY OF UNIQUE LABELS FROM DIFFERENT FORMATS OF TARGET

738 HELPER FUNCTIONS

- GENEVENSICES GENERATOR TO CREATE NPACKS OF SLICES GOING UP TO N USED IN SKLEARN
- DECOMPOSITIONDICTLEARNING ANDSKLEARNCLUSTERKMEANS
- SAFEMASK HELPER FUNCTION TO CONVERT A MASK TO THE FORMAT EXPECTED BY THE NUMPY ARRAY OR SCIPY SPARSE MATRIX ON WHICH TO USE IT SPARSE MATRICES SUPPORT INTEGER INDICES ONLY WHILE NUMPY ARRAYS SUPPORT BOTH BOOLEAN MASKS AND INTEGER INDICES
- SAFESQR HELPER FUNCTION FOR UNIFIED SQUARING 2 OF ARRAYLIKES MATRICES AND SPARSE MATRICES

739 HASH FUNCTIONS

- MURMURHASH332 PROVIDES A PYTHON WRAPPER FOR THE MURMURHASH3X8632 C NON CRYPTOGRAPHIC HASH FUNCTION THIS HASH FUNCTION IS SUITABLE FOR IMPLEMENTING LOOKUP TABLES BLOOM FILTERS COUNT MIN SKETCH FEATURE HASHING AND IMPLICITLY DEFINED SPARSE RANDOM PROJECTIONS

```
FROM SKLEARNUTILS IMPORT MURMURHASH332
MURMURHASH332SOME FEATURE SEED0 384616559
TRUE
MURMURHASH332SOME FEATURE SEED0 POSITIVE TRUE 3910350737
TRUE
```

THESKLEARNUTILSMURMURHASH MODULE CAN ALSO BE “CIMPORTED” FROM OTHER CYTHON MODULES SO AS TO BENEFIT FROM THE HIGH PERFORMANCE OF MURMURHASH WHILE SKIPPING THE OVERHEAD OF THE PYTHON INTERPRETER

7310 WARNINGS AND EXCEPTIONS

- DEPRECATED DECORATOR TO MARK A FUNCTION OR CLASS AS DEPRECATED
- SKLEARNEXCEPTIONSCONVERGENCEWARNING CUSTOM WARNING TO CATCH CONVERGENCE PROBLEMS USED

INSKLEARNCOVARIANCEGRAPHICALASSO

74 HOW TO OPTIMIZE FOR SPEED

THE FOLLOWING GIVES SOME PRACTICAL GUIDELINES TO HELP YOU WRITE EFFICIENT CODE FOR THE SCIKITLEARN PROJECT NOTE WHILE IT IS ALWAYS USEFUL TO PROFILE YOUR CODE SO AS TO CHECK PERFORMANCE ASSUMPTIONS IT IS ALSO HIGHLY RECOMMENDED TO REVIEW THE LITERATURE TO ENSURE THAT THE IMPLEMENTED ALGORITHM IS THE STATE OF THE ART FOR THE TASK BEFORE INVESTING INTO COSTLY IMPLEMENTATION OPTIMIZATION TIMES AND TIMES HOURS OF EFFORTS INVESTED IN OPTIMIZING COMPLICATED IMPLEMENTATION DETAILS HAVE BEEN RENDERED IRRELEVANT BY THE SUBSEQUENT DISCOVERY OF SIMPLE ALGORITHMIC TRICKS OR BY USING ANOTHER ALGORITHM ALTOGETHER THAT IS BETTER SUITED TO THE PROBLEM

2436 CHAPTER 7 DEVELOPER’S GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

THE SECTION A SIMPLE ALGORITHMIC TRICK WARM RESTARTS GIVES AN EXAMPLE OF SUCH A TRICK

741 PYTHON CYTHON OR CC

IN GENERAL THE SCIKITLEARN PROJECT EMPHASIZES THE READABILITY OF THE SOURCE CODE TO MAKE IT EASY FOR THE PROJECT USERS TO DIVE INTO THE SOURCE CODE SO AS TO UNDERSTAND HOW THE ALGORITHM BEHAVES ON THEIR DATA BUT ALSO FOR EASE OF MAINTAINABILITY BY THE DEVELOPERS

WHEN IMPLEMENTING A NEW ALGORITHM IS THUS RECOMMENDED TO START IMPLEMENTING IT IN PYTHON USING NUMPY AND SCIPY BY TAKING CARE OF AVOIDING LOOPING CODE USING THE VECTORIZED IDIOMS OF THOSE LIBRARIES IN PRACTICE THIS MEANS TRYING TO REPLACE ANY NESTED FOR LOOPS BY CALLS TO EQUIVALENT NUMPY ARRAY METHODS THE GOAL IS TO AVOID THE CPU WASTING TIME IN THE PYTHON INTERPRETER RATHER THAN CRUNCHING NUMBERS TO FIT YOUR STATISTICAL MODEL IT’S GENERALLY A GOOD IDEA TO CONSIDER NUMPY AND SCIPY PERFORMANCE TIPS [HTTPSCIPYGITHUBIOOLDWIKIPAGESPERFORMANCETIPS](http://scipy.github.io/oldwiki/pagess/performance_tips)

SOMETIMES HOWEVER AN ALGORITHM CANNOT BE EXPRESSED EFFICIENTLY IN SIMPLE VECTORIZED NUMPY CODE IN THIS CASE THE RECOMMENDED STRATEGY IS THE FOLLOWING

1PROFILE THE PYTHON IMPLEMENTATION TO FIND THE MAIN BOTTLENECK AND ISOLATE IT IN A DEDICATED MODULE LEVEL FUNCTION THIS FUNCTION WILL BE REIMPLEMENTED AS A COMPILED EXTENSION MODULE

2 IF THERE EXISTS A WELL MAINTAINED BSD OR MIT CC IMPLEMENTATION OF THE SAME ALGORITHM THAT IS NOT TOO BIG YOU CAN WRITE A CYTHON WRAPPER FOR IT AND INCLUDE A COPY OF THE SOURCE CODE OF THE LIBRARY IN THE SCIKITLEARN SOURCE TREE THIS STRATEGY IS USED FOR THE CLASSES SVMLINEARSVC SVM SVC AND LINEAR MODEL LOGISTIC REGRESSION WRAPPERS FOR LIBLINEAR AND LIBSVM

3 OTHERWISE WRITE AN OPTIMIZED VERSION OF YOUR PYTHON FUNCTION USING CYTHON DIRECTLY THIS STRATEGY IS USED FOR THE LINEAR MODEL ELASTIC NET AND LINEAR MODEL SGD CLASSIFIER CLASSES FOR INSTANCE

4MOVE THE PYTHON VERSION OF THE FUNCTION IN THE TESTS AND USE IT TO CHECK THAT THE RESULTS OF THE COMPILED EXTENSION ARE CONSISTENT WITH THE GOLD STANDARD EASY TO DEBUG PYTHON VERSION

5 ONCE THE CODE IS OPTIMIZED NOT SIMPLE BOTTLENECK SPOTTABLE BY PROFILING CHECK WHETHER IT IS POSSIBLE TO HAVE COARSE GRAINED PARALLELISM THAT IS AMENABLE TO MULTIPROCESSING BY USING THE JOBLIB PARALLEL CLASS

WHEN USING CYTHON USE EITHER

- PYTHON SETUPPY BUILD EXT I
- PYTHON SETUPPY INSTALL

TO GENERATE C FILES YOU ARE RESPONSIBLE FOR ADDING C CPP EXTENSIONS ALONG WITH BUILD PARAMETERS IN EACH SUBMODULE SETUPPY

CC GENERATED FILES ARE EMBEDDED IN DISTRIBUTED STABLE PACKAGES THE GOAL IS TO MAKE IT POSSIBLE TO INSTALL SCIKITLEARN STABLE VERSION ON ANY MACHINE WITH PYTHON NUMPY SCIPY AND CC COMPILER

742 PROFILING PYTHON CODE

IN ORDER TO PROFILE PYTHON CODE WE RECOMMEND TO WRITE A SCRIPT THAT LOADS AND PREPARE YOU DATA AND THEN USE THE IPYTHON INTEGRATED PROFILER FOR INTERACTIVELY EXPLORING THE RELEVANT PART FOR THE CODE

SUPPOSE WE WANT TO PROFILE THE NON NEGATIVE MATRIX FACTORIZATION MODULE OF SCIKITLEARN LET US SETUP A NEW IPYTHON SESSION AND LOAD THE DIGITS DATASET AND AS IN THE RECOGNIZING HANDWRITTEN DIGITS EXAMPLE

```
IN 1 FROM SKLEARN DECOMPOSITION IMPORT NMF
IN 2 FROM SKLEARN DATASETS IMPORT LOADDIGITS
IN 3 X LOADDIGITS DATA
```

74 HOW TO OPTIMIZE FOR SPEED 2437

SCIKITLEARN USER GUIDE RELEASE 0213

BEFORE STARTING THE PROFILING SESSION AND ENGAGING IN TENTATIVE OPTIMIZATION ITERATIONS IT IS IMPORTANT TO MEASURE THE TOTAL EXECUTION TIME OF THE FUNCTION WE WANT TO OPTIMIZE WITHOUT ANY KIND OF PROFILER OVERHEAD AND SAVE IT SOMEWHERE FOR LATER REFERENCE

IN 4 TIMEIT NMFNCOMPONENTS16 TOL1E2FITX  
1 LOOPS BEST OF 3 17 S PER LOOP

TO HAVE A LOOK AT THE OVERALL PERFORMANCE PROFILE USING THE PRUN MAGIC COMMAND

IN 5 PRUN L NMFPY NMFNCOMPONENTS16 TOL1E2FITX  
14496 FUNCTION CALLS IN1682 CPU SECONDS  
ORDERED BY INTERNAL TIME  
LIST REDUCED FROM90 TO 9 DUE TO RESTRICTION NMFPY  
NCALLS TOTTIME PERCALL CUMTIME PERCALL FILENAMESLINENOFUNCTION  
36 0609 0017 1499 0042 NMFPY151NLSSUBPROBLEM  
1263 0157 0000 0157 0000 NMFPY18POS  
1 0053 0053 1681 1681 NMFPY352FITTRANSFORM  
673 0008 0000 0057 0000 NMFPY28NORM  
1 0006 0006 0047 0047 NMFPY42INITIALIZENMF  
36 0001 0000 0010 0000 NMFPY36SPARSENESS  
30 0001 0000 0001 0000 NMFPY23NEG  
1 0000 0000 0000 0000 NMFPY337INIT  
1 0000 0000 1681 1681 NMFPY461FIT

THETOTTIME COLUMN IS THE MOST INTERESTING IT GIVES TO TOTAL TIME SPENT EXECUTING THE CODE OF A GIVEN FUNCTION  
IGNORING THE TIME SPENT IN EXECUTING THE SUBFUNCTIONS THE REAL TOTAL TIME LOCAL CODE SUBFUNCTION CALLS IS GIVEN BY THECUMTIME COLUMN

NOTE THE USE OF THE L NMFPY THAT RESTRICTS THE OUTPUT TO LINES THAT CONTAINS THE “NMFPY” STRING THIS IS USEFUL TO HAVE A QUICK LOOK AT THE HOTSPOT OF THE NMF PYTHON MODULE ITSELF IGNORING ANYTHING ELSE

HERE IS THE BEGINNING OF THE OUTPUT OF THE SAME COMMAND WITHOUT THE L NMFPY FILTER

IN 5 PRUN NMFNCOMPONENTS16 TOL1E2FITX  
16159 FUNCTION CALLS IN1840 CPU SECONDS  
ORDERED BY INTERNAL TIME  
NCALLS TOTTIME PERCALL CUMTIME PERCALL FILENAMESLINENOFUNCTION  
2833 0653 0000 0653 0000 NUMPYCOREDOTBLASDOT  
46 0651 0014 1636 0036 NMFPY151NLSSUBPROBLEM  
1397 0171 0000 0171 0000 NMFPY18POS  
2780 0167 0000 0167 0000 METHOD SUM OF NUMPYNDARRAY  
↪OBJECTS  
1 0064 0064 1840 1840 NMFPY352FITTRANSFORM  
1542 0043 0000 0043 0000 METHOD FLATTEN OF NUMPYNDARRAY  
↪OBJECTS  
337 0019 0000 0019 0000 METHOD ALL OF NUMPYNDARRAY  
↪OBJECTS  
2734 0011 0000 0181 0000 FROMNUMERICPY1185SUM  
2 0010 0005 0010 0005 NUMPYLINALGLAPACKLITEDGESDD  
748 0009 0000 0065 0000 NMFPY28NORM

THE ABOVE RESULTS SHOW THAT THE EXECUTION IS LARGELY DOMINATED BY DOT PRODUCTS OPERATIONS DELEGATED TO BLAS HENCE THERE IS PROBABLY NO HUGE GAIN TO EXPECT BY REWRITING THIS CODE IN CYTHON OR CC IN THIS CASE OUT OF THE 17S TOTAL EXECUTION TIME ALMOST 07S ARE SPENT IN COMPILED CODE WE CAN CONSIDER OPTIMAL BY REWRITING THE REST OF THE PYTHON

2438 CHAPTER 7 DEVELOPER’S GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

CODE AND ASSUMING WE COULD ACHIEVE A 1000 BOOST ON THIS PORTION WHICH IS HIGHLY UNLIKELY GIVEN THE SHALLOWNESS OF THE PYTHON LOOPS WE WOULD NOT GAIN MORE THAN A 24X SPEEDUP GLOBALLY

HENCE MAJOR IMPROVEMENTS CAN ONLY BE ACHIEVED BY ALGORITHMIC IMPROVEMENTS IN THIS PARTICULAR EXAMPLE EG TRYING TO FIND OPERATION THAT ARE BOTH COSTLY AND USELESS TO AVOID COMPUTING THEN RATHER THAN TRYING TO OPTIMIZE THEIR IMPLEMENTATION

IT IS HOWEVER STILL INTERESTING TO CHECK WHAT’S HAPPENING INSIDE THE NLSSUBPROBLEM FUNCTION WHICH IS THE HOTSPOT IF WE ONLY CONSIDER PYTHON CODE IT TAKES AROUND 100 OF THE ACCUMULATED TIME OF THE MODULE IN ORDER TO BETTER UNDERSTAND THE PROFILE OF THIS SPECIFIC FUNCTION LET US INSTALL LINEPROFILER AND WIRE IT TO IPYTHON

PIP INSTALL LINEPROFILER

- UNDER IPYTHON 013 FIRST CREATE A CONFIGURATION PROFILE

IPYTHON PROFILE CREATE

THEN REGISTER THE LINEPROFILER EXTENSION IN IPYTHONPROFILEDEFAULTIPYTHONCONFIGPY

CTERMINALIPYTHONAPPEXTENSIONSAPPENDLINEPROFILER

CINTERACTIVESHELLAPPEXTENSIONSAPPENDLINEPROFILER

THIS WILL REGISTER THE LPRUN MAGIC COMMAND IN THE IPYTHON TERMINAL APPLICATION AND THE OTHER FRONTENDS SUCH AS QTCONSOLE AND NOTEBOOK

NOW RESTART IPYTHON AND LET US USE THIS NEW TOY

IN 1 FROM SKLEARNDATASETS IMPORT LOADDIGITS

IN 2 FROM SKLEARNDECOMPOSITIONNMF IMPORT NLSSUBPROBLEM NMF

IN 3 X LOADDIGITSDATA

IN 4 LPRUN F NLSSUBPROBLEM NMFNCOMPONENTS16 TOL1E2FITX

TIMER UNIT 1E06 S

FILE SKLEARNDECOMPOSITIONNMF.PY

FUNCTION NLSSUBPROBLEM AT LINE 137

TOTAL TIME 173153 S

LINE HITS TIME PER HIT TIME LINE CONTENTS

137 DEFNLSSUBPROBLEM(V, W, H, INIT)

↳TOL, MAXITER

138 NONNEGATIVE LEAST SQUARE

↳SOLVER

170

171 48 5863 1221 03 IFH, INIT, 0, ANY

172 RAISE ValueError('ERROR: NEGATIVE

↳VALUES IN H, INIT PASSED TO NLS SOLVER

173

174 48 139 29 00 H, INIT

175 48 112141 23363 58 WTV, NPDOTWT, V

176 48 16144 3363 08 WTW, NPDOTWT, W

177

178 VALUES JUSTIFIED IN THE PAPER

179 48 144 30 00 ALPHA, 1

180 48 113 24 00 BETA, 0.1

74 HOW TO OPTIMIZE FOR SPEED 2439

```
SCIKITLEARN USER GUIDE RELEASE 0213
181 638 1880 29 01 FORNITERINRANGE1 MAXITER
↩→ 1
182 638 195133 3059 102 GRAD NPDOTWTW H WTV
183 638 495761 7771 259 PROJGRADIENT NORMGRADNP
↩→ LOGICALORGRAD 0 H 0
184 638 2449 38 01 IFPROJGRADIENT TOL
185 48 130 27 00 BREAK
186
187 1474 4474 30 02 FORINNERITER INRANGE1
↩→ 20
188 1474 83833 569 44 HN H ALPHA GRAD
189 HN NPWHEREHN 0
↩→ HN 0
190 1474 194239 1318 101 HN POSHN
191 1474 48858 331 25 D HN H
192 1474 150407 1020 78 GRADD NPSUMGRAD D
193 1474 515390 3497 269 DQD NPSUMNPDOTWTW
↩→ DD
```

BY LOOKING AT THE TOP VALUES OF THE TIME COLUMN IT IS REALLY EASY TO PINPOINT THE MOST EXPENSIVE EXPRESSIONS THAT WOULD DESERVE ADDITIONAL CARE

743 MEMORY USAGE PROFILING

YOU CAN ANALYZE IN DETAIL THE MEMORY USAGE OF ANY PYTHON CODE WITH THE HELP OF MEMORYPROFILER FIRST INSTALL THE LATEST VERSION

PIP INSTALL U MEMORYPROFILER

THEN SETUP THE MAGICS IN A MANNER SIMILAR TO LINEPROFILER

- UNDER IPYTHON 011 FIRST CREATE A CONFIGURATION PROFILE

IPYTHON PROFILE CREATE

THEN REGISTER THE EXTENSION IN IPYTHONPROFILEDEFAULTIPYTHONCONFIGPY ALONGSIDE THE LINE PROFILER

CTERMINALIPYTHONAPPEXTENSIONSAPPENDMEMORYPROFILER

CINTERACTIVESHELLAPPEXTENSIONSAPPENDMEMORYPROFILER

THIS WILL REGISTER THE MEMIT ANDMPRUN MAGIC COMMANDS IN THE IPYTHON TERMINAL APPLICATION AND THE OTHER FRONTENDS SUCH AS QTCONSOLE AND NOTEBOOK

MPRUN IS USEFUL TO EXAMINE LINEBYLINE THE MEMORY USAGE OF KEY FUNCTIONS IN YOUR PROGRAM IT IS VERY SIMILAR TO LPRUN DISCUSSED IN THE PREVIOUS SECTION FOR EXAMPLE FROM THE MEMORYPROFILEREXAMPLES DIRECTORY

IN 1FROM EXAMPLE IMPORT MYFUNC

IN 2 MPRUN F MYFUNC MYFUNC

FILENAME EXAMPLEPY

LINE MEM USAGE INCREMENT LINE CONTENTS

3 PROFILE

4 597 MB 000 MB DEFMYFUNC

2440 CHAPTER 7 DEVELOPER’S GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

5 1361 MB 764 MB A 1 106  
6 16620 MB 15259 MB B 2 2107  
7 1361 MB 15259 MB DELB  
8 1361 MB 000 MB RETURNA

ANOTHER USEFUL MAGIC THAT MEMORYPROFILER DEFINES ISMEMIT WHICH IS ANALOGOUS TO TIMEIT IT CAN BE USED AS FOLLOWS

IN 1 IMPORT NUMPY AS NP  
IN 2 MEMIT NPZEROS1E7  
MAXIMUM OF 3 76402344 MB PER LOOP

FOR MORE DETAILS SEE THE DOCSTRINGS OF THE MAGICS USING MEMIT ANDMPRUN

744 PERFORMANCE TIPS FOR THE CYTHON DEVELOPER

IF PROFILING OF THE PYTHON CODE REVEALS THAT THE PYTHON INTERPRETER OVERHEAD IS LARGER BY ONE ORDER OF MAGNITUDE OR MORE THAN THE COST OF THE ACTUAL NUMERICAL COMPUTATION EG FOR LOOPS OVER VECTOR COMPONENTS NESTED EVALUATION OF CONDITIONAL EXPRESSION SCALAR ARITHMETIC IT IS PROBABLY ADEQUATE TO EXTRACT THE HOTSPOT PORTION OF THE CODE AS A STANDALONE FUNCTION IN A PYX FILE ADD STATIC TYPE DECLARATIONS AND THEN USE CYTHON TO GENERATE A C PROGRAM SUITABLE TO BE COMPILED AS A PYTHON EXTENSION MODULE

THE OFFICIAL DOCUMENTATION AVAILABLE AT HTTPDOCSCYTHONORG CONTAINS A TUTORIAL AND REFERENCE GUIDE FOR DEVELOPING SUCH A MODULE IN THE FOLLOWING WE WILL JUST HIGHLIGHT A COUPLE OF TRICKS THAT WE FOUND IMPORTANT IN PRACTICE ON THE EXISTING CYTHON CODEBASE IN THE SCIKITLEARN PROJECT

TODO HTML REPORT TYPE DECLARATIONS BOUND CHECKS DIVISION BY ZERO CHECKS MEMORY ALIGNMENT DIRECT BLAS CALLS

- [HTTPSWWWYOUTUBECOMWATCHVGMVKIQGOW8](https://www.youtube.com/watch?v=GMVKIQGOW8)
- [HTTPCONFERENCESCIPYORGPROCEEDINGSSCIPY2009PAPER1](http://conferences.cipy.org/proceedings/scipy2009/paper1)
- [HTTPCONFERENCESCIPYORGPROCEEDINGSSCIPY2009PAPER2](http://conferences.cipy.org/proceedings/scipy2009/paper2)

USING OPENMP

SINCE SCIKITLEARN CAN BE BUILT WITHOUT OPENMP SUPPORT IT'S NECESSARY TO PROTECT EACH DIRECT CALL TO OPENMP THIS CAN BE DONE USING THE FOLLOWING SYNTAX

```
IMPORTING OPENMP
IF SKLEARNOPENMPSUPPORTED
    CALLING OPENMP
IF SKLEARNOPENMPSUPPORTED
    MAXTHREADS OPENMPOMPGETMAXTHREADS
ELSE
    MAXTHREADS 1
```

NOTE PROTECTING THE PARALLEL LOOP PRANGE IS ALREADY DONE BY CYTHON

74 HOW TO OPTIMIZE FOR SPEED 2441

SCIKITLEARN USER GUIDE RELEASE 0213

745 PROFILING COMPILED EXTENSIONS

WHEN WORKING WITH COMPILED EXTENSIONS WRITTEN IN CC WITH A WRAPPER OR DIRECTLY AS CYTHON EXTENSION THE DEFAULT PYTHON PROFILER IS USELESS WE NEED A DEDICATED TOOL TO INTROSPECT WHAT’S HAPPENING INSIDE THE COMPILED EXTENSION IT SELF

USING YEP AND GPERFTOOLS

EASY PROFILING WITHOUT SPECIAL COMPILATION OPTIONS USE YEP

- [HTTPSPYIORGPROJECTYEP](http://pyio.org/project/yep)
- [HTTPFABIANPNETBLOG2011APROFILERFORPYTHONEXTENSIONS](http://fabianp.net/blog/2011/aprofilerforpythonextensions)

USING GPROF

IN ORDER TO PROFILE COMPILED PYTHON EXTENSIONS ONE COULD USE GPROF AFTER HAVING RECOMPILED THE PROJECT WITH GCC PG AND USING THE PYTHONDBG VARIANT OF THE INTERPRETER ON DEBIAN UBUNTU HOWEVER THIS APPROACH REQUIRES TO ALSO HAVENUMPY ANDSCIPY RECOMPILED WITH PG WHICH IS RATHER COMPLICATED TO GET WORKING

FORTUNATELY THERE EXIST TWO ALTERNATIVE PROFILERS THAT DON’T REQUIRE YOU TO RECOMPILE EVERYTHING

USING VALGRIND CALLGRIND KCACHEGRIND

KCACHEGRIND

YEP CAN BE USED TO CREATE A PROFILING REPORT KCACHEGRIND PROVIDES A GRAPHICAL ENVIRONMENT TO VISUALIZE THIS REPORT

RUN YEP TO PROFILE SOME PYTHON SCRIPT

PYTHON M YEP C MYFILEPY

OPEN MYFILEPYCALLGRIN WITH KCACHEGRIND

KCACHEGRIND MYFILEPYPROF

NOTEYEP CAN BE EXECUTED WITH THE ARGUMENT LINES ORLTO COMPILE A PROFILING REPORT ‘LINE BY LINE’

746 MULTICORE PARALLELISM USING JOBLIBPARALLEL

TODO GIVE A SIMPLE TEASER EXAMPLE HERE

CHECKOUT THE OFFICIAL JOBLIB DOCUMENTATION

- [HTTPSJOBLIBREADTHEDOCSIO](https://joblib.readthedocs.io)

747 A SIMPLE ALGORITHMIC TRICK WARM RESTARTS

SEE THE GLOSSARY ENTRY FOR WARMSTART

2442 CHAPTER 7 DEVELOPER’S GUIDE



SCIKITLEARN USER GUIDE RELEASE 0213

75 ADVANCED INSTALLATION INSTRUCTIONS

THERE ARE DIFFERENT WAYS TO GET SCIKITLEARN INSTALLED

- INSTALL AN OFFICIAL RELEASE THIS IS THE BEST APPROACH FOR MOST USERS IT WILL PROVIDE A STABLE VERSION AND PREBUILD PACKAGES ARE AVAILABLE FOR MOST PLATFORMS
  - INSTALL THE VERSION OF SCIKITLEARN PROVIDED BY YOUR OPERATING SYSTEM OR PYTHON DISTRIBUTION THIS IS A QUICK OPTION FOR THOSE WHO HAVE OPERATING SYSTEMS THAT DISTRIBUTE SCIKITLEARN IT MIGHT NOT PROVIDE THE LATEST RELEASE VERSION
  - BUILDING THE PACKAGE FROM SOURCE THIS IS BEST FOR USERS WHO WANT THE LATESTANDGREATEST FEATURES AND AREN'T AFRAID OF RUNNING BRANDNEW CODE THIS DOCUMENT DESCRIBES HOW TO BUILD FROM SOURCE
- NOTE IF YOU WISH TO CONTRIBUTE TO THE PROJECT YOU NEED TO INSTALL THE LATEST DEVELOPMENT VERSION

751 INSTALLING NIGHTLY BUILDS

THE CONTINUOUS INTEGRATION SERVERS OF THE SCIKITLEARN PROJECT BUILD TEST AND UPLOAD WHEEL PACKAGES FOR THE MOST RECENT PYTHON VERSION ON A NIGHTLY BASIS TO HELP USERS TEST BLEEDING EDGE FEATURES OR BUG FIXES

PIP INSTALL PRE F HTTPSSKLEARNNIGHTLYSCDN8SECURERAXCDNCOM SCIKITLEARN

752 BUILDING FROM SOURCE

IN THE VAST MAJORITY OF CASES BUILDING SCIKITLEARN FOR DEVELOPMENT PURPOSES CAN BE DONE WITH

PIP INSTALL CYTHON PYTEST FLAKE8

THEN IN THE MAIN REPOSITORY

PIP INSTALL EDITABLE

PLEASE READ BELOW FOR DETAILS AND MORE ADVANCED INSTRUCTIONS

DEPENDENCIES

SCIKITLEARN REQUIRES

- PYTHON 35
- NUMPY 111
- SCIPY 017
- JOBLIB 011

NOTE FOR INSTALLING ON PYPY PYPY3V510 NUMPY 1140 AND SCIPY 110 ARE REQUIRED FOR PYPY ONLY INSTALLA

TION INSTRUCTIONS WITH PIP APPLY

BUILDING SCIKITLEARN ALSO REQUIRES

- CYTHON 0285

75 ADVANCED INSTALLATION INSTRUCTIONS 2443

SCIKITLEARN USER GUIDE RELEASE 0213

• OPENMP

NOTE IT IS POSSIBLE TO BUILD SCIKITLEARN WITHOUT OPENMP SUPPORT BY SETTING THE SKLEARNNOOPENMP ENVIRONMENT VARIABLE BEFORE CYTHONIZATION THIS IS NOT RECOMMENDED SINCE IT WILL FORCE SOME ESTIMATORS TO RUN IN SEQUENTIAL MODE AND THEIR NJOBS PARAMETER WILL BE IGNORED

RUNNING TESTS REQUIRES

• PYTEST 330

SOME TESTS ALSO REQUIRE PANDAS

RETRIEVING THE LATEST CODE

WE USE GIT FOR VERSION CONTROL AND GITHUB FOR HOSTING OUR MAIN REPOSITORY

YOU CAN CHECK OUT THE LATEST SOURCES WITH THE COMMAND

`GIT CLONE GITGITHUBCOMSCIKITLEARNSCIKITLEARNGIT`

IF YOU WANT TO BUILD A STABLE VERSION YOU CAN `GIT CHECKOUT VERSION` TO GET THE CODE FOR THAT PARTICULAR

VERSION OR DOWNLOAD AN ZIP ARCHIVE OF THE VERSION FROM GITHUB

ONCE YOU HAVE ALL THE BUILD REQUIREMENTS INSTALLED SEE BELOW FOR DETAILS YOU CAN BUILD AND INSTALL THE PACKAGE IN THE FOLLOWING WAY

IF YOU RUN THE DEVELOPMENT VERSION IT IS CUMBERSOME TO REINSTALL THE PACKAGE EACH TIME YOU UPDATE THE SOURCES

THEREFORE IT'S RECOMMENDED THAT YOU INSTALL IN EDITABLE MODE WHICH ALLOWS YOU TO EDIT THE CODE INPLACE THIS BUILDS

THE EXTENSION IN PLACE AND CREATES A LINK TO THE DEVELOPMENT DIRECTORY SEE THE PIP DOCS

`PIP INSTALL EDITABLE`

NOTE THIS IS FUNDAMENTALLY SIMILAR TO USING THE COMMAND `PYTHON SETUPPY DEVELOP` SEE THE `SETUPTOOL`

DOCS IT IS HOWEVER PREFERRED TO USE PIP

NOTE YOU WILL HAVE TO RERUN

`PIP INSTALL EDITABLE`

EVERY TIME THE SOURCE CODE OF A COMPILED EXTENSION IS CHANGED FOR INSTANCE WHEN SWITCHING BRANCHES OR PULLING

CHANGES FROM UPSTREAM COMPILED EXTENSIONS ARE CYTHON FILES ENDING IN `PYX` OR `RPXD`

ON UNIXLIKE SYSTEMS YOU CAN EQUIVALENTLY TYPE `MAKE` IN FROM THE `TOplevel` FOLDER HAVE A LOOK AT THE `MAKEFILE`

FOR ADDITIONAL UTILITIES

MAC OSX

THE DEFAULT C COMPILER `APPLECLANG` ON MAC OSX DOES NOT DIRECTLY SUPPORT OPENMP THE FIRST SOLUTION TO BUILD

SCIKITLEARN IS TO INSTALL ANOTHER C COMPILER SUCH AS `GCC` OR `LLVMCLANG` ANOTHER SOLUTION IS TO ENABLE OPENMP SUPPORT

ON THE DEFAULT `APPLECLANG` IN THE FOLLOWING WE PRESENT HOW TO CONFIGURE THIS SECOND OPTION

YOU FIRST NEED TO INSTALL THE OPENMP LIBRARY

2444 CHAPTER 7 DEVELOPER'S GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

BREW INSTALL LIBOMP

THEN YOU NEED TO SET THE FOLLOWING ENVIRONMENT VARIABLES

EXPORT CCUSRBINCLANG

EXPORT CXXUSRBINCLANG

EXPORT CPPFLAGSCPPFLAGS XPREPROCESSOR FOPENMP

EXPORT CFLAGSCFLAGS IUSRLOCALOPTLIBOMPINCLUDE

EXPORT CXXFLAGSCXXFLAGS IUSRLOCALOPTLIBOMPINCLUDE

EXPORT LDFLAGSLDFLAGS LUSRLOCALOPTLIBOMPLIB LOMP

EXPORT DYLDLIBRARYPATHUSRLOCALOPTLIBOMPLIB

FINALLY YOU CAN BUILD THE PACKAGE USING THE STANDARD COMMAND

FREEBSD

THE CLANG COMPILER INCLUDED IN FREEBSD 120 AND 112 BASE SYSTEMS DOES NOT INCLUDE OPENMP SUPPORT YOU NEED TO INSTALL THEOPENMP LIBRARY FROM PACKAGES OR PORTS

SUDO PKG INSTALL OPENMP

THIS WILL INSTALL HEADER FILES IN USRLOCALINCLUDE AND LIBS INUSRLOCALLIB SINCE THESE DIRECTORIES ARE NOT SEARCHED BY DEFAULT YOU CAN SET THE ENVIRONMENT VARIABLES TO THESE LOCATIONS

EXPORT CFLAGSCFLAGS IUSRLOCALINCLUDE

EXPORT CXXFLAGSCXXFLAGS IUSRLOCALINCLUDE

EXPORT LDFLAGSLDFLAGS LUSRLOCALLIB LOMP

EXPORT DYLDLIBRARYPATHUSRLOCALLIB

FINALLY YOU CAN BUILD THE PACKAGE USING THE STANDARD COMMAND

FOR THE UPCOMMING FREEBSD 121 AND 113 VERSIONS OPENMP WILL BE INCLUDED IN THE BASE SYSTEM AND THESE STEPS WILL NOT BE NECESSARY

753 INSTALLING BUILD DEPENDENCIES

LINUX

INSTALLING FROM SOURCE WITHOUT CONDA REQUIRES YOU TO HAVE INSTALLED THE SCIKITLEARN RUNTIME DEPENDENCIES PYTHON DEVELOPMENT HEADERS AND A WORKING CC COMPILER UNDER DEBIANBASED OPERATING SYSTEMS WHICH INCLUDE UBUNTU

SUDO APTGET INSTALL BUILDESSENTIAL PYTHON3DEV PYTHON3SETUPTOOLS

PYTHON3PIP

AND THEN

PIP3 INSTALL NUMPY SCIPY CYTHON

NOTE IN ORDER TO BUILD THE DOCUMENTATION AND RUN THE EXAMPLE CODE CONTAINS IN THIS DOCUMENTATION YOU WILL NEED MATPLOTLIB

PIP3 INSTALL MATPLOTLIB

75 ADVANCED INSTALLATION INSTRUCTIONS 2445

SCIKITLEARN USER GUIDE RELEASE 0213

WHEN PRECOMPILED WHEELS ARE NOT AVAILAIBLE FOR YOUR ARCHITECTURE YOU CAN INSTALL THE SYSTEM VERSIONS

SUDO APTGET INSTALL CYTHON3 PYTHON3NUMPY PYTHON3SCIPY PYTHON3MATPLOTLIB

ON RED HAT AND CLONES EG CENTOS INSTALL THE DEPENDENCIES USING

SUDO YUM Y INSTALL GCC GCCC PYTHONDEVEL NUMPY SCIPY

NOTE TO USE A HIGH PERFORMANCE BLAS LIBRARY EG OPENBLAS SEE SCIPY INSTALLATION INSTRUCTIONS

WINDOWS

TO BUILD SCIKITLEARN ON WINDOWS YOU NEED A WORKING CC COMPILER IN ADDITION TO NUMPY SCIPY AND SETUPTOOLS

THE BUILDING COMMAND DEPENDS ON THE ARCHITECTURE OF THE PYTHON INTERPRETER 32BIT OR 64BIT YOU CAN CHECK THE ARCHITECTURE BY RUNNING THE FOLLOWING IN CMD ORPOWERSHELL CONSOLE

PYTHON C IMPORT STRUCT PRINTSTRUCTCALCSIZEP 8

THE ABOVE COMMANDS ASSUME THAT YOU HAVE THE PYTHON INSTALLATION FOLDER IN YOUR PATH ENVIRONMENT VARIABLE

YOU WILL NEED BUILD TOOLS FOR VISUAL STUDIO 2017

FOR 64BIT PYTHON CONFIGURE THE BUILD ENVIRONMENT WITH

SET DISTUTILSUSESDK1

CPROGRAM FILES X86MICROSOFT VISUAL

↳STUDIO2017BUILDTOOLSVCAUXILIARYBUILD VCVARALLBAT X64

AND BUILD SCIKITLEARN FROM THIS ENVIRONMENT

PYTHON SETUPPY INSTALL

REPLACEX64 BYX86 TO BUILD FOR 32BIT PYTHON

BUILDING BINARY PACKAGES AND INSTALLERS

THEWHL PACKAGE AND EXE INSTALLERS CAN BE BUILT WITH

PIP INSTALL WHEEL

PYTHON SETUPPY BDISTWHEEL BDISTWININST B DOCLOGOSSCIKITLEARNLOGOBMP

THE RESULTING PACKAGES ARE GENERATED IN THE DIST FOLDER

USING AN ALTERNATIVE COMPILER

IT IS POSSIBLE TO USE MINGW A PORT OF GCC TO WINDOWS OS AS AN ALTERNATIVE TO MSVC FOR 32BIT PYTHON NOT THAT EXTENSIONS BUILT WITH MINGW32 CAN BE REDISTRIBUTED AS REUSABLE PACKAGES AS THEY DEPEND ON GCC RUNTIME LIBRARIES TYPICALLY NOT INSTALLED ON ENDUSERS ENVIRONMENT

TO FORCE THE USE OF A PARTICULAR COMPILER PASS THE COMPILER FLAG TO THE BUILD STEP

2446 CHAPTER 7 DEVELOPER'S GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213  
PYTHON SETUPPY BUILD COMPILERMYCOMPILER INSTALL  
WHERE MYCOMPILER SHOULD BE ONE OF MINGW32 ORMSVC  
76 MAINTAINER COREDEVELOPER INFORMATION  
761 BEFORE A RELEASE  
1 UPDATE AUTHORS TABLE  
CD BUILDTOOLS MAKE AUTHORS CD  
AND COMMIT  
2 CONFIRM ANY BLOCKERS TAGGED FOR THE MILESTONE ARE RESOLVED AND THAT OTHER ISSUES TAGGED FOR THE MILESTONE CAN BE POSTPONED  
3 ENSURE THE CHANGE LOG AND COMMITS CORRESPOND WITHIN REASON AND THAT THE CHANGE LOG IS REASONABLY WELL CURATED SOME TOOLS FOR THESE TASKS INCLUDE  
• MAINTTOOLSSORTWHAT'SNEW.PY CAN PUT WHAT'S NEW ENTRIES INTO SECTIONS  
• THE MAINTTOOLSWHAT'SMISSING.SH SCRIPT MAY BE USED TO IDENTIFY PULL REQUESTS THAT WERE MERGED BUT LIKELY MISSING FROM WHAT'S NEW  
PREPARING A BUGFIX RELEASE  
SINCE ANY COMMITS TO A RELEASED BRANCH EG 0999X WILL AUTOMATICALLY UPDATE THE WEB SITE DOCUMENTATION IT IS BEST TO DEVELOP A BUGFIX RELEASE WITH A PULL REQUEST IN WHICH 0999X IS THE BASE IT ALSO ALLOWS YOU TO KEEP TRACK OF ANY TASKS TOWARDS RELEASE WITH A TO DO LIST  
MOST DEVELOPMENT OF THE BUG FIX RELEASE AND ITS DOCUMENTATION SHOULD HAPPEN IN MASTER TO AVOID ASYNCHRONY TO SELECT COMMITS FROM MASTER FOR USE IN THE BUG FIX VERSION 09993 YOU CAN USE  
GIT CHECKOUT B RELEASE09993 MASTER  
GIT REBASE I 0999X  
THEN PICK THE COMMITS FOR RELEASE AND RESOLVE ANY ISSUES AND CREATE A PULL REQUEST WITH 0999X AS BASE ADD A COMMIT UPDATING SKLEARNVERSION ADDITIONAL COMMITS CAN BE CHERRY-PICKED INTO THE RELEASE09993 BRANCH WHILE PREPARING THE RELEASE  
762 MAKING A RELEASE  
1 UPDATE DOCS  
• EDIT THE DOCWHAT'SNEW.RST FILE TO ADD RELEASE TITLE AND COMMIT STATISTICS YOU CAN RETRIEVE COMMIT STATISTICS WITH  
GIT SHORTLOG S 09993 CUT F2 SORT IGNORECASE TR N  
↩SED S GS  
• UPDATE THE RELEASE DATE IN WHAT'SNEW.RST  
• EDIT THE DOCINDEX.RST TO CHANGE THE 'NEWS' ENTRY OF THE FRONT PAGE  
76 MAINTAINER COREDEVELOPER INFORMATION 2447

SCIKITLEARN USER GUIDE RELEASE 0213

- NOTE THAT THESE CHANGES SHOULD BE MADE IN MASTER AND CHERRY PICKED INTO THE RELEASE BRANCH
- 2 ON THE BRANCH FOR RELEASING UPDATE THE VERSION NUMBER IN SKLEARNINITPY THE VERSION VARIABLE BY REMOVINGDEVONLY WHEN READY TO RELEASE ON MASTER INCREMENT THE VERSION IN THE SAME PLACE WHEN BRANCHING FOR RELEASE

3 CREATE THE TAG AND PUSH IT

- GIT TAG A 0999
- GIT PUSH GITGITHUBCOMSCIKITLEARNSCIKITLEARNGIT TAGS

4 CREATE THE SOURCE TARBALL

- WIPE CLEAN YOUR REPO
- GIT CLEAN XFD
- GENERATE THE TARBALL
- PYTHON SETUPPY SDIST

THE RESULT SHOULD BE IN THE DIST FOLDER WE WILL UPLOAD IT LATER WITH THE WHEELS CHECK THAT YOU CAN INSTALL IT IN A NEW VIRTUALENV AND THAT THE TESTS PASS

5 UPDATE THE DEPENDENCY VERSIONS AND SET BUILDCOMMIT VARIABLE TO THE RELEASE TAG AT

HTTPSGITHUBCOMMACPYTHONSCIKITLEARNWHEELS

ONCE THE CI HAS COMPLETED SUCCESSFULLY COLLECT THE GENERATED BINARY WHEEL PACKAGES AND UPLOAD THEM TO PYPI BY RUNNING THE FOLLOWING COMMANDS IN THE SCIKITLEARN SOURCE FOLDER CHECKED OUT AT THE RELEASE TAG

- PIP INSTALL U WHEELHOUSEUPLOADER TWINE
- PYTHON SETUPPY FETCHARTIFACTS

6 CHECK THE CONTENT OF THE DIST FOLDER IT SHOULD CONTAIN ALL THE WHEELS ALONG WITH THE SOURCE TARBALL “SCIKIT LEARNXXXTARGZ”

MAKE SURE THAT YOU DO NOT HAVE DEVELOPER VERSIONS OR OLDER VERSIONS OF THE SCIKITLEARN PACKAGE IN THAT FOLDER UPLOAD EVERYTHING AT ONCE TO HTTPSPYPIORG

TWINE UPLOAD DIST

7 FOR MAJORMINOR NOT BUGFIX RELEASE UPDATE THE SYMLINK FOR STABLE IN HTTPSGITHUBCOMSCIKITLEARN SCIKITLEARNGITHUBIO

- CD TMP
- GIT CLONE DEPTH 1 NOCHECKOUT GITGITHUBCOMSCIKITLEARNSCIKITLEARN

↪GITHUBIOGIT

CD SCIKITLEARNGITHUBIO

ECHO STABLE GITINFOSPARSECHECKOUT

GIT CHECKOUT MASTER

LN SF 0999 STABLE

GIT PUSH ORIGIN MASTER

THE FOLLOWING GITHUB CHECKLIST MIGHT BE HELPFUL IN A RELEASE PR

UPDATE NEWS ANDWHATS NEW DATE IN MASTER AND RELEASE BRANCH

CREATE TAG

UPDATE DEPENDENCIES ANDRELEASE TAG AT HTTPSGITHUBCOMMACPYTHONSCIKIT

↪LEARNWHEELS

2448 CHAPTER 7 DEVELOPER’S GUIDE

SCIKITLEARN USER GUIDE RELEASE 0213

TWINE THE WHEELS TO PYPI WHEN THATS GREEN

HTTPSGITHUBCOMSCIKITLEARNSCIKITLEARNRELEASES DRAFT

CONFIRM BOT DETECTED AT HTTPSGITHUBCOMCONDAFORGESCIKITLEARNFEEDSTOCK

↩→ANDWAITFORMERGE

HTTPSGITHUBCOMSCIKITLEARNSCIKITLEARNRELEASES PUBLISH

ANNOUNCE ON MAILING LIST

REGENERATE DASH DOCS HTTPSGITHUBCOMKAPELIDASHUSERCONTRIBUTIONSTREE

↩→MASTERDOCSETSSCIKIT

763 THE SCIKITLEARNORG WEB SITE

THE SCIKITLEARN WEB SITE HTTPSCIKITLEARNORG IS HOSTED AT GITHUB BUT SHOULD RARELY BE UPDATED MANUALLY BY PUSHING TO THE HTTPSGITHUBCOMSCIKITLEARNSCIKITLEARNGITHUBIO REPOSITORY MOST UPDATES CAN BE MADE BY PUSHING TO MASTER FOR DEV OR A RELEASE BRANCH LIKE 099X FROM WHICH CIRCLE CI BUILDS AND UPLOADS THE DOCUMENTATION AUTOMATICALLY

764 TRAVIS CRON JOBS

FROM HTTPSDOCSTRAVISCICOMUSERCRONJOBS TRAVIS CI CRON JOBS WORK SIMILARLY TO THE CRON UTILITY THEY RUN BUILDS AT REGULAR SCHEDULED INTERVALS INDEPENDENTLY OF WHETHER ANY COMMITS WERE PUSHED TO THE REPOSITORY CRON JOBS ALWAYS FETCH THE MOST RECENT COMMIT ON A PARTICULAR BRANCH AND BUILD THE PROJECT AT THAT STATE CRON JOBS CAN RUN DAILY WEEKLY OR MONTHLY WHICH IN PRACTICE MEANS UP TO AN HOUR AFTER THE SELECTED TIME SPAN AND YOU CANNOT SET THEM TO RUN AT A SPECIFIC TIME

FOR SCIKITLEARN CRON JOBS ARE USED FOR BUILDS THAT WE DO NOT WANT TO RUN IN EACH PR AS AN EXAMPLE THE BUILD WITH THE DEV VERSIONS OF NUMPY AND SCIPIY IS RUN AS A CRON JOB MOST OF THE TIME WHEN THIS NUMPYDEV BUILD FAIL IT IS RELATED TO A NUMPY CHANGE AND NOT A SCIKITLEARN ONE SO IT WOULD NOT MAKE SENSE TO BLAME THE PR AUTHOR FOR THE TRAVIS FAILURE THE DEFINITION OF WHAT GETS RUN IN THE CRON JOB IS DONE IN THE TRAVISYML CONFIG FILE EXACTLY THE SAME WAY AS THE OTHER TRAVIS JOBS WE USE A IF TYPE CRON FILTER IN ORDER FOR THE BUILD TO BE RUN ONLY IN CRON JOBS

THE BRANCH TARGETED BY THE CRON JOB AND THE FREQUENCY OF THE CRON JOB IS SET VIA THE WEB UI AT HTTPSWWWTRAVISCIORGSCIKITLEARNSCIKITLEARNSETTINGS

765 EXPERIMENTAL FEATURES

THESKLEARNEXPERIMENTAL MODULE WAS INTRODUCED IN 021 AND CONTAINS EXPERIMENTAL FEATURES ESTIMATORS THAT ARE SUBJECT TO CHANGE WITHOUT DEPRECATION CYCLE

TO CREATE AN EXPERIMENTAL MODULE YOU CAN JUST COPY AND MODIFY THE CONTENT OF ENABLEHISTGRADIENTBOOSTINGPY OR ENABLEITERATIVEIMPUTERPY

NOTE THAT THE PUBLIC IMPORT PATH MUST BE TO A PUBLIC SUBPACKAGE LIKE SKLEARNENSEMBLE ORSKLEARNIMPUTE NOT JUST APY MODULE ALSO THE PRIVATE EXPERIMENTAL FEATURES THAT ARE IMPORTED MUST BE IN A SUBMODULESUBPACKAGE OF THE PUBLIC SUBPACKAGE EG SKLEARNENSEMBLEHISTGRADIENTBOOSTING ORSKLEARNIMPUTE

ITERATIVEPY THIS IS NEEDED SO THAT PICKLES STILL WORK IN THE FUTURE WHEN THE FEATURES AREN'T EXPERIMENTAL ANYMORE

PLEASE ALSO WRITE BASIC TESTS FOLLOWING THOSE IN TESTENABLEHISTGRADIENTBOOSTINGPY

MAKE SURE EVERY USERFACING CODE YOU WRITE EXPLICITLY MENTIONS THAT THE FEATURE IS EXPERIMENTAL AND ADD A NOQA COMMENT TO AVOID PEP8RELATED WARNINGS

TO USE THIS EXPERIMENTAL FEATURE WE NEED TO EXPLICITLY ASK FOR IT

FROM SKLEARNEXPERIMENTAL IMPORT ENABLEHISTGRADIENTBOOSTING NOQA

FROM SKLEARNENSEMBLE IMPORT HISTGRADIENTBOOSTINGREGRESSOR

76 MAINTAINER COREDEVELOPER INFORMATION 2449

SCIKITLEARN USER GUIDE RELEASE 0213  
FOR THE DOCS TO RENDER PROPERLY PLEASE ALSO IMPORT ENABLEMYEXPERIMENTALFEATURE INDOCCONFPY  
ELSE SPHINX WON'T BE ABLE TO IMPORT THE CORRESPONDING MODULES NOTE THAT USING FROM SKLEARNEXPERIMENTAL  
IMPORTDOES NOT WORK  
NOTE THAT SOME EXPERIMENTAL CLASSES FUNCTIONS ARE NOT INCLUDED IN THE SKLEARNEXPERIMENTAL MODULE  
SKLEARNDATASETSFETCHOPENML  
2450 CHAPTER 7 DEVELOPER'S GUIDE



BIBLIOGRAPHY

M2012 "MACHINE LEARNING A PROBABILISTIC PERSPECTIVE" MURPHY K P CHAPTER 1443 PP 492493 THE MIT PRESS 2012

RW2006 CARL EDUARD RASMUSSEN AND CHRISTOPHER KI WILLIAMS "GAUSSIAN PROCESSES FOR MACHINE LEARNING" MIT PRESS 2006 LINK TO AN OFFICIAL COMPLETE PDF VERSION OF THE BOOK HERE

B1999 L BREIMAN "PASTING SMALL VOTES FOR CLASSIFICATION IN LARGE DATABASES AND ONLINE" MACHINE LEARNING 361 85103 1999

B1996 L BREIMAN "BAGGING PREDICTORS" MACHINE LEARNING 242 123140 1996

H1998 T HO "THE RANDOM SUBSPACE METHOD FOR CONSTRUCTING DECISION FORESTS" PATTERN ANALYSIS AND MACHINE INTELLIGENCE 208 832844 1998

LG2012 G LOUPPE AND P GEURTS "ENSEMBLES ON RANDOM PATCHES" MACHINE LEARNING AND KNOWLEDGE DISCOVERY IN DATABASES 346361 2012

B2001 12 BREIMAN "RANDOM FORESTS" MACHINE LEARNING 451 532 2001

B1998 12 BREIMAN "ARCING CLASSIFIERS" ANNALS OF STATISTICS 1998

L2014 G LOUPPE "UNDERSTANDING RANDOM FORESTS FROM THEORY TO PRACTICE" PHD THESIS U OF LIEGE 2014

FS1995 Y FREUND AND R SCHAPIRE "A DECISIONTHEORETIC GENERALIZATION OF ONLINE LEARNING AND AN APPLICATION TO BOOSTING" 1997

ZZRH2009 J ZHU H ZOU S ROSSET T HASTIE "MULTICLASS ADABOOST" 2009

D1997 8 DRUCKER "IMPROVING REGRESSORS USING BOOSTING TECHNIQUES" 1997

HTF T HASTIE R TIBSHIRANI AND J FRIEDMAN "ELEMENTS OF STATISTICAL LEARNING ED 2" SPRINGER 2009

VEB2009 VINH EPPS AND BAILEY 2009 "INFORMATION THEORETIC MEASURES FOR CLUSTERINGS COMPARISON" PROCEEDINGS OF THE 26TH ANNUAL INTERNATIONAL CONFERENCE ON MACHINE LEARNING ICML '09 DOI10114515533741553511 ISBN 9781605585161

VEB2010 VINH EPPS AND BAILEY 2010 "INFORMATION THEORETIC MEASURES FOR CLUSTERINGS COMPARISON VARIANTS PROPERTIES NORMALIZATION AND CORRECTION FOR CHANCE" JMLR HTTPJMLRCSAILMITEDUPAPERSVOLUME11 VINH10AVINH10APDF

YAT2016 YANG ALGESHEIMER AND TESSONE 2016 "A COMPARATIVE ANALYSIS OF COMMUNITY DETECTION ALGORITHMS ON ARTIFICIAL NETWORKS" SCIENTIFIC REPORTS 6 30750 DOI101038SREP30750

B2011 IDENTIFICATION AND CHARACTERIZATION OF EVENTS IN SOCIAL MEDIA HILA BECKER PHD THESIS

MRL09 "ONLINE DICTIONARY LEARNING FOR SPARSE CODING" J MAIRAL F BACH J PONCE G SAPIRO 2009

JEN09 "STRUCTURED SPARSE PRINCIPAL COMPONENT ANALYSIS" R JENATTON G OBOZINSKI F BACH 2009 2451

SCIKITLEARN USER GUIDE RELEASE 0213  
R45F14345C0001 12 BREIMAN "RANDOM FORESTS" MACHINE LEARNING 451 532 2001  
RF91CAB2DC4271 12 BREIMAN "RANDOM FORESTS" MACHINE LEARNING 451 532 2001  
RF91CAB2DC4272 P GEURTS D ERNST AND L WEHENKEL "EXTREMELY RANDOMIZED TREES" MACHINE LEARNING 631 342 2006  
RC8F28BFAD63F1 P GEURTS D ERNST AND L WEHENKEL "EXTREMELY RANDOMIZED TREES" MACHINE LEARNING 631 342 2006  
RA7D0C8995FBC1 P GEURTS D ERNST AND L WEHENKEL "EXTREMELY RANDOMIZED TREES" MACHINE LEARNING 631 342 2006  
GUYON2015 I GUYON K BENNETT G CAWLEY HJ ESCALANTE S ESCALERA TK HO N MACIÀ B RAY M SAEED  
AR STATNIKOV E VIEGAS DESIGN OF THE 2015 CHALEARN AUTOML CHALLENGE IJCNN 2015  
MOSLEY2013 L MOSLEY A BALANCED APPROACH TO THE MULTICLASS IMBALANCE PROBLEM IJCV 2010  
KELLEHER2015 JOHN D KELLEHER BRIAN MAC NAMEE AOIFE D'ARCY FUNDAMENTALS OF MACHINE LEARNING FOR PREDICTIVE DATA ANALYTICS ALGORITHMS WORKED EXAMPLES AND CASE STUDIES 2015  
URBANOWICZ2015 URBANOWICZ RJ MOORE JH EXSTRACS 20 DESCRIPTION AND EVALUATION OF A SCALABLE LEARNING CLASSIFIER SYSTEM EVOL INTEL 2015 8 89  
MANNING2008 CD MANNING P RAGHAVAN H SCHÜTZE INTRODUCTION TO INFORMATION RETRIEVAL 2008  
EVERINGHAM2010 M EVERINGHAM L VAN GOOL CKI WILLIAMS J WINN A ZISSERMAN THE PASCAL VISUAL OBJECT CLASSES VOC CHALLENGE IJCV 2010  
DAVIS2006 J DAVIS M GOADRICH THE RELATIONSHIP BETWEEN PRECISIONRECALL AND ROC CURVES ICML 2006  
FLACH2015 PA FLACH M KULL PRECISIONRECALLGAIN CURVES PR ANALYSIS DONE RIGHT NIPS 2015  
HTF2009 T HASTIE R TIBSHIRANI AND J FRIEDMAN THE ELEMENTS OF STATISTICAL LEARNING SECOND EDITION SECTION 10132 SPRINGER 2009  
MOL2019 C MOLNAR INTERPRETABLE MACHINE LEARNING SECTION 51 2019  
NQY18 J NOTHMAN H QIN AND R YURCHAK 2018 "STOP WORD LISTS IN FREE OPENSOURCE SOFTWARE PACKAGES" IN PROC WORKSHOP FOR NLP OPEN SOURCE SOFTWARE  
RR2007 "RANDOM FEATURES FOR LARGESCALE KERNEL MACHINES" RAHIMI A AND RECHT B ADVANCES IN NEURAL INFORMATION PROCESSING 2007  
LS2010 "RANDOM FOURIER APPROXIMATIONS FOR SKEWED MULTIPLICATIVE HISTOGRAM KERNELS" RANDOM FOURIER APPROXIMATIONS FOR SKEWED MULTIPLICATIVE HISTOGRAM KERNELS LECTURE NOTES FOR COMPUTER SCIENCDD DAGM  
VZ2010 "EFFICIENT ADDITIVE KERNELS VIA EXPLICIT FEATURE MAPS" VEDALDI A AND ZISSERMAN A COMPUTER VISION AND PATTERN RECOGNITION 2010  
VVZ2010 "GENERALIZED RBF FEATURE MAPS FOR EFFICIENT DETECTION" VEMPATI S AND VEDALDI A AND ZISSERMAN A AND JAWAHAR CV 2010  
R57CF438D70601 OBTAINING CALIBRATED PROBABILITY ESTIMATES FROM DECISION TREES AND NAIVE BAYESIAN CLASSIFIERS B ZADROZNY C ELKAN ICML 2001  
R57CF438D70602 TRANSFORMING CLASSIFIER SCORES INTO ACCURATE MULTICLASS PROBABILITY ESTIMATES B ZADROZNY C ELKAN KDD 2002  
R57CF438D70603 PROBABILISTIC OUTPUTS FOR SUPPORT VECTOR MACHINES AND COMPARISONS TO REGULARIZED LIKELIHOOD METHODS J PLATT 1999  
R57CF438D70604 PREDICTING GOOD PROBABILITIES WITH SUPERVISED LEARNING A NICULESCUMIZIL R CARUANA ICML 2005  
2452 BIBLIOGRAPHY

SCIKITLEARN USER GUIDE RELEASE 0213

R2C55E37003FE1 ANKERST MIHAEL MARKUS M BREUNIG HANSPETER KRIEGEL AND JÖRG SANDER “OPTICS ORDERING POINTS TO IDENTIFY THE CLUSTERING STRUCTURE” ACM SIGMOD RECORD 28 NO 2 1999 4960

R2C55E37003FE2 SCHUBERT ERICH MICHAEL GERTZ “IMPROVING THE CLUSTER STRUCTURE EXTRACTED FROM OPTICS PLOTS” PROC OF THE CONFERENCE “LERNEN WISSEN DATEN ANALYSEN” LWDA 2018 318329

1 ANKERST MIHAEL MARKUS M BREUNIG HANSPETER KRIEGEL AND JÖRG SANDER “OPTICS ORDERING POINTS TO IDENTIFY THE CLUSTERING STRUCTURE” ACM SIGMOD RECORD 28 NO 2 1999 4960

R68AE096DA0E41 ROUSSEEUW PJ VAN DRIESSEN K “A FAST ALGORITHM FOR THE MINIMUM COVARIANCE DETERMINANT ESTIMATOR” TECHNOMETRICS 413 212 1999

RVD A FAST ALGORITHM FOR THE MINIMUM COVARIANCE DETERMINANT ESTIMATOR 1999 AMERICAN STATISTICAL ASSOCIATION AND THE AMERICAN SOCIETY FOR QUALITY TECHNOMETRICS

RVDRIESSEN A FAST ALGORITHM FOR THE MINIMUM COVARIANCE DETERMINANT ESTIMATOR 1999 AMERICAN STATISTICAL ASSOCIATION AND THE AMERICAN SOCIETY FOR QUALITY TECHNOMETRICS

R9F63E655F7BDROUSEEUW1984 P J ROUSSEEUW LEAST MEDIAN OF SQUARES REGRESSION J AM STAT ASS 79871 1984

R9F63E655F7BDROUSSEEUW A FAST ALGORITHM FOR THE MINIMUM COVARIANCE DETERMINANT ESTIMATOR 1999 AMERICAN STATISTICAL ASSOCIATION AND THE AMERICAN SOCIETY FOR QUALITY TECHNOMETRICS

R9F63E655F7BDBUTLERDAVIES R W BUTLER P L DAVIES AND M JHUN ASYMPTOTICS FOR THE MINIMUM COVARIANCE DETERMINANT ESTIMATOR THE ANNALS OF STATISTICS 1993 V OL 21 NO 3 13851400

RVD A FAST ALGORITHM FOR THE MINIMUM COVARIANCE DETERMINANT ESTIMATOR 1999 AMERICAN STATISTICAL ASSOCIATION AND THE AMERICAN SOCIETY FOR QUALITY TECHNOMETRICS

RVDRIESSEN A FAST ALGORITHM FOR THE MINIMUM COVARIANCE DETERMINANT ESTIMATOR 1999 AMERICAN STATISTICAL ASSOCIATION AND THE AMERICAN SOCIETY FOR QUALITY TECHNOMETRICS

1 DHILLON I S 2001 AUGUST COCLUSTERING DOCUMENTS AND WORDS USING BIPARTITE SPECTRAL GRAPH PARTITIONING IN PROCEEDINGS OF THE SEVENTH ACM SIGKDD INTERNATIONAL CONFERENCE ON KNOWLEDGE DISCOVERY AND DATA MINING PP 269274 ACM

1 KLUGER Y BASRI R CHANG J T GERSTEIN M 2003 SPECTRAL BICLUSTERING OF MICROARRAY DATA COCLUSTERING GENES AND CONDITIONS GENOME RESEARCH 134 703716

1 I GUYON “DESIGN OF EXPERIMENTS FOR THE NIPS 2003 VARIABLE SELECTION BENCHMARK” 2003

1 J FRIEDMAN “MULTIVARIATE ADAPTIVE REGRESSION SPLINES” THE ANNALS OF STATISTICS 19 1 PAGES 167 1991

2 L BREIMAN “BAGGING PREDICTORS” MACHINE LEARNING 24 PAGES 123140 1996

1 J FRIEDMAN “MULTIVARIATE ADAPTIVE REGRESSION SPLINES” THE ANNALS OF STATISTICS 19 1 PAGES 167 1991

2 L BREIMAN “BAGGING PREDICTORS” MACHINE LEARNING 24 PAGES 123140 1996

1 J FRIEDMAN “MULTIVARIATE ADAPTIVE REGRESSION SPLINES” THE ANNALS OF STATISTICS 19 1 PAGES 167 1991

2 L BREIMAN “BAGGING PREDICTORS” MACHINE LEARNING 24 PAGES 123140 1996

1 10 ZHU H ZOU S ROSSET T HASTIE “MULTICLASS ADABOOST” 2009

1 T HASTIE R TIBSHIRANI AND J FRIEDMAN “ELEMENTS OF STATISTICAL LEARNING ED 2” SPRINGER 2009

1 G CELEUX M EL ANBARI JM MARIN C P ROBERT “REGULARIZATION IN REGRESSION COMPARING BAYESIAN AND FREQUENTIST METHODS IN A POORLY INFORMATIVE SITUATION” 2009

1 S MARSLAND “MACHINE LEARNING AN ALGORITHMIC PERSPECTIVE” CHAPTER 10 2009 HTTPSEATMASSEYACNZ PERSONALSRMARSLANDCODE10LLEPY

R33E4EC8C4AD51 Y FREUND R SCHAPIRE “A DECISIONTHEORETIC GENERALIZATION OF ONLINE LEARNING AND AN APPLICATION TO BOOSTING” 1995

BIBLIOGRAPHY 2453

SCIKITLEARN USER GUIDE RELEASE 0213

R33E4EC8C4AD52 10 ZHU H ZOU S ROSSET T HASTIE “MULTICLASS ADABOOST” 2009

R0C261B7DEE9D1 Y FREUND R SCHAPIRE “A DECISIONTHEORETIC GENERALIZATION OF ONLINE LEARNING AND AN APPLICATION TO BOOSTING” 1995

R0C261B7DEE9D2 8 DRUCKER “IMPROVING REGRESSORS USING BOOSTING TECHNIQUES” 1997

RB1846455D0E51 L BREIMAN “PASTING SMALL VOTES FOR CLASSIFICATION IN LARGE DATABASES AND ONLINE” MACHINE LEARNING 361 85103 1999

RB1846455D0E52 L BREIMAN “BAGGING PREDICTORS” MACHINE LEARNING 242 123140 1996

RB1846455D0E53 T HO “THE RANDOM SUBSPACE METHOD FOR CONSTRUCTING DECISION FORESTS” PATTERN ANALYSIS AND MACHINE INTELLIGENCE 208 832844 1998

RB1846455D0E54 G LOUPPE AND P GEURTS “ENSEMBLES ON RANDOM PATCHES” MACHINE LEARNING AND KNOWLEDGE DISCOVERY IN DATABASES 346361 2012

R4D113BA76FC01 L BREIMAN “PASTING SMALL VOTES FOR CLASSIFICATION IN LARGE DATABASES AND ONLINE” MACHINE LEARNING 361 85103 1999

R4D113BA76FC02 L BREIMAN “BAGGING PREDICTORS” MACHINE LEARNING 242 123140 1996

R4D113BA76FC03 T HO “THE RANDOM SUBSPACE METHOD FOR CONSTRUCTING DECISION FORESTS” PATTERN ANALYSIS AND MACHINE INTELLIGENCE 208 832844 1998

R4D113BA76FC04 G LOUPPE AND P GEURTS “ENSEMBLES ON RANDOM PATCHES” MACHINE LEARNING AND KNOWLEDGE DISCOVERY IN DATABASES 346361 2012

RD7AE0A2AE6881 LIU FEI TONY TING KAI MING AND ZHOU ZHIHUA “ISOLATION FOREST” DATA MINING 2008

ICDM’08 EIGHTH IEEE INTERNATIONAL CONFERENCE ON

RD7AE0A2AE6882 LIU FEI TONY TING KAI MING AND ZHOU ZHIHUA “ISOLATIONBASED ANOMALY DETECTION” ACM TRANSACTIONS ON KNOWLEDGE DISCOVERY FROM DATA TKDD 61 2012 3

R6E47E53BACBD1 P GEURTS D ERNST AND L WEHENKEL “EXTREMELY RANDOMIZED TREES” MACHINE LEARNING 631 342 2006

R6E47E53BACBD2 MOOSMANN F AND TRIGGS B AND JURIE F “FAST DISCRIMINATIVE VISUAL CODEBOOKS USING RANDOMIZED CLUSTERING FORESTS” NIPS 2007

R1B90AC3CA370YATES2011 R BAEZAYATES AND B RIBEIRONETO 2011 MODERN INFORMATION RETRIEVAL ADDISON WESLEY PP 6874

R1B90AC3CA370MRS2008 CD MANNING P RAGHAVAN AND H SCHÜTZE 2008 INTRODUCTION TO INFORMATION RETRIEVAL CAMBRIDGE UNIVERSITY PRESS PP 118120

RE310F679C81E1 GUYON I WESTON J BARNHILL S VAPNIK V “GENE SELECTION FOR CANCER CLASSIFICATION USING SUPPORT VECTOR MACHINES” MACH LEARN 4613 389-422 2002

R6F4D61CEB4111 GUYON I WESTON J BARNHILL S VAPNIK V “GENE SELECTION FOR CANCER CLASSIFICATION USING SUPPORT VECTOR MACHINES” MACH LEARN 4613 389-422 2002

1 MUTUAL INFORMATION ON WIKIPEDIA

2 A KRASKOV H STOGBAUER AND P GRASSBERGER “ESTIMATING MUTUAL INFORMATION” PHYS REV E 69 2004

3 B C ROSS “MUTUAL INFORMATION BETWEEN DISCRETE AND CONTINUOUS DATA SETS” PLOS ONE 92 2014

4 L F KOZACHENKO N N LEONENKO “SAMPLE ESTIMATE OF THE ENTROPY OF A RANDOM VECTOR PROBL PEREDACHI INF 232 1987 916

1 MUTUAL INFORMATION ON WIKIPEDIA

2 A KRASKOV H STOGBAUER AND P GRASSBERGER “ESTIMATING MUTUAL INFORMATION” PHYS REV E 69 2004

2454 BIBLIOGRAPHY

SCIKITLEARN USER GUIDE RELEASE 0213

3 B C ROSS “MUTUAL INFORMATION BETWEEN DISCRETE AND CONTINUOUS DATA SETS” PLOS ONE 92 2014

4 L F KOZACHENKO N N LEONENKO “SAMPLE ESTIMATE OF THE ENTROPY OF A RANDOM VECTOR” PROBL PEREDACHI INF 232 1987 916

RCD31B817A31E1 STEF VAN BUUREN KARIN GROOTHUISOUDSHOORN 2011 “MICE MULTIVARIATE IMPUTATION BY CHAINED EQUATIONS IN R” JOURNAL OF STATISTICAL SOFTWARE 45 167

RCD31B817A31E2 S F BUCK 1960 “A METHOD OF ESTIMATION OF MISSING VALUES IN MULTIVARIATE DATA SUITABLE FOR USE WITH AN ELECTRONIC COMPUTER” JOURNAL OF THE ROYAL STATISTICAL SOCIETY 222 302306

RE4616EF910FB1 PETER J HUBER ELVEZIO M RONCHETTI ROBUST STATISTICS CONCOMITANT SCALE ESTIMATES PG 172

RE4616EF910FB2 ART B OWEN 2006 A ROBUST HYBRID OF LASSO AND RIDGE REGRESSION HTTPSTATWEBSTANFORDEDU OWENREPORTSHHUPDF

R80CE5B25CF9D1 HTTPSENWIKIPEDIAORGWIKIRANSAC

R80CE5B25CF9D2 HTTPSWWWSRICOMSITESDEFAULTFILESPUBLICATIONSRRANSACPUBLICATIONPDF

R80CE5B25CF9D3 HTTPWWWBMAORGBMVC2009PAPERSPAPER355PAPER355PDF

1 “LEAST ANGLE REGRESSION” EFRON ET AL HTTPSTATWEBSTANFORDEDUTIBSFTPLARSPDF

2 WIKIPEDIA ENTRY ON THE LEASTANGLE REGRESSION

3 WIKIPEDIA ENTRY ON THE LASSO

1 “LEAST ANGLE REGRESSION” EFRON ET AL HTTPSTATWEBSTANFORDEDUTIBSFTPLARSPDF

2 WIKIPEDIA ENTRY ON THE LEASTANGLE REGRESSION

3 WIKIPEDIA ENTRY ON THE LASSO

R7F4D308F50541 TENENBAUM JB DE SILVA V LANGFORD JC A GLOBAL GEOMETRIC FRAMEWORK FOR NONLINEAR DIMENSIONALITY REDUCTION SCIENCE 290 5500

R62E36DD1B0561 ROWEIS S SAUL L NONLINEAR DIMENSIONALITY REDUCTION BY LOCALLY LINEAR EMBEDDING SCIENCE 2902323 2000

R62E36DD1B0562 DONOHO D GRIMES C HESSIAN EIGENMAPS LOCALLY LINEAR EMBEDDING TECHNIQUES FOR HIGH DIMENSIONAL DATA PROC NATL ACAD SCI U S A 1005591 2003

R62E36DD1B0563 ZHANG Z WANG J MLLE MODIFIED LOCALLY LINEAR EMBEDDING USING MULTIPLE WEIGHTS HTTPCITESEERXISTPSUEDUVIEWDOC SUMMARYDOI101170382

R62E36DD1B0564 ZHANG Z ZHA H PRINCIPAL MANIFOLDS AND NONLINEAR DIMENSIONALITY REDUCTION VIA TANGENT SPACE ALIGNMENT JOURNAL OF SHANGHAI UNIV 8406 2004

1 ROWEIS S SAUL L NONLINEAR DIMENSIONALITY REDUCTION BY LOCALLY LINEAR EMBEDDING SCIENCE 2902323 2000

2 DONOHO D GRIMES C HESSIAN EIGENMAPS LOCALLY LINEAR EMBEDDING TECHNIQUES FOR HIGHDIMENSIONAL DATA PROC NATL ACAD SCI U S A 1005591 2003

3 ZHANG Z WANG J MLLE MODIFIED LOCALLY LINEAR EMBEDDING USING MULTIPLE WEIGHTS HTTPCITESEERXIST PSUEDUVIEWDOC SUMMARYDOI101170382

4 ZHANG Z ZHA H PRINCIPAL MANIFOLDS AND NONLINEAR DIMENSIONALITY REDUCTION VIA TANGENT SPACE ALIGNMENT JOURNAL OF SHANGHAI UNIV 8406 2004

1 WIKIPEDIA ENTRY FOR THE AVERAGE PRECISION

1 BRODERSEN KH ONG CS STEPHAN KE BUHMANN JM 2010 THE BALANCED ACCURACY AND ITS POSTERIOR DISTRIBUTION PROCEEDINGS OF THE 20TH INTERNATIONAL CONFERENCE ON PATTERN RECOGNITION 312124

2 JOHN D KELLEHER BRIAN MAC NAMEE AOIFE D'ARCY 2015 FUNDAMENTALS OF MACHINE LEARNING FOR PREDICTIVE DATA ANALYTICS ALGORITHMS WORKED EXAMPLES AND CASE STUDIES

BIBLIOGRAPHY 2455

SCIKITLEARN USER GUIDE RELEASE 0213

1 WIKIPEDIA ENTRY FOR THE BRIER SCORE

1 J COHEN 1960 "A COEFFICIENT OF AGREEMENT FOR NOMINAL SCALES" EDUCATIONAL AND PSYCHOLOGICAL MEASUREMENT 2013746 DOI101177001316446002000104

2 R ARTSTEIN AND M POESIO 2008 "INTERCODER AGREEMENT FOR COMPUTATIONAL LINGUISTICS" COMPUTATIONAL LINGUISTICS 344555596

3 WIKIPEDIA ENTRY FOR THE COHEN'S KAPPA

1 WIKIPEDIA ENTRY FOR THE CONFUSION MATRIX WIKIPEDIA AND OTHER REFERENCES MAY USE A DIFFERENT CONVENTION FOR AXES

1 WIKIPEDIA ENTRY FOR THE F1SCORE

1 R BAEZAYATES AND B RIBEIRONETO 2011 MODERN INFORMATION RETRIEVAL ADDISON WESLEY PP 327328

2 WIKIPEDIA ENTRY FOR THE F1SCORE

1 GRIGORIOS TSOUMAKAS IOANNIS KATAKIS MULTILABEL CLASSIFICATION AN OVERVIEW INTERNATIONAL JOURNAL OF DATA WAREHOUSING MINING 33 113 JULYSEPTEMBER 2007

2 WIKIPEDIA ENTRY ON THE HAMMING DISTANCE

1 WIKIPEDIA ENTRY ON THE HINGE LOSS

2 KOBAYASHI YORAM SINGER ON THE ALGORITHMIC IMPLEMENTATION OF MULTICLASS KERNELBASED VECTOR MACHINES JOURNAL OF MACHINE LEARNING RESEARCH 2 2001 265292

3 L1 AND L2 REGULARIZATION FOR MULTICLASS HINGE LOSS MODELS BY ROBERT C MOORE JOHN DENERO

1 WIKIPEDIA ENTRY FOR THE JACCARD INDEX

1 BALDI BRUNAK CHAUVIN ANDERSEN AND NIELSEN 2000 ASSESSING THE ACCURACY OF PREDICTION ALGORITHMS FOR CLASSIFICATION AN OVERVIEW

2 WIKIPEDIA ENTRY FOR THE MATTHEWS CORRELATION COEFFICIENT

3 GORODKIN 2004 COMPARING TWO KCATEGORY ASSIGNMENTS BY A KCATEGORY CORRELATION COEFFICIENT

4 JURMAN RICCADONNA FURLANELLO 2012 A COMPARISON OF MCC AND CEN ERROR MEASURES IN MULTICLASS PREDICTION

1 WIKIPEDIA ENTRY FOR THE PRECISION AND RECALL

2 WIKIPEDIA ENTRY FOR THE F1SCORE

3 DISCRIMINATIVE METHODS FOR MULTILABELED CLASSIFICATION ADVANCES IN KNOWLEDGE DISCOVERY AND DATA MINING 2004 PP 2230 BY SHANTANU GODBOLE SUNITA SARAWAGI

1 WIKIPEDIA ENTRY FOR THE RECEIVER OPERATING CHARACTERISTIC

2 FAWCETT T AN INTRODUCTION TO ROC ANALYSISJ PATTERN RECOGNITION LETTERS 2006 278861874

3 ANALYZING A PORTION OF THE ROC CURVE MCCLISH 1989

1 WIKIPEDIA ENTRY FOR THE RECEIVER OPERATING CHARACTERISTIC

2 FAWCETT T AN INTRODUCTION TO ROC ANALYSISJ PATTERN RECOGNITION LETTERS 2006 278861874

1 WIKIPEDIA ENTRY ON THE COEFFICIENT OF DETERMINATION

1 TSOUMAKAS G KATAKIS I VLAHAVAS I 2010 MINING MULTILABEL DATA IN DATA MINING AND KNOWLEDGE DISCOVERY HANDBOOK PP 667685 SPRINGER US

1 TSOUMAKAS G KATAKIS I VLAHAVAS I 2010 MINING MULTILABEL DATA IN DATA MINING AND KNOWLEDGE DISCOVERY HANDBOOK PP 667685 SPRINGER US

2456 BIBLIOGRAPHY

SCIKITLEARN USER GUIDE RELEASE 0213

1 VINH EPPS AND BAILEY 2010 INFORMATION THEORETIC MEASURES FOR CLUSTERINGS COMPARISON VARIANTS PROPERTIES NORMALIZATION AND CORRECTION FOR CHANCE JMLR

2 WIKIPEDIA ENTRY FOR THE ADJUSTED MUTUAL INFORMATION

HUBERT1985 L HUBERT AND P ARABIE COMPARING PARTITIONS JOURNAL OF CLASSIFICATION 1985 HTTPSLINKSPRINGERCOM ARTICLE1010072FBF01908075

WK HTTPSENWIKIPEDIAORGWIKIRANDINDEXADJUSTEDRANDINDEX

1 T CALINSKI AND J HARABASZ 1974 "A DENDRITE METHOD FOR CLUSTER ANALYSIS" COMMUNICATIONS IN STATISTICS

1 DAVIES DAVID L BOULDIN DONALD W 1979 "A CLUSTER SEPARATION MEASURE" IEEE TRANSACTIONS ON PATTERN ANALYSIS AND MACHINE INTELLIGENCE PAMI1 2 224227

1 ANDREW ROSENBERG AND JULIA HIRSCHBERG 2007 VMEASURE A CONDITIONAL ENTROPYBASED EXTERNAL CLUSTER EVALUA TION MEASURE

1 E B FOWKLES AND C L MALLOWS 1983 "A METHOD FOR COMPARING TWO HIERARCHICAL CLUSTERINGS" JOURNAL OF THE AMERICAN STATISTICAL ASSOCIATION

2 WIKIPEDIA ENTRY FOR THE FOWLKESMALLOWS INDEX

1 ANDREW ROSENBERG AND JULIA HIRSCHBERG 2007 VMEASURE A CONDITIONAL ENTROPYBASED EXTERNAL CLUSTER EVALUA TION MEASURE

1 PETER J ROUSSEEUW 1987 "SILHOUETTES A GRAPHICAL AID TO THE INTERPRETATION AND VALIDATION OF CLUSTER ANALYSIS" COMPUTATIONAL AND APPLIED MATHEMATICS 20 5365

2 WIKIPEDIA ENTRY ON THE SILHOUETTE COEFFICIENT

1 PETER J ROUSSEEUW 1987 "SILHOUETTES A GRAPHICAL AID TO THE INTERPRETATION AND VALIDATION OF CLUSTER ANALYSIS" COMPUTATIONAL AND APPLIED MATHEMATICS 20 5365

2 WIKIPEDIA ENTRY ON THE SILHOUETTE COEFFICIENT

1 ANDREW ROSENBERG AND JULIA HIRSCHBERG 2007 VMEASURE A CONDITIONAL ENTROPYBASED EXTERNAL CLUSTER EVALUA TION MEASURE

R16529824BFF21 BISHOP CHRISTOPHER M 2006 "PATTERN RECOGNITION AND MACHINE LEARNING" V OL 4 NO 4 NEW YORK SPRINGER

R16529824BFF22 HAGAI ATTIAS 2000 "A VARIATIONAL BAYESIAN FRAMEWORK FOR GRAPHICAL MODELS" IN ADVANCES IN NEURAL INFORMATION PROCESSING SYSTEMS 12

R16529824BFF23 BLEI DAVID M AND MICHAEL I JORDAN 2006 "VARIATIONAL INFERENCE FOR DIRICHLET PROCESS MIX TURES" BAYESIAN ANALYSIS 11

R2EDDAEEC08491 "SOLVING MULTICLASS LEARNING PROBLEMS VIA ERRORCORRECTING OUTPUT CODES" DIETTERICH T BAKIRI G JOURNAL OF ARTIFICIAL INTELLIGENCE RESEARCH 2 1995

R2EDDAEEC08492 "THE ERROR CODING METHOD AND PICTS" JAMES G HASTIE T JOURNAL OF COMPUTATIONAL AND GRAPH ICAL STATISTICS 7 1998

R2EDDAEEC08493 "THE ELEMENTS OF STATISTICAL LEARNING" HASTIE T TIBSHIRANI R FRIEDMAN J PAGE 606 SECOND EDITION 2008

RCA479BB498411 BREUNIG M M KRIEGEL H P NG R T SANDER J 2000 MAY LOF IDENTIFYING DENSITY BASED LOCAL OUTLIERS IN ACM SIGMOD RECORD

RF9B6BAEE82291 J GOLDBERGER G HINTON S ROWEIS R SALAKHUTDINOV "NEIGHBOURHOOD COMPONENTS ANALYSIS" ADVANCES IN NEURAL INFORMATION PROCESSING SYSTEMS 17 513520 2005 HTTPWWWCSNYUEDUROWEISPAPERS NCANIPSPDF

BIBLIOGRAPHY 2457

SCIKITLEARN USER GUIDE RELEASE 0213

RF9B6BAEE82292 WIKIPEDIA ENTRY ON NEIGHBORHOOD COMPONENTS ANALYSIS [HTTPS://ENWIKIPEDIA.ORG/WIKI/NEIGHBOURHOODCOMPONENTS/ANALYSIS](https://en.wikipedia.org/wiki/Neighborhood_components_analysis)

RF3E1504535DE1 IK YEO AND RA JOHNSON "A NEW FAMILY OF POWER TRANSFORMATIONS TO IMPROVE NORMALITY OR SYMMETRY" BIOMETRIKA 874 PP954959 2000

RF3E1504535DE2 GEP BOX AND DR COX "AN ANALYSIS OF TRANSFORMATIONS" JOURNAL OF THE ROYAL STATISTICAL SOCIETY B 26 211252 1964

1 IK YEO AND RA JOHNSON "A NEW FAMILY OF POWER TRANSFORMATIONS TO IMPROVE NORMALITY OR SYMMETRY" BIOMETRIKA 874 PP954959 2000

2 GEP BOX AND DR COX "AN ANALYSIS OF TRANSFORMATIONS" JOURNAL OF THE ROYAL STATISTICAL SOCIETY B 26 211252 1964

R0FECF191E4B81 PING LI T HASTIE AND K W CHURCH 2006 "VERY SPARSE RANDOM PROJECTIONS" [HTTPS://WEBSTANFORD.EDU/HASTIE/PAPERS/PING\\_KDD06\\_RPPDF](https://web.stanford.edu/hastie/papers/ping_kdd06_rppdf)

R0FECF191E4B82 D ACHLIOPTAS 2001 "DATABASE-FRIENDLY RANDOM PROJECTIONS" [HTTPS://USERS/SOE/UCSC/EDU/OPTAS/PAPERS/JL.PDF](https://users.soe.ucsc.edu/optas/papers/jl.pdf)

1 [HTTPS://ENWIKIPEDIA.ORG/WIKI/JOHNSONE28093LINDENSTRAUSS/LEMMA](https://en.wikipedia.org/wiki/Johnstone_28093_Lindenstrauss_lemma)

2 SANJOY DASGUPTA AND ANUPAM GUPTA 1999 "AN ELEMENTARY PROOF OF THE JOHNSON-LINDENSTRAUSS LEMMA" [HTTP://CITSEER.IST.PSU.EDU/VIEWDOC/SUMMARY/DOI1011453654](http://citeseer.ist.psu.edu/viewdoc/summary?doi=10.1.1.453.654)

R20C70293EF721 LIBSVM: A LIBRARY FOR SUPPORT VECTOR MACHINES

R20C70293EF722 PLATT JOHN 1999 "PROBABILISTIC OUTPUTS FOR SUPPORT VECTOR MACHINES AND COMPARISON TO REGULARIZED LIKELIHOOD METHODS"

RB1EC977CD3071 [HTTPS://ENWIKIPEDIA.ORG/WIKI/DECISIONTREE/LEARNING](https://en.wikipedia.org/wiki/Decision_tree_learning)

RB1EC977CD3072 L BREIMAN J FRIEDMAN R OLSHEN AND C STONE "CLASSIFICATION AND REGRESSION TREES" WADSWORTH BELMONT CA 1984

RB1EC977CD3073 T HASTIE R TIBSHIRANI AND J FRIEDMAN "ELEMENTS OF STATISTICAL LEARNING" SPRINGER 2009

RB1EC977CD3074 L BREIMAN AND A CUTLER "RANDOM FORESTS" [HTTPS://WWW.STAT.BERKELEY.EDU/BREIMANRANDOMFORESTSCCHOME.HTM](https://www.stat.berkeley.edu/~breiman/randomforestscchome.htm)

RA37B7E3ADB191 [HTTPS://ENWIKIPEDIA.ORG/WIKI/DECISIONTREE/LEARNING](https://en.wikipedia.org/wiki/Decision_tree_learning)

RA37B7E3ADB192 L BREIMAN J FRIEDMAN R OLSHEN AND C STONE "CLASSIFICATION AND REGRESSION TREES" WADSWORTH BELMONT CA 1984

RA37B7E3ADB193 T HASTIE R TIBSHIRANI AND J FRIEDMAN "ELEMENTS OF STATISTICAL LEARNING" SPRINGER 2009

RA37B7E3ADB194 L BREIMAN AND A CUTLER "RANDOM FORESTS" [HTTPS://WWW.STAT.BERKELEY.EDU/BREIMANRANDOMFORESTSCCHOME.HTM](https://www.stat.berkeley.edu/~breiman/randomforestscchome.htm)

RDD99A0224C6E1 P GEURTS D ERNST AND L WEHENKEL "EXTREMELY RANDOMIZED TREES" MACHINE LEARNING 631 342 2006

R4939D63D5A491 P GEURTS D ERNST AND L WEHENKEL "EXTREMELY RANDOMIZED TREES" MACHINE LEARNING 631 342 2006

1 WIKIPEDIA ENTRY FOR THE JACCARD INDEX

2458 BIBLIOGRAPHY























































































