Contents

1	Anoop index:	2
2	experimental0	2
3	How Do I Get Started? 3.1 articles-read, worth rereading	2 2 3 3 3
1	https://machinelearningmastery.com/python-machine-learning- 4.1 pre requisite https://machinelearningmastery.com/gentle-in- 4.2 1 - install	ntroduction-to-the-bias-v 5 5 5 5 7
	 4.4 3 - Load csv 4.5 4 - use pandas.DataFrame helper functions to describe data with statistics 4.6 5 - Basic Data Visualization 4.7 6 - Preprocessing data 4.7.1 skip for now, come back later 4.8 7 - Resampling 4.9 Google Developers intro to ml 4.9.1 1 hello apple or orange using decision tree 4.9.2 2 Decision Tree visualization 	14 15 16 16 17 18 18 18
	4.10 7 - Resampling	19 validation/ 19 20 22

4.11.1 Classification Metrics
4.11.2 Regression Metrics
4.12 9 - Spot checking
4.12.1 https://machinelearningmastery.com/spot-check-classification-machine-learningmastery.com/spot-c
4.12.2 https://machinelearningmastery.com/spot-check-regression-machine-learningmastery.com/spot
4.12.3 How to Develop a Reusable Framework to Spot-Check
Algorithms in Python:
4.13 10 - Model selection by comparison between models 27
4.14 11 - Improve accuracy by tuning the selected algorithm 30

1 Anoop index:

How Do I Get Started? Applied Machine Learning Process Linear Algebra Statistical Methods Understand Machine Learning Algorithms Python Machine Learning (scikit-learn) Code Algorithm from Scratch (Python) Introduction to Time Series Forecasting (Python) XGBoost in Python (Stochastic Gradient Boosting) Deep Learning (Keras) Long Short-Term Memory (LSTM) Deep Learning for Natural Language Processing (NLP) Deep Learning for Time Series Forecasting

2 experimental0

chrome history suggestions based on a neural network read to get ideas for projects: https://machinelearningmastery.com/self-study-machine-learning-projects/maintain a ml daily blog? how to get computers to talk to each other using socket programming use free time with phone to study google maps of Bangalore like guttapalli buy toothbrush, bigger bulb for room sabji kitne ki hai app

3 How Do I Get Started?

3.1 articles-read, worth rereading

Applied Machine Learning Process

Intermediate: Python Ecosystem.

\$\$ Step 4: Practice on Datasets. Select datasets to work on and practice

the process.

Practice Machine Learning with Small In-Memory Datasets

Tour of Real-World Machine Learning Problems \$\$ Work on Machine Learning Problems That Matter To You

\$\$ Step 5: Build a Portfolio. Gather results and demonstrate your skills.

Builda Machine Learning Port folio

Get Paid To Apply Machine Learning \$\$ Machine Learning For Money For more on this top-down approach, see:

The Machine Learning Mastery Method

Machine Learning for Programmers

Many of my students have used this approach to go on and do well in Kaggle competitions and get jobs as Machine Learning Engineers and Data Scientists.

3.2 notes

3.2.1 intro

sckit-learn is higher level than numpy & scipy machine learning is a subset of artificial intelligence artificial learning is a more consistent name for machine learning

3.2.2 some key word definitions:

Model: A machine learning model can be a mathematical representation of a real-world process. To generate a machine learning model you will need to provide training data to a machine learning algorithm to learn from.

Algorithm: Machine Learning algorithm is the hypothesis set that is taken at the beginning before the training starts with real-world data. When we say Linear Regression algorithm, it means a set of functions that define similar characteristics as defined by Linear Regression and from those set of functions we will choose one function that fits the most by the training data.

Training: While training for machine learning, you pass an algorithm with training data. The learning algorithm finds patterns in the training data such that the input parameters correspond to the target. The output of the training process is a machine learning model which you can then use to make predictions. This process is also called learning.

Regression: Regression techniques are used when the output is real-valued based on continuous variables. For example, any time series data. This technique involves fitting a line.

Classification: In classification, you will need to categorize data into predefined classes. For example, an email can either be spam or not spam.

Target: The target is whatever the output of the input variables. It could be the individual classes that the input variables maybe mapped to in case of a classification problem or the output value range in a regression problem. If the training set is considered then the target is the training output values that will be considered.

Feature: Features are individual independent variables that act as the input in your system. Prediction models use features to make predictions. New features can also be obtained from old features using a method known as feature engineering. More simply, you can consider one column of your data set to be one feature. Sometimes these are also called attributes. And the number of features are called dimensions.

Label: Labels are the final output. You can also consider the output classes to be the labels. When data scientists speak of labeled data, they mean groups of samples that have been tagged to one or more labels.

Overfitting: An important consideration in machine learning is how well the approximation of the target function that has been trained using training data, generalizes to new data. Generalization works best if the signal or the sample that is used as the training data has a high signal to noise ratio. If that is not the case, generalization would be poor and we will not get good predictions. A model is overfitting if it fits the training data too well and there is a poor generalization of new data.

Regularization: Regularization is the method to estimate a preferred complexity of the machine learning model so that the model generalizes and the over-fit/under-fit problem is avoided. This is done by adding a penalty on the different parameters of the model thereby reducing the freedom of the model.

Parameter and Hyper-Parameter: Parameters are configuration variables that can be thought to be internal to the model as they can be estimated from the training data. Algorithms have mechanisms to optimize parameters. On the other hand, hyperparameters cannot be estimated from the training data. Hyperparameters of a model are set and tuned depending on a combination of some heuristics and the experience and domain knowledge of the data scientist.

4 https://machinelearningmastery.com/python-machine-learning-mini-

14Lessons

 $1.1 \quad ext{pre requisite https://machinelearningmastery.com/gentle-introduction-to-the-presentation}$

Bias Error Bias are the simplifying assumptions made by a model to make the target function easier to learn. Low Bias: Suggests less assumptions about the form of the target function. High-Bias: Suggests more assumptions about the form of the target function.

Examples of low-bias machine learning algorithms include: Decision Trees, k-Nearest Neighbors and Support Vector Machines. Examples of high-bias machine learning algorithms include: Linear Regression, Linear Discriminant Analysis and Logistic Regression.

Variance Error Variance is the amount that the estimate of the target function will change if different training data was used.

Examples of low-variance machine learning algorithms include: Linear Regression, Linear Discriminant Analysis and Logistic Regression. Examples of high-variance machine learning algorithms include: Decision Trees, k-Nearest Neighbors and Support Vector Machines.

 $Fig~1.~bulls-eye~visualise~ \verb|http://scott.fortmann-roe.com/docs/BiasVariance.| html|$

skill of the model, a score with a high variance = (that may change a lot based on the data used to fit the model), or a high bias = (such as an overestimate of the skill of the model).

4.2 1 - install

4.2.1 next time try using this tutorial https://sourabhbajaj.com/mac-setup/Python/numpy.html

4.2.2 make a new virtualenv

pwd

Use :session: property to speed up org-babel when possible to use the same session

source ~/.bashrc
#mkvirtualenv mlm
workon mlm
which python

```
(pyvenv-workon "mlm2")
```

4.2.3 install https://stackoverflow.com/questions/26319762/how-to-install-scipy-stack-w

pip install numpy brew install gcc pip install scipy brew install freetype pip install matplotlib pip install nose pip install pandas pip install sympy pip install ipython[all] brew install pyqt brew install qt brew install sip #after this edit the 2 scripts

4.2.4 check if properly isntalled using .__version _ after import

using python snippets inside orgmode https://orgmode.org/worg/org-contrib/babel/languages/ob-doc-python.html installed python-mode from package-list-packages for emacs

```
(pyvenv-workon "mlm")
```

Fixed some errors using: pip install nose pyparsing python-dateutil pep8

```
import sys
import scipy
import numpy
import matplotlib
import pandas
import sklearn

print(f'python: {sys.version}')
print(f'scipy: {scipy.__version__}')
print(f'numpy: {numpy.__version__}')
print(f'matplotlib {matplotlib.__version__}')
print(f'pandas {pandas.__version__}')
print(f'sklearn {sklearn.__version__}')
```

4.3 2 - python, pandas, numpy, mathplotlib veeery basics

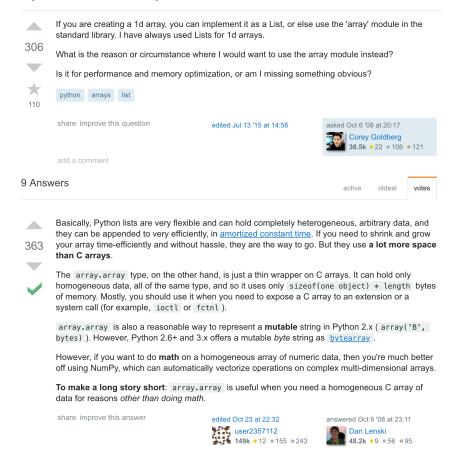
4.3.1 python, also refer in-y-minutes file for future reference

```
if 1>2:
    print("wtf")
else:
    print("ok")
```

```
try:
    # Use "raise" to raise an error
    raise IndexError("This is an index error")
except IndexError as e:
    print("its indexerror")
                         # Pass is just a no-op. Usually you would do recovery here.
except (TypeError, NameError):
    print("its typeerror or nameerror")
                         # Multiple exceptions can be handled together, if required.
else:
                         # Optional clause to the try/except block. Must follow all ex
    print("All good!") # Runs only if the code in try raises no exceptions
                         # Execute under all circumstances
finally:
    print("We can clean up resources here")
def accepts_variable_number_of_arguments(*args):
    print(type(args))
    print(args)
accepts_variable_number_of_arguments(1,2,3)
def accepts_variable_number_of_keyword_arguments(**kwargs):
    print(type(kwargs))
    print(kwargs)
accepts_variable_number_of_keyword_arguments(name="anoop", work="code")
def accepts_both_args_and_kwargs(*args, **kwargs):
    print(args)
    print("----")
    print(kwargs)
accepts_both_args_and_kwargs(1,2,3,a="4",b="5",last0="6")
```

4.3.2 numpy basics

Python List vs. Array - when to use?



numpy official tutorial https://docs.scipy.org/doc/numpy-1.15.
0/user/quickstart.html

```
import numpy as np
a = np.arange(15).reshape(3,5)
print(a)
a.shape
a.ndim
a.dtype.name
#dir(a)
a.size
```

```
type(a)
 b = np.array([6,7,8])
 type(b)
 np.zeros([2,3])
 np.arange(15)
 np.linspace(0,9, 19)
 from numpy import pi
 x = np.linspace(0, 2*pi, 5)
 np.sin(x)
  #2 decimal places
 np.around(np.sin(x), decimals=2)
 A = np.array([[1,2],[3,4]])
  I = np.array([[1,0],[0,1]])
 elementwise = A * I
 matrix_product = A @ I
 print(elementwise, "\n", matrix_product)
 a = np.ones(3, dtype=np.int32)
 b = np.linspace(0, 1, 3)
 c = a + b
 print(a,b,c)
 c.dtype.name
 c*1j
 np.ones(1)
 my_e = np.exp(np.ones(1))
 from numpy import e
 print(e, my_e)
 d = np.exp(c*1j)
 d.dtype
  \# exp, \sin etc are called numpy universal functions
  # multidimensional array
c = np.array([
   Γ
^^I[ 0, 1, 2],
```

```
^^I[ 10, 12, 13]
    ],
    Г
^^I[100,101,102],
^^I[110,112,113]
])
c.shape
# Visualize0 2,2,3 as you traverse from the topmost bracket to the inner ones
for i in c.flat:
    print(i, end=" // ")
print("\n")
id(c) #id is unique identifier of an object in python
   axis in numpy
   file: /ml/flipshope/screenshots0/Screenshot 2018-12-11 at 6.06.43 PM.png
   matplotlib python is not installed as a framework error, solution: https:
//stackoverflow.com/questions/34977388/matplotlib-runtimeerror-python-is-not-installed
   Above is a hacky solution I need to switch away from virtualenv & vir-
tualenvwrapper and move to venv entirely Also, reddit recommends to avoid
the pyenv or any wrapper around venv && strongly recommends to use venv
directly venv ships by default with python >= 3.3 https://matplotlib.
org/faq/osx_framework.html https://news.ycombinator.com/item?id=
18612590 https://news.ycombinator.com/item?id=18247512:) andybak
54 days ago [-] In case this scares any new users, I've used nothing more than
pip and virtualenv for several years with no issues of note.
import numpy as np
import matplotlib.pyplot as plt
#import matplotlib
#matplotlib.use('TkAqq')
#import matplotlib.pyplot as plt
def mandelbrot( h,w, maxit=20 ):
    """Returns an image of the Mandelbrot fractal of size (h,w)."""
    y,x = np.ogrid[-1.4:1.4:h*1j, -2:0.8:w*1j]
    c = x+y*1j
    z = c
    divtime = maxit + np.zeros(z.shape, dtype=int)
```

```
for i in range(maxit):
^{1}z = z**2 + c
^{\text{ldiverge}} = z*np.conj(z) > 2**2
                                               # who is diverging
^^Idiv_now = diverge & (divtime==maxit) # who is diverging now
^^Idivtime[div_now] = i
                                            # note when
^{\text{lz[diverge]}} = 2
                                            # avoid diverging too much
    return divtime
plt.imshow(mandelbrot(400,400))
plt.show()
import sys
print(sys.path)
   real-python tutorial https://realpython.com/numpy-array-programming/
   When it comes to computation, there are really three concepts that lend
NumPy its power: Vectorization Broadcasting Indexing
import numpy as np
arr = np.arange(36).reshape(3,4,3)
arr
11 11 11
visualize0:-
00 01 02 03 04 05 06 07 08 09 10 11
12 13 14 15 16 17 18 19 20 21 22 23
24 25 26 27 28 29 30 31 32 33 34 35
00 01 02
03 04 05
06 07 08
09 10 11
[#3 items
[#4 items
IJ
[]
[]
```

```
[]
J...
]
11 11 11
a = np.array([2,3,4])
b = 2
b_broadcasted = np.broadcast(a,b)
print(list(b_broadcasted))
# rest of this tutorial seemed a bit advanced, skip for now, come back later
   Note Pandas is a library built on top of NumPy
   todo: switch to jupyter notebook instead of emacs: https://github.
com/millejoh/emacs-ipython-notebook http://millejoh.github.io/emacs-ipython-notebook/
https://www.youtube.com/watch?v=wtVF5cMhBjghttps://news.ycombinator.
com/item?id=9728143 https://github.com/gregsexton/ob-ipython
4.3.3 matplotlib basics
matplotlib, is written in pure Python and is heavily dependent on NumPy
   Matplotlib is conceptually divided into three parts: pylab interface (sim-
ilar to MATLAB) pylab tutorial, this is the matplotlib.pyplot import Mat-
plotlib frontend or API artist tutorial backends drawing devices or renderers
   Lets learn pyplot https://matplotlib.org/users/pyplot_tutorial.
html#pyplot-tutorial
   venv basics to switch from virtualenv
   python3 -m venv mlm2 source mlm2/bin/activate
   pyvenv-deactivate
import matplotlib
#matplotlib.use('MacOSX')
import matplotlib.pyplot as plt
plt.plot([0, 1, 4, 9, 16])
#plt.show()
plt.plot([0, 1, 4, 9, 16], 'ro')
#plt.show()
plt.plot([0.5], 'y.')
plt.show()
```

4.3.4 Pandas -basics:

print(p.shape)

```
Todo:Official beginner tutorial: https://pandas.pydata.org/pandas-docs/
stable/10min.html Todo: Intermediate - Julia Evans - https://jvns.
ca/blog/2013/12/22/cooking-with-pandas/ "take a real dataset or three,
play around with it, and learn how to use pandas along the way."
   Panda Series & DataFrames: https://medium.freecodecamp.org/series-and-dataframe-in-pyt
#Series and DataFrames
import pandas as pd
x1 = pd.Series([6,3,4,6])
x = pd.Series([6,3,4,6], index=['a','b','c','d'])
y = pd.Series(3, index=['a', 'b', 'c', 'd'])
#DataFrames
import numpy as np
dates = pd.date_range('20181201', periods = 8)
my_narray = np.random.randn(8,3)
list('ABC')
df = pd.DataFrame(index = dates, data = my_narray, columns = ['A','B','C'])
df_absolute = df.apply(abs)
   So pandas is kinda like an excel sheet0
import numpy as np
np.arange(4)
ma = np.arange(4).reshape((2,2))
import pandas
p = pandas.DataFrame(ma)
print(ma[1,0])
print(p[1])
print(p[1][0] == ma[1][0])
```

Todo:https://pandas.pydata.org/pandas-docs/stable/dsintro.html#dsintro

4.4 3 - Load csv

```
Datasets
   Work with csv using python's csv module
import csv
with open('iris.csv') as csv_file:
    for line in csv_file:
^^Iprint(line)
^{\hat{I}pass}
    11 11 11
    csv_reader = csv.reader(csv_file)
    for line in csv_reader:
^^I#print(line)
^^Ipass
my_fieldnames = ("sepal_length", "sepal_width", "petal_length", "petal_width", "class"
with open('iris.csv') as csv_file:
    csv_dict_reader = csv.DictReader(csv_file, fieldnames=my_fieldnames)
    for line in csv_dict_reader:
^^I#print(line)
^^Ipass
with open('test_writeout.csv', mode='w') as out_file:
    csv_writer = csv.writer(out_file)
    csv_writer.writerow(["row", "1"])
    csv_writer.writerow(["row", "2"])
    csv_dict_writer = csv.DictWriter(out_file, fieldnames = my_fieldnames)
    csv_dict_writer.writeheader()
    csv_dict_writer.writerow({"sepal_length": 1, "sepal_width": 2, "petal_length": 3,
   Work with csv using numpy
import numpy as np
with open("numpy_loadtxt_input.txt") as input_file:
```

https://realpython.com/python-csv/https://github.com/jbrownlee/

```
for line in input_file:
^^Iprint(line)
    HHHH
    my_nparray = np.loadtxt(input_file, delimiter=" ")
    print(my_nparray)
    print(my_nparray.dtype)
   Work with csv usign pandas
import pandas
df = pandas.read_csv('pandas_read_csv.csv')
print(df)
df2 = pandas.read_csv('pandas_read_csv.csv', parse_dates=['Hire Date'], index_col='Nam
print(df2)
my_col_names = ("Name", "Hired_on", "Salary", "sick_days_remaining")
df3 = pandas.read_csv('pandas_read_csv.csv', header=None, names=my_col_names,
parse_dates=['Hired_on'])
print(df3)
df3.to_csv('pandas_to_csv.csv')
4.5 4 - use pandas.DataFrame helper functions to describe
     data with statistics
# Scatter Plot Matrix
import matplotlib.pyplot as plt
import pandas
my_col_names = ("num_preg", "plasma_glucose", "blood_pressure", "triceps_skin_thicknes
df = pandas.read_csv("https://raw.githubusercontent.com/jbrownlee/Datasets/master/pima
#print(df)
from pandas.plotting import scatter_matrix
scatter_matrix(df)
plt.show()
df.corr()
```

4.6 5 - Basic Data Visualization

Using pandas and matplotlib together

```
# Scatter Plot Matrix
import matplotlib.pyplot as plt
import pandas
my_fieldnames = ("sepal_length", "sepal_width", "petal_length", "petal_width", "class"
data = pandas.read_csv('iris.csv', names=my_fieldnames)
#print(data)
se = data.loc[data["class"] == "Iris-setosa"]
ve = data.loc[data["class"] == "Iris-versicolor"]
vi = data.loc[data["class"] == "Iris-virginica"]
#se.hist()
#ve.hist()
#vi.hist()
se.plot(kind="box")
from pandas.plotting import scatter_matrix
#scatter_matrix(se)
plt.savefig("img/lesson5.png")
return "img/lesson5.png"
     6 - Preprocessing data
Standardize numerical data (e.g. mean of 0 and standard deviation of 1)
using the scale and center options.
   Simple example:
import numpy
narray = numpy.arange(0,4).reshape(2,2)
import pandas
df = pandas.DataFrame(narray, columns=("c1", "c2"))
from sklearn.preprocessing import StandardScaler
```

```
scaler = StandardScaler()
scaler.fit(df) #calculates and store mean & standard-deviation for each column
print(scaler.transform(df))
type(scaler.transform(df))
df_standardized = pandas.DataFrame(scaler.transform(df))
   Complex example:
import matplotlib.pyplot as plt
import pandas
#df is my short for DataFrame
my_col_names = ("num_preg", "plasma_glucose", "blood_pressure", "triceps_skin_thicknes
df = pandas.read_csv("https://raw.githubusercontent.com/jbrownlee/Datasets/master/pima
#print(df)
import numpy
array = df.values
X = array[:, :-1]
Y = array[:, -1:]
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler().fit(X)
X_standardized = scaler.transform(X)
#printout
numpy.set_printoptions(precision=2)
print(X_standardized)
df_standardized = pandas.DataFrame(X_standardized)
df_standardized.describe() #you can see that mean=0, sdev=1 for each column
```

4.7.1 skip for now, come back later

Normalize numerical data (e.g. to a range of 0-1) using the range option. Explore more advanced feature engineering such as Binarizing.

4.8 7 - Resampling

statistical methods called resampling methods are used to split your training dataset up into subsets, some are used to train the model and others are held back and used to estimate the accuracy of the model on unseen data.

4.9 Google Developers intro to ml

complete google developers course, its a pre-requisite for this lesson as per me https://www.youtube.com/playlist?list=PLOU2XLYxmsIIuiBfYad6rFYQU_jL2ryal

4.9.1 1 hello apple or orange using decision tree

```
from sklearn import tree
features_original = [[140, "smooth"], [130, "smooth"], [150, "bumpy"], [170, "bumpy"]]
labels_original = ["apple", "apple", "orange", "orange"]

#smooth=1 // bumpy=0
#orange=1 // apple=0
features = [[140, 1], [130, 1], [150, 0], [170, 0]] #weight, texture of fruit
labels = [0, 0, 1, 1]

#train a classifier:
clf = tree.DecisionTreeClassifier() #instantiate an empty box of rules
clf = clf.fit(features, labels) #learning algorithm fills the above box with rules

#print(clf.predict([[150,0]]))
print(clf.predict([[200,0]]))
```

4.9.2 2 Decision Tree visualization

```
from sklearn.datasets import load_iris
from sklearn import tree
iris = load_iris()
print(dir(iris))

print(iris.feature_names) #in this example data_names is more suitable
print(iris.data[0])

print(iris.target_names)
```

```
print(iris.target)
#Resampling:
test_ids = [0,50,100]
import numpy as np
#training data
train_target = np.delete(iris.target, test_ids)
train_data = np.delete(iris.data, test_ids, axis=0)
#testing data
test_target = iris.target[test_ids]
test_data = iris.data[test_ids]
clf = tree.DecisionTreeClassifier()
clf.fit(train_data, train_target)
predicted_target = clf.predict(test_data)
print(f"Reality: {test_data} features is {test_target}")
print(f"Prediction: {test_data} has been predicted as {predicted_target}")
#Visualize: https://medium.com/@rnbrown/creating-and-visualizing-decision-trees-with-p
from sklearn.externals.six import StringIO
from sklearn.tree import export_graphviz
import pydotplus
dot_data = StringIO() #StringIO() behaves like a file
export_graphviz(clf, out_file=dot_data, feature_names=iris.feature_names, class_names=
graph = pydotplus.graph_from_dot_data(dot_data.getvalue())
graph.write_png("img/iris.png")
return "img/iris.png"
4.10
     7 - Resampling
```

4.10.1 https://machinelearningmastery.com/k-fold-cross-validation/

Shuffle the dataset randomly. Split the dataset into k groups For each unique group: Take the group as a hold out or test data set Take the remaining groups as a training data set Fit a model on the training set and evaluate it on the test set Retain the evaluation score and discard the model Summarize the skill of the model using the sample of model evaluation scores

k-fold cross validation leave one out cross validation = n-fold cross validation

Train/Test Split: = 1-fold cross validation (Taken to one extreme, k may be set to 1 such that a single train/test split is created to evaluate the model.)

To summarize, there is a bias-variance trade-off associated with the choice of k in k-fold cross-validation. Typically, given these considerations, one performs k-fold cross-validation using k=5 or k=10, as these values have been shown empirically to yield test error rate estimates that suffer neither from excessively high bias nor from very high variance.

KFold() scikit-learn class can be used

```
from sklearn.model_selection import KFold
import numpy

array = [0.1, 0.2, 0.3, 0.4, 0.5, 0.6]
na = numpy.array(array)

my_kfold = KFold(n_splits=3, shuffle=True, random_state=6)

#split method generates the indices
for train, test in my_kfold.split(na):
    print("indices for train and test respectively", train, test)
    print("actual train and test data respectively", na[train], na[test])
    print()
```

4.10.2 Example = apply 10-fold to diabetes data

print(array)

use scikit-learn to estimate the accuracy of the Logistic Regression algorithm on the Pima Indians onset of diabetes dataset using 10-fold cross validation.

Emacs help: file: /ml/flipshope/screenshots0/Screenshot 2018-12-13 at 12.57.01 PM.png I used decision tree to fit this data

```
import pandas as pd
my_names = ['preg', 'plas', 'pres', 'skin', 'test', 'mass', 'pedi', 'age', 'class']
df = pd.read_csv("https://raw.githubusercontent.com/jbrownlee/Datasets/master/pima-ind
array = df.values
```

```
features = array[:,:-1]
is_diabetic = array[:,-1:]
from sklearn.model_selection import cross_val_score
from sklearn.model_selection import KFold
my_kfold = KFold(n_splits=10, random_state=6)
#my_kfold.split(features)
11 11 11
X = []
y = []
clf = tree.DecisionTreeClassifier()
clf.fit(X, y)
print("prediction:", clf.predict(X_test), " // actual:", y_test)
from sklearn import tree
clf = tree.DecisionTreeClassifier()
results = cross_val_score(clf, features, is_diabetic, cv=my_kfold)
print(f"Mean accuracy: {results.mean()}, Sdev: {results.std()}")
Example given on website use LogisticRegression
# Evaluate using Cross Validation
from pandas import read_csv
from sklearn.model_selection import KFold
from sklearn.model_selection import cross_val_score
from sklearn.linear_model import LogisticRegression
url = "https://raw.githubusercontent.com/jbrownlee/Datasets/master/pima-indians-diabete
names = ['preg', 'plas', 'pres', 'skin', 'test', 'mass', 'pedi', 'age', 'class']
dataframe = read_csv(url, names=names)
array = dataframe.values
X = array[:,0:8]
Y = array[:,8]
kfold = KFold(n_splits=10, random_state=7)
model = LogisticRegression()
results = cross_val_score(model, X, Y, cv=kfold)
print(f"Mean accuracy: {results.mean()}, Sdev: {results.std()}")
```

4.11 8 - Algorithm evaluation metrics

stopped this lesson in between and skipped for now https://machinelearningmastery.com/metrics-evaluate-machine-learning-algorithms-python/

Contents of this lesson: You will learn about 3 classification metrics: Accuracy. Logarithmic Loss. Area Under ROC Curve.

Also 2 convenience methods for classification prediction results: Confusion Matrix. Classification Report.

And 3 regression metrics: Mean Absolute Error. Mean Squared Error. \mathbb{R}^2 .

Classification metrics -for problems like diabetic or not (binary classification 0 or 1 is 'class')

Regression metrics -for problems like boston house pricing (continuous price metric is 'class')

All recipes evaluate the same algorithms, Logistic Regression for classification and Linear Regression for the regression problems

4.11.1 Classification Metrics

Apply Logistic Regression to diabetes problem and watch how each algo evaluation recipe works keeping the algo constant

Recipes:

1. Classification Accuracy.

from sklearn.model_selection import KFold

```
my_kfold = KFold(n_splits=10, random_state=6)

from sklearn.linear_model import LogisticRegression
my_estimator_model = LogisticRegression()

# scoring="name_of_algorithm_evaluation_metric0" refer to https://scikit-learn.org
results = cross_val_score(my_estimator_model, X, y, cv=my_kfold, scoring="accuracy
print(f"mean accuracy of diabetes prediction mu = {results.mean()}, sdev = {results.mean()}
```

2. Logarithmic Loss. Some evaluation metrics (like mean squared error) are naturally descending scores (the smallest score is best) and as such are reported as negative by the cross_{valscore}() function. This is important to note, because some scores will be reported as negative that by definition can never be negative.

from sklearn.model_selection import cross_val_score

The theory of this I havent yet understood, skip for now http://wiki.fast.ai/index.php/Log_Loss

```
scoring = "neg_log_loss"
#use above value in model_selection.cross_val_score()
```

- 3. Area Under ROC Curve.
- 4. Confusion Matrix.
- 5. Classification Report. Making use of convenience feature in sklearn from sklearn.metrics import classification_report

4.11.2 Regression Metrics

Mean Absolute Error. Mean Squared Error. R².

4.12 9 - Spot checking

4.12.1 https://machinelearningmastery.com/spot-check-classification-machine-learning-

We will spot check 6 classification algorithms

2 Linear Machine Learning Algorithms: Logistic Regression, lr ;;;; NOTE: ALTHOUGH ITS A CLASSIFICATION PROBLEM, THIS MODEL STILL HAS *'Regression'* IN ITS NAME Linear Discriminant Analysis, lda

4 Nonlinear Machine Learning Algorithms: K-Nearest Neighbors, knn Naive Bayes, nb Decision Trees (or Classification and Regression Trees, CART), in this case we will use classification tree Support Vector Machines, sym

Lets try out on pima-indians-diabetes Lets suppose that mean() accuracy on 10-fold cross validation represents the performance of each algorithm

DOUBT A note about sklearn.model_{selection.crossvalscore}() function: If no scoring is specified, the estimator(classifier) passed should have a 'score' method https://github.com/scikit-learn/scikit-learn/blob/14031f6/sklearn/model_selection/_validation.py#L36

"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."

 $\label{lem:https://stackoverflow.com/questions/42825714/what-is-the-score-function-formula-of-skleptons and the state of the state of$

```
import pandas
feature_names = ["pregnant", "plasma_glucose", "blood_pressure",
^^I^^I "skin_fold_thickness", "serum_insulin", "bmi", "pedigree", "age", "class"]
df = pandas.read_csv("https://raw.githubusercontent.com/jbrownlee/Datasets/master/pima
array = df.values
features = array[:, :-1]
is_diabetic = array[:, -1]
from sklearn.model_selection import KFold
from sklearn.model_selection import cross_val_score
from sklearn.linear_model import LogisticRegression
my_kfold = KFold(n_splits=10, random_state=6)
classifier_lr = LogisticRegression()
results_logistic_regression = cross_val_score(classifier_lr, features, is_diabetic, cv
print("lr: ", results_logistic_regression.mean(), results_logistic_regression.std())
from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
classifier_lda = LinearDiscriminantAnalysis()
results_linear_discriminant_analysis = cross_val_score(classifier_lda, features, is_dia
```

print("lda: ", results_linear_discriminant_analysis.mean(), results_linear_discriminant

```
from sklearn.neighbors import KNeighborsClassifier
classifier_knn = KNeighborsClassifier()
results_knn = cross_val_score(classifier_knn, features, is_diabetic, cv=my_kfold)
print("knn: ", results_knn.mean(), results_knn.std())
from sklearn.naive_bayes import GaussianNB
classifier_nb = GaussianNB()
results_nb = cross_val_score(classifier_nb, features, is_diabetic, cv=my_kfold)
print("nb: ", results_nb.mean(), results_nb.std())
from sklearn.tree import DecisionTreeClassifier
classifier_dt = DecisionTreeClassifier()
results_dt = cross_val_score(classifier_dt, features, is_diabetic, cv=my_kfold)
print("dt: ", results_dt.mean(), results_dt.std())
from sklearn.svm import SVC
classifier_svc = SVC()
results_svc = cross_val_score(classifier_svc, features, is_diabetic, cv=my_kfold)
print("svc: ", results_svc.mean(), results_svc.std())
   Note:python submodule doubt https://stackoverflow.com/questions/
12229580/python-importing-a-sub-package-or-sub-module Try it out
yourself
import sklearn
dir(sklearn)
import sklearn.discriminant_analysis
dir(sklearn) #now more things are shown, why?
```

4.12.2 https://machinelearningmastery.com/spot-check-regression-machine-learning-algorithms.

We will spot check 7 regression algorithms

- **4 Linear Machine Learning Algorithms:** Linear Regression Ridge Regression LASSO Linear Regression Elastic Net Regression
- **3 Nonlinear Machine Learning Algorithms:** K-Nearest Neighbors Classification and Regression Trees, in this case we will use regression tree Support Vector Machines

Lets try out on boston-house-pricing Lets suppose that the negative mean squared error measures on 10-fold cross validation represents the performance of each algorithm

```
import pandas
data_names = ("crim", "zn", "indus", "chas", "nox", "rm", "age", "dis", "rad", "tax",
df = pandas.read_csv("https://raw.githubusercontent.com/anoopemacs/Datasets/master/hous
#chas is binary, rest are all continuous features
array = df.values
features = array[:, :-1] #X
house_value = array[:, -1] #y
from sklearn.model_selection import KFold
my_kfold = KFold(n_splits=10, random_state=6)
from sklearn.model_selection import cross_val_score
#LINEAR REGRESSION MODELS x 4
from sklearn.linear_model import LinearRegression
estimator_lr = LinearRegression()
results_lr = cross_val_score(estimator_lr, features, house_value, scoring="neg_mean_sq"
print("lr: ", results_lr.mean())
from sklearn.linear_model import Ridge
estimator_rr = Ridge()
results_rr = cross_val_score(estimator_rr, features, house_value, scoring="neg_mean_sq"
print("rr: ", results_rr.mean())
from sklearn.linear_model import Lasso
estimator_lasso = Lasso()
results_lasso = cross_val_score(estimator_lasso, features, house_value, scoring="neg_m
print("lasso: ", results_lasso.mean())
from sklearn.linear_model import ElasticNet
estimator_elasticnet = ElasticNet()
results_elasticnet = cross_val_score(estimator_elasticnet, features, house_value, score
print("elasticnet: ", results_elasticnet.mean())
#NON LINEAR MODELS x 3
from sklearn.neighbors import KNeighborsRegressor
estimator_knn = KNeighborsRegressor()
results_knn = cross_val_score(estimator_knn, features, house_value, scoring="neg_mean_
```

```
print("knn: ", results_knn.mean())

from sklearn.tree import DecisionTreeRegressor
    estimator_cart = DecisionTreeRegressor()
    results_cart = cross_val_score(estimator_cart, features, house_value, scoring="neg_mean
    print("cart: ", results_cart.mean())

from sklearn.svm import SVR
    estimator_svr = SVR()
    results_svr = cross_val_score(estimator_svr, features, house_value, scoring="neg_mean_print("svr: ", results_svr.mean())
```

4.12.3 How to Develop a Reusable Framework to Spot-Check Algorithms in Python:

https://machinelearningmastery.com/spot-check-machine-learning-algorithms-in-python/ This is an important lesson, indirectly I will also learn how to write/maintain my own python module as well But skip for now as I want to come back to topics like this in a second revision

4.13 10 - Model selection by comparison between models

Advice to beginners from JSAT java ml package author edward raff: http://jsatml.blogspot.com/2014/10/beginner-advice-on-learning-to.html His book recommendations:

The Elements of Statistical Learning This is a very common book for people to learn and get, especially since it is free. In my opinion, the book isn't particularly great in any area but its not particularly bad at anything either. Though if you are more of a visual learner this book has lots of graphs and diagrams to try and give intuitions of how things work / what's going on. If your goal is to implement for learning, recreating their visual results would be an excellent exercise.

Machine learning: a Probabilistic Perspective This is currently my favorite ML book. Murphy does an excellent job explain the algorithms, and relating them to each other to help foster a deeper understanding. Though the learning curve is a bit wonky at times, and the later chapters dont have quite as much lower level details its overall excellent. For implementing in particular, much more refined (and explained) pseudo code is present for many of the algorithms and for many of the chapters and algorithms goes through the math needed to develop these algorithms. Occasionally some

implementation considerations are discussed. While not often, it is more than most books on Machine Learning.

https://machinelearningmastery.com/compare-machine-learning-algorithms-python-sciking-python-py

Lets compare the following on binary classification of pima-indians-diabetes

dataset LR: Logistic Regression LDA: Linear Discriminant Analysis KNN:

K-Nearest Neighbors CART: Classification and Regression Trees NB: Naive

Bayes SVM: Support Vector Machines

NOTE: trick to use to find dir(sklearn) doesnt give all submodules, use linux ls: ls Users/anoop/ml/flipshope/mlm2/lib/python3.7/site-packages/sklearn

```
#hide warnings
def warn(*args, **kwargs):
    pass
import warnings
warnings.warn = warn
import pandas
data_names = ["pregnant", "plasma_glucose", "blood_pressure",
         "skin_fold_thickness", "serum_insulin", "bmi", "pedigree", "age", "class"]
df = pandas.read_csv("https://raw.githubusercontent.com/jbrownlee/Datasets/master/pima-
print(df.shape)
array = df.values
features = array[:, :-1]
is_diabetic = array[:, -1]
#print(features)
#print(is_diabetic)
from sklearn.model_selection import KFold
my_kfold = KFold(n_splits=10, random_state=6)
from sklearn.model_selection import cross_val_score
models = []
from sklearn.linear_model import LogisticRegression
```

```
models.append(("LR", LogisticRegression()))
from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
models.append(("LDA", LinearDiscriminantAnalysis()))
from sklearn.neighbors import KNeighborsClassifier
models.append(("KNN", KNeighborsClassifier()))
from sklearn.tree import DecisionTreeClassifier
models.append(("CART", DecisionTreeClassifier()))
from sklearn.naive_bayes import GaussianNB
models.append(("NB", GaussianNB()))
from sklearn.svm import SVC
models.append(("SVM", SVC()))
results = []
model_names = []
for (name, model) in models:
    #print(name, model)
^^Iresult = cross_val_score(model, features, is_diabetic, cv=my_kfold, scoring='accura
^^I#print(name, ": ", result.mean(), " // ", result.std())
^^Iresults.append(result)
^^Imodel_names.append(name)
    except Exception:
^^Iprint(name)
^^Icontinue
print(model_names)
print(results)
import numpy
from matplotlib import pyplot
pyplot.boxplot(results, labels=model_names)
#pyplot.show()
pyplot.savefig("img/model-selection-by-comparison.png")
```

```
return "img/model-selection-by-comparison.png"
```

import pandas

From the above fig it seems like for this problem, at first glance it seems like LR and LDA sound promising for further evaluation

4.14 11 - Improve accuracy by tuning the selected algorithm

2 ways: Tune the parameters of an algorithm using a grid search that you specify. Tune the parameters of an algorithm using a random search.

```
data_names = ["pregnant", "plasma_glucose", "blood_pressure",
         "skin_fold_thickness", "serum_insulin", "bmi", "pedigree", "age", "class"]
df = pandas.read_csv("https://raw.githubusercontent.com/jbrownlee/Datasets/master/pima-
array = df.values
X = array[:, :-1]
y = array[:, -1]
from sklearn.linear_model import Ridge
classifier = Ridge()
from sklearn.model_selection import GridSearchCV
my_param_grid = {"alpha" : [1, 0.1, 0.01, 0.001, 0.0001, 0]}
grid = GridSearchCV(classifier, param_grid = my_param_grid)
grid.fit(X, y)
print(grid.best_score_)
print(grid.best_estimator_.alpha)
************************************
from sklearn.model_selection import RandomizedSearchCV
classifier2 = Ridge()
from scipy.stats import randint
my_param_dist = {"alpha": randint(0,10)} #doubt, can alpha be > 1, also wtf is scipy.s
random_search = RandomizedSearchCV(estimator=classifier2, param_distributions = my_param_distributions)
random_search.fit(X, y)
print(random_search.best_estimator_)
```