

LHC Package Documentation

January 25, 2021

<code>add_rbf_features</code>	<i>Add $n_centroids$ RBF features to X</i>
-------------------------------	---

Description

Add $n_centroids$ RBF features to X

Usage

```
add_rbf_features(X, s, n_centroids, Xi = NULL)
```

Arguments

<code>X</code>	covariate matrix
<code>s</code>	median pairwise distance of points in X
<code>n_centroids</code>	number of RBF features to add
<code>Xi</code>	matrix of reference points to calculate RBF w.r.t

Value

X augmented covariate matrix

AMS_data

AMS data object class

Description

This reference class object is used to store the AMS metric of a classification model at different decision thresholds. AMS is a performance measure which includes the sample weightings and is defined by the Higgs Boson Kaggle Competition. A vector of true sample classifications (0 or 1), a vector of estimated probabilities from a model, and a vector of scaled sample weights are needed to initialise.

Fields

`y` A vector of true sample classifications (0 or 1),
`prob` A vector of the samples estimated probabilities from a model
`weights` A vector of scaled sample weights.
`thresholds` A vector of 30 decision thresholds.
`ams` A vector of the AMS metric at each threshold.
`max_ams` maximum ams
`max_thresh` threshold of maximum ams

Methods

`calc_ams()` Calculate the AMS at each thresholds.
`initialize(y, prob, weights)` Provide true sample labels, estimated probabilities, and sample weights. A vector of decision thresholds is initialised.
`plot_ams()` Plot AMS against threshold.

ams_metric

Calculate AMS metric

Description

Calculate AMS metric

Usage

```
ams_metric(y, y_pred, weights, Ns = 691.9886, Nb = 410999.8)
```

ams_threshold	<i>Calculate optimal threshold</i>
---------------	------------------------------------

Description

Calculate optimal threshold

Usage

```
ams_threshold(amss)
```

Arguments

amss	a list of AMS_data objects
------	----------------------------

Value

optimal decision threshold considering all curves#

average_auc	<i>Calculate average AUC</i>
-------------	------------------------------

Description

Calculate average AUC

Usage

```
average_auc(rocs)
```

Arguments

rocs	a list of ROC_curve objects
------	-----------------------------

Value

average AUC

avg_median_pairwise_distance

Calculate average of median pairwise distances for between all adjacent points

Description

Calculate average of median pairwise distances for between all adjacent points

Usage

avg_median_pairwise_distance(X)

Arguments

X covariate matrix

Value

s median pairwise distance

backtrack_linesearch *Backtrack linesearch*

Description

Backtracking linesearch to find "optimal" step size

Usage

backtrack_linesearch(f, gradf, x, deltax, alpha, beta)

Arguments

f	function to minimising
gradf	gradient of f
x	parameter to optimising over
deltax	newton step
alpha	linesearch parameter
beta	linesearch update parameter

Value

stop to take

calc_K	<i>Compute the kernel matrix over the (training) set X</i>
--------	--

Description

Compute the kernel matrix over the (training) set X

Usage

```
calc_K(X, ckernel)
```

Arguments

X	covariate matrix (nxd)
ckernel	kernel function with args (x_i, x_j)

Value

K kernel matrix (nxn)

decide	<i>Thresholding function</i>
--------	------------------------------

Description

Thresholding function

Usage

```
decide(p, thresh = 0.5)
```

Arguments

p	vector of probabilities
thresh	threshold over which we assign output of 1

fit_l1_logreg	<i>Fit L1 constrained logistic regression</i>
---------------	---

Description

Fit L1 constrained logistic regression

Usage

```
fit_l1_logreg(X, y, C = 1)
```

generate_blues	<i>Generate a set of blues</i>
----------------	--------------------------------

Description

Generate a set of blues

Usage

generate_blues(n)

Arguments

n number of colours to generate

generate_colours	<i>Generate distinct colours</i>
------------------	----------------------------------

Description

Generate distinct colours

Usage

generate_colours(ncolours)

Arguments

ncolours number of colours we want

Value

vector of hex colours

get_const_features	<i>find colnames for columns that are constant (e.g. all 1, -999, NA etc)</i>
--------------------	---

Description

find colnames for columns that are constant (e.g. all 1, -999, NA etc)

Usage

```
get_const_features(X)
```

Arguments

X	matrix of covariates
---	----------------------

Value

list of column names

get_model_idx	<i>helper function to get the index corresponding to the model built on folds $l \neq k$ and for jet number j 1,2,3, ordering columns by fold and then with nesting on j i.e. first six cols are $k=1, j=1,2,3$; $k=2, j=1,2,3$; etc</i>
---------------	---

Description

helper function to get the index corresponding to the model built on folds $l \neq k$ and for jet number j 1,2,3, ordering columns by fold and then with nesting on j i.e. first six cols are $k=1, j=1,2,3$; $k=2, j=1,2,3$; etc

Usage

```
get_model_idx(j, k, K)
```

Arguments

j	this jet group
k	this fold
K	number of folds

get_rbf_centroids	<i>Get reference points for RBF centroids</i>
-------------------	---

Description

Get reference points for RBF centroids

Usage

```
get_rbf_centroids(X, n_centroids, idx = NULL)
```

Arguments

X	covariate matrix
n_centroids	number of RBF centroids
idx	[Optional] location of RBF centroid reference points

Value

list of Xi (centroid points) and idx (location of them)

get_subset_idx	<i>Select indices of rows in x which correspond to values of labels, which can be multiple elements</i>
----------------	---

Description

Select indices of rows in x which correspond to values of labels, which can be multiple elements

Usage

```
get_subset_idx(x, labels)
```

Arguments

x	vector of labels
labels	reference labels to compare to e.g. x <- c("t", "b", "t", "v", "b"); labels <- c("t", "b") get_subset_idx(x, labels) = T, T, T, F, T

get_valid_cols	<i>Gets the columns from header that aren't in features_to_rm</i>
----------------	---

Description

Gets the columns from header that aren't in features_to_rm

Usage

```
get_valid_cols(header, features_to_rm, j)
```

Arguments

header	list of column names e.g. names(X)
features_to_rm	list of feature names we want to remove
j	jet group

Value

index of columns to retain

idx_higgs_mass	<i>Get boolean index for rows with missing/or not missing (depending on G/j) Higgs mass</i>
----------------	---

Description

Get boolean index for rows with missing/or not missing (depending on G/j) Higgs mass

Usage

```
idx_higgs_mass(X, j, G)
```

Arguments

X	matrix of covariates
j	jet group
G	number of jet groups

Value

vector of bools

idx_jet_cat	<i>Get boolean vector of rows with j=0,1 or 2+</i>
-------------	--

Description

Get boolean vector of rows with j=0,1 or 2+

Usage

```
idx_jet_cat(nj, j)
```

Arguments

nj	Vector of number of jets for each point
j	jet group

import_data	<i>Import raw LHC data</i>
-------------	----------------------------

Description

This function provides a standard way to load in the LHC dataset for our analysis pipeline. Undefined data (-999s) are replaced with NAs, columns are split into variables, labels and supplementary data.

Usage

```
import_data(filepath = "atlas-higgs-challenge-2014-v2.csv")
```

Arguments

filepath	str location of csv
----------	---------------------

Value

named list of X, y, w, kaggle_w, kaggle_s, e_id, nj

interior_point_fit	<i>Interior point fit</i>
--------------------	---------------------------

Description

x contains lambda (dual var) as well

Usage

```
interior_point_fit(
    f,
    dualf,
    gradf,
    Hf,
    x,
    m,
    mu = 10,
    eps = 1e-06,
    eps_feas = 1e-06
)
```

Arguments

f	objective function
dualf	dual objective function
gradf	residual vector
Hf	Hessian for residual
x	primal-dual point
m	number of inequality constraints
mu	interior-point step parameter
eps	tolerance for problem
eps_feas	tolerance for feasibility of primal-dual points

invert_angle_sign	<i>Invert Angle Sign</i>
-------------------	--------------------------

Description

Uses the sign of the pseudorapidity of the tau particle to modify the sign of the pseudorapidity of the leptons and jets on the basis that the interaction should be invariant to rotations of pi about the beam (z) axis.

$$\eta(\theta) = -\log \tan \frac{\theta}{2}$$

$$\eta(\pi - \theta) = -\eta(\theta)$$

Usage

```
invert_angle_sign(X)
```

inv_model_idx	<i>does the inverse procedure to get_model_idx, s.t. if idx <- get_model_idx(j, k, K) then (j,k) <- inv_model_idx(idx) (if R output tuples)</i>
---------------	---

Description

does the inverse procedure to get_model_idx, s.t. if idx <- get_model_idx(j, k, K) then (j,k) <- inv_model_idx(idx) (if R output tuples)

Usage

```
inv_model_idx(idx, K)
```

Arguments

idx	model index
K	number of folds

Value

numeric pair of j and k

kernel_svm	<i>Kernel SVM</i>
------------	-------------------

Description

Fit a kernel support vector machine for binary classification.

Usage

```
kernel_svm(X, y, C, ckernel)
```

Arguments

X	An nxd matrix with samples as rows and features as columns.
y	A length-n vector of -1s and 1s indicating true sample classes.
C	Regularisation parameter.
ckernel	Kernel function with hyperparameters set

l1_logistic_reg	<i>L1 constrained logistic regression</i>
-----------------	---

Description

L1 constrained logistic regression

Usage

```
l1_logistic_reg(X, y, C = 1)
```

Arguments

X	covariate matrix
y	response vector
C	[Optional] regularisation parameter, defaults to 1

Value

b vector of coefficients

lin_kernel	<i>Define linear kernel</i>
------------	-----------------------------

Description

Define linear kernel

Usage

```
lin_kernel(x_i, x_j)
```

Arguments

x_i	point in R^d
x_j	point in R^d

logisticf	<i>Calculate logistic function</i>
-----------	------------------------------------

Description

Calculate logistic function

Usage

```
logisticf(x)
```

Arguments

x float

Value

logisticf(x)

logistic_model	<i>Logistic model object class</i>
----------------	------------------------------------

Description

This reference class object fits a binary classification model, using the `logistic_reg` function. The model can be used to predict the classes of new samples. The sample classes must be 0 or 1, and the prediction returns the estimated probabilities that each sample is class 1. A decision threshold should be subsequently.

Fields

X An $n \times d$ matrix with samples as rows and features as columns.
y A length- n vector of 0s and 1s indicating true sample classes.
coeffs A length- d vector of model coefficients.
lambda regularization parameter

Methods

`initialize(X, y, lambda = 1e-06)` Provide X and y and the coeffs field will be calculated using `logistic_reg`
`predict(X_test)` Provide a matrix of new samples and a vector of $P(y=1)$ is returned

logistic_reg	<i>Logistic regression</i>
--------------	----------------------------

Description

Fit a logistic regression model by IRWLS. This is the same method that glm uses when family="binomial".

Usage

```
logistic_reg(X, y, lambda = 0)
```

Arguments

X	covariate matrix
y	response vector
lambda	[Optional] L2 regularisation parameter, defaults to 0.
r	[Optional] weight vector

Value

b vector of coefficients

logit	<i>Calculate logit function</i>
-------	---------------------------------

Description

Calculate logit function

Usage

```
logit(p)
```

Arguments

p	float in [0,1]
---	----------------

Value

logit(p)

<code>newton_step</code>	<i>Newton step</i>
--------------------------	--------------------

Description

Compute newton step

Usage

`newton_step(grad, H)`

Arguments

<code>grad</code>	gradient vector
<code>H</code>	Hessian matrix

Value

step to take

<code>pairwise_distance</code>	<i>Calculate distance for each row of X0 and X1</i>
--------------------------------	---

Description

Calculate distance for each row of X0 and X1

Usage

`pairwise_distance(X0, X1)`

Arguments

<code>X0</code>	covariate matrix
<code>X1</code>	covariate matrix

Value

vector of distances

partition_data	<i>Partition data into (random) folds for cross-validation.</i>
----------------	---

Description

Partition data into (random) folds for cross-validation.

Usage

```
partition_data(n, k, random = FALSE)
```

Arguments

n	number of rows
k	number of folds
random	flag to choose whether to randomly select

Value

ind vector of integers denoting the OOS fold each row belongs to Returns vector of indices denoting the OOS index, i.e. for rows with I_i=1, those are OOS for i=1

permute_matrix	<i>Cyclic permutation of rows of X by r rows</i>
----------------	--

Description

Cyclic permutation of rows of X by r rows

Usage

```
permute_matrix(X, r = 1)
```

Arguments

X	covariate matrix
r	number of rows to permute (default=1)

plot_amss	<i>define a function to plot multiple ams objects on the same axes</i>
-----------	--

Description

define a function to plot multiple ams objects on the same axes

Usage

```
plot_amss(amss, title = NULL, info = "", min.max = TRUE, scale = 0.8)
```

Arguments

amss	list of ams objects
title	str title to give plot, if null a default title is generated
info	additional info to add to the default title
scale	controls plot cex arguments to size text

plot_distributions	<i>Density plots of variables</i>
--------------------	-----------------------------------

Description

Given a matrix plot the density of the listed variables using ggplot2 facet_wrap.

Usage

```
plot_distributions(X, variables = NULL, labels = NULL)
```

Arguments

X	An nxd matrix with samples as rows and features as columns.
variables	Optional vector of column names to be plotted
lables	Optional vector of class labels to view distributions by class

Details

If the matrix X is more than 10,000 samples, a random 10,000 samples will be selected to keep the amount of data plotted reasonable.

plot_rocs	<i>Plot Multiple ROC curves</i>
-----------	---------------------------------

Description

Plot Multiple ROC curves

Usage

```
plot_rocs(rocs, title = NULL, info = "", scale = 0.8)
```

Arguments

rocs	a list of ROC_curve objects
title	str title to give plot, if null a default title is generated
info	additional info to add to the default title
scale	controls plot cex arguments to size text

poly_kernel	<i>Define polynomial kernel</i>
-------------	---------------------------------

Description

Define polynomial kernel

Usage

```
poly_kernel(x_i, x_j, b)
```

Arguments

x_i	point in R^d
x_j	point in R^d
b	order of polynomial

Value

scalar of $(1+x^T x)^b$

poly_transform	<i>Polynomial transform</i>
----------------	-----------------------------

Description

Run a b-degree polynomial transform on the columns of X (of order b)

Usage

```
poly_transform(X, b = 2)
```

Arguments

X	matrix of covariates
b	order of polynomial

Value

augmented matrix of covariates

rbf_feature	<i>Compute single RBF feature at some centroid i in idx (or xi in Xi)</i>
-------------	---

Description

Compute single RBF feature at some centroid i in idx (or xi in Xi)

Usage

```
rbf_feature(X, s, idx = NULL, xi = NULL)
```

Arguments

X	covariate matrix
s	median pairwise distance of points in X
idx	[Optional] location of reference centroid
xi	[Optional] reference centroid

rbf_kernel	<i>RBF kernel</i>
------------	-------------------

Description

RBf kernel

Usage

```
rbf_kernel(x_i, x_j, sigma)
```

Arguments

x_i	point in R^d
x_j	point in R^d
sigma	bandwidth hyperparameter

reduce_features	<i>Reduce feature space dimensionality by exploiting redundancy</i>
-----------------	---

Description

Reduce feature space dimensionality by exploiting redundancy

Usage

```
reduce_features(X)
```

Arguments

X	matrix of covariates
---	----------------------

Value

X augmented matrix of covariates

ROC_curve	<i>ROC curve object class</i>
-----------	-------------------------------

Description

This reference class object is used to plot a Receiver Operating Characteristic curve. An ROC curve is a performance measure of a classification model, created by plotting the true positive rate (TPR) against the false positive rate (FPR) as the decision threshold is varied. The object finds suitable thresholds, calculates FPR and TPR at each, and can calculate the Area Under the Curve (AUC). A vector of true sample classifications (0 or 1) and a vector of estimated probabilities from a model are needed to initialise.

Fields

thresholds A vector of 30 decision thresholds.
 FP A vector of the false positive rate at each threshold.
 TP A vector of the true positive rate at each threshold.
 auc A numeric that is the area under the ROC curve.

Methods

calc_auc() If the AUC has not already be calculated, this calls the calculation.
 initialize(y, prob) Provide sample labels and probabilites, and the FPR and TPR are calculated at 30 decision thresholds
 plot_curve() Plot the ROC curve.

save_fig	<i>Wrap figure saving</i>
----------	---------------------------

Description

Wrap figure saving

Usage

```
save_fig(plot_func, filepath, filetype = pdf)
```

Arguments

plot_func	partially called function of no arguments to generate plot
filepath	string file path to save to
filetype	type to save as (function name)

scale_dat	<i>Scale features</i>
-----------	-----------------------

Description

This define a function to scale features of a matrix with reference to another matrix useful because you can normalise X_{train} , and apply the same transformation to X_{test} . Use with caution if either dataset contain extreme outliers.

Usage

```
scale_dat(X, ref, na.rm = FALSE, add.intercept = TRUE)
```

Arguments

X	matrix of covariates.
ref	matrix of covariates from which to calculate mu and sd.
na.rm	a logical to indicate if NAs should be stripped in mean and standard deviation computations.
add.intercept	a logical to indicate if column of 1s should be added to the output (Intercept)

Value

augmented matrix of covariates, standardized and an intercept column

set_features_to_rm	<i>create list of feature names we want to omit based on jet group, or constant values/missing</i>
--------------------	--

Description

create list of feature names we want to omit based on jet group, or constant values/missing

Usage

```
set_features_to_rm(X, G, kI, nj)
```

Arguments

X	matrix of covariates
G	number of jet groups
kI	fold indices (test label)
nj	Vector of number of jets for each point

Value

nested list of column names

svm

*Soft margin SVM***Description**

Fit a soft margin support vector machine for binary classification.

Usage

```
svm(X, y, C = 1)
```

Arguments

X	An nxd matrix with samples as rows and features as columns.
y	A length-n vector of -1s and 1s indicating true sample classes.
C	Regularisation parameter.

trig_kernel

*Define trigonometric kernel***Description**

Define trigonometric kernel

Usage

```
trig_kernel(x_i, x_j, b = 0)
```

Arguments

x_i	point in R^d
x_j	point in R^d
b	order of polynomial

tuned_kernel	<i>Function factory to partially call kernel function to return the kernel function with it's hyperparameters set</i>
--------------	---

Description

Function factory to partially call kernel function to return the kernel function with it's hyperparameters set

Usage

```
tuned_kernel(ckernel, ...)
```

Arguments

ckernel	kernel function with args (x_i,x_j,hyper)
---------	---

Value

function with args (x_i,x_j)

Index

add_rbf_features, [1](#)
AMS_data, [2](#)
ams_metric, [2](#)
ams_threshold, [3](#)
average_auc, [3](#)
avg_median_pairwise_distance, [4](#)

backtrack_linesearch, [4](#)

calc_K, [5](#)

decide, [5](#)

fit_l1_logreg, [5](#)

generate_blues, [6](#)
generate_colours, [6](#)
get_const_features, [7](#)
get_model_idx, [7](#)
get_rbf_centroids, [8](#)
get_subset_idx, [8](#)
get_valid_cols, [9](#)

idx_higgs_mass, [9](#)
idx_jet_cat, [10](#)
import_data, [10](#)
interior_point_fit, [11](#)
inv_model_idx, [12](#)
invert_angle_sign, [11](#)

kernel_svm, [12](#)

l1_logistic_reg, [13](#)
lin_kernel, [13](#)
logistic_model, [14](#)
logistic_reg, [15](#)
logisticf, [14](#)
logit, [15](#)

newton_step, [16](#)

pairwise_distance, [16](#)

partition_data, [17](#)
permute_matrix, [17](#)
plot_amss, [18](#)
plot_distributions, [18](#)
plot_rocs, [19](#)
poly_kernel, [19](#)
poly_transform, [20](#)

rbf_feature, [20](#)
rbf_kernel, [21](#)
reduce_features, [21](#)
ROC_curve, [22](#)

save_fig, [22](#)
scale_dat, [23](#)
set_features_to_rm, [23](#)
svm, [24](#)

trig_kernel, [24](#)
tuned_kernel, [25](#)