

Package ‘QDRS’

August 26, 2020

Type Package

Title Quantitative Disease Risk Score

Version 0.1.0

Author Danqing Xu

Maintainer Danqing Xu <xud@pstat.ucsb.edu>

Description Functions for calculating quantitative disease risk scores: (1) Phenotype Risk Score (PheRS), (2) a spectral approach (Eigen), (3) principal component approaches (PCA) that include individual PCs and a linear combination of multiple PCs (LPC), (4), Non-negative Matrix Factorization with Rank 1 (NMF1), and (5) Latent Variable Score based on an unsupervised multivariate mixed model framework (LVS).

License GPL-3

Encoding UTF-8

LazyData true

Imports MASS,
NMF,
AssocTests,
gllvm,
PRROC,
ggplot2,
ggpubr

RoxygenNote 7.1.1

R topics documented:

curve.palettes	2
dist.plot	2
EHR	3
eigen.score	4
example.scores	5
LPC	5
LVS.fit	6
LVS.score	7
multiple.pr.curve	8
multiple.roc.curve	9
NMF1	10
pairwise.auc	10

pairwise.upr	11
pairwise.wilcox	12
PC	12
PheRS	13
predictQDRS	14
rankOne.R	15
Index	16

curve.palettes	<i>Curve Palettes</i>
----------------	-----------------------

Description

A vector of colors for performance curve

Usage

data(curve.palettes)

Format

An object of class character of length 10.

Examples

```
data(curve.palettes)
## Not run:
curve.palettes

## End(Not run)
```

dist.plot	<i>Distribution Plots of QDRS</i>
-----------	-----------------------------------

Description

Generate a pdf file with a density plot, a boxplot and a prevalence plot for the resulting scores. These three plots show the distribution of scores by disease status for visualization of risk stratification.

Usage

```
dist.plot(
  disease,
  output.date = NULL,
  score.mat,
  score.name,
  group,
  cutoff = c(seq(0.05, 0.45, by = 0.05), seq(0.5, 1, by = 0.01)),
  unknown.show = FALSE
)
```

Arguments

disease	A string of disease name.
output.date	An output date. The default is system date (today).
score.mat	A resulting score matrix with multiple types of scores.
score.name	The score of interest.
group	A grouping vector.
cutoff	A vector of cutoff points to set the bins of prevalence plot. The default vector has 60 bins.
unknown.show	A logical value indicates whether unknowns should be used for plotting.

Value

It produces a pdf file of plots and returns the list of three distribution plots.

Examples

```
## Not run:
group1 = example.scores$group
group1[group1 != "Control"] = "Case"
set.seed(830)
na.ind = sample(1:length(group), size = 1000)
group1[na.ind] = NA
res = dist.plot(disease = "Disease",
  score.mat = example.scores$score.mat,
  score.name = "LPC",
  group = group1,
  unknown.show = T)

## End(Not run)
```

EHR

An example data set

Description

It contains `sample.set` a set of 95 features for 5,000 subjects, `sample.group` a vector of group with two levels "Case" and "Control", and `training` a logical vector that indicates the usage for training.

Usage

```
data(EHR)
```

Format

An object of class `list` of length 3.

Examples

```
data(EHR)
## Not run:
# check the number of cases and controls
table(EHR$sample.group)
# check the prevalence of features
colMeans(EHR$sample.set)

## End(Not run)
```

eigen.score

*Eigen Score***Description**

Compute Eigen weights and scores. This method is unsupervised, and assigns weights that are proportional to the balanced accuracies (the average between the sensitivity and the specificity) of the input feature.

Usage

```
eigen.score(X, training, scale = TRUE)
```

Arguments

X	The original data set that include training and test sets. It should be a matrix of numbers.
training	A logical or index vector to indicate whether the subject belongs to the training set.
scale	A logical value to indicate whether the input features need to be scaled.

Value

It returns a list of following components:

weights	The Eigen weights for input features.
scores	The resulting Eigen scores for the whole set.

References

Iuliana Ionita-Laza, Kenneth McCallum, Bin Xu, and Joseph D Buxbaum. A spectral approach integrating functional genomic annotations for coding and noncoding variants. *Nature genetics*, 48(2):214, 2016.

Examples

```
## Not run:
data(EHR)
res1 <- eigen.score(EHR$sample.set, EHR$training)

## End(Not run)
```

example.scores

An example resulting score set

Description

It contains `score.mat` a matrix of five types of QDRSs for 5,000 subjects, `score.names` a vector of score names, `group` a vector of group with two levels "Case" and "Control", and `group.levels` the unique grouping values.

Usage

```
data(example.scores)
```

Format

An object of class `list` of length 4.

Examples

```
data(example.scores)
## Not run:
# check the number of cases and controls
table(example.scores$.group)

## End(Not run)
```

LPC

Linear Combination of Principal Components

Description

Computes LPC weights and scores based on an almost unsupervised method that combines multiple PCs and only requires weak labels to help select the signs of individual PCs.

Usage

```
LPC(X, group, training, scale = TRUE)
```

Arguments

<code>X</code>	The original data set that include training and test sets. It should be a matrix of numbers.
<code>group</code>	A vector that indicate cases ("Case") and controls ("Control"). NA is allowed, and means that the observation is not used in individual PC sign determination.
<code>training</code>	A logical or index vector to indicate whether the subject belongs to the training set.
<code>scale</code>	A logical value to indicate whether the input features need to be scaled.

Value

It returns a list of following components:

lpc.n	The number of significant PCs (eigenvalues) suggested by Tracy-Widom test, a vector of the sign of individual PC.
pc.sign	A vector of the sign of individual PC.
weights	The LPC weights for input features.
scores	The resulting LPC scores for the whole set.

Examples

```
## Not run:
data(EHR)
res1 <- LPC(X = EHR$sample.set, group = EHR$sample.group, training = EHR$training)

## End(Not run)
```

LVS.fit

Fit a Latent Variable Model

Description

LVS.fit computes a generalized linear latent variable model for multivariate bernoulli data, which assumes one latent variable.

Usage

```
LVS.fit(X, family, starting.choice = "random", p.seed = 3080)
```

Arguments

X	The training data set. It should be a matrix of numbers.
starting.choice	Starting values can be generated by fitting model without latent variables, and applying factorial analysis to residuals to get starting values for latent variables and their coefficients (starting.choice = "res"). Another options are to use zeros as a starting values (starting.choice = "zero") or initialize starting values for latent variables with (n x 1) matrix. Defaults to "random".
p.seed	A random seed for model fitting.
family	distribution function for responses. Options are poisson(link = "log"), "negative.binomial" (with log link), binomial(link = "probit") , zero inflated poisson ("ZIP"), gaussian(link = "identity"), "gamma" (with log link), and "ordinal".

Value

An object of class "gllvm" includes the following components.

LVS.score

*Latent Variable Score***Description**

Computes Latent Variable Scores based on an unsupervised multivariate mixed model framework for multivariate bernoulli data that assumes one latent variable.

Usage

```
LVS.score(
  X,
  Y = NULL,
  family = "binomial",
  starting.choice = "random",
  p.seed = 3080
)
```

Arguments

X	A training data set. It should be a matrix of numbers.
Y	A new data set. It should be a matrix of numbers. Default is NULL, then X will be used in prediction.
family	distribution function for responses. Options are poisson(link = "log"), "negative.binomial" (with log link), binomial(link = "probit") , zero inflated poisson ("ZIP"), gaussian(link = "identity"), "gamma" (with log link), and "ordinal".
starting.choice	Starting values can be generated by fitting model without latent variables, and applying factorial analysis to residuals to get starting values for latent variables and their coefficients (starting.choice = "res"). Another options are to use zeros as a starting values (starting.choice = "zero") or initialize starting values for latent variables with (n x 1) matrix. Defaults to "random".
p.seed	A random seed for model fitting.

Value

lvs The resulting LVSs for the new set Y.

References

Jenni Niku, Wesley Brooks, Riki Herliansyah, Francis KC Hui, Sara Taskinen, and David I Warton. Efficient estimation of generalized linear latent variable models. PloS one, 14(5), 2019.

Examples

```
## Not run:
data(EHR)
sample.set = EHR$sample.set
res1 <- LVS.score(X = sample.set[seq(1,50),])

## End(Not run)
```

multiple.pr.curve	<i>Multiple PR curves</i>
-------------------	---------------------------

Description

Generate a pdf file with multiple plots for several group level contrast. Each plot consists of PR curves for different quantitative disease scores, and the corresponding AUPRs for performance comparison.

Usage

```
multiple.pr.curve(
  disease,
  output.date = NULL,
  score.mat,
  score.names,
  group,
  group.levels,
  legend.pos = "topright",
  pairs.sub = NULL
)
```

Arguments

disease	A string of disease name.
output.date	An output date. The default is system date (today).
score.mat	A resulting score matrix with multiple types of scores.
score.names	A vector that indicates the scores of interest for comparison.
group	A grouping vector.
group.levels	A vector of grouping levels that reflects the severity from low to high. The typical first element represents control or healthy status.
legend.pos	A vector of legend position. The default is "bottomright" for all plots.
pairs.sub	An index vector for subset of group level pairs, e.g. the pairs contrasting control vs. other case stages.

Value

It produces a pdf file of plots.

Examples

```
## Not run:
multiple.pr.curve(disease = "Disease",
  score.mat = example.scores$score.mat,
  score.names = example.scores$score.names,
  group = example.scores$group,
  group.levels = example.scores$group.levels,
  legend.pos = "topright", pairs.sub = 1:5)

## End(Not run)
```

multiple.roc.curve	<i>Multiple ROC curves</i>
--------------------	----------------------------

Description

Generate a pdf file with multiple plots for several group level contrast. Each plot consists of ROC curves for different quantitative disease scores, and the corresponding AUROCs for performance comparison.

Usage

```
multiple.roc.curve(
  disease,
  output.date = NULL,
  score.mat,
  score.names,
  group,
  group.levels,
  legend.pos = "bottomright",
  pairs.sub = NULL
)
```

Arguments

disease	A string of disease name.
output.date	An output date. The default is system date (today).
score.mat	A resulting score matrix with multiple types of scores.
score.names	A vector that indicates the scores of interest for comparison.
group	A grouping vector.
group.levels	A vector of grouping levels that reflects the severity from low to high. The typical first element represents control or healthy status.
legend.pos	A vector of legend position. The default is "bottomright" for all plots.
pairs.sub	An index vector for subset of group level pairs, e.g. the pairs contrasting control vs. other case stages.

Value

It produces a pdf file of plots.

Examples

```
## Not run:
multiple.roc.curve(disease = "Disease",
  score.mat = example.scores$score.mat,
  score.names = example.scores$score.names,
  group = example.scores$group,
  group.levels = example.scores$group.levels,
  legend.pos = "bottomright", pairs.sub = 1:5)

## End(Not run)
```

NMF1	<i>Non-negative Matrix Factorization with Rank 1 Score</i>
------	--

Description

Computes NMF1 scores with factor loadings constrained to be non-negative.

Usage

```
NMF1(X, p.seed = 123)
```

Arguments

X	The original data set. It should be a matrix of numbers.
p.seed	A random seed for model fitting.

Value

It returns a list of following components:

weights	The NMF1 weights for input features.
scores	The resulting NMF1 scores for the whole set.

Examples

```
## Not run:
data(EHR)
sample.set = EHR$sample.set
res1 <- NMF1(X = sample.set[seq(1,50),])

## End(Not run)
```

pairwise.auc	<i>Pairwise Area Under the ROC Curve</i>
--------------	--

Description

Calculates AUROC for pairwise levels of group contrast .

Usage

```
pairwise.auc(x, g, format = "l")
```

Arguments

x	A response vector.
g	A grouping vector or factor.
format	Format of result. The default is "l", long format. The other option is square format.

Value

A data frame with group 1, group 2, and the AUROC with comparison between them.

Examples

```
## Not run:
data(EHR)
x = c(rnorm(2500), runif(2500))
res = pairwise.auc(x, EHR$sample.group)

## End(Not run)
```

pairwise.upr	<i>Pairwise Area Under the PR Curve</i>
--------------	---

Description

Calculates AUPRC for pairwise group levels contrast.

Usage

```
pairwise.upr(x, g, format = "l")
```

Arguments

x	A response vector.
g	A grouping vector or factor.
format	Format of result. The default is "l", long format. The other option is square format.

Value

A data frame with group 1, group 2, and the AUPRC with comparison between them.

Examples

```
## Not run:
data(EHR)
x = c(rnorm(2500), runif(2500))
res = pairwise.upr(x, EHR$sample.group)

## End(Not run)
```

pairwise.wilcox

Pairwise Wilcoxon Rank Sum Tests

Description

Calculates pairwise comparisons between group levels with corrections for multiple testing.

Usage

```
pairwise.wilcox(
  x,
  g,
  alternative = "two.sided",
  p.adjust.method = "bonferroni"
)
```

Arguments

x	A response vector.
g	A grouping vector or factor.
alternative	A character string specifying the alternative hypothesis, must be one of "two.sided" (default), "greater" or "less". You can specify just the initial letter.
p.adjust.method	Method for adjusting p values (can be abbreviated): "holm", "hochberg", "hommel", "bonferroni", "BH", "BY", "fdr", "none". The default is "bonferroni".

Value

A data frame with group 1, group 2, and the Wilcoxon Rank Sum test p value of comparison between them.

Examples

```
## Not run:
data(EHR)
x = c(rnorm(2500), runif(2500))
res = pairwise.wilcox(x, EHR$sample.group)

## End(Not run)
```

PC

Individual Principal Components with Selected Signs

Description

Computes individual PC weights and scores. The approach only requires weak labels to help select the signs of individual PCs.

Usage

```
PC(X, group, training, scale = TRUE, pc.num)
```

Arguments

X	The original data set that include training and test sets. It should be a matrix of numbers.
group	A grouping vector or factor that indicate cases ("Case") and controls ("Control"). NA is allowed, and means that the observation is not used in individual PC sign determination.
training	A logical or index vector to indicate whether the subject belongs to the training set.
scale	A logical value to indicate whether the input features need to be scaled.
pc.num	The vector of desired individual PCs.

Value

It returns a list of following components:

weights	The selected individual PCs' weights for input features.
scores	The resulting individual PC scores for the whole set.

Examples

```
## Not run:
data(EHR)
res1 <- PC(X = EHR$sample.set,
  group = EHR$sample.group,
  training = EHR$training, pc.num = 1:2)

## End(Not run)
```

PheRS

Phenotype Risk Score

Description

Computes Phenotype Risk Scores based on a approach that was proposed in the context of rare Mendelian phenotypes.

Usage

```
PheRS(X, feature.prevalence = NULL, group = NULL)
```

Arguments

<code>X</code>	The original data set that include training and test sets. It should be a matrix of binary numbers.
<code>feature.prevalence</code>	A vector of feature prevalence.
<code>group</code>	A vector that indicate cases ("Case") and controls ("Control"). NA is allowed. The default is NULL, meaning all observations will be used in prevalence computation if the feature prevalence vector is not provided by the user. Otherwise, only the controls will be used.

Value

It returns a list of following components:

<code>weights</code>	The PheRS weights for input features.
<code>scores</code>	The resulting PheRSs for the whole set.

References

Lisa Bastarache, Jacob J Hughey, Scott Hebring, Joy Marlo, Wanke Zhao, Wanting T Ho, Sara L Van Driest, Tracy L McGregor, Jonathan D Mosley, Quinn S Wells, et al. Phenotype risk scores identify patients with unrecognized mendelian disease patterns. *Science*, 359(6381):1233– 1239, 2018.

Examples

```
## Not run:
data(EHR)
res1 <- PheRS(X = EHR$sample.set, group = EHR$sample.group)
res2 <- PheRS(X = EHR$sample.set, group = NULL)

## End(Not run)
```

predictQDRS

Predict Quantitative Disease Risk Score

Description

Predicts Quantitative Disease Risk Scores using weights derived from training set.

Usage

```
predictQDRS(Y, weights)
```

Arguments

<code>Y</code>	A new data set for prediction.
<code>weights</code>	A matrix of weights.

Value

It returns a matrix of resulting quantitative disease risk scores.

Examples

```
## Not run:
data(EHR)
res1 <- PC(X = EHR$sample.set,
  group = EHR$sample.group,
  training = EHR$training, pc.num = 1:2)
res2 <- predictQDRS(Y = EHR$sample.set[seq(1,50),], weights = res1$weights)

## End(Not run)
```

rankOne.R

*R Matrix with Rank One for Eigen Approach***Description**

Computes an estimated rank-one matrix R with unit-norm eigenvector v . Up to a sign ambiguity, the entries of v are proportional to the balanced accuracies (the average between the sensitivity and the specificity) of the input features.

Usage

```
rankOne.R(Qmat)
```

Arguments

`Qmat` The correlation matrix Q .

Value

`Rmat` An estimate of the rank-one R matrix

Index

*Topic **datasets**

curve.palettes, [2](#)

EHR, [3](#)

example.scores, [5](#)

curve.palettes, [2](#)

dist.plot, [2](#)

EHR, [3](#)

eigen.score, [4](#)

example.scores, [5](#)

LPC, [5](#)

LVS.fit, [6](#)

LVS.score, [7](#)

multiple.pr.curve, [8](#)

multiple.roc.curve, [9](#)

NMF1, [10](#)

pairwise.auc, [10](#)

pairwise.upr, [11](#)

pairwise.wilcox, [12](#)

PC, [12](#)

PheRS, [13](#)

predictQDRS, [14](#)

rankOne.R, [15](#)