# Package 'GSAR'

## November 7, 2025

140vember 7, 2023
Type Package
Title Gene Set Analysis in R
<b>Version</b> 1.45.0
<b>Date</b> 2025-5-12
Author Yasir Rahmatallah <pre><pre><pre><pre><pre><pre><pre><pre></pre></pre></pre></pre></pre></pre></pre></pre>
Maintainer Yasir Rahmatallah <pre></pre>
<b>Depends</b> R ( $>= 3.0.1$ ), igraph ( $>= 0.7.1$ )
Imports stats, graphics
Suggests MASS, GSVAdata, ALL, tweeDEseqCountData, GSEABase, annotate, org.Hs.eg.db, Biobase, genefilter, hgu95av2.db, edgeR, BiocStyle
LazyData yes
biocViews Software, StatisticalMethod, DifferentialExpression
<b>Description</b> Gene set analysis using specific alternative hypotheses. Tests for differential expression, scale and net correlation structure.
License GPL (>=2)
git_url https://git.bioconductor.org/packages/GSAR
git_branch devel
git_last_commit 25443bb
git_last_commit_date 2025-10-29
Repository Bioconductor 3.23
Date/Publication 2025-11-06
Contents
GSAR-package ADtest AggrFtest CVMtest findMST2

2 GSAR-package

GSAR-package		Gene Set Analysis in R	
Index			40
	WWtest		38
	RMDtest		33
	RKStest		31
	RCVMtest		29
	RADtest		26
	radial.ranking		25
	plotMST2.pathway		22
	p53DataSet		21
	MDtest		19
	KStest		17
	HDP.ranking		16
	GSNCAtest		13
	findMST2.PPI		11

#### Description

Package GSAR provides a set of statistical methods for self-contained gene set analysis. It consists of two-sample multivariate nonparametric statistical methods to test a null hypothesis against specific alternative hypotheses, such as differences in shift (functions KStest, MDtest, ADtest, and CVMtest), scale (functions RKStest, RMDtest, RADtest, RCVMtest, and AggrFtest) or correlation structure (function GSNCAtest) between two conditions. It also offers a graphical visualization tool for correlation networks to examine the change in the net correlation structure of a gene set between two conditions (function plotMST2.pathway). The visualization scheme is based on the minimum spanning trees (MSTs). Function findMST2 is used to find the unioin of the first and second MSTs. The same tool works as well for protein-protein interaction (PPI) networks to highlight the most essential interactions among proteins and reveal fine network structure as was already shown in Zybailov et. al. 2016. Function findMST2.PPI is used to find the unioin of the first and second MSTs of PPI networks. Gene set analysis methods available in this package were proposed in Rahmatallah et. al. 2012, Rahmatallah et. al. 2014, and Rahmatallah and Glazko 2024. The performance of different methods was tested using simulation and real gene expression data. These methods can be applied to RNA-Seq count data given that proper normalization is used. Proper normalization must take into account both the within-sample differences (mainly gene length) and between-samples differences (library size). However, because the count data often follows the negative binomial distribution, special attention should be paid to applying the variance tests (RKStest, RMDtest, RADtest, RCVMtest, and AggrFtest). The variance of the negative binomial distribution is proportional to it's mean and multivariate tests of variance designed specifically for RNA-seq count data are under-explored.

## Author(s)

Yasir Rahmatallah <yrahmatallah@uams.edu>, Galina Glazko <gvglazko@uams.edu> Maintainer: Yasir Rahmatallah <yrahmatallah@uams.edu>, Galina Glazko <gvglazko@uams.edu> ADtest 3

#### References

Rahmatallah Y. and Glazko G. (2024) Gene Set Analysis: improving data interpretability with new differential variance tests. 09 September 2024, PREPRINT (Version 1) available at Research Square [https://doi.org/10.21203/rs.3.rs-4888767/v1].

Rahmatallah Y., Emmert-Streib F. and Glazko G. (2014) Gene sets net correlations analysis (GSNCA): a multivariate differential coexpression test for gene sets. Bioinformatics **30**, 360–368.

Rahmatallah Y., Emmert-Streib F. and Glazko G. (2012) Gene set analysis for self-contained tests: complex null and specific alternative hypotheses. Bioinformatics **28**, 3073–3080.

Zybailov B., Byrd A., Glazko G., Rahmatallah Y. and Raney K. (2016) Protein-protein interaction analysis for functional characterization of helicases. Methods, **108**, 56–64.

#### See Also

igraph.

ADtest	Multivariate Anderson-Darling Test of Means
Abecse	munivariate macroon Darting Test of means

## **Description**

Performs two-sample nonparametric multivariate test of means based on the minimum spanning tree (MST) and Anderson-Darling statistic. It tests the null hypothesis that a set of features has the same mean in two conditions versus different means.

## Usage

```
ADtest(object, group, nperm=1000, pvalue.only=TRUE)
```

#### **Arguments**

object	a numeric matrix with columns and rows respectively corresponding to samples and features.
group	a numeric vector indicating group associations for samples. Possible values are 1 and 2. $$
nperm	number of permutations used to estimate the null distribution of the test statistic. If not given, a default value 1000 is used.
pvalue.only	logical. If TRUE (default), the p-value is returned. If FALSE a list of length three containing the observed statistic, the vector of permuted statistics, and the p-value is returned.

4 ADtest

#### **Details**

This function tests the null hypothesis that a set of features has no shift between two conditions. It performs a two-sample nonparametric multivariate test based on the minimum spanning tree (MST) and Anderson-Darling statistic as proposed in Rahmatallah and Glazko (2024). The MST of the weighted undirectional graph created from the samples is found. The nodes of the MST are ranked based on their position in the MST. The MST is rooted at the node with largest geodisic distance and then nodes are ranked in the High Directed Preorder (HDP) traversal of the tree (Rahmatallah et. al. 2012). The Anderson-Darling statistic can be defined as

$$A = \frac{1}{n_1 n_2} \sum_{i=1}^{N-1} \frac{(r_i N - i n_1)^2}{i(N-i)}$$

where  $r_i$  is the number of nodes (samples) from condition 1 which ranked lower than i,  $1 \le i \le N$ , N is the total number of samples, and  $n_1$  and  $n_2$  are respectively the number of samples in groups 1 and 2. The Anderson-Darling statistic put more emphasis on the tails of the deviation between the empirical cumulative distribution functions (CDFs) of samples of two phenotypes in the MST. The performance of this test under different alternative hypotheses was examind in Rahmatallah and Glazko (2024). The null distribution of the test statistic is estimated by permuting sample labels nperm times and calculating the test statistic for each. P-value is calculated as

$$p.value = \frac{\sum_{k=1}^{nperm} I\left[A_k \ge A_{obs}\right] + 1}{nperm + 1}$$

where  $A_k$  is the test statistic for permutation k,  $A_{obs}$  is the observed test statistic, and I is the indicator function.

#### Value

When pvalue.only=TRUE (default), function ADtest returns the p-value indicating the attained significance level. When pvalue.only=FALSE, function ADtest produces a list of length 3 with the following components:

statistic the value of the observed test statistic.

perm. stat numeric vector of the resulting test statistic for nperm random permutations of

sample labels.

p.value p-value indicating the attained significance level.

## Note

This function invokes function HDP. ranking which does not work properly if there is any node in the MST with more than 26 links. However, this situation is almost impossible for a dataset composed of a few hundreds or less of samples.

#### Author(s)

Yasir Rahmatallah and Galina Glazko

AggrFtest 5

#### References

Rahmatallah Y. and Glazko G. (2025) Improving data interpretability with new differential sample variance tests. BMC Bioinformatics **26**, 103.

Rahmatallah Y., Emmert-Streib F. and Glazko G. (2012) Gene set analysis for self-contained tests: complex null and specific alternative hypotheses. Bioinformatics **28**, 3073–3080.

#### See Also

```
KStest, MDtest, CVMtest, RADtest, RKStest, RMDtest, RCVMtest, HDP.ranking.
```

#### **Examples**

```
## generate a feature set of length 20 in two conditions
## each condition has 20 samples
## use multivariate normal distribution
library(MASS)
ngenes <- 20
nsamples <- 40
## let the mean vector have zeros of length 20 in both conditions
zero_vector <- array(0,c(1,ngenes))</pre>
## set the covariance matrix to be an identity matrix for both conditions
cov_mtrx <- diag(ngenes)</pre>
gp <- mvrnorm(nsamples, zero_vector, cov_mtrx)</pre>
## apply a mean shift of 3 to half of the features under condition 1
gp[1:20,1:10] <- gp[1:20,1:10] + 3
dataset \leftarrow aperm(gp, c(2,1))
## first 20 samples belong to condition 1
## second 20 samples belong to condition 2
pvalue <- ADtest(object=dataset, group=c(rep(1,20),rep(2,20)))</pre>
```

AggrFtest

Aggregated F-Test of Variance Using Fisher's Probability Combining Method

## **Description**

Performs two-sample nonparametric test of variance. The univariate F-test is used for every gene in the gene set and the resulted p-values are aggregated together using Fisher's probability combining method and used as the test statistic. The null distribution of the test statistic is estimated by permuting sample labels and calculating the test statistic for a large number of times. This statistic tests the null hypothesis that none of the genes shows significant difference in variance between two conditions against the alternative hypothesis that at least one gene shows significant difference in variance between two conditions according to the F-test.

## Usage

```
AggrFtest(object, group, nperm=1000, pvalue.only=TRUE)
```

#### **Arguments**

object a numeric matrix with columns and rows respectively corresponding to samples

and features (genes).

group a numeric vector indicating group associations for samples. Possible values are

1 and 2.

nperm a numeric value indicating the number of permutations used to estimate the null

distribution of the test statistic. If not given, a default value 1000 is used.

pyalue.only logical. If TRUE (default), the p-value is returned. If FALSE a list of length three

containing the observed statistic, the vector of permuted statistics, and the p-

value is returned.

#### **Details**

This function tests the null hypothesis that none of the genes in a gene set shows a significant difference in variance between two conditions according to the F-test against the alternative hypothesis that at least one gene shows significant difference in variance according to the F-test. It performs a two-sample nonparametric test of variance by using the univariate F-test for every gene in a set, adjust for multiple testing using the Benjamini and Hochberg method (also known as FDR) as shown in Benjamini and Hochberg (1995), and then aggregates the obtained adjusted p-values using Fisher's probability combining method to get a test statistic (T) for the gene set

$$T = -2\sum_{i=1}^{p} \log_e(p_i)$$

where  $p_i$  is the adjusted p-value of the univariate F-test for gene i. The null distribution of the test statistic is estimated by permuting sample labels nperm times and calculating the test statistic T for each. P-value is calculated as

$$p.value = \frac{\sum_{k=1}^{nperm} I\left[T_k \ge T_{obs}\right] + 1}{nperm + 1}$$

where  $T_k$  is the test statistic for permutation k,  $T_{obs}$  is the observed test statistic, and I is the indicator function.

## Value

When pvalue.only=TRUE (default), function AggrFtest returns the p-value indicating the attained significance level. When pvalue.only=FALSE, function AggrFtest produces a list of length 3 with the following components:

statistic the value of the observed test statistic.

perm. stat numeric vector of the resulting test statistic for nperm random permutations of

sample labels.

p.value p-value indicating the attained significance level.

## Author(s)

Yasir Rahmatallah and Galina Glazko

CVMtest 7

#### References

Benjamini Y. and Hochberg Y. (1995) Controlling the false discovery rate: a practical and powerful approach to multiple testing. Journal of the Royal Statistical Society Series B **57**, 289–300.

#### See Also

```
RKStest, RMDtest, RADtest, RCVMtest.
```

#### **Examples**

```
## generate a feature set of length 20 in two conditions
## each condition has 20 samples
## use multivariate normal distribution
library(MASS)
ngenes <- 20
nsamples <- 40
## let the mean vector have zeros of length 20 for both conditions
zero_vector <- array(0,c(1,ngenes))</pre>
## set the covariance matrix to be an identity matrix for condition 1
cov_mtrx <- diag(ngenes)</pre>
gp1 <- mvrnorm((nsamples/2), zero_vector, cov_mtrx)</pre>
## set some scale difference in the covariance matrix for condition 2
cov_mtrx <- cov_mtrx*3</pre>
gp2 <- mvrnorm((nsamples/2), zero_vector, cov_mtrx)</pre>
## combine the data of two conditions into one dataset
gp <- rbind(gp1,gp2)</pre>
dataset <- aperm(gp, c(2,1))</pre>
## first 20 samples belong to group 1
## second 20 samples belong to group 2
pvalue <- AggrFtest(object=dataset, group=c(rep(1,20),rep(2,20)))</pre>
```

**CVMtest** 

Multivariate Cramer-Von Mises of Means

## Description

Performs two-sample nonparametric multivariate test of means based on the minimum spanning tree (MST) and Cramer-Von Mises statistic. It tests the null hypothesis that a set of features has the same mean in two conditions versus different means.

#### Usage

```
CVMtest(object, group, nperm=1000, pvalue.only=TRUE)
```

8 CVMtest

#### **Arguments**

object a numeric matrix with columns and rows respectively corresponding to samples

and features.

group a numeric vector indicating group associations for samples. Possible values are

1 and 2.

nperm number of permutations used to estimate the null distribution of the test statistic.

If not given, a default value 1000 is used.

pvalue.only logical. If TRUE (default), the p-value is returned. If FALSE a list of length three

containing the observed statistic, the vector of permuted statistics, and the p-

value is returned.

#### **Details**

This function tests the null hypothesis that a set of features has no shift between two conditions. It performs a two-sample nonparametric multivariate test based on the minimum spanning tree (MST) and Cramer-Von Mises statistic as proposed in Rahmatallah and Glazko (2024). The MST of the weighted undirectional graph created from the samples is found. The nodes of the MST are ranked based on their position in the MST. The MST is rooted at the node with largest geodisic distance (rank 1) and then nodes are ranked in the High Directed Preorder (HDP) traversal of the tree (Rahmatallah et. al. 2012). The Cramer-Von Mises statistic can be defined as

$$C = \frac{n_1 n_2}{N^2} \left\{ \sum_{i=1}^{n_1} \left[ \frac{r_i}{n_1} - i \left( \frac{1}{n_1} + \frac{1}{n_2} \right) \right]^2 + \sum_{j=1}^{n_2} \left[ \frac{s_j}{n_2} - j \left( \frac{1}{n_1} + \frac{1}{n_2} \right) \right]^2 \right\}$$

where  $r_i$  and  $s_j$  are respectively the number of samples from conditions 1 and 2 which ranked lower than i,  $1 \le i \le N$ ,  $n_1$  and  $n_2$  are respectively the number of samples in groups 1 and 2, and N is the total number of samples. The performance of this test under different alternative hypotheses was examind in Rahmatallah and Glazko (2024). The null distribution of the test statistic is estimated by permuting sample labels nperm times and calculating the test statistic for each. P-value is calculated as

$$p.value = \frac{\sum_{k=1}^{nperm} I[C_k \ge C_{obs}] + 1}{nperm + 1}$$

where  $C_k$  is the test statistic for permutation k,  $C_{obs}$  is the observed test statistic, and I is the indicator function.

#### Value

When pvalue.only=TRUE (default), function ADtest returns the p-value indicating the attained significance level. When pvalue.only=FALSE, function ADtest produces a list of length 3 with the following components:

statistic the value of the observed test statistic.

perm. stat numeric vector of the resulting test statistic for nperm random permutations of

sample labels.

p.value p-value indicating the attained significance level.

findMST2

#### Note

This function invokes function HDP.ranking which does not work properly if there is any node in the MST with more than 26 links. However, this situation is almost impossible for a dataset composed of a few hundreds or less of samples.

#### Author(s)

Yasir Rahmatallah and Galina Glazko

#### References

Rahmatallah Y. and Glazko G. (2025) Improving data interpretability with new differential sample variance tests. BMC Bioinformatics **26**, 103.

Rahmatallah Y., Emmert-Streib F. and Glazko G. (2012) Gene set analysis for self-contained tests: complex null and specific alternative hypotheses. Bioinformatics **28**, 3073–3080.

#### See Also

```
RCVMtest, KStest, MDtest, ADtest, RKStest, RMDtest, RADtest, HDP.ranking.
```

#### **Examples**

```
## generate a feature set of length 20 in two conditions
## each condition has 20 samples
## use multivariate normal distribution
library(MASS)
ngenes <- 20
nsamples <- 40
## let the mean vector have zeros of length 20 in both conditions
zero_vector <- array(0,c(1,ngenes))</pre>
## set the covariance matrix to be an identity matrix for both conditions
cov_mtrx <- diag(ngenes)</pre>
gp <- mvrnorm(nsamples, zero_vector, cov_mtrx)</pre>
## apply a mean shift of 3 to half of the features under condition 1
gp[1:20,1:10] \leftarrow gp[1:20,1:10] + 3
dataset \leftarrow aperm(gp, c(2,1))
## first 20 samples belong to condition 1
## second 20 samples belong to condition 2
pvalue <- CVMtest(object=dataset, group=c(rep(1,20),rep(2,20)))</pre>
```

findMST2

Union of the First and Second Minimum Spanning Trees

#### Description

Find the union of the first and second minimum spanning trees.

10 findMST2

#### Usage

findMST2(object, cor.method="pearson", min.sd=1e-3, return.MST2only=TRUE)

#### **Arguments**

object a numeric matrix with columns and rows respectively corresponding to samples

and features.

cor.method a character string indicating which correlation coefficient is to be computed.

Possible values are "pearson" (default), "spearman" and "kendall".

min.sd the minimum allowed standard deviation for any feature. If any feature has a

standard deviation smaller than min. sd the execution stops and an error message

is returned.

return.MST2only

logical. If TRUE (default), an object of class igraph containing the MST2 is returned. If FALSE, a list of length three containing objects of class igraph is returned. The first and second objects are the first and second MSTs, respectively.

The third is the union of the first and second, MST2.

#### **Details**

This function produces the union of the first and second minimum spanning trees (MSTs) as an object of class igraph (check package igraph for details). It can as well return the first and second minimum spanning trees when return.MST2only is FALSE (default). It starts by calculating the correlation (coexpression) matrix and using it to obtain a weighting matrix for a complete graph using the equation  $w_{ij} = 1 - |r_{ij}|$  where  $r_{ij}$  is the correlation between features i and j and  $w_{ij}$  is the weight of the link between vertices (nodes) i and j in the graph G(V, E).

For the graph G(V, E) where V is the set of vertices and E is the set of edges, the first MST is defined as the acyclic subset  $T_1 \subseteq E$  that connects all vertices in V and whose total length  $\sum_{i,j \in T_1} d(v_i, v_j)$  is minimal (Rahmatallah et. al. 2014). The second MST is defined as the MST of the reduced graph  $G(V, E - T_1)$ . The union of the first and second MSTs is denoted as MST2.

It was shown in Rahmatallah et. al. 2014 that MST2 can be used as a graphical visualization tool to highlight the most highly correlated genes in the correlation network. A gene that is highly correlated with all the other genes tends to occupy a central position and has a relatively high degree in the MST2 because the shortest paths connecting the vertices of the first and second MSTs tend to pass through the vertex corresponding to this gene. In contrast, a gene with low intergene correlations most likely occupies a non-central position in the MST2 and has a degree of 2.

In rare cases, a feature may have a constant or nearly constant level across the samples. This results in a zero or a tiny standard deviation. Such case produces an error in command cor used to compute the correlations between features. To avoid this situation, standard deviations are checked in advance and if any is found below the minimum limit min.sd (default is 1e-3), the execution stops and an error message is returned indicating the the number of feature causing the problem (if only one the index of that feature is given too).

#### Value

When return.MST2only=TRUE (default), function findMST2 returns an object of class igraph representing the MST2. If return.MST2only=FALSE, function findMST2 returns a list of length 3 with the following components:

findMST2.PPI

MST2 an object of class igraph containing the union of the first and second MSTs.

first.mst an object of class igraph containing the first MST. second.mst an object of class igraph containing the second MST.

#### Author(s)

Yasir Rahmatallah and Galina Glazko

#### References

Rahmatallah Y., Emmert-Streib F. and Glazko G. (2014) Gene sets net correlations analysis (GSNCA): a multivariate differential coexpression test for gene sets. Bioinformatics **30**, 360–368.

#### See Also

```
GSNCAtest, plotMST2.pathway.
```

## Examples

```
## generate a dataset of 20 features and 20 samples
## use multivariate normal distribution with different covariance matrices
library(MASS)
ngenes <- 20
nsamples <- 20
zero_vector <- array(0,c(1,ngenes))</pre>
## create a covariance matrix with high off-diagonal elements
## for the first 5 features and low for the remaining 15 features
cov_mtrx <- diag(ngenes)</pre>
cov_mtrx[!diag(ngenes)] <- 0.1</pre>
mask <- diag(ngenes/4)</pre>
mask[!diag(ngenes/4)] <- 0.6</pre>
cov_mtrx[1:(ngenes/4),1:(ngenes/4)] <- mask</pre>
gp <- mvrnorm(nsamples, zero_vector, cov_mtrx)</pre>
dataset <- aperm(gp, c(2,1))</pre>
## findMST2 returns a list of length 3
## trees[[1]] is an object of class igraph containing the MST2
trees <- findMST2(dataset)</pre>
```

findMST2.PPI

Union of the First and Second Minimum Spanning Trees for PPI Networks

## Description

Find the union of the first and second minimum spanning trees for protein-protein interaction (PPI) networks.

## Usage

```
findMST2.PPI(object, return.MST2only=TRUE)
```

12 findMST2.PPI

## **Arguments**

object an object of class igraph representing the PPI network.

return.MST2only

logical. If TRUE (default), an object of class igraph containing the MST2 is returned. If FALSE, a list of length three containing objects of class igraph is returned. The first and second objects are the first and second MSTs, respectively.

The third is the union of the first and second, MST2.

#### **Details**

This function produces the union of the first and second minimum spanning trees (MSTs) as an igraph object (check package igraph for details). It can as well return the first and second minimum spanning trees when return. MST2 only is FALSE.

For the graph G(V, E) where V is the set of vertices and E is the set of edges, the first MST is defined as the acyclic subset  $T_1 \subseteq E$  that connects all vertices in V and whose total length  $\sum_{i,j \in T_1} d(v_i,v_j)$ is minimal (Rahmatallah et. al. 2014). The second MST is defined as the MST of the reduced graph  $G(V, E - T_1)$ . The union of the first and second MSTs is denoted as MST2.

It was shown in Zybailov et. al. 2016 that MST2 can be informative as a graphical visualization tool in deciphering the properties of protein-protein interaction (PPI) networks by highlighting the minimum set of essential interactions among proteins. Most influential proteins with many interactions tend to occupy central position and have relatively high connectivity degree in the MST2 because the shortest paths connecting the vertices of the first and second MSTs tend to pass through the verteces corresponding to these proteins. In contrast, proteins with few interactions most likely occupy non-central positions in the MST2 and have a degree of 2.

#### Value

If return. MST2only=TRUE (default), function findMST2. PPI returns an object of class igraph representing the MST2. If return. MST2only=FALSE, function findMST2. PPI returns a list of length 3 with the following components:

MST2 an object of class igraph containing the union of the first and second MSTs.

an object of class igraph containing the first MST. first.mst an object of class igraph containing the second MST. second.mst

#### Author(s)

Yasir Rahmatallah and Galina Glazko

## References

Zybailov B., Byrd A., Glazko G., Rahmatallah Y. and Raney K. (2016) Protein-protein interaction analysis for functional characterization of helicases. Methods, 108, 56–64.

#### See Also

GSNCAtest, plotMST2.pathway.

GSNCAtest 13

## **Examples**

```
## generate a random undirected graph with power-law
## distribution degree where minimum degree is 4 and
## maximum degree is 100
set.seed(123)
degs <- sample(c(4:100), 100, replace=TRUE, prob=c(4:100)^-2)
if(floor(sum(degs)/2) != (sum(degs)/2)) degs[1] <- degs[1] + 1
randomGraph <- sample_degseq(degs, method="vl")
## find MST2 of the random graph and highlight vertices
## with degree greater than 10 with red color
mst2.ppi <- findMST2.PPI(object=randomGraph, return.MST2only=TRUE)
degs <- degree(mst2.ppi)
ind <- which(degs > 10)
V(mst2.ppi)$color <- "yellow"
V(mst2.ppi)$color[ind] <- "red"</pre>
```

GSNCAtest

Gene Sets Net Correlations Analysis

#### **Description**

Performs Gene Sets Net Correlation Analysis (GSNCA) test to detect differentially coexpressed gene sets.

## Usage

```
GSNCAtest(object, group, nperm=1000, cor.method="pearson", check.sd=TRUE,
min.sd=1e-3, max.skip=10, pvalue.only=TRUE)
```

## **Arguments**

object	a numeric matrix with columns and rows respectively corresponding to samples and features.
group	a numeric vector indicating group associations for samples. Possible values are $1\ \mathrm{and}\ 2.$
nperm	number of permutations used to estimate the null distribution of the test statistic. If not given, a default value 1000 is used.
cor.method	a character string indicating which correlation coefficient is to be computed. Possible values are "pearson" (default), "spearman" and "kendall".
check.sd	logical. Should the standard deviations of features checked for small values before the intergene correlations are computed? Default is TRUE (recommended).
min.sd	the minimum allowed standard deviation for any feature. If any feature has a standard deviation smaller than min. sd the execution stops and an error message is returned.
max.skip	maximum number of skipped random permutations which yield any feature with a standard deviation less than min.sd.

14 GSNCAtest

pvalue.only

logical. If TRUE (default), the p-value is returned. If FALSE a list of length three containing the observed statistic, the vector of permuted statistics, and the p-value is returned.

#### **Details**

This function performs the Gene Sets Net Correlations Analysis (GSNCA), a two-sample nonparametric multivariate differential coexpression test that accounts for the correlation structure between features (genes). The test assigns weight factors to features under one condition and adjust these weights simultaneously such that equality is achieved between each feature's weight and the sum of its weighted absolute correlations with other features in the feature set. The problem is solved as an eigenvector problem with a unique solution (see Rahmatallah et. al. 2014 for details). The test statistic  $w_{GSNCA}$  is given by the first norm between the scaled weight vectors  $w^{(1)}$  and  $w^{(2)}$  (each vector is multiplied by its norm) between two conditions

$$w = \sum_{i=1}^{p} |w_i^{(1)} - w_i^{(2)}|$$

This test statistic tests the null hypothesis that w=0 against the alternative that w does not equal to zero. The performance of this test was thoroughly examind in Rahmatallah et. al. (2014). The null distribution of the test statistic is estimated by permuting sample labels nperm times and calculating the test statistic for each. P-value is calculated as

$$p.value = \frac{\sum_{k=1}^{nperm} I[W_k \ge W_{obs}] + 1}{nperm + 1}$$

where  $W_k$  is the test statistic for permutation k,  $W_{obs}$  is the observed test statistic, and I is the indicator function.

In the case of RNA-seq count data, some non-expressed genes may have zero counts across the samples under one or two conditions. Such situation results in zero or tiny standard deviation for one or more features. Such case produces an error in command cor used to compute the correlation coefficients between features. To avoid this situation, standard deviations are checked in advance when check.sd is TRUE (default) and if any is found below the minimum limit min.sd (default is 1e-3), the execution stops and an error message is returned indicating the number of feature causing the problem (if only one the index of that feature is given too). If a feature has nearly a constant level for some samples under both conditions, permuting sample labels may group such samples under one condition and produce a standard deviation smaller than min.sd. To allow the test to skip such permutations without causing excessive delay, we set an upper limit for the number of allowed skips by the argument max.skip (default is 10). If the upper limit is exceeded, an error message is returned. Allowing this skipping may or may not solve the issue depending on the proportion of samples causing the problem in the feature set.

If the user is certain that the tested feature sets contain no feature with nearly equal levels over many samples (such as the case with microarrays), the checking stage for tiny standard deviations can be skipped by setting check.sd to FALSE in order to reduce the execution time.

GSNCAtest 15

#### Value

When pvalue.only=TRUE (default), function GSNCAtest returns the p-value indicating the attained significance level. When pvalue.only=FALSE, function GSNCAtest produces a list of length 3 with the following components:

statistic the value of the observed test statistic.

perm. stat numeric vector of the resulting test statistic for nperm random permutations of

sample labels.

p.value p-value indicating the attained significance level.

#### Author(s)

Yasir Rahmatallah and Galina Glazko

#### References

Rahmatallah Y., Emmert-Streib F. and Glazko G. (2014) Gene sets net correlations analysis (GSNCA): a multivariate differential coexpression test for gene sets. Bioinformatics **30**, 360–368.

#### See Also

```
findMST2, plotMST2.pathway.
```

#### **Examples**

```
## generate a feature set of length 20 in two conditions
## each condition has 20 samples
## use multivariate normal distribution with different covariance matrices
library(MASS)
ngenes <- 20
nsamples <- 40
zero_vector <- array(0,c(1,ngenes))</pre>
## create a covariance matrix with low off-diagonal elements
cov_mtrx1 <- diag(ngenes)</pre>
cov_mtrx1[!diag(ngenes)] <- 0.1</pre>
## create a covariance matrix with high off-diagonal elements
## for the first 5 features and low for the rest 15 features
cov_mtrx2 <- diag(ngenes)</pre>
cov_mtrx2[!diag(ngenes)] <- 0.1</pre>
mask <- diag(ngenes/4)</pre>
mask[!diag(ngenes/4)] <- 0.6</pre>
cov_mtrx2[1:(ngenes/4),1:(ngenes/4)] <- mask</pre>
gp1 <- mvrnorm((nsamples/2), zero_vector, cov_mtrx1)</pre>
gp2 <- mvrnorm((nsamples/2), zero_vector, cov_mtrx2)</pre>
gp <- rbind(gp1,gp2)</pre>
dataset <- aperm(gp, c(2,1))</pre>
## first 20 samples belong to group 1
## second 20 samples belong to group 2
pvalue <- GSNCAtest(object=dataset, group=c(rep(1,20),rep(2,20)))</pre>
```

16 HDP.ranking

HDP.ranking

High Directed Preorder Ranking of MST

#### **Description**

Rank nodes in an object of class igraph (see package igraph for the definition of class igraph) containing a minimum spanning tree (MST) according to the High Directed Preorder traversal of the tree.

## Usage

HDP.ranking(object)

#### **Arguments**

object

object of class igraph that consists of a minimum spanning tree.

#### **Details**

Rank nodes in an object of class igraph (see package igraph) containing a minimum spanning tree (MST). The MST is rooted at a node with the largest geodesic distance and the rest of the nodes are ranked according to the high directed preorder (HDP) traversal of the tree (Friedman and Rafsky 1979).

## Value

Numeric vector giving the node ranks according to HDP traversal of the MST.

#### Note

This function does not work properly if there is any node in the MST with more than 26 links. However, this situation is almost impossible for a dataset composed of a few hundreds or less of samples.

## Author(s)

Yasir Rahmatallah and Galina Glazko

#### References

Rahmatallah Y., Emmert-Streib F. and Glazko G. (2012) Gene set analysis for self-contained tests: complex null and specific alternative hypotheses. Bioinformatics **28**, 3073–3080.

Friedman J. and Rafsky L. (1979) Multivariate generalization of the Wald-Wolfowitz and Smirnov two-sample tests. Ann. Stat. 7, 697–717.

#### See Also

radial.ranking, KStest, MDtest.

KStest 17

#### **Examples**

```
## generate random data using normal distribution
## generate 20 features in 20 samples
object <- matrix(rnorm(400),20,20)
objt <- aperm(object, c(2,1))
## calculate the weight matrix
Wmat <- as.matrix(dist(objt, method = "euclidean", diag = TRUE, upper = TRUE, p = 2))
## create a weighted undirectional graph from the weight matrix
gr <- graph_from_adjacency_matrix(Wmat, weighted = TRUE, mode = "undirected")
## find the minimum spanning tree
MST <- mst(gr)
HDP.ranks <- HDP.ranking(MST)</pre>
```

**KStest** 

Multivariate Kolmogorov-Smirnov Test of Means

#### **Description**

Performs two-sample nonparametric multivariate test of means based on the minimum spanning tree (MST) and Kolmogorov-Smirnov statistic. It tests the null hypothesis that a set of features has the same mean in two conditions versus different means.

#### Usage

```
KStest(object, group, nperm=1000, pvalue.only=TRUE)
```

## **Arguments**

object	a numeric matrix with columns and rows respectively corresponding to samples and features.
group	a numeric vector indicating group associations for samples. Possible values are 1 and 2.
nperm	number of permutations used to estimate the null distribution of the test statistic. If not given, a default value 1000 is used.
pvalue.only	logical. If TRUE (default), the p-value is returned. If FALSE a list of length three

logical. If TRUE (default), the p-value is returned. If FALSE a list of length three containing the observed statistic, the vector of permuted statistics, and the p-

value is returned.

#### **Details**

This function tests the null hypothesis that a set of features has no shift between two conditions. It performs a two-sample nonparametric multivariate test based on the minimum spanning tree (MST) and Kolmogorov-Smirnov statistic as proposed by Friedman and Rafsky (1979). The MST of the weighted undirectional graph created from the samples is found. The nodes of the MST are ranked based on their position in the MST. The MST is rooted at the node with largest geodisic distance (rank 1) and then nodes are ranked in the High Directed Preorder (HDP) traversal of the tree (Rahmatallah et. al. 2012). The quantity  $d_i = (r_i/n_1) - (s_i/n_2)$  is calculated where  $r_i(s_i)$ 

18 KStest

is the number of nodes (samples) from condition 1(2) which ranked lower than i,  $1 \le i \le N$  and N is the total number of samples. The Kolmogorov-Smirnov statistic is given by the maximum absolute difference  $D = \sqrt{\frac{n_1 n_2}{n_1 + n_2}} max |d_i|$ . The performance of this test under different alternative hypotheses was thoroughly examind in Rahmatallah et. al. (2012). The null distribution of the test statistic is estimated by permuting sample labels nperm times and calculating the test statistic for each. P-value is calculated as

$$p.value = \frac{\sum_{k=1}^{nperm} I[D_k \ge D_{obs}] + 1}{nperm + 1}$$

where  $D_k$  is the test statistic for permutation k,  $D_{obs}$  is the observed test statistic, and I is the indicator function.

#### Value

When pvalue.only=TRUE (default), function KStest returns the p-value indicating the attained significance level. When pvalue.only=FALSE, function KStest produces a list of length 3 with the following components:

statistic the value of the observed test statistic.

perm. stat numeric vector of the resulting test statistic for nperm random permutations of

sample labels.

p.value p-value indicating the attained significance level.

## Note

This function invokes function HDP.ranking which does not work properly if there is any node in the MST with more than 26 links. However, this situation is almost impossible for a dataset composed of a few hundreds or less of samples.

#### Author(s)

Yasir Rahmatallah and Galina Glazko

#### References

Rahmatallah Y., Emmert-Streib F. and Glazko G. (2012) Gene set analysis for self-contained tests: complex null and specific alternative hypotheses. Bioinformatics **28**, 3073–3080.

Friedman J. and Rafsky L. (1979) Multivariate generalization of the Wald-Wolfowitz and Smirnov two-sample tests. Ann. Stat. 7, 697–717.

#### See Also

MDtest, ADtest, CVMtest, WWtest, RKStest, RMDtest, RADtest, RCVMtest, HDP.ranking.

MDtest 19

#### **Examples**

```
## generate a feature set of length 20 in two conditions
## each condition has 20 samples
## use multivariate normal distribution
library(MASS)
ngenes <- 20
nsamples <- 40
## let the mean vector have zeros of length 20 in both conditions
zero_vector <- array(0,c(1,ngenes))</pre>
## set the covariance matrix to be an identity matrix for both conditions
cov_mtrx <- diag(ngenes)</pre>
gp <- mvrnorm(nsamples, zero_vector, cov_mtrx)</pre>
## apply a mean shift of 3 to half of the features under condition 1
gp[1:20,1:10] <- gp[1:20,1:10] + 3
dataset <- aperm(gp, c(2,1))</pre>
## first 20 samples belong to condition 1
## second 20 samples belong to condition 2
pvalue <- KStest(object=dataset, group=c(rep(1,20),rep(2,20)))</pre>
```

MDtest

Multivariate Mean Deviation Test of Means

## Description

Performs two-sample nonparametric multivariate test of means based on the minimum spanning tree (MST). It calculates the mean deviation between the cumulative distribution functions (CDFs) of sample ranks in two conditions. It tests the null hypothesis that a set of features has the same mean in two conditions versus different means.

## Usage

```
MDtest(object, group, nperm=1000, pvalue.only=TRUE)
```

#### **Arguments**

nperm

object	a numeric matrix with columns and rows respectively corresponding to samples and features.
group	a numeric vector indicating group associations for samples. Possible values are 1 and 2.

number of permutations used to estimate the null distribution of the test statistic.

If not given, a default value 1000 is used.

pvalue.only logical. If TRUE (default), the p-value is returned. If FALSE a list of length three

containing the observed statistic, the vector of permuted statistics, and the p-

value is returned.

#### **Details**

This function tests the null hypothesis that a set of features has no difference in mean (shift) between two conditions. It performs a two-sample nonparametric multivariate test by ranking samples based on the minimum spanning tree (MST) as proposed by Friedman and Rafsky (1979). The MST of the weighted undirectional graph created from the samples is found. The nodes of the MST are ranked based on their position in the MST. The MST is rooted at the node with largest geodisic distance and then nodes are ranked in the High Directed Preorder (HDP) traversal of the tree (Rahmatallah et. al. 2012). The mean deviation between the cumulative distribution functions (CDFs) of sample ranks in two conditions is calculated. The null distribution of the test statistic is estimated by permuting sample labels nperm times and calculating the test statistic for each. P-value is calculated as

$$p.value = \frac{\sum_{k=1}^{nperm} I\left[|D_k| \ge |D_{obs}|\right] + 1}{nperm + 1}$$

where  $D_k$  is the test statistic for permutation k,  $D_{obs}$  is the observed test statistic, and I is the indicator function. This statistic was introduced for a single-sample version of Gene Set Enrichment Analysis (ssGSEA) in Barbie et al. (2009) to estimate enrichment scores for gene sets. It was repurposed in package GSAR to test if the mean deviation between the empirical CDFs of sample ranks of two groups in the MST is significant.

#### Value

When pvalue.only=TRUE (default), function MDtest returns the p-value indicating the attained significance level. When pvalue.only=FALSE, function MDtest produces a list of length 3 with the following components:

statistic the value of the observed test statistic.

perm. stat numeric vector of the resulting test statistic for nperm random permutations of

sample labels.

p. value p-value indicating the attained significance level.

#### Note

This function invokes function HDP.ranking which does not work properly if there is any node in the MST with more than 26 links. However, this situation is almost impossible for a dataset composed of a few hundreds or less of samples.

## Author(s)

Yasir Rahmatallah and Galina Glazko

#### References

Rahmatallah Y., Emmert-Streib F. and Glazko G. (2012) Gene set analysis for self-contained tests: complex null and specific alternative hypotheses. Bioinformatics **28**, 3073–3080.

Barbie D., Tamayo P., Boehm J., et al. (2009) Systematic RNA interference reveals that oncogenic KRAS-driven cancers require TBK1. Nature **462**, 108–112.

Friedman J. and Rafsky L. (1979) Multivariate generalization of the Wald-Wolfowitz and Smirnov two-sample tests. Ann. Stat. 7, 697–717.

p53DataSet 21

#### See Also

KStest, ADtest, CVMtest, WWtest, RKStest, RMDtest, RADtest, RCVMtest, HDP.ranking.

#### **Examples**

```
## generate a feature set of length 20 in two conditions
## each condition has 20 samples
## use multivariate normal distribution
library(MASS)
ngenes <- 20
nsamples <- 40
## let the mean vector have zeros of length 20 both conditions
zero_vector <- array(0,c(1,ngenes))</pre>
## set the covariance matrix to be an identity matrix for both conditions
cov_mtrx <- diag(ngenes)</pre>
gp <- mvrnorm(nsamples, zero_vector, cov_mtrx)</pre>
## apply a mean shift of 3 to half of the features under condition 2
gp[1:20,1:10] \leftarrow gp[1:20,1:10] + 3
dataset \leftarrow aperm(gp, c(2,1))
## first 20 samples belong to group 1
## second 20 samples belong to group 2
pvalue <- MDtest(object=dataset, group=c(rep(1,20),rep(2,20)))</pre>
```

p53DataSet

p53 Dataset of the NCI-60 Cell Lines

#### **Description**

A matrix of gene expression profiles for a processed version of the p53 dataset obtained from the NCI-60 cell lines using the hgu95av2 microarray platform.

## Usage

```
data(p53DataSet)
```

#### **Format**

A matrix of 8655 rows and 50 columns where rows correspond to genes and columns correspond to samples. Gene symbol identifiers are used for rows. Column names indicate the class of the samples (wild type p53 or mutated p53) with the first 17 column names starting with WT1 and ending with WT17 and next 33 column names starting with MUT1 and ending with MUT33.

#### Details

p53 is a major tumor suppressor protein. The p53 dataset comprises 50 samples of NCI-60 cell lines differentiated based on the status of the TP53 gene: 17 cell lines carrying wild type (WT) TP53 and 33 cell lines carrying mutated (MUT) TP53 (Olivier et. al. 2002, Subramanian et. al. 2005). Transcriptional profiles obtained from microarrays of platform hgu95av2 were obtained from the available datasets at the GSEA Broad Institute's website.

22 plotMST2.pathway

Probe level intensities were quantile normalized and transformed to the log scale using log2(1 + intensity). Probes originally had Affymetrix identifiers which are mapped to unique gene symbol identifiers. Probes without mapping to entrez and gene symbol identifiers were discarded. Probes with duplicate intensities were assessed and the probe with the largest absolute value of t-statistic between WT and MUT conditions was selected as the gene match. Genes were assigned gene symbol identifiers and columns were assigned names indicating weither they belong to WT or MUT condition. The columns were sorted such that the first 17 columns are WT samples and the next 33 columns are the MUT samples. p53DataSet was used in the analysis presented in Rahmatallah et. al. 2014.

#### Source

Broad Institute (http://www.broadinstitute.org/gsea/datasets.jsp)

#### References

Rahmatallah Y., Emmert-Streib F. and Glazko G. (2014) Gene sets net correlations analysis (GSNCA): a multivariate differential coexpression test for gene sets. Bioinformatics **30**, 360–368.

Subramanian A., Tamayo P., Mootha V., Mukherjee S., Ebert B., Gillette M., Paulovich A., Pomeroy S., Golub T., Lander E. and Mesirov J. (2005) Gene set enrichment analysis: A knowledge-based approach for interpreting genome-wide expression profiles. Proc. Natl. Acad. Sci. **102**, 15545–15550.

Olivier M., Eeles R., Hollstein M., Khan M., Harris C. and Hainaut P. (2002) The IARC TP53 database: new online mutation analysis and recommendations to users. Hum. Mutat. **19**, 607–614.

## **Examples**

```
data(p53DataSet)
dim(p53DataSet)
```

plotMST2.pathway

Plot MST2 for a pathway in two conditions

#### **Description**

This is a wrapper function which uses function findMST2 to find the union of the first and second minimum spanning trees (or MST2) of the correlation network for a feature set (pathway) under two conditions. It plots the MST2 of the correlation network of the feature set under both conditions side-by-side and highlights hub nodes to facilitate a visual comparison.

## Usage

```
plotMST2.pathway(object, group, name=NULL, cor.method="pearson",
min.sd=1e-3, legend.size=1, leg.x=-0.8, leg.y=1.5, return.weights=FALSE,
group1.name="Group 1", group2.name="Group 2", label.size=1,
label.color="black", label.dist=0.5, vertex.size=8, vertex.label.font=1,
edge.width=1)
```

plotMST2.pathway 23

#### **Arguments**

object a numeric matrix with columns and rows respectively corresponding to samples and features. Gene names are provided to this function as the rownames of this a numeric vector indicating group associations for samples. Possible values are group 1 and 2. name an optional character string giving the name of the feature set (gene set). If given, the name will be displayed at the top of the plot. cor.method a character string indicating which correlation coefficient is to be computed. Possible values are "pearson" (default), "spearman" and "kendall". Default value is "pearson". a numeric value indicating the minimum allowed standard deviation for any min.sd feature. If any feature has a standard deviation smaller than min.sd then the execution stops and an error message is returned. Default value is 1e-3. legend.size an optional numeric value controlling the relative font size of the legend to the default font size. Default is 1. leg.x a numeric value indicating the amount of horizontal shift of the legend box to allow better positioning in the plot. a numeric value indicating the amount of vertical shift of the legend box to allow leg.y better positioning in the plot. return.weights logical. Default value is FALSE. If the weight factors assigned to the genes by the GSNCA method are desired, setting this parameter to TRUE returns the weight factors in a matrix with 2 columns (for class 1 and class 2) and number of rows equal to the number of genes in the gene set. If the rownames of object are provided, then they will be used as rownames for the returned matrix. If the rownames of object are abscent, node labels will be set to as.character(c(1:nrow(object))). group1.name an optional character string to be presented as the given name for class 1 in the plot. Default value is "Group 1" an optional character string to be presented as the given name for class 2 in the group2.name plot. Default value is "Group 2" label.size a numeric value passed to argument vertex.label.cex in command plot.igraph to specify the vertex label size. Default value is 1. label.color a character string specifying the color of vertex labels. Default value is "black". label.dist a numeric value passed to argument vertex.label.dist in command plot.igraph to specify the distance between vertex labels and the centers of vertices. Default value is 0.5. a numeric value passed to argument vertex.size in command plot.igraph to specvertex size ify the vertex size. Default value is 8. vertex.label.font a numeric value passed to argument vertex.label.font in command plot.igraph to specify the used font type. Default value is 1. edge.width a numeric value passed to argument edge. width in command plot.igraph to specify the edge width in the plot.

24 plotMST2.pathway

#### **Details**

This is a wrapper plotting function for the convenience of users. It uses function findMST2 to find the union of the first and second minimum spanning trees (or MST2) of the correlation network for a feature set (pathway) under two conditions and plots them side-by-side. It also lists the hub nodes and their weight factors (w) under each condition (see Rahmatallah et. al. 2014 for details). The range in which weight factors fall is indicated by the node colors defined in the legend. Weight factor have values mostly ranging between 0.5 (low coexpression) and 1.5 (high coexpression). To allow the users more control over plotting parameters and to present different feature sets appropriately, two optional arguments were introduced: legend.size and label.size. Node lables will be the names of the features in the set, i.e. rownames(object). If the rownames attribute is not set for object, node labels will be set to as.character(c(1:nrow(object))).

The weight factors, inferred from the Gene Sets Net Correlations Analysis (GSNCA) method (see GSNCAtest), correlate to some extent with genes centralities in the MST2: genes with large weights are placed near the center of the MST2, and genes with small weights are placed on the periphery (Rahmatallah et. al. 2014). Adopting network terminology, a gene with the largest weight is a hub gene, coexpressed with most of the other genes in a pathway (see findMST2). Therefore, MST2 is a convenient graphical visualization tool to examine the pathways tested by the GSNCA method (see GSNCAtest).

The correlation (coexpression) network is obtained using the weight matrix W with elements  $w_{ij} = 1 - |r_{ij}|$  where  $r_{ij}$  is the correlation between features i and j and  $w_{ij}$  is the weight of the link between vertices (nodes) i and j in the network. The correlation coefficient used is indicated by the argument cor.method with three possible values: "pearson" (default), "spearman" and "kendall".

In some cases (especially for RNA-Seq count data), a feature (or more) may have a constant or nearly constant level across the samples in one or both conditions. This results in a zero or a tiny standard deviation. Such case produces an error in command cor used to compute the correlation coefficients between features. To avoid this situation, standard deviations are checked in advance and if any is found below the minimum limit min.sd (default is 1e-3), the execution stops and an error message is returned indicating the number of feature causing the problem (if only one the index of that feature is given too).

## Note

This function is suitable for a feature set of roughly 80 features or less. It works for feature sets with larger number of features but the placements of nodes and their labels in the plot will be too crowded for a useful visual presentation.

#### Author(s)

Yasir Rahmatallah and Galina Glazko

#### References

Rahmatallah Y., Emmert-Streib F. and Glazko G. (2014) Gene sets net correlations analysis (GSNCA): a multivariate differential coexpression test for gene sets. Bioinformatics **30**, 360–368.

radial.ranking 25

#### See Also

```
findMST2, GSNCAtest.
```

#### **Examples**

```
## generate a feature set of length 20 in two conditions
## each condition has 20 samples
## use multivariate normal distribution with different covariance matrices
library(MASS)
ngenes <- 20
nsamples <- 40
zero_vector <- array(0,c(1,ngenes))</pre>
## create a covariance matrix with low off-diagonal elements
cov_mtrx1 <- diag(ngenes)</pre>
cov_mtrx1[!diag(ngenes)] <- 0.1</pre>
## create a covariance matrix with high off-diagonal elements
## for the first 5 features and low for the rest 15 features
cov_mtrx2 <- diag(ngenes)</pre>
cov_mtrx2[!diag(ngenes)] <- 0.1</pre>
mask <- diag(ngenes/4)</pre>
mask[!diag(ngenes/4)] <- 0.6
cov_mtrx2[1:(ngenes/4),1:(ngenes/4)] <- mask</pre>
gp1 <- mvrnorm((nsamples/2), zero_vector, cov_mtrx1)</pre>
gp2 <- mvrnorm((nsamples/2), zero_vector, cov_mtrx2)</pre>
gp <- rbind(gp1,gp2)</pre>
dataset <- aperm(gp, c(2,1))
## first 20 samples belong to group 1
## second 20 samples belong to group 2
## since rowname(object)=NULL, node labels will be automatically
## set to as.character(c(1:nrow(object)))
plotMST2.pathway(object=dataset, group=c(rep(1,20),rep(2,20)),
name="Example Pathway")
```

radial.ranking

Radial Ranking of MST

## **Description**

Rank vertices in an object of class igraph (see package igraph for the definition of class igraph) that consists of a minimum spanning tree (MST) or the union of multiple MSTs radially such that vertices with higher depth and distance from the centroid are given higher ranks.

## Usage

```
radial.ranking(object)
```

#### **Arguments**

object

object of class igraph that consists of a minimum spanning tree or the union of multiple spanning trees.

26 RADtest

#### **Details**

Rank nodes in an object of class igraph (see package igraph) that consists of a minimum spanning tree (MST) or the union of multiple MSTs radially. The MST is rooted at the node of smallest geodesic distance (centroid) and nodes with largest depths from the root are assigned higher ranks. Hence, ranks are increasing radially from the root of the MST (Friedman and Rafsky 1979).

#### Value

Numeric vector giving the radial node ranks in the MST or union of MSTs.

#### Author(s)

Yasir Rahmatallah and Galina Glazko

#### References

Rahmatallah Y., Emmert-Streib F. and Glazko G. (2012) Gene set analysis for self-contained tests: complex null and specific alternative hypotheses. Bioinformatics **28**, 3073–3080.

Friedman J. and Rafsky L. (1979) Multivariate generalization of the Wald-Wolfowitz and Smirnov two-sample tests. Ann. Stat. 7, 697–717.

#### See Also

```
HDP.ranking, RKStest, RMDtest.
```

## **Examples**

```
## generate random data using normal distribution
## generate 20 features in 20 samples
object <- matrix(rnorm(400),20,20)
objt <- aperm(object, c(2,1))
## calculate the weight matrix
Wmat <- as.matrix(dist(objt, method = "euclidean", diag = TRUE, upper = TRUE, p = 2))
## create a weighted undirectional graph from the weight matrix
gr <- graph_from_adjacency_matrix(Wmat, weighted = TRUE, mode = "undirected")
## find the minimum spanning tree
MST <- mst(gr)
radial.ranks <- radial.ranking(MST)</pre>
```

**RADtest** 

Multivariate Radial Anderson-Darling Test of Variance

## Description

Performs two-sample nonparametric multivariate test of variance based on the minimum spanning tree (MST) and Anderson-Darling statistic. It tests the null hypothesis that a set of features has the same scale in two conditions versus different scales.

RADtest 27

#### Usage

RADtest(object, group, mst.order=1, nperm=1000, pvalue.only=TRUE)

#### **Arguments**

object a numeric matrix with columns and rows respectively corresponding to samples

and features.

group a numeric vector indicating group associations for samples. Possible values are

1 and 2.

mst.order numeric value to indicate the consideration of the union of the first mst.order

MSTs. Default value is 1. Maximum allowed value is 5.

nperm number of permutations used to estimate the null distribution of the test statistic.

If not given, a default value 1000 is used.

pvalue.only logical. If TRUE (default), the p-value is returned. If FALSE a list of length three

containing the observed statistic, the vector of permuted statistics, and the p-

value is returned.

#### **Details**

This function tests the null hypothesis that a set of features has the same scale in two conditions. It performs a two-sample nonparametric multivariate test based on the minimum spanning tree (MST) and Anderson-Darling statistic as proposed by Rahmatallah and Glazko (2024). The MST of the weighted undirectional graph created from the samples is found. The nodes of the MST are ranked based on their position in the MST. The MST is rooted at the node with smallest geodisic distance and nodes with higher depths from the root are assigned higher ranks (radial ranking). The Anderson-Darling statistic can be defined as

$$A = \frac{1}{n_1 n_2} \sum_{i=1}^{N-1} \frac{(r_i N - i n_1)^2}{i(N-i)}$$

where  $r_i$  is the number of nodes (samples) from condition 1 which ranked lower than i,  $1 \le i \le N$ ,  $n_1$  and  $n_2$  are respectively the number of samples in groups 1 and 2, and N is the total number of samples. The performance of this test under different alternative hypotheses was examind in Rahmatallah and Glazko (2024). The null distribution of the test statistic is estimated by permuting sample labels nperm times and calculating the test statistic for each. P-value is calculated as

$$p.value = \frac{\sum_{k=1}^{nperm} I\left[A_k \ge A_{obs}\right] + 1}{nperm + 1}$$

where  $A_k$  is the test statistic for permutation k,  $A_{obs}$  is the observed test statistic, and I is the indicator function.

#### Value

When pvalue.only=TRUE (default), function RKStest returns the p-value indicating the attained significance level. When pvalue.only=FALSE, function RKStest produces a list of length 3 with the following components:

28 RADtest

statistic the value of the observed test statistic.

perm. stat numeric vector of the resulting test statistic for nperm random permutations of sample labels.

p.value p-value indicating the attained significance level.

#### Note

The variance of both the Poisson and negative Bionomial distributions, used to model count data, is a function of their mean. Therefore, using the radial Anderson-Darling test (RADtest) to detect pathways with differential variance for RNA-Seq counts is not recommended without proper data normalization.

#### Author(s)

Yasir Rahmatallah and Galina Glazko

#### References

Rahmatallah Y. and Glazko G. (2025) Improving data interpretability with new differential sample variance tests. BMC Bioinformatics **26**, 103.

Rahmatallah Y., Emmert-Streib F. and Glazko G. (2012) Gene set analysis for self-contained tests: complex null and specific alternative hypotheses. Bioinformatics **28**, 3073–3080.

#### See Also

RKStest, RMDtest, RCVMtest, ADtest, KStest, MDtest, CVMtest, radial.ranking.

## Examples

```
## generate a feature set of length 20 in two conditions
## each condition has 20 samples
## use multivariate normal distribution
library(MASS)
ngenes <- 20
nsamples <- 40
## let the mean vector have zeros of length 20 for both conditions
zero_vector <- array(0,c(1,ngenes))</pre>
## set the covariance matrix to be an identity matrix for condition 1
cov_mtrx <- diag(ngenes)</pre>
gp1 <- mvrnorm((nsamples/2), zero_vector, cov_mtrx)</pre>
## set some scale difference in the covariance matrix for condition 2
cov_mtrx <- cov_mtrx*3</pre>
gp2 <- mvrnorm((nsamples/2), zero_vector, cov_mtrx)</pre>
## combine the data of two conditions into one dataset
gp <- rbind(gp1,gp2)</pre>
dataset <- aperm(gp, c(2,1))</pre>
## first 20 samples belong to group 1
## second 20 samples belong to group 2
pvalue <- RADtest(object=dataset, group=c(rep(1,20),rep(2,20)))</pre>
```

RCVMtest 29

RCVMtest

Multivariate Radial Cramer-Von Mises Test of Variance

#### Description

Performs two-sample nonparametric multivariate test of variance based on the minimum spanning tree (MST) and Cramer-Von Mises statistic. It tests the null hypothesis that a set of features has the same scale in two conditions versus different scales.

#### Usage

RCVMtest(object, group, mst.order=1, nperm=1000, pvalue.only=TRUE)

#### **Arguments**

_	
object	a numeric matrix with columns and rows respectively corresponding to samples and features.
group	a numeric vector indicating group associations for samples. Possible values are 1 and 2.
mst.order	numeric value to indicate the consideration of the union of the first mst.order MSTs. Default value is 1. Maximum allowed value is 5.
nperm	number of permutations used to estimate the null distribution of the test statistic. If not given, a default value 1000 is used.
pvalue.only	logical. If TRUE (default), the p-value is returned. If FALSE a list of length three containing the observed statistic, the vector of permuted statistics, and the p-value is returned.

#### **Details**

This function tests the null hypothesis that a set of features has the same scale in two conditions. It performs a two-sample nonparametric multivariate test based on the minimum spanning tree (MST) and Cramer-Von Mises statistic as proposed by Rahmatallah and Glazko (2024). The MST of the weighted undirectional graph created from the samples is found. The nodes of the MST are ranked based on their position in the MST. The MST is rooted at the node with smallest geodisic distance and nodes are assigned ranks according to their distance from the root (radial ranking) in the MST (Rahmatallah et. al. 2012). The Cramer-Von Mises statistic can be defined as

$$C = \frac{n_1 n_2}{N^2} \left\{ \sum_{i=1}^{n_1} \left[ \frac{r_i}{n_1} - i \left( \frac{1}{n_1} + \frac{1}{n_2} \right) \right]^2 + \sum_{j=1}^{n_2} \left[ \frac{s_j}{n_2} - j \left( \frac{1}{n_1} + \frac{1}{n_2} \right) \right]^2 \right\}$$

where  $r_i$  and  $s_j$  are respectively the number of samples from conditions 1 and 2 which ranked lower than  $i, 1 \le i \le N$ ,  $n_1$  and  $n_2$  are respectively the number of samples in groups 1 and 2, and N is the total number of samples. The performance of this test under different alternative hypotheses was examind in Rahmatallah and Glazko (2024). The null distribution of the test statistic is estimated by

30 RCVMtest

permuting sample labels nperm times and calculating the test statistic for each. P-value is calculated as

$$p.value = \frac{\sum_{k=1}^{nperm} I\left[C_k \ge C_{obs}\right] + 1}{nperm + 1}$$

where  $C_k$  is the test statistic for permutation k,  $C_{obs}$  is the observed test statistic, and I is the indicator function.

#### Value

When pvalue.only=TRUE (default), function RKStest returns the p-value indicating the attained significance level. When pvalue.only=FALSE, function RKStest produces a list of length 3 with the following components:

statistic the value of the observed test statistic.

perm. stat numeric vector of the resulting test statistic for nperm random permutations of

sample labels.

p. value p-value indicating the attained significance level.

#### Note

The variance of both the Poisson and negative Bionomial distributions, used to model count data, is a function of their mean. Therefore, using the radial Anderson-Darling test (RADtest) to detect pathways with differential variance for RNA-Seq counts is not recommended without proper data normalization.

## Author(s)

Yasir Rahmatallah and Galina Glazko

#### References

Rahmatallah Y. and Glazko G. (2025) Improving data interpretability with new differential sample variance tests. BMC Bioinformatics **26**, 103.

Rahmatallah Y., Emmert-Streib F. and Glazko G. (2012) Gene set analysis for self-contained tests: complex null and specific alternative hypotheses. Bioinformatics **28**, 3073–3080.

## See Also

```
CVMtest, RKStest, RMDtest, RADtest, KStest, MDtest, ADtest, radial.ranking.
```

## Examples

```
## generate a feature set of length 20 in two conditions
## each condition has 20 samples
## use multivariate normal distribution
library(MASS)
ngenes <- 20
nsamples <- 40</pre>
```

RKStest 31

```
## let the mean vector have zeros of length 20 for both conditions
zero_vector <- array(0,c(1,ngenes))
## set the covariance matrix to be an identity matrix for condition 1
cov_mtrx <- diag(ngenes)
gp1 <- mvrnorm((nsamples/2), zero_vector, cov_mtrx)
## set some scale difference in the covariance matrix for condition 2
cov_mtrx <- cov_mtrx*3
gp2 <- mvrnorm((nsamples/2), zero_vector, cov_mtrx)
## combine the data of two conditions into one dataset
gp <- rbind(gp1,gp2)
dataset <- aperm(gp, c(2,1))
## first 20 samples belong to group 1
## second 20 samples belong to group 2
pvalue <- RCVMtest(object=dataset, group=c(rep(1,20),rep(2,20)))</pre>
```

RKStest

Multivariate Radial Kolmogorov-Smirnov Test of Variance

## **Description**

Performs two-sample nonparametric multivariate test of variance based on the minimum spanning tree (MST) and Kolmogorov-Smirnov statistic. It tests the null hypothesis that a set of features has the same scale in two conditions versus different scales.

#### Usage

```
RKStest(object, group, mst.order=1, nperm=1000, pvalue.only=TRUE)
```

#### **Arguments**

object	a numeric matrix with columns and rows respectively corresponding to samples and features.
group	a numeric vector indicating group associations for samples. Possible values are 1 and 2. $$
mst.order	numeric value to indicate the consideration of the union of the first mst.order MSTs. Default value is 1. Maximum allowed value is 5.
nperm	number of permutations used to estimate the null distribution of the test statistic. If not given, a default value 1000 is used.
pvalue.only	logical. If TRUE (default), the p-value is returned. If FALSE a list of length three containing the observed statistic, the vector of permuted statistics, and the p-value is returned.

#### **Details**

This function tests the null hypothesis that a set of features has the same scale in two conditions. It performs a two-sample nonparametric multivariate test based on the minimum spanning tree (MST) and Kolmogorov-Smirnov statistic as proposed by Friedman and Rafsky (1979). The MST of the weighted undirectional graph created from the samples is found. The nodes of the MST are ranked based on their position in the MST. The MST is rooted at the node with smallest geodisic distance (rank 1) and nodes with higher depths from the root are assigned higher ranks. The quantity  $d_i = (r_i/n_1) - (s_i/n_2)$  is calculated where  $r_i(s_i)$  is the number of nodes (samples) from condition 1(2) which ranked lower than  $i, 1 \le i \le N$  and N is the total number of samples. The Kolmogorov-Smirnov statistic is given by the maximum absolute difference  $D = \sqrt{\frac{n_1 n_2}{n_1 + n_2} max |d_i|}$ . The performance of this test under different alternative hypotheses was thoroughly examind in Rahmatallah et. al. (2012). The null distribution of the test statistic is estimated by permuting sample labels nperm times and calculating the test statistic for each. P-value is calculated as

$$p.value = \frac{\sum_{k=1}^{nperm} I[D_k \ge D_{obs}] + 1}{nperm + 1}$$

where  $D_k$  is the test statistic for permutation k,  $D_{obs}$  is the observed test statistic, and I is the indicator function.

#### Value

When pvalue.only=TRUE (default), function RKStest returns the p-value indicating the attained significance level. When pvalue.only=FALSE, function RKStest produces a list of length 3 with the following components:

statistic the value of the observed test statistic.

perm. stat numeric vector of the resulting test statistic for nperm random permutations of

sample labels.

p. value p-value indicating the attained significance level.

## Note

The variance of both the Poisson and negative Bionomial distributions, used to model count data, is a function of their mean. Therefore, using the radial Kolmogorov-Smirnov test (RKStest) to detect pathways with differential variance for RNA-Seq counts is not recommended without proper data normalization.

#### Author(s)

Yasir Rahmatallah and Galina Glazko

#### References

Rahmatallah Y., Emmert-Streib F. and Glazko G. (2012) Gene set analysis for self-contained tests: complex null and specific alternative hypotheses. Bioinformatics **28**, 3073–3080.

Friedman J. and Rafsky L. (1979) Multivariate generalization of the Wald-Wolfowitz and Smirnov two-sample tests. Ann. Stat. 7, 697–717.

RMDtest 33

#### See Also

RMDtest, RADtest, RCVMtest, WWtest, KStest, MDtest, ADtest, CVMtest.

#### **Examples**

```
## generate a feature set of length 20 in two conditions
## each condition has 20 samples
## use multivariate normal distribution
library(MASS)
ngenes <- 20
nsamples <- 40
## let the mean vector have zeros of length 20 for both conditions
zero_vector <- array(0,c(1,ngenes))</pre>
## set the covariance matrix to be an identity matrix for condition 1
cov_mtrx <- diag(ngenes)</pre>
gp1 <- mvrnorm((nsamples/2), zero_vector, cov_mtrx)</pre>
## set some scale difference in the covariance matrix for condition 2
cov_mtrx <- cov_mtrx*3
gp2 <- mvrnorm((nsamples/2), zero_vector, cov_mtrx)</pre>
## combine the data of two conditions into one dataset
gp <- rbind(gp1,gp2)</pre>
dataset <- aperm(gp, c(2,1))</pre>
## first 20 samples belong to group 1
## second 20 samples belong to group 2
pvalue <- RKStest(object=dataset, group=c(rep(1,20),rep(2,20)))</pre>
```

RMDtest

Multivariate Radial Mean Deviation Test of Variance

#### **Description**

Performs two-sample nonparametric multivariate test of variance based on the minimum spanning tree (MST). It calculates the mean deviation between the cumulative distribution functions (CDFs) of sample ranks in two conditions. It tests the null hypothesis that a set of features has the same variance (scale) in two conditions versus different variances.

#### Usage

```
RMDtest(object, group, mst.order=1, nperm=1000, pvalue.only=TRUE)
```

#### **Arguments**

object	a numeric matrix with columns and rows respectively corresponding to samples and features.
group	a numeric vector indicating group associations for samples. Possible values are 1 and 2.
mst.order	numeric value to indicate the consideration of the union of the first mst.order MSTs. Default value is 1. Maximum allowed value is 5.

34 RMDtest

nperm number of permutations used to estimate the null distribution of the test statistic.

If not given, a default value 1000 is used.

pvalue.only logical. If TRUE (default), the p-value is returned. If FALSE a list of length three

containing the observed statistic, the vector of permuted statistics, and the p-

value is returned.

#### **Details**

This function tests the null hypothesis that a set of features has the same scale in two conditions. It performs a two-sample nonparametric multivariate test based on the minimum spanning tree (MST) as proposed by Friedman and Rafsky (1979). The MST of the weighted undirectional graph created from the samples is found. The nodes of the MST are ranked based on their position in the MST. The MST is rooted at the node with smallest geodisic distance (rank 1) and nodes with higher depths from the root are assigned higher ranks. The mean deviation between the cumulative distribution functions (CDFs) of sample ranks in two conditions is calculated. The null distribution of the test statistic is estimated by permuting sample labels nperm times and calculating the test statistic for each. P-value is calculated as

$$p.value = \frac{\sum_{k=1}^{nperm} I\left[|D_k| \ge |D_{obs}|\right] + 1}{nperm + 1}$$

where  $D_k$  is the test statistic for permutation k,  $D_{obs}$  is the observed test statistic, and I is the indicator function. This statistic was introduced for a single-sample version of Gene Set Enrichment Analysis (ssGSEA) in Barbie et al. (2009) to estimate enrichment scores for gene sets. It was repurposed in package GSAR to test if the mean deviation between the empirical CDFs of sample ranks of two groups in the MST is significant.

#### Value

When pvalue.only=TRUE (default), function RMDtest returns the p-value indicating the attained significance level. When pvalue.only=FALSE, function RMDtest produces a list of length 3 with the following components:

statistic the value of the observed test statistic.

perm. stat numeric vector of the resulting test statistic for nperm random permutations of

sample labels.

p. value p-value indicating the attained significance level.

## Note

The variance of both the Poisson and negative Bionomial distributions, used to model count data, is a function of their mean. Therefore, using the radial mean deviation test (RMDtest) to detect pathways with differential variance for RNA-Seq counts is not recommended without proper data normalization.

#### Author(s)

Yasir Rahmatallah and Galina Glazko

TestGeneSets 35

#### References

Rahmatallah Y., Emmert-Streib F. and Glazko G. (2012) Gene set analysis for self-contained tests: complex null and specific alternative hypotheses. Bioinformatics **28**, 3073–3080.

Barbie D., Tamayo P., Boehm J., et al. (2009) Systematic RNA interference reveals that oncogenic KRAS-driven cancers require TBK1. Nature **462**, 108–112.

Friedman J. and Rafsky L. (1979) Multivariate generalization of the Wald-Wolfowitz and Smirnov two-sample tests. Ann. Stat. 7, 697–717.

#### See Also

```
RKStest, RADtest, RCVMtest, WWtest, KStest, MDtest, ADtest, CVMtest.
```

#### **Examples**

```
## generate a feature set of length 20 in two conditions
## each condition has 20 samples
## use multivariate normal distribution
library(MASS)
ngenes <- 20
nsamples <- 40
## let the mean vector have zeros of length 20 for both conditions
zero_vector <- array(0,c(1,ngenes))</pre>
## set the covariance matrix to be an identity matrix for condition 1
cov_mtrx <- diag(ngenes)</pre>
gp1 <- mvrnorm((nsamples/2), zero_vector, cov_mtrx)</pre>
## set some scale difference in the covariance matrix for condition 2
cov_mtrx <- cov_mtrx*3</pre>
gp2 <- mvrnorm((nsamples/2), zero_vector, cov_mtrx)</pre>
## combine the data of two conditions into one dataset
gp <- rbind(gp1,gp2)</pre>
dataset \leftarrow aperm(gp, c(2,1))
## first 20 samples belong to group 1
## second 20 samples belong to group 2
pvalue <- RMDtest(object=dataset, group=c(rep(1,20),rep(2,20)))</pre>
```

TestGeneSets

Test a List of Gene Sets Using a Specific Statistical Method

## **Description**

A wrapper function that invokes a specific statistical method from the ones available in package GSAR (see Rahmatallah and Glazko 2024, Rahmatallah et. al. 2014, and Rahmatallah et. al. 2012 for details) to test a list of gene sets in a sequential order and returns results in a list object.

#### Usage

```
TestGeneSets(object, group, geneSets=NULL, min.size=10, max.size=500,
test=NULL, nperm=1000, mst.order=1, pvalue.only=TRUE)
```

36 TestGeneSets

## **Arguments**

object	a numeric matrix with columns and rows respectively corresponding to samples and features.
group	a numeric vector indicating group associations for samples. Possible values are 1 and 2.
geneSets	a list of character vectors providing the identifiers of features to be considered in each gene set.
min.size	a numeric value indicating the minimum allowed gene set size. Default value is 10.
max.size	a numeric value indicating the maximum allowed gene set size. Default value is 500.
test	a character parameter indicating which statistical method to use for testing the gene sets. Must be one of "GSNCAtest", "WWtest", "KStest", "MDtest", "ADtest" ("CVMtest", "RKStest", "RMDtest", "RADtest", or "RCVMtest".
nperm	number of permutations used to estimate the null distribution of the test statistic. If not given, a default value 1000 is used.
mst.order	numeric value to indicate the consideration of the union of the first mst.order MSTs when "RKStest", "RMDtest", "RADtest", or "RCVMtest" are used. Default value is 1. Maximum allowed value is 5.
pvalue.only	logical. If TRUE (default), the p-value is returned. If FALSE a list of length three containing the observed statistic, the vector of permuted statistics, and the p-value is returned.

#### **Details**

This is a wrapper function that facilitates the use of any statistical method in package GSAR for multiple gene sets that are provided in a list object. The function filters out any gene that is abscent in the considered data (input parameter object) from the gene sets and discard any set that is too small in size (has less than min.size genes) or too large (has more than max.size genes). The function performs the specified method for all the remaining gene sets in a sequential order and return results in a list object.

#### Value

A list object of length equals the length of the provided gene set list. When pvalue.only=TRUE (default), each item in the returned list by function TestGeneSets consists of a numeric p-value indicating the attained significance level obtained by the specified method. When pvalue.only=FALSE, each item in the returned list is a list of length 3 with the following components:

statistic the value of the observed test statistic.

perm.stat numeric vector of the resulting test statistic for nperm random permutations of sample labels.

p.value p-value indicating the attained significance level.

## Author(s)

Yasir Rahmatallah and Galina Glazko

TestGeneSets 37

#### References

Rahmatallah Y. and Glazko G. (2024) Gene Set Analysis: improving data interpretability with new differential variance tests. 09 September 2024, PREPRINT (Version 1) available at Research Square [https://doi.org/10.21203/rs.3.rs-4888767/v1].

Rahmatallah Y., Emmert-Streib F. and Glazko G. (2014) Gene sets net correlations analysis (GSNCA): a multivariate differential coexpression test for gene sets. Bioinformatics **30**, 360–368.

Rahmatallah Y., Emmert-Streib F. and Glazko G. (2012) Gene set analysis for self-contained tests: complex null and specific alternative hypotheses. Bioinformatics **28**, 3073–3080.

#### See Also

```
GSNCAtest, WWtest, MDtest, KStest, RKStest, RMDtest.
```

## **Examples**

```
## generate a feature set of size 50 in two conditions
## where each condition has 20 samples
## use multivariate normal distribution
library(MASS)
ngenes <- 50
nsamples <- 40
## let the mean vector have zeros of length 50 for both conditions
zero_vector <- array(0,c(1,ngenes))</pre>
## set the covariance matrix to be an identity matrix for both conditions
cov_mtrx <- diag(ngenes)</pre>
gp <- mvrnorm(nsamples, zero_vector, cov_mtrx)</pre>
## apply a mean shift of 5 to the first 10 features under condition 1
gp[1:20,1:10] \leftarrow gp[1:20,1:10] + 5
dataset <- aperm(gp, c(2,1))</pre>
## assign a unique identifier to each gene
rownames(dataset) <- as.character(c(1:ngenes))</pre>
## first 20 samples belong to condition 1
## second 20 samples belong to condition 2
sample.labels \leftarrow c(rep(1,20),rep(2,20))
## construct 3 named gene sets such that they respectively consist of
## genes 1 to 20, 11 to 40, and 31 to 50. Notice that gene sets
## can have intersections and can be of different sizes
## Sine only the first 10 genes have a significant difference between
## the two conditions the only the first gene set (set1) returns a
## small p-value when KStest is selected
geneSets <- list("set1"=as.character(c(1:20)), "set2"=as.character(c(11:40)),</pre>
"set3"=as.character(c(31:40)))
results <- TestGeneSets(object=dataset, group=sample.labels,</pre>
geneSets=geneSets, test="KStest")
```

38 WWtest

WWtest

Multivariate Generalization of the Wald-Wolfowitz Runs Test

## **Description**

Performs two-sample nonparametric multivariate generalization of the Wald-Wolfowitz runs test based on the minimum spanning tree (MST). It tests the alternative hypothesis that a set of features has different distributions in two conditions against the null hypothesis of having the same distribution.

#### Usage

WWtest(object, group, nperm=1000, pvalue.only=TRUE)

#### **Arguments**

object	a numeric matrix with columns and rows respectively corresponding to samples and features.
group	a numeric vector indicating group associations for samples. Possible values are 1 and 2. $$
nperm	number of permutations used to estimate the null distribution of the test statistic. If not given, a default value 1000 is used.
pvalue.only	logical. If TRUE (default), the p-value is returned. If FALSE a list of length three containing the observed statistic, the vector of permuted statistics, and the p-value is returned.

#### **Details**

This function tests the alternative hypothesis that a set of features has different distributions in two conditions against the null hypothesis of having the same distribution. It performs the two-sample nonparametric multivariate generalization of the Wald-Wolfowitz runs test based on the minimum spanning tree (MST) as proposed by Friedman and Rafsky (1979). The performance of this test under different alternative hypotheses was thoroughly examind in Rahmatallah et. al. (2012). The null distribution of the test statistic is estimated by permuting sample labels nperm times and calculating the test statistic for each. P-value is calculated as

$$p.value = \frac{\sum_{k=1}^{nperm} I\left[W_k \le W_{obs}\right] + 1}{nperm + 1}$$

where  $W_k$  is the test statistic for permutation k,  $W_{obs}$  is the observed test statistic, and I is the indicator function.

WWtest 39

#### Value

When pvalue.only=TRUE (default), function WWtest returns the p-value indicating the attained significance level. When pvalue.only=FALSE, function WWtest produces a list of length 3 with the following components:

statistic the value of the observed test statistic.

perm. stat numeric vector of the resulting test statistic for nperm random permutations of

sample labels.

p.value p-value indicating the attained significance level.

## Author(s)

Yasir Rahmatallah and Galina Glazko

#### References

Rahmatallah Y., Emmert-Streib F. and Glazko G. (2012) Gene set analysis for self-contained tests: complex null and specific alternative hypotheses. Bioinformatics **28**, 3073–3080.

Friedman J. and Rafsky L. (1979) Multivariate generalization of the Wald-Wolfowitz and Smirnov two-sample tests. Ann. Stat. 7, 697–717.

#### See Also

KStest, RKStest, MDtest, RMDtest, ADtest, RADtest, CVMtest, RCVMtest.

## Examples

```
## generate a feature set of length 20 in two conditions
## each condition has 20 samples
## use multivariate normal distribution
library(MASS)
ngenes <- 20
nsamples <- 40
## let the mean vector have zeros of length 20 for condition 1
zero_vector <- array(0,c(1,ngenes))</pre>
## let the mean vector have 2s of length 20 for condition 2
mu_vector <- array(2,c(1,ngenes))</pre>
## set the covariance matrix to be an identity matrix
cov_mtrx <- diag(ngenes)</pre>
gp1 <- mvrnorm((nsamples/2), zero_vector, cov_mtrx)</pre>
gp2 <- mvrnorm((nsamples/2), mu_vector, cov_mtrx)</pre>
## combine the data of two conditions into one dataset
gp <- rbind(gp1,gp2)</pre>
dataset \leftarrow aperm(gp, c(2,1))
## first 20 samples belong to group 1
## second 20 samples belong to group 2
pvalue <- WWtest(object=dataset, group=c(rep(1,20),rep(2,20)))</pre>
```

## **Index**

* arith	GSAR-package, 2
HDP.ranking, 16	
radial.ranking,25	ADtest, 2, 3, 9, 18, 21, 28, 30, 33, 35, 39
* datasets	AggrFtest, 2, 5
p53DataSet, 21	CVMtest, 2, 5, 7, 18, 21, 28, 30, 33, 35, 39
* dplot	Cyritest, 2, 3, 7, 10, 21, 20, 30, 33, 33, 39
findMST2, 9	findMST2, 2, 9, 15, 22, 24, 25
findMST2.PPI, 11	findMST2.PPI, 2, 11
plotMST2.pathway, 22	, ,
* graphs	GSAR (GSAR-package), 2
findMST2,9	GSAR-package, 2
findMST2.PPI, 11	GSNCAtest, 2, 11, 12, 13, 24, 25, 37
plotMST2.pathway, 22	
* multivariate	HDP.ranking, 4, 5, 9, 16, 18, 20, 21, 26
ADtest, 3	
CVMtest, 7	igraph, 3, 10, 12, 16, 25, 26
GSNCAtest, 13	KStest, 2, 5, 9, 16, 17, 21, 28, 30, 33, 35, 37,
KStest, 17	39
MDtest, 19	37
RADtest, 26	MDtest, 2, 5, 9, 16, 18, 19, 28, 30, 33, 35, 37,
RCVMtest, 29	39
RKStest, 31	
RMDtest, 33	p53DataSet, 21
TestGeneSets, 35	plotMST2.pathway, 2, 11, 12, 15, 22
WWtest, 38	
* nonparametric	radial.ranking, 16, 25, 28, 30
ADtest, 3	RADtest, 2, 5, 7, 9, 18, 21, 26, 30, 33, 35, 39
AggrFtest, 5	RCVMtest, 2, 5, 7, 9, 18, 21, 28, 29, 33, 35, 39
CVMtest, 7	RKStest, 2, 5, 7, 9, 18, 21, 26, 28, 30, 31, 35,
GSNCAtest, 13	37, 39
KStest, 17	RMDtest, 2, 5, 7, 9, 18, 21, 26, 28, 30, 33, 33,
MDtest, 19	37, 39
RADtest, 26	TestGeneSets, 35
RCVMtest, 29	restdenesets, 33
RKStest, 31	WWtest, 18, 21, 33, 35, 37, 38
RMDtest, 33	
TestGeneSets, 35	
WWtest, 38	
* package	