

Package ‘GSAR’

October 15, 2023

Type Package

Title Gene Set Analysis in R

Version 1.34.0

Date 2023-2-14

Author Yasir Rahmatallah <yrahmatallah@uams.edu>, Galina Glazko <gvg1azko@uams.edu>

Maintainer Yasir Rahmatallah <yrahmatallah@uams.edu>, Galina Glazko <gvg1azko@uams.edu>

Depends 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

Description 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 RELEASE_3_17

git_last_commit 8c3f047

git_last_commit_date 2023-04-25

Date/Publication 2023-10-15

R topics documented:

GSAR-package	2
AggrFtest	3
findMST2	5
findMST2.PPI	7
GSNCAtest	9
HDP.ranking	11

KStest	13
MDtest	15
p53DataSet	17
plotMST2.pathway	18
radial.ranking	21
RKStest	22
RMDtest	24
TestGeneSets	26
WWtest	29

Index	31
--------------	-----------

GSAR-package

Gene Set Analysis in R

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` and `MDtest`), scale (functions `RKStest`, `RMDtest`, 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 union 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 union of the first and second MSTs of PPI networks. Some of the methods available in this package were proposed in Rahmatallah et. al. 2014 and Friedman and Rafsky 1979. The performance of different methods available in this package was thoroughly tested using simulated data and microarray datasets in Rahmatallah et. al. 2012 and Rahmatallah et. al. 2014. These methods can also 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 or sequencing depth). However, because the count data often follows the negative binomial distribution, special attention should be paid to applying the variance tests (`RKStest`, `RMDtest`, and `AggrFtest`). The variance of the negative binomial distribution is proportional to its mean and multivariate tests of variance designed specifically for RNA-seq count data are virtually unavailable. The performance of variance tests in this package with count data highly depends on the used normalization and remains currently 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>

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.
- 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.
- 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](#).

AggrFtest	<i>Aggregated F-Test of Variance Using Fisher's Probability Combining Method</i>
-----------	--

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.
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 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

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

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:

<code>statistic</code>	the value of the observed test statistic.
<code>perm.stat</code>	numeric vector of the resulting test statistic for <code>nperm</code> random permutations of sample labels.
<code>p.value</code>	p-value indicating the attained significance level.

Author(s)

Yasir Rahmatallah and Galina Glazko

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](#).

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))
## 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 <- AggrFtest(object=dataset, group=c(rep(1,20),rep(2,20)))

```

findMST2

Union of the First and Second Minimum Spanning Trees

Description

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

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 <code>igraph</code> containing the MST2 is returned. If FALSE, a list of length three containing objects of class <code>igraph</code> 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:

<code>MST2</code>	an object of class <code>igraph</code> containing the union of the first and second MSTs.
<code>first.mst</code>	an object of class <code>igraph</code> containing the first MST.
<code>second.mst</code>	an object of class <code>igraph</code> 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))
## 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)
cov_mtrx[!diag(ngenes)] <- 0.1
mask <- diag(ngenes/4)
mask[!diag(ngenes/4)] <- 0.6
cov_mtrx[1:(ngenes/4),1:(ngenes/4)] <- mask
gp <- mvrnorm(nsamples, zero_vector, cov_mtrx)
dataset <- aperm(gp, c(2,1))
## findMST2 returns a list of length 3
## trees[[1]] is an object of class igraph containing the MST2
trees <- findMST2(dataset)
```

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)
```

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.MST2only` 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 vertices 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:

<code>MST2</code>	an object of class <code>igraph</code> containing the union of the first and second MSTs.
<code>first.mst</code>	an object of class <code>igraph</code> containing the first MST.
<code>second.mst</code>	an object of class <code>igraph</code> containing the 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](#).

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="v1")
## 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"
```


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

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

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

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.

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:

<code>statistic</code>	the value of the observed test statistic.
<code>perm.stat</code>	numeric vector of the resulting test statistic for <code>nperm</code> random permutations of sample labels.
<code>p.value</code>	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))
## create a covariance matrix with low off-diagonal elements
cov_mtrx1 <- diag(ngenes)
cov_mtrx1[!diag(ngenes)] <- 0.1
## 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)
cov_mtrx2[!diag(ngenes)] <- 0.1
mask <- diag(ngenes/4)
mask[!diag(ngenes/4)] <- 0.6
cov_mtrx2[1:(ngenes/4),1:(ngenes/4)] <- mask
gp1 <- mvrnorm((nsamples/2), zero_vector, cov_mtrx1)
gp2 <- mvrnorm((nsamples/2), zero_vector, cov_mtrx2)
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 <- GSNCAtest(object=dataset, group=c(rep(1,20),rep(2,20)))
```

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](#).

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)
```

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 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 unidirectional 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 $r_i(s_i)$ is the number of nodes (samples) from condition 1(2) which ranked lower than i , $1 \leq i \leq N$ 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

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

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:

<code>statistic</code>	the value of the observed test statistic.
<code>perm.stat</code>	numeric vector of the resulting test statistic for <code>nperm</code> random permutations of sample labels.
<code>p.value</code>	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`, `WWtest`, `RKStest`, `RMDtest`, `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))
## set the covariance matrix to be an identity matrix for both conditions
cov_mtrx <- diag(ngenes)
gp <- mvrnorm(nsamples, zero_vector, cov_mtrx)
## 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))
## 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)))
```

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

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 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 (rank 1) 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

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

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 `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:

<code>statistic</code>	the value of the observed test statistic.
<code>perm.stat</code>	numeric vector of the resulting test statistic for <code>nperm</code> random permutations of sample labels.
<code>p.value</code>	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

[KStest](#), [WWtest](#), [RKStest](#), [RMDtest](#), [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))
## set the covariance matrix to be an identity matrix for both conditions
cov_mtrx <- diag(ngenes)
gp <- mvrnorm(nsamples, zero_vector, cov_mtrx)
## apply a mean shift of 3 to half of the features under condition 2
gp[1:20,1:10] <- gp[1:20,1:10] + 3
dataset <- 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)))
```

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.

Probe level intensities were quantile normalized and transformed to the log scale using $\log_2(1 + \text{intensity})$. 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 whether 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)
```

Arguments

<code>object</code>	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 matrix.
<code>group</code>	a numeric vector indicating group associations for samples. Possible values are 1 and 2.
<code>name</code>	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.
<code>cor.method</code>	a character string indicating which correlation coefficient is to be computed. Possible values are “pearson” (default), “spearman” and “kendall”. Default value is “pearson”.
<code>min.sd</code>	a numeric value indicating the minimum allowed standard deviation for any feature. If any feature has a standard deviation smaller than <code>min.sd</code> then the execution stops and an error message is returned. Default value is 1e-3.
<code>legend.size</code>	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.
leg.y	a numeric value indicating the amount of vertical shift of the legend box to allow 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 absent, node labels will be set to <code>as.character(c(1:nrow(object)))</code> .
group1.name	an optional character string to be presented as the given name for class 1 in the plot. Default value is "Group 1"
group2.name	an optional character string to be presented as the given name for class 2 in the plot. Default value is "Group 2"
label.size	a numeric value passed to argument <code>vertex.label.cex</code> in command <code>plot.igraph</code> 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 <code>vertex.label.dist</code> in command <code>plot.igraph</code> to specify the distance between vertex labels and the centers of vertices. Default value is 0.5.
vertex.size	a numeric value passed to argument <code>vertex.size</code> in command <code>plot.igraph</code> to specify the vertex size. Default value is 8.
vertex.label.font	a numeric value passed to argument <code>vertex.label.font</code> in command <code>plot.igraph</code> to specify the used font type. Default value is 1.
edge.width	a numeric value passed to argument <code>edge.width</code> in command <code>plot.igraph</code> to specify the edge width in the plot.

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 labels 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.

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))
## create a covariance matrix with low off-diagonal elements
cov_mtrx1 <- diag(ngenes)
cov_mtrx1[!diag(ngenes)] <- 0.1
## 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)
```

```
cov_mtrx2[!diag(ngenes)] <- 0.1
mask <- diag(ngenes/4)
mask[!diag(ngenes/4)] <- 0.6
cov_mtrx2[1:(ngenes/4),1:(ngenes/4)] <- mask
gp1 <- mvrnorm((nsamples/2), zero_vector, cov_mtrx1)
gp2 <- mvrnorm((nsamples/2), zero_vector, cov_mtrx2)
gp <- rbind(gp1, gp2)
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 <code>igraph</code> that consists of a minimum spanning tree or the union of multiple spanning trees.
--------	---

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)
```

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 , $1 \leq i \leq N$ 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 examined 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

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

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:

<code>statistic</code>	the value of the observed test statistic.
<code>perm.stat</code>	numeric vector of the resulting test statistic for <code>nperm</code> random permutations of sample labels.
<code>p.value</code>	p-value indicating the attained significance level.

Note

The variance of both the Poisson and negative Binomial 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.

See Also

[RMDtest](#), [WWtest](#), [MDtest](#), [KStest](#).

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))
## 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 <- RKStest(object=dataset, group=c(rep(1,20),rep(2,20)))
```

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

<code>object</code>	a numeric matrix with columns and rows respectively corresponding to samples and features.
<code>group</code>	a numeric vector indicating group associations for samples. Possible values are 1 and 2.
<code>mst.order</code>	numeric value to indicate the consideration of the union of the first <code>mst.order</code> MSTs. Default value is 1. Maximum allowed value is 5.
<code>nperm</code>	number of permutations used to estimate the null distribution of the test statistic. If not given, a default value 1000 is used.
<code>pvalue.only</code>	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

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

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 `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:

<code>statistic</code>	the value of the observed test statistic.
<code>perm.stat</code>	numeric vector of the resulting test statistic for <code>nperm</code> random permutations of sample labels.
<code>p.value</code>	p-value indicating the attained significance level.

Note

The variance of both the Poisson and negative Binomial 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

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

[RKStest](#), [WWtest](#), [MDtest](#), [KStest](#).

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))
## 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 <- RMDtest(object=dataset, group=c(rep(1,20),rep(2,20)))
```

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 et. al. 2014, Rahmatallah et. al. 2012 and Friedman and Rafsky 1979 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)
```

Arguments

<code>object</code>	a numeric matrix with columns and rows respectively corresponding to samples and features.
<code>group</code>	a numeric vector indicating group associations for samples. Possible values are 1 and 2.
<code>geneSets</code>	a list of character vectors providing the identifiers of features to be considered in each gene set.
<code>min.size</code>	a numeric value indicating the minimum allowed gene set size. Default value is 10.
<code>max.size</code>	a numeric value indicating the maximum allowed gene set size. Default value is 500.
<code>test</code>	a character parameter indicating which statistical method to use for testing the gene sets. Must be one of "GSNCAtest", "WWtest", "KStest", "MDtest", "RKStest", "RMDtest".
<code>nperm</code>	number of permutations used to estimate the null distribution of the test statistic. If not given, a default value 1000 is used.
<code>mst.order</code>	numeric value to indicate the consideration of the union of the first <code>mst.order</code> MSTs when "RKStest" or "RMDtest" are used. Default value is 1. Maximum allowed value is 5.
<code>pvalue.only</code>	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 absent 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:

<code>statistic</code>	the value of the observed test statistic.
<code>perm.stat</code>	numeric vector of the resulting test statistic for <code>nperm</code> 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.

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

[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))

## set the covariance matrix to be an identity matrix for both conditions
cov_mtx <- diag(ngenes)
gp <- mvrnorm(nsamples, zero_vector, cov_mtx)
## apply a mean shift of 5 to the first 10 features under condition 1
gp[1:20,1:10] <- gp[1:20,1:10] + 5
dataset <- aperm(gp, c(2,1))
## assign a unique identifier to each gene
rownames(dataset) <- as.character(c(1:ngenes))
## first 20 samples belong to condition 1
## second 20 samples belong to condition 2
sample.labels <- 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
## Since 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)),
"set3"=as.character(c(31:40)))
results <- TestGeneSets(object=dataset, group=sample.labels,
geneSets=geneSets, test="KStest")
```

 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 examined 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

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

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

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:

<code>statistic</code>	the value of the observed test statistic.
<code>perm.stat</code>	numeric vector of the resulting test statistic for <code>nperm</code> random permutations of sample labels.
<code>p.value</code>	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](#).

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))
## let the mean vector have 2s of length 20 for condition 2
mu_vector <- array(2,c(1,ngenes))
## set the covariance matrix to be an identity matrix
cov_mtrx <- diag(ngenes)
gp1 <- mvrnorm((nsamples/2), zero_vector, cov_mtrx)
gp2 <- mvrnorm((nsamples/2), mu_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 <- WWtest(object=dataset, group=c(rep(1,20),rep(2,20)))
```

Index

- * **arith**
 - HDP.ranking, 11
 - radial.ranking, 21
 - * **datasets**
 - p53DataSet, 17
 - * **dplot**
 - findMST2, 5
 - findMST2.PPI, 7
 - plotMST2.pathway, 18
 - * **graphs**
 - findMST2, 5
 - findMST2.PPI, 7
 - plotMST2.pathway, 18
 - * **multivariate**
 - GSNCAtest, 9
 - KStest, 13
 - MDtest, 15
 - RKStest, 22
 - RMDtest, 24
 - TestGeneSets, 26
 - WWtest, 29
 - * **nonparametric**
 - AggrFtest, 3
 - GSNCAtest, 9
 - KStest, 13
 - MDtest, 15
 - RKStest, 22
 - RMDtest, 24
 - TestGeneSets, 26
 - WWtest, 29
 - * **package**
 - GSAR-package, 2
- AggrFtest, 2, 3
- findMST2, 2, 5, 11, 18–20
- findMST2.PPI, 2, 7
- GSAR (GSAR-package), 2
- GSAR-package, 2
- GSNCAtest, 2, 6, 8, 9, 19, 20, 28
- HDP.ranking, 11, 14, 16, 22
- igraph, 3, 6, 7, 11, 12, 21
- KStest, 2, 12, 13, 16, 24, 26, 28, 30
- MDtest, 2, 12, 14, 15, 24, 26, 28, 30
- p53DataSet, 17
- plotMST2.pathway, 2, 6, 8, 11, 18
- radial.ranking, 12, 21
- RKStest, 2, 4, 14, 16, 22, 22, 26, 28, 30
- RMDtest, 2, 4, 14, 16, 22, 24, 24, 28, 30
- TestGeneSets, 26
- WWtest, 14, 16, 24, 26, 28, 29