

Package ‘BioTIP’

October 14, 2021

Type Package

Title BioTIP: An R package for characterization of Biological Tipping-Point

Version 1.6.0

Author Zhezhen Wang, Andrew Goldstein, Yuxi Sun, Biniam Feleke, Qier An, Antonio Feliciano, Xinnan Yang

Maintainer Yuxi (Jennifer) Sun <ysun11@uchicago.edu>, Zhezhen Wang <zhezhen@uchicago.edu>, and X Holly Yang <xyang2@uchicago.edu>

Description Adopting tipping-point theory to transcriptome profiles to unravel disease regulatory trajectory.

Depends R (>= 3.6)

License GPL-2

Encoding UTF-8

LazyData true

Roxygen list(markdown = TRUE)

RoxygenNote 7.1.0

Suggests knitr, markdown, base, rmarkdown, ggplot2

VignetteBuilder knitr

biocViews Sequencing, RNASeq, GeneExpression, Transcription, Software

URL <https://github.com/xyang2uchicago/BioTIP>

Imports igraph, cluster, psych, stringr, GenomicRanges, Hmisc, MASS

git_url <https://git.bioconductor.org/packages/BioTIP>

git_branch RELEASE_3_13

git_last_commit 42fb23f

git_last_commit_date 2021-05-19

Date/Publication 2021-10-14

R topics documented:

| | |
|---------------------------------|-----------|
| avg.cor.shrink | 2 |
| cod | 4 |
| cor.shrink | 4 |
| gencode | 6 |
| getBiotypes | 6 |
| getCluster_methods | 8 |
| getCTS | 9 |
| getIc | 10 |
| getIc.new | 12 |
| getMaxMCImember | 13 |
| getMaxStats | 15 |
| getMCI | 16 |
| getMCI_inner | 18 |
| getNetwork | 19 |
| getReadthrough | 20 |
| GSE6136_cli | 21 |
| GSE6136_matrix | 22 |
| ILEF | 22 |
| intron | 23 |
| optimize.sd_selection | 23 |
| plotBar_MCI | 25 |
| plotIc | 26 |
| plotMaxMCI | 27 |
| plot_Ic_Simulation | 28 |
| plot_MCI_Simulation | 30 |
| plot_SS_Simulation | 31 |
| sd_selection | 32 |
| simulationMCI | 34 |
| simulation_Ic | 35 |
| simulation_Ic_sample | 36 |
| Index | 39 |

| | |
|----------------|--|
| avg.cor.shrink | <i>Estimation of average values of correlation</i> |
|----------------|--|

Description

This function takes in one (or two) matrix X (rows are genes, columns are samples) (or Y). It then calculates the average pairwise correlation between genes or samples. This method uses the method outlined by Schafer and Strimmer in "A Shrinkage Approach to Large-Scale Covariance Matrix Estimation and Implications for Functional Genomics" (2005)

Usage

```
avg.cor.shrink(
  X,
  Y = NULL,
  MARGIN = c(1, 2),
  shrink = TRUE,
  abs = FALSE,
  target = c("zero", "average", "half")
)
```

Arguments

| | |
|--------|---|
| X | A $G1 \times S$ matrix of counts. Rows correspond to genes, columns correspond to samples. |
| Y | A $G2 \times S$ matrix of counts. Rows correspond to genes, columns correspond to samples. By default is NULL. |
| MARGIN | An integer indicate whether the rows (1, genes) or the columns (2, samples) to be calculated for pairwise correlation. |
| shrink | A flag specifying whether to shrink the correlation or not. |
| abs | A flag specifying whether to take the absolute value before taking the average (used for gene-gene correlations, not sample-sample correlations) |
| target | A character choose among ('zero', 'average', 'half'), indicating whether to shrink towards zero (for gene-gene correlations), shrink towards their empirical common average, one or 0.5 (for sample-sample correlations). |

Value

The average pairwise correlation between genes or samples.

Author(s)

Andrew Goldstein <andrewgoldstein@uchicago.edu>; Xinan H Yang <xyang2@uchicago.edu>

References

Schafer and Strimmer (2005) "A Shrinkage Approach to Large-Scale Covariance Matrix Estimation and Implications for Functional Genomics"

Examples

```
## Generating a data X as coming from a multivariate normal distribution
## with 10 highly correlated variables, roughly simulating correlated genes.
M = matrix(.9, nrow = 10, ncol = 10)
diag(M) = 1
mu = rnorm(10)
X = MASS::mvrnorm(1000, mu, M)
dim(X) #1000 10
```

```
## Mean value of standard pairwise correlation between 1000 genes
# cor_tX = cor(t(X))
# mean(abs(cor_tX[upper.tri(cor_tX, diag = FALSE)])) # 0.9150228

## Calculating estimated pairwise correlation between 1000 genes
avg.cor.shrink(X, MARGIN=1,shrink = TRUE, targe='zero') # 0.8287838

M = matrix(.9, nrow = 20, ncol = 10)
diag(M) = 1
Y = rbind(M,X)
dim(Y) #1020 10
avg.cor.shrink(X, Y, MARGIN=1,shrink = TRUE, targe='zero') #0.8197959
```

| | |
|-----|--------------------|
| cod | <i>cod dataset</i> |
|-----|--------------------|

Description

A subset of gencode_gr extracted as: `cod <- subset(gencode_gr, biotype == 'protein_coding')`

Usage

cod

Format

A dataframe with 3 data columns and 1 metadata column.

seqnames chromosome names (chr21)

ranges rangeschromosome ranges on the genome (10906201-11029719)

strand specific strand of the genomic location (+,-,*)

| | |
|------------|---|
| cor.shrink | <i>Estimation of correlation matrix</i> |
|------------|---|

Description

This function takes in one (or two) matrix X (rows are genes, columns are samples) (or Y). It then calculates the average pairwise correlation between genes or samples. This method uses the method outlined by Schafer and Strimmer (2005).

Usage

```
cor.shrink(
  X,
  Y = NULL,
  MARGIN = c(1, 2),
  shrink = TRUE,
  target = c("zero", "average", "half")
)
```

Arguments

| | |
|--------|---|
| X | A $G_1 \times S$ matrix of counts. Rows correspond to genes, columns correspond to samples. |
| Y | A $G_2 \times S$ matrix of counts. Rows correspond to genes, columns correspond to samples. By default is NULL. |
| MARGIN | An integer indicate whether the rows (1, genes) or the columns (2, samples) to be calculated for pairwise correlation. |
| shrink | A flag specifying whether to shrink the correlation or not. |
| target | A character choose among ('zero', 'average', 'half'), indicating whether to shrink towards zero (for gene-gene correlations), shrink towards their empirical common average, one or 0.5 (for sample-sample correlations). |

Value

The pairwise correlation between genes or samples. If $Y == \text{NULL}$, a $G_1 \times G_1$ matrix is returned; otherwise, a $G_1 \times G_2$ matrix is returned.

Author(s)

Andrew Goldstein <andrewgoldstein@uchicago.edu>

References

Schafer and Strimmer (2005) "A Shrinkage Approach to Large-Scale Covariance Matrix Estimation and Implications for Functional Genomics"

Examples

```
require(MASS)
## Generating a data X as coming from a multivariate normal distribution
## with 10 highly correlated variables, roughly simulating correlated genes.
M = matrix(.9, nrow = 10, ncol = 10)
diag(M) = 1
mu = rnorm(10)
X = MASS::mvrnorm(500, mu, M)
dim(X) #500 10

## Calculating pairwise correlation among 500 correlated genes
cor_tX = cor(t(X))
```

```

mean(abs(cor_tX[upper.tri(cor_tX, diag = FALSE)])) # 0.9468098

## Calculating estimated pairwise correlation among 500 correlated genes
cor.matrix <- cor.shrink(X, MARGIN=1, shrink = TRUE, targe='zero') # 0.8287838
dim(cor.matrix)  #[1] 500 500
mean(upper.tri(cor.matrix, diag=FALSE))  # 0.499

## Calculating stimated pairwise correlation among 500 correlated genes
## and additional 100 random genes
Y = rbind(X, matrix(rnorm(300*10), nrow = 300))
dim(Y)  #800 10
cor.matrix <- cor.shrink(X, Y, MARGIN=1, shrink = TRUE, targe='zero')
dim(cor.matrix)  #[1] 500 800
mean(upper.tri(cor.matrix, diag=FALSE))  # 0.6868

```

| | |
|---------|---|
| gencode | <i>A chr21 data from GENCODE GRCh37</i> |
|---------|---|

Usage

gencode

Format

A data frame of chr21 with 2619444 rows and 9 variables:

Rows Rows, include chromosome numbers

Columns Columns include seqname, source, feature, start , end, score, strand, frame and attribute

cod_gr cod_gr is a subset of sample data 'gencode_gr'

| | |
|-------------|--------------------------------------|
| getBiotypes | <i>Assigning Transcript Biotypes</i> |
|-------------|--------------------------------------|

Description

The purpose of the getBiotypes() function is to class both coding and noncoding transcripts into biotypes using the most recent GENCODE annotations. This tool can also be used to define potential lncRNAs, given an available genome transcriptome assembly (a gtf file) or any genomic loci of interest.

Usage

```
getBiotypes(full_gr, gencode_gr, intron_gr = NULL, minoverlap = 1L)
```

Arguments

| | |
|------------|--|
| full_gr | A GRanges object which contains either coding or noncoding transcripts. Each GRanges objects' columns requires a unique identifications. For further details refer to the GRanges package. |
| gencode_gr | A GRanges object containing a human Chr21 GENCODE reference annotation. A metadata column, "biotype", describes the transcript type. |
| intron_gr | A GRanges object containing the coordinates of non-coding transcripts. |
| minoverlap | An IRanges argument which detects minimum overlap between two IRanges objects. For more information about minoverlap argument refer to the IRanges package. |

Details

For details of findOverlaps, type.partialOverlap, type.50Overlap type.toPlot, queryhits, and subjecthits see GenomicRanges <https://www.bioconductor.org/packages/release/bioc/html/GenomicRanges.html>, IRanges <https://www.bioconductor.org/packages/release/bioc/html/IRanges.html>, and BiocManager <http://bioconductor.org/install/index.html>.

Value

A GRanges object that returns classified transcriptome biotypes.

Note

Replace the PATH_FILE when loading your data locally.

Author(s)

Zhezhen Wang and Biniam Feleke

Source

Reference GRCh37 genome https://www.gencodegenes.org/human/release_25lift37.html for details on gtf format visit ensemble <https://useast.ensembl.org/info/website/upload/gff.html>

References

Wang, Z.Z., J. M. Cunningham and X. H. Yang (2018). 'CisPi: a transcriptomic score for disclosing cis-acting disease-associated lincRNAs.' *Bioinformatics* 34(17): 664-670', PMID: 30423099'

Examples

```
# Input datasets from our package's data folder
library(GenomicRanges)
data("gencode")
data("intron")
data("ILEF")
```

```
# Converting datasets to GRanges object
gencode_gr = GRanges(gencode)
ILEF_gr = GRanges(ILEF)
cod_gr = GRanges(cod)
intron_gr= GRanges(intron)

# Filtering non-coding transcripts
getBiotypes(ILEF_gr, gencode_gr, intron_gr)

## Not run: getBiotypes(intron_gr)
```

getCluster_methods *Clustering Network Nodes*

Description

This function runs over all states which are grouped samples. For each state, this function splits the correlation network generated from the function [getNetwork](#) into several sub-networks (which we called 'module'). The network nodes will be defined by the end-user. For transcriptome analysis, network nodes can be the expressed transcripts. The outputs of this function include the module IDs and node IDs per module.

Usage

```
getCluster_methods(
  igrapl,
  method = c("rw", "hcm", "km", "pam", "natural"),
  cutoff = NULL
)
```

Arguments

- | | |
|--------|---|
| igrapl | A list of numerical matrices or a list of igraph objects. The list of igraph objects can be the output from the getNetwork function. |
| method | A mathematical clustering model for analyzing network nodes. Default is a random walk ('rw'). A method could be 'rw', 'hcm', 'km', 'pam', or 'natural', where: <ul style="list-style-type: none"> • rw: random walk using cluster_walktrap function in igraph package. 'igrapl' has to be a list of igraph. • hcm: hierarchical clustering using function hclust) and dist, using method 'complete'. • km and pam: k-medoids or PAM algorithm using KMedoids. • natrual: if nodes are disconnected, they may naturally cluster and form sub-networks. |
| cutoff | A numeric value, default is NULL. For each method it means: <ul style="list-style-type: none"> • rw: the number of steps needed, see cluster_walktrap for more detail. If "cutoff" is not assigned, default of 4 will be used. |

- hcm, km and pam: number of clusters wanted. No default assigned.
- natural: does not use this parameter.

Value

When method=rw: A list of `communities` objects of R package `igraph`, whose length is the length of the input object `igraphL`. These `communities` objects can be used for visualization when being assigned to the `'mark.groups'` parameter of the `plot.igraph` function of the `igraph` package. Otherwise this function returns a list of vectors, whose length is the length of the input object `igraphL`. The names of each vector are the pre-selected transcript IDs by the function `sd_selection`. Each vector, whose length is the number of pre-selected transcript in a state, contains the module IDs.

Author(s)

Zhezhen Wang <zhezhen@uchicago.edu>

Examples

```
test = list('state1' = matrix(sample(1:10, 6), 3, 3), 'state2' =
matrix(sample(1:10, 6), 3, 3), 'state3' = matrix(sample(1:10, 6), 3, 3))
#assign colnames and rownames to the matrix

for(i in names(test)){
  colnames(test[[i]]) = 1:3
  row.names(test[[i]]) = 1:3}

#using 'rw' or 'natural' method
igraphL <- getNetwork(test, fdr=1)
#[1] "state1:3 nodes"
#[1] "state2:3 nodes"
#[1] "state3:3 nodes"

c1 <- getCluster_methods(igraphL)

#using 'km', 'pam' or 'hcm'
c1 <- getCluster_methods(test, method = 'pam', cutoff=2)
```

getCTS

Obtain the identified BioTiP and its length

Description

getCTS obtains the identified BioTiP and its length based off of MCI scores.

Usage

```
getCTS(maxMCI, maxMCIms)
```

Arguments

- `maxMCI` A numeric vector, whose length is the number of states. This parameter is the maximum MCI score of each state, and it can be obtained from the output of `getMaxStats`. Names need to be included in names of `maxMCImS`.
- `maxMCImS` A list of character vectors per state. The vectors are network nodes (e.g. transcript ids). This parameter is the second element of the output of the function `getMaxMCI`.

Value

A character vector, in which the elements are the unique IDs of the network nodes of the BioTiP.

Author(s)

Antonio Feliciano y Pleyto and Zhezhen Wang <zhezhen@uchicago.edu>

Examples

```
maxMCI <- c(a = 2.56, b = 8.52, c = 2.36, d = 4.81, e = 5.26)
maxMCImS <- list(a = c("A100", "A293", "C403"),
                 b = c("B853", "D826", "A406"),
                 c = c("J198", "D103", "B105"),
                 d = c("K529", "D385", "E358"),
                 e = c("J019", "U926", "N824"))
identical(names(maxMCI), names(maxMCImS))
# TRUE
getCTS(maxMCI, maxMCImS)
# "Length: 3"
# "B853" "D826" "A406"
```

getIc

Get Index for Critical transition (Ic score)

Description

Retrieve Ic scores (Pearson correlation of genes / Pearson correlation of samples) for the identified critical transition state

Usage

```
getIc(
  counts,
  sampleL,
  genes,
  output = c("Ic", "PCCg", "PCCs"),
  fun = c("cor", "BioTIP"),
  shrink = TRUE,
```

```

use = c("everything", "all.obs", "complete.obs", "na.or.complete",
"pairwise.complete.obs")
)

```

Arguments

| | |
|---------|---|
| counts | A numeric matrix or data frame. The rows and columns represent unique transcript IDs (geneID) and sample names, respectively. |
| sampleL | A list of vectors, whose length is the number of states. Each vector gives the sample names in a state. Note that the vector s (sample names) has to be among the column names of the R object 'df'. |
| genes | A character vector consisting of unique CTS gene ids. This can be obtained from getMaxMCI member |
| output | A string. Please select from 'Ic', 'PCCg', or 'PCCs'. Uses 'Ic' by default. 'PCCg' is the PCC between genes (numerator) and 'PCCs' is PCC between samples (denominator) |
| fun | An optional character string indicating the R function to calculate correlations for all possible pairs of columns of a matrix. When using "BioTIP", The method is modified to ignore missing values, analogous to how <code>cor(X, use = "pairwise.complete.obs")</code> works. Note that the "BioTIP" option only function together with <code>shrink = TRUE</code> . |
| shrink | A flag specifying whether to shrink the correlation or not. This approach uses the method outlined by Schafer and Strimmer in "A Shrinkage Approach to Large-Scale Covariance Matrix Estimation and Implications for Functional Genomics" (2005) Comparing to <code>fun='cor'</code> , the 'BioTIP' method without shrinkage is modified to ignore missing values, analogous to how <code>cor(X, use = "pairwise.complete.obs")</code> works. |
| use | An optional character string, when <code>fun=="cor"</code> , it gives a method for computing covariances in the presence of missing values. This must be (an abbreviation of) one of the strings "everything", "all.obs", "complete.obs", "na.or.complete", or "pairwise.complete.obs". |

Value

A list of numeric values, whose length and names are inherited from `sampleL`

Author(s)

Zhezhen Wang <zhezhen@uchicago.edu>; Xinan H Yang <xyang2@uchicago.edu>

References

- Schafer and Strimmer (2005) "A Shrinkage Approach to Large-Scale Covariance Matrix Estimation and Implications for Functional Genomics"
- M. Mojtahedi et al., Cell Fate Decision as High-Dimensional Critical State Transition. PLoS Biol 14, e2000640 (2016).

Examples

```

counts = matrix(sample(1:100, 27), 3, 9)
colnames(counts) = 1:9
row.names(counts) = c('loci1', 'loci2', 'loci3')
cli = cbind(1:9, rep(c('state1', 'state2', 'state3'), each = 3))
colnames(cli) = c('samples', 'group')
samplesL <- split(cli[, 1], f = cli[, 'group'])
CTS = c('loci1', 'loci2')

## Comparing the results with an estimated correlation matrix with that without estimation.
Ic = getIc(counts, samplesL, CTS, fun='cor')
Ic.2 = getIc(counts, samplesL, CTS, fun='BioTIP', shrink=FALSE)
BioTIP = getIc(counts, samplesL, CTS, fun='BioTIP')

```

| | |
|-----------|--|
| getIc.new | <i>Index of criticality Scoring System with estimated correlation, an updated Ic-score</i> |
|-----------|--|

Description

This function calculates the BioTIP score on a given data matrix X (or two matrixes X and Y). It can also calculate the I_c score, if desired.

This approach uses the method outlined by Schafer and Strimmer in "A Shrinkage Approach to Large-Scale Covariance Matrix Estimation and Implications for Functional Genomics" (2005)

This approach is modified to ignore missing values, analogous to how `cor(X, use = "pairwise.complete.obs")` works.

The gene-gene correlations are shrunk towards 0, whereas the sample-sample correlations are shrunk towards their empirical average.

Usage

```

getIc.new(
  X,
  method = c("BioTIP", "Ic"),
  PCC_sample.target = c("average", "zero", "half"),
  output = c("IndexScore", "PCCg", "PCCs")
)

```

Arguments

| | |
|--------|--|
| X | A G x S matrix of counts. Rows correspond to genes, columns correspond to samples. |
| method | A flag specifying whether to calculate the BioTIP score or the I_c score |

`PCC_sample.target` A character choose among ('average', 'zero', 'half'), indicating whether to shrink PCC towards towards their empirical common average, zero or 0.5 (for sample-sample correlations).

`output` A string. Please select from 'IndexScore', 'PCCg', or 'PCCs'. Uses 'IndexScore' by default. 'PCCg' is the PCC between genes (numerator) and 'PCCs' is PCC between samples (denominator).

Value

A value containing the shrunk BioTIP or non-shrunk I_c score

Author(s)

Andrew Goldstein <andrewgoldstein@uchicago.edu>

Examples

```
## Generating a data X as coming from a multivariate normal distribution
## with 10 highly correlated variables, roughly simulating correlated genes.
M = matrix(.9, nrow = 10, ncol = 10)
diag(M) = 1
mu = rnorm(10)
X = MASS::mvrnorm(1000, mu, M)
dim(X) #1000 10

## Calculating pairwise correlation between 1000 genes; then the mean value
## in two ways, respectively
cor_tX = cor(t(X))
mean(abs(cor_tX[upper.tri(cor_tX, diag = FALSE)])) # 0.9150228

getIc.new(X, method = "Ic", output = 'PCCg') # 0.9150228
getIc.new(X, method = "BioTIP", output = 'PCCg') # 0.8287838

## Using the Index of critical scoring system, in two ways, respectively
(newscore = getIc.new(X, method = "BioTIP"))
(oldscore = getIc.new(X, method = "Ic"))
```

getMaxMCI member

Identifying the 'Biomodule'

Description

This function reports the 'biomodule', which is the module with the maximum Module Critical Index (MCI) scores for each state. Each state can have multiple modules (groups of subnetworks derived from the function [getCluster_methods](#)). This function runs over all states.

Usage

```
getMaxMCImember(membersL, MCI1, minsize = 1)
```

Arguments

| | |
|----------|---|
| membersL | A list of integer vectors with unique ids as names. Each vector represents the cluster number assign to that unique id. The length of this list is equal to the number of states in the study. This can be the first element of the output from function <code>getMCI</code> or the output from <code>getCluster_methods</code> , see Examples for more detail. |
| MCI1 | A list of numeric vectors with unique cluster numbers as names. Each vector represents the MCI scores of that module. This can be the second element of the output from function <code>getMCI</code> . |
| minsize | A numerical value of the minimum module size (the number of transcripts in a cluster) to output for downstream analysis. |

Value

A nested list whose length is the length of the input object `membersL`. Each internal list contains two objects: one object is the vector of biomodule IDs across states, and the other object is a list of transcript IDs (each defines the biomodule per state) across states.

Author(s)

Zhezhen Wang <zhezhen@uchicago.edu>

Examples

```
#1st option: get the input directly from getMCI function
test = list('state1' = matrix(sample(1:10, 6), 4, 3),
           'state2' = matrix(sample(1:10, 6), 4, 3),
           'state3' = matrix(sample(1:10, 6), 4, 3))

# assign colnames and rownames to the matrix
for(i in names(test)){
  colnames(test[[i]]) = 1:3
  row.names(test[[i]]) = c('g1', 'g2', 'g3', 'g4')}

cluster = list(c(1, 2, 2, 1), c(1, 2, 3, 1), c(2, 2, 1, 1))
names(cluster) = names(test)
for(i in names(cluster)){
  names(cluster[[i]]) = c('g1', 'g2', 'g3', 'g4')}

membersL_noweight <- getMCI(cluster, test)
maxMCIm <- getMaxMCImember(membersL_noweight[[1]], membersL_noweight[[2]], min =3)
#The same as
maxMCIm <- getMaxMCImember(cluster, membersL_noweight[[2]], min =2)

## case1: using 'rw' method by default
igraphL <- getNetwork(test, fdr=1)
```

```

cl <- getCluster_methods(igraphL)
## make sure every element in list cl is a \code{communities} object
sapply(cl, class)
##      state1      state2      state3
## "communities" "communities" "communities"

## If there is(are) state(s) that is(are) empty which will not be a communities object(s),
## please manually remove that state(s).

cl = cl[which(sapply(cl, class) == 'communities')]

## and then run
library(igraph)
cluster = lapply(cl, membership)
maxCIms <- getMaxMCImember(cluster, membersL_noweight[[2]], min =2)

## or run function 'getMCI' and use the 1st option
membersL_noweight <- getMCI(cl, test)

## case2: using methods other than the default
cl <- getCluster_methods(test, method = "pam", cutoff = 2)
## check to make sure membersL_noweight[[2]] has values and run
maxCIms <- getMaxMCImember(cl, membersL_noweight[[2]], min =2)

```

getMaxStats

Get the cluster index and network nodes of biomodule

Description

This function retrieves the cluster index and network-node ids for the identified biomodule (that shows the maximum MCI score) at each state in the study.

Usage

```
getMaxStats(membersL, idx)
```

Arguments

| | |
|----------|--|
| membersL | A two-layer nested list of character or numeric values, any one out of the five elements output by the function getMCI . |
| idx | A vector of integers that are cluster ids of the biomodule (the module with the highest MCI score) per state. This is the first element of the result from getMaxMCImember . |

Value

A list describing the biomodule of each state, corresponding to one of the five elements (members, MCI, Sd, PCC, and PCCo) outputted by the function [getMCI](#). The class of the vector depends on the class of the input parameter membersL.

Author(s)

Zhezhen Wang <zhezhen@uchicago.edu>

See Also

[getMCI](#)

Examples

```
test = list('state1' = matrix(sample(1:10, 6), 4, 3),
           'state2' = matrix(sample(1:10, 6), 4, 3),
           'state3' = matrix(sample(1:10, 6), 4, 3))

# assign colnames and rownames to the matrix
for(i in names(test)){
  colnames(test[[i]]) = 1:3
  row.names(test[[i]]) = c('g1', 'g2', 'g3', 'g4')
}

cluster = list(c(1, 2, 2, 1), c(1, 2, 3, 1), c(2, 2, 1, 1))
names(cluster) = names(test)
for(i in names(cluster)){
  names(cluster[[i]]) = c('g1', 'g2', 'g3', 'g4')
}

membersL_noweight <- getMCI(cluster, test)
idx = c(1, 2, 1)
names(idx) = names(membersL_noweight[['sd']])
selectedSD = getMaxStats(membersL_noweight[['sd']], idx)
```

getMCI

Calculating MCI Scores

Description

This function calculates a module critical index (MCI) score for each module per state within a dataset. Each module is a cluster of transcripts generated from the function [getCluster_methods](#). Note that a dataset should contains three or more states (samples in groups).

Usage

```
getMCI(
  groups,
  countsL,
  adjust.size = FALSE,
  fun = c("cor", "BioTIP"),
  df = NULL
)
```


Arguments

| | |
|-------------|--|
| groups | A list of elements whose length is the member of states. The elements could be either be vectors or communities object of the R package <code>igraph</code> . If a vector, it is the output of the function <code>getCluster_methods</code> . The names of each vector are the pre-selected transcript IDs generated by the function <code>sd_selection</code> . Each vector, whose length is the number of pre-selected transcripts in a state, contains the module IDs. If a communities object, it can be obtained by <code>getCluster_methods</code> using the "rw" method. It is also an output of the function <code>sd_selection</code> . |
| countsl | A list of x numeric count matrices or x data frame, where x is the number of states. |
| adjust.size | A Boolean value indicating if MCI score should be adjusted by module size (the number of transcripts in the module) or not. Default FALSE. This parameter is not recommended for fun=BioTIP. |
| fun | A character chosen between ("cor", "BioTIP"), indicating where an adjusted correlation matrix will be used to calculate the MCI score. |
| df | NULL or a numeric matrix or data frame, where rows and columns represent unique transcript IDs (geneID) and sample names, respectively. Used only when fun='BioTIP'. By default is NULL, estimating the correlation among selected genes. Otherwise, estimating the correlation among all genes in the df, ensuring cross state comparison. |

Value

A list of five elements (members, MCI, Sd, PCC, and PCCo). Each of element is a two-layer nested list whose length is the length of the input object groups. Each internal nested list is structured according to the number of modules identified in that state.

- members: vectors of unique ids
- MCI: the MCI score
- sd: standard deviation
- PCC: Mean of pairwise Pearson Correlation Coefficient calculated among the loci in a module.
- PCCo: Mean of pairwise Pearson Correlation Coefficient calculated between the loci in a module and the loci outside that module but inside the same state.

Author(s)

Zhezhen Wang <zhezhen@uchicago.edu>; Xinan H Yang <xyang2@uchicago.edu>

Examples

```
test = list('state1' = matrix(sample(1:10, 6), 4, 3), 'state2' =
matrix(sample(1:10, 6), 4, 3), 'state3' = matrix(sample(1:10, 6), 4, 3))

## Assign colnames and rownames to the matrix
for(i in names(test)){
  colnames(test[[i]]) = 1:3
  row.names(test[[i]]) = c('g1', 'g2', 'g3', 'g4')}
```

```

cluster = list(c(1, 2, 2, 1), c(1, 2, 3, 1), c(2, 2, 1, 1))
names(cluster) = names(test)
for(i in names(cluster)){
  names(cluster[[i]]) = c('g1', 'g2', 'g3', 'g4')}

membersL_noweight <- getMCI(cluster, test, fun='cor')
names(membersL_noweight)
## [1] "members" "MCI"      "sd"      "PCC"      "PCCo"

```

getMCI_inner

Calculating MCI Score for randomly selected

Description

This function calculates random MCI score, allowing an estimation of correlation matrix using the Schafer-Strimmer Method for the PCC_in component in the MCI score.

Usage

```

getMCI_inner(
  members,
  countsL,
  adjust.size,
  fun = c("cor", "BioTIP"),
  PCC_gene.target = "zero",
  M = NULL
)

```

Arguments

| | |
|-----------------|---|
| members | An integer that is the length of genes in the CTS (critical transition signal). |
| countsL | A list of subset of dat matrix. The list length is the number of states. Each subset of matrix gives the genes in rows and samples in column. |
| adjust.size | A boolean value indicating if MCI score should be adjust by module size (the number of transcripts in the module) or not. Default FALSE. |
| fun | A character chosen between ("cor", "BioTIP"), indicating where an adjusted correlation matrix will be used to calculated the MCI score. |
| PCC_gene.target | A character 'zero' indicating that gene-gene correlation matrix will be shrunk. towards zero, used only for fun='BioTIP'. |
| M | is the overall shrunk correlation matrix, used only for fun='BioTIP'. |

Value

A vector recording one MCI score per state.

Author(s)

Zhezhen Wang <zhezhen@uchicago.edu>; Xinan H Yang <xyang2@uchicago.edu>

getNetwork

Building Networks of Nodes

Description

This function builds one correlation network for each state (sample group) and runs across all states. The network nodes are defined by the context of the input dataset. For transcriptomic network analysis, network nodes can be the expressed transcript IDs and network links can be the correlation coefficients. Using the Pearson Correlation Coefficient (PCC) analysis, this function assembles a correlation network of nodes (e.g., co-expressed transcripts) for each state using the R package igraph.

Usage

```
getNetwork(optimal, fdr = 0.05)
```

Arguments

| | |
|---------|---|
| optimal | A list of x numeric data frames, where x is the number of states studied. Each data frame consists of loci with high standard deviations. This object can be obtained through <code>sd_selection</code> function. |
| fdr | A numeric cutoff value for a Pearson Correlation Coefficient (PCC) analysis. Default is 0.05. Transcripts are linked into a network if their correlations meet this PCC-significance criterion. |

Value

A list of igraph objects whose length is the length of the input object `optimal`. Each object is a network of correlated nodes whose PCCs meet the significant criteria based on the false discovery rate (FDR) control. The length of the list is the number of states with PCC networks. If no PCC meets the significant criteria in a state, the state will be deleted from the output.

Author(s)

Zhezhen Wang <zhezhen@uchicago.edu>; Xinan H Yang <xyang2@uchicago.edu>

Examples

```
test = list('state1' = matrix(sample(1:10, 6), 2, 3),
           'state2' = matrix(sample(1:10, 6), 2, 3),
           'state3' = matrix(sample(1:10, 6), 2, 3))

for(i in names(test)){
  colnames(test[[i]]) = 1:3
}
```

```
row.names(test[[i]]) = 1:2}

igraphL <- getNetwork(test, fdr=1)
#[1] "state1:2 nodes"
#[1] "state2:2 nodes"
#[1] "state3:2 nodes"
```

| | |
|----------------|-----------------------------------|
| getReadthrough | <i>Overlapping Coding Regions</i> |
|----------------|-----------------------------------|

Description

The `getReadthrough()` function is used to find long transcripts that cover more than two coding regions for gene regions of interest.

Usage

```
getReadthrough(gr, cod_gr)
```

Arguments

| | |
|---------------------|---|
| <code>gr</code> | A GRanges object that shows the start and end loci on genome. |
| <code>cod_gr</code> | A GRanges object containing coding regions. |

Details

For details of `findOverlaps`, `type.partialOverlap`, `type.50Overlap` `type.toPlot`, `queryhits`, `readthrough` and `subjecthits` see, `GenomicRanges` <https://www.bioconductor.org/packages/release/bioc/html/GenomicRanges.html>, `IRanges` <https://www.bioconductor.org/packages/release/bioc/html/IRanges.html>, and `BiocManager` <http://bioconductor.org/install/index.html>.

Value

A GRanges object that returns overlapping regions of the classified transcript biotypes.

Note

Replace the `path_file` when loading data locally to the data directory.

Author(s)

Zhezhen Wang and Biniam Feleke

Source

Reference GRCh37 genome https://www.gencodegenes.org/human/release_25lift37.html. For details on gtf format visit ensemble <https://useast.ensembl.org/info/website/upload/gff.html>.

References

Wang, Z. Z., J. M. Cunningham and X. H. Yang (2018). 'CisPi: a transcriptomic score for disclosing cis-acting disease-associated lincRNAs.' *Bioinformatics*34(17): 664-670'

Examples

```
#First Load datasets and libraries
library(GenomicRanges)
data("gencode")
data("ILEF")
data("cod")

# Assigning datasets a GRanges object
gencode_gr = GRanges(gencode)
ILEF_gr = GRanges(ILEF)
cod_gr = GRanges(cod)

getReadthrough(ILEF_gr, cod_gr)

## Not run: getReadthrough(cod_gr)
```

GSE6136_cli

GSE6136 cli dataset

Description

A gene annotation samples with their corresponding geneId extracted from GENCODE database. A python script was then run to extcat GSE6136_cli dataset. The script is included in the R/data_raw folder.

Usage

```
GSE6136_cli
```

Format

A dataframe with 22690 columns and 27 rows column.

GSM142398-GSM142423 Names of gene interest

Rows Summary of GSM genes

Source

<https://www.encodegenes.org/human/>

| | |
|----------------|-------------------------------|
| GSE6136_matrix | <i>GSE6136 matrix dataset</i> |
|----------------|-------------------------------|

Description

A gene annotation samples with their corresponding geneId extracted from GENCODE database. A python script was then run to extcat GSE6136_cli dataset. The script is included in the R/data_raw folder.

Usage

GSE6136_matrix

Format

A dataframe with 22690 columns and 27 rows column.

GSM142398-GSM142423 Names of gene interest

ID_REF Reference ID of the target genes

Source

<https://www.genecodegenes.org/human/>

| | |
|------|---|
| ILEF | <i>Chromosome ranges of chr21 dataset</i> |
|------|---|

Description

A dataset containing chromosomes in the genome regions of interest for 137 chromosome. The variables are as follows:

Usage

ILEF

Format

A data frame of chr21 with 137 rows and 4 variables:

seqnames chromosome names (chr1,chr6,chr21)

start gene read start position (167684657,167729508)

end end of gene read position (15710335,43717938)

width width of gene

strand specific strand of the genomic location (+,-,*)

Row.names name of the data rows(A1BG,vawser)

| | |
|--------|--|
| intron | <i>Coding transcriptome in chr21 dataset</i> |
|--------|--|

Description

A dataset containing chromosomes in the genome regions of interest. The variables are as follows:

Usage

```
intron
```

Format

A data frame with 659327 rows and 5 variables:

seqnames chromosome names (chr1,chrM,chr21)
ranges chromosome ranges on the genome(167684657–167729508)
strand specific strand of the genomic location (+,-,*)
name internal assigned names(uc001aaa.3_intron_0_0_chr1_12228_f)
score score not used in this data set(0)

| | |
|-----------------------|-------------------------------------|
| optimize.sd_selection | <i>Optimization of sd selection</i> |
|-----------------------|-------------------------------------|

Description

The `optimize.sd_selection` filters a multi-state dataset based on a cutoff value for standard deviation per state and optimizes. By default, a cutoff value of 0.01 is used. Suggested if each state contains more than 10 samples.

Usage

```
optimize.sd_selection(  
  df,  
  samplesL,  
  B = 100,  
  percent = 0.8,  
  times = 0.8,  
  cutoff = 0.01,  
  method = c("other", "reference", "previous", "itself", "longitudinal reference"),  
  control_df = NULL,  
  control_samplesL = NULL  
)
```

Arguments

| | |
|------------------|---|
| df | A dataframe of numerics. The rows and columns represent unique transcript IDs (geneID) and sample names, respectively. |
| samplesL | A list of n vectors, where n equals to the number of states. Each vector gives the sample names in a state. Note that the vectors (sample names) has to be among the column names of the R object 'df'. |
| B | An integer indicating number of times to run this optimization, default 1000. |
| percent | A numeric value indicating the percentage of samples will be selected in each round of simulation. |
| times | A numeric value indicating the percentage of B times a transcript need to be selected in order to be considered a stable signature. |
| cutoff | A positive numeric value. Default is 0.01. If < 1, automatically goes to select top x percentage transcripts using the a selecting method (which is either the reference, other or previous stage), e.g. by default it will select top 1 percentage of the transcripts. |
| method | Selection of methods from reference, other, previous, default uses other. Partial match enabled. <ul style="list-style-type: none"> • itself, or longitudinal reference. Some specific requirements for each option: • reference, the reference has to be the first. • previous, make sure samplesL is in the right order from benign to malign. • itself, make sure the cutoff is smaller than 1. • longitudinal reference make sure control_df and control_samplesL are not NULL. The row numbers of control_df is the same as df and all transcript in df are also in control_df. |
| control_df | A count matrix with unique loci as row names and samples names of control samples as column names, only used for method longitudinal reference. |
| control_samplesL | A list of characters with stages as names of control samples, required for method 'longitudinal reference'. |

Value

A list of dataframe of filtered transcripts with the highest standard deviation are selected from df based on a cutoff value assigned. The resulting dataframe represents a subset of the raw input df.

Author(s)

Zhezhen Wang <zhezhen@uchicago.edu>

See Also

[sd_selection](#)

Examples

```

counts = matrix(sample(1:100, 30), 2, 30)
colnames(counts) = 1:30
row.names(counts) = paste0('loci', 1:2)
cli = cbind(1:30, rep(c('state1', 'state2', 'state3'), each = 10))
colnames(cli) = c('samples', 'group')
samplesL <- split(cli[, 1], f = cli[, 'group'])
test_sd_selection <- optimize.sd_selection(counts, samplesL, B = 3, cutoff = 0.01)

```

plotBar_MCI

plot MCI barplots

Description

A barplot of MCI for all clusters in all states.

Usage

```

plotBar_MCI(
  MCI1,
  ylim = NULL,
  nr = 1,
  nc = NULL,
  order = NULL,
  minsize = 3,
  states = NULL
)

```

Arguments

| | |
|---------|---|
| MCI1 | A list can be obtained through getMCI. |
| ylim | A vector needed if the output barplots need to be on the same y scale. |
| nr | The number of rows to plot. |
| nc | The number of columns to plot, default length(groups). |
| order | A character vector of the order of the barplot. Default is NULL which uses the input list order. |
| minsize | A non-negative numeric value of minimum size allowed for a cluster. |
| states | A character of the state names to be shown on the plot. Default is NULL, assign this if you want to show all states including states without nodes. |

Value

Return a barplot of MCI scores across states.

Author(s)

Zhezhen Wang <zhezhen@uchicago.edu>

References

Chen L, Liu R, Liu Z, Li M & Aihara K (2012) Detecting early-warning signals for sudden deterioration of complex diseases by dynamical network biomarkers *Scientific Reports* 2:342

Examples

```
test = list('state1' = matrix(sample(1:10, 6), 4, 3),
  'state2' = matrix(sample(1:10, 6), 4, 3),
  'state3' = matrix(sample(1:10, 6), 4, 3))
# assign colnames and rownames to the matrix
for(i in names(test)){
  colnames(test[[i]]) = 1:3
  row.names(test[[i]]) = c('g1', 'g2', 'g3', 'g4')
}

cluster = list(c(1, 2, 2, 1), c(1, 2, 3, 1), c(2, 2, 1, 1))
names(cluster) = names(test)
for(i in names(cluster)){
  names(cluster[[i]]) = c('g1', 'g2', 'g3', 'g4')
}
membersL_noweight <- getMCI(cluster, test)
plotBar_MCI(membersL_noweight)
```

plotIc

plot a line plot of Ic scores for each state.

Description

plot a line plot with Ic score for each state

Usage

```
plotIc(
  Ic,
  las = 0,
  order = NULL,
  ylab = "Ic",
  col = "black",
  main = NULL,
  add = FALSE,
  ylim = NULL,
  lty = 1:5,
  lwd = 1
)
```

Arguments

| | |
|-------|---|
| Ic | A vector with names of states. If order is not assigned, then plot by the order of this vector. |
| las | Numeric in 0, 1, 2, 3; the style of axis labels. Default is 0, meaning labels are parallel. (link to http://127.0.0.1:21580/library/graphics/html/par.html). |
| order | A vector of state names in the customized order to be plotted, setting to NULL by default. |
| ylab | titles y axes, as in plot. |
| col | vector of colors. Colors are used cyclically. |
| main | A character vector. The title of the plot. Default is NULL. |
| add | logical. If TRUE, plots are added to current one. This is inherited from matplot . |
| ylim | An integer vector of length 2. Default is NULL. |
| lty | An vector of line types. This is also inherited from matplot . |
| lwd | An integer of line widths. This is also inherited from matplot . |

Value

Return a line plot of Ic score across states.

Author(s)

Zhezhen Wang <zhezhen@uchicago.edu>

Examples

```
Ic = c('state3' = 3.4, 'state1' = 5.6, 'state2' = 2)
plotIc(Ic, order = c('state1', 'state2', 'state3'))
```

plotMaxMCI

Plot the Maximized MCI per State

Description

This function generates a line plot over multiple states with the maximum MCI score per state. The module size (i.e., number of network nodes) is specified at each state in parentheses.

Usage

```
plotMaxMCI(maxMCIms, MCI1, las = 0, order = NULL, states = NULL)
```

Arguments

| | |
|---------|---|
| maxMCIm | A list of 2 elements. The 1st element is an integer vector of module ids whose names are the state names. The 2nd element is a list of character vectors per state. The vectors are network nodes (e.g. transcript ids). This parameter can be obtained by running function getMaxMCI . |
| MCI1 | A list of numeric vectors whose names are unique cluster ids. Each vector represents the MCI scores of modules in a state. This can be the second element of the output from the function getMCI . |
| las | Numeric in 0, 1, 2, 3; the style of axis labels. Default is 0, meaning labels are parallel. See getMCI for more detail. |
| order | A vector of state names in the customized order to be plotted, setting to NULL by default. |
| states | A character vector of state names that will be shown on the plot, setting to NULL by default. Assign this if you want to show all states, including states with no resulted modules. This parameter will overwrite the parameter 'order'. |

Value

Returns a line plot of maximum MCI scores across the states

Author(s)

Zhezhen Wang <zhezhen@uchicago.edu>

Examples

```
maxMCIm = list(c(state1 = 1, state2 = 2, state3 = 1),
  c(list(state1 = c('g1', 'g2', 'g3'),
    state2 = c('g3', 'g5'),
    state3 = c('g2', 'g6'))))

MCI1 = list(state1=c('1' = 8.84, '2' = 6.4),
  state2 = c('1' =NA, '2' = 9.5, '3' = NA),
  state3 = c('1' = 2.3, '2' = 1.4))

plotMaxMCI(maxMCIm, MCI1)
```

plot_Ic_Simulation *Line or boxplot of an observed and its simulated scores*

Description

Generate a line (or box) plot of Ic score and simulated Ic scores, with three horizontal lines: the min, max and 2*(max-min) value of the state of interests, or all values.

Usage

```
plot_Ic_Simulation(  
  Ic,  
  simulation,  
  las = 0,  
  ylim = NULL,  
  order = NULL,  
  main = NULL,  
  ylab = "Ic",  
  fun = c("matplot", "boxplot"),  
  which2point = NULL  
)
```

Arguments

| | |
|------------|---|
| Ic | A vector with names of states. If order is not assigned, then plot by the order of this vector. |
| simulation | A numeric matrix of Ic scores in which rows are states and columns are numbers of simulated times. It can be obtained from simulation_Ic |
| las | Numeric in 0, 1, 2, 3; the style of axis labels. Default is 0, meaning labels are parallel. (link to http://127.0.0.1:21580/library/graphics/html/par.html). |
| ylim | An integer vector of length 2. Default is NULL. |
| order | Characters of names of Ic to be plotted in a desired order. Default is NULL. |
| main | A character vector. The title of the plot. Default is NULL. |
| ylab | titles y axes, as in plot. |
| fun | A character choose between ('matplot', 'boxplot'), indicating plot type. |

Value

Return a plot of the observed Ic (red) and simulated Ic (grey) scores per states.

Author(s)

Zhezhen Wang <zhezhen@uchicago.edu>; Xinan H Yang <xyang2@uchicago.edu>

Examples

```
sim = matrix(sample(1:10, 9), 3, 3)  
row.names(sim) = paste0('state', 1:3)  
Ic = c('state1' = 3.4, 'state2' = 5.6, 'state3' = 2)  
plot_Ic_Simulation(Ic, sim)
```

plot_MCI_Simulation *Plot observed and simulated MCI Scores*

Description

Box plots of observed (red) and simulated MCI scores by bootstrapping genes B times, with three horizontal lines: the min, max and 2*(max-min) value of the state of interests, or all all values.

Usage

```
plot_MCI_Simulation(
  MCI,
  simulation,
  las = 0,
  order = NULL,
  ylim = NULL,
  main = NULL,
  which2point = NULL,
  ...
)
```

Arguments

| | |
|-------------|--|
| MCI | A named vector of max CI scores per state, can be obtained from function getMaxStats . |
| simulation | A matrix state * number of simulated times, can be obtained from function simulationMCI . |
| las | Numeric in 0, 1, 2, 3; the style of axis labels. Default is 0, meaning labels are parallel. (link to http://127.0.0.1:21580/library/graphics/html/par.html) |
| order | A vector of state names in the customized order to be plotted, set to NULL by default. |
| ylim | An integer vector of length 2. Default is NULL. |
| main | A character vector. The title of the plot. Default is NULL. |
| which2point | A character (or integer) which state's values were used to set up the three horizontal lines. by default is NULL, indicating the values of all states will be used. |
| ... | Other parameters passed to this function |

Value

Return a box plot of MCI(red) and simulated MCI(grey) scores per state.

Author(s)

Zhezhen Wang <zhezhen@uchicago.edu>

Examples

```
MCI = c(1:3); names(MCI) = c('a', 'b', 'c')
simMCI = matrix(sample(1:100, 9), 3, 3)
row.names(simMCI) = names(MCI)
plot_MCI_Simulation(MCI, simMCI)
```

| | |
|--------------------|---|
| plot_SS_Simulation | <i>Density plot the leading two distance between any two states from random scores of all states in a system.</i> |
|--------------------|---|

Description

Generate a density plot of Ic score (orBioTIP score) from a simulation, which is the distance between the first-largest and the second-largest random scores. This is an alternative method to estimate the significance of an observed BioTIP (or Ic) score in a system. This measurement makes more sense to evaluate random scores of sample-label shuffling, in which the nature sample-sample correlation within a phenotypic state (or cell subpopulation) was removed.

Usage

```
plot_SS_Simulation(
  Ic,
  simulation,
  las = 0,
  xlim = NULL,
  ylim = NULL,
  order = NULL,
  main = "1st max - 2nd max",
  ylab = "1st max - 2nd max"
)
```

Arguments

| | |
|-------------|---|
| Ic | A vector with names of states. If order is not assigned, then plot by the order of this vector. |
| simulation | A numeric matrix of Ic scores in which rows are states and columns are numbers of simulated times. It can be obtained from simulation_Ic |
| las | Numeric in 0, 1, 2, 3; the style of axis labels. Default is 0, meaning labels are parallel. (link to http://127.0.0.1:21580/library/graphics/html/par.html). |
| xlim | An integer vector of length 2. Default is NULL. |
| ylim | An integer vector of length 2. Default is NULL. |
| order | Characters of names of Ic to be plotted in a desired order. Default is NULL. |
| main | A character vector. The title of the plot. Default is NULL. |
| ylab | titles y axes, as in plot. |
| which2point | A character (or integer) which state's values were used to set up the three horizontal lines. by default is NULL, indicating the values of all states will be used. |

Value

Return a plot of the observed Ic (red) and simulated Ic (grey) scores per state.

Author(s)

Xinan H Yang <xyang2@uchicago.edu>

Examples

```
sim = matrix(sample(1:10, 9), 3, 3)
row.names(sim) = paste0('state', 1:3)
Ic = c('state1' = 3.4, 'state2' = 15.6, 'state3' = 2)
plot_SS_Simulation(Ic, sim)
```

sd_selection

Selecting Highly Oscillating Transcripts

Description

sd_selection pre-selects highly oscillating transcripts from the input dataset df. The dataset must contain multiple sample groups (or 'states'). For each state, the function filters the dataset using a cutoff value for standard deviation. The default cutoff value is 0.01 (i.e., higher than the top 1 percentage standard deviation).

Usage

```
sd_selection(
  df,
  samplesL,
  cutoff = 0.01,
  method = c("other", "reference", "previous", "itself", "longitudinal reference"),
  control_df = NULL,
  control_samplesL = NULL
)
```

Arguments

| | |
|----------|--|
| df | A numeric matrix or data frame. The rows and columns represent unique transcript IDs (geneID) and sample names, respectively. |
| samplesL | A list of vectors, whose length is the number of states. Each vector gives the sample names in a state. Note that the vectors (sample names) has to be among the column names of the R object 'df'. |
| cutoff | A positive numeric value. Default is 0.01. If < 1, automatically selects top x transcripts using the a selecting method (which is either the reference, other stages or previous stage), e.g. by default it will select top 1 percentage of the transcripts. |

| | |
|------------------|---|
| method | <p>Selection of methods from reference,other, previous, default uses other. Partial match enabled.</p> <ul style="list-style-type: none"> • itself, or longitudinal reference. Some specific requirements for each option: • reference, the reference has to be the first. • previous, make sure sampleL is in the right order from benign to malign. • itself, make sure the cutoff is smaller than 1. • longitudinal reference make sure control_df and control_samplesL are not NULL. The row numbers of control_df is the same as df and all transcripts in df are also in control_df. |
| control_df | A count matrix with unique loci as row names and samples names of control samples as column names, only used for method longitudinal reference |
| control_samplesL | A list of characters with stages as names of control samples, required for method 'longitudinal reference' |

Value

sd_selection() A list of data frames, whose length is the number of states. The rows in each data frame are the filtered transcripts with highest standard deviation selected from df and based on an assigned cutoff value. Each resulting data frame represents a subset of the raw input df, with the sample ID of the same state in the column.

Author(s)

Zhezhen Wang <zhezhen@uchicago.edu>

See Also

[optimize.sd_selection](#)

Examples

```
counts = matrix(sample(1:100, 18), 2, 9)
colnames(counts) = 1:9
row.names(counts) = c('loci1', 'loci2')
cli = cbind(1:9, rep(c('state1', 'state2', 'state3'), each = 3))
colnames(cli) = c('samples', 'group')
samplesL <- split(cli[, 1], f = cli[, 'group'])
test_sd_selection <- sd_selection(counts, samplesL, 0.01)
```

simulationMCI *Get MCI Scores for randomly selected genes*

Description

This function gets the MCI scores for randomly selected features (e.g. transcript ids),

Usage

```
simulationMCI(  
  len,  
  samplesL,  
  df,  
  adjust.size = FALSE,  
  B = 1000,  
  fun = c("cor", "BioTIP")  
)
```

Arguments

| | |
|-------------|---|
| len | An integer that is the length of genes in the CTS (critical transition signal). |
| samplesL | A list of vectors, whose length is the number of states. Each vector gives the sample names in a state. Note that the vector <i>s</i> (sample names) has to be among the column names of the R object 'df'. |
| df | A numeric matrix or dataframe of numerics, factor or character. The rows and columns represent unique transcript IDs (geneID) and sample names, respectively |
| adjust.size | A boolean value indicating if MCI score should be adjust by module size (the number of transcripts in the module) or not. Default FALSE. |
| B | An integer, setting the permutation with B runs. Default is 1000. |
| fun | A character chosen between ("cor", "BioTIP"), indicating where an adjusted correlation matrix will be used to calculated the MCI score. |

Value

A numeric matrix indicating the MCI scores of permutation. The dimension (row X column) of this matrix is the length of samplesL * B.

Author(s)

Zhezhen Wang <zhezhen@uchicago.edu>; Xinan H Yang <xyang2@uchicago.edu>

Examples

```

counts = matrix(sample(1:100, 18), 3, 9)
colnames(counts) = 1:9
row.names(counts) = c('loci1', 'loci2', 'loci3')
cli = cbind(1:9, rep(c('state1', 'state2', 'state3'), each = 3))
colnames(cli) = c('samples', 'group')
samplesL <- split(cli[, 1], f = cli[, 'group'])
simMCI = simulationMCI(2, samplesL, counts, B=2)
simMCI
#           [,1]      [,2]
#state1  2.924194  2.924194
#state2 20.877138 20.877138
#state3  2.924194  2.924194

```

| | |
|---------------|--|
| simulation_Ic | <i>Calculating random Index of Critical transition (Ic scores) for randomly-selected genes</i> |
|---------------|--|

Description

Simulating Ic scores for x randomly selected samples, where x should be the same as the length of identified critical-transition signal (CTS) (e.g., number of genes) and B is self-defined running times.

Usage

```

simulation_Ic(
  obs.x,
  sampleL,
  counts,
  B = 1000,
  fun = c("cor", "BioTIP"),
  shrink = TRUE,
  use = c("everything", "all.obs", "complete.obs", "na.or.complete",
         "pairwise.complete.obs"),
  output = c("Ic", "PCCg", "PCCs")
)

```

Arguments

| | |
|---------|--|
| obs.x | An integer, length of identified CTS. |
| sampleL | A list of vectors, whose length is the number of states. Each vector gives the sample names in a state. Note that the vector s (sample names) has to be among the column names of the R object 'df'. |
| counts | A numeric matrix or dataframe in which columns are samples and rows are transcripts. Each row needs to have a unique row name (i.e. transcript ID). |

| | |
|--------|---|
| B | An integer, setting the permutation with B runs. Default is 1000. |
| fun | An optional character string indicating the R function to calculate correlations for all possible pairs of columns of a matrix. When using "BioTIP", The method is modified to ignore missing values, analogous to how <code>cor(X, use = "pairwise.complete.obs")</code> works. Note that the "BioTIP" option only function together with <code>shrink = TRUE</code> . |
| shrink | A flag specifying whether to shrink the correlation or not. This approach uses the method outlined by Schafer and Strimmer in "A Shrinkage Approach to Large-Scale Covariance Matrix Estimation and Implications for Functional Genomics" (2005) Comparing to <code>fun='cor'</code> , the 'BioTIP' method without shrinkage is modified to ignore missing values, analogous to how <code>cor(X, use = "pairwise.complete.obs")</code> works. |
| use | An optional character string, when <code>fun=="cor"</code> , it gives a method for computing covariances in the presence of missing values. This must be (an abbreviation of) one of the strings "everything", "all.obs", "complete.obs", "na.or.complete", or "pairwise.complete.obs". |
| output | A string. Please select from 'Ic', 'PCCg', or 'PCCs'. Uses 'Ic' by default. 'PCCg' is the PCC between genes (numerator) and 'PCCs' is PCC between samples (denominator) |

Value

A matrix of y rows and B columns where y is the length of `sampleL` and B is self-defined. Each column is a set of Ic scores calculated for each state

Author(s)

Zhezhen Wang <zhezhen@uchicago.edu>

Examples

```
counts = matrix(sample(1:100, 27), 3, 9)
colnames(counts) = 1:9
row.names(counts) = c('loci1', 'loci2', 'loci3')
cli = cbind(1:9, rep(c('state1', 'state2', 'state3'), each = 3))
colnames(cli) = c('samples', 'group')
samplesL <- split(cli[, 1], f = cli[, 'group'])
simulation_Ic(2, samplesL, counts, B =3, fun="BioTIP", shrink=TRUE)
```

`simulation_Ic_sample` *Calculating (and plot) random Ic scores (Mojtahedi et al. 2016) based on shuffling sample labelling.*

Description

Run B times of sample-label shuffling to calculate the Ic score, where x should be the same as the length of identified BioTIP and B is self-defined.

Usage

```
simulation_Ic_sample(
  counts,
  sampleNo,
  Ic = NULL,
  genes,
  B = 1000,
  ylim = NULL,
  main = "simulation of samples",
  fun = c("cor", "BioTIP"),
  shrink = TRUE,
  use = c("everything", "all.obs", "complete.obs", "na.or.complete",
    "pairwise.complete.obs"),
  output = c("Ic", "PCCg", "PCCs"),
  plot = FALSE
)
```

Arguments

| | |
|----------|---|
| counts | A numeric matrix or data frame. The rows and columns represent unique transcript IDs (geneID) and sample names, respectively. |
| sampleNo | An integer of sample size at the tipping-point state. |
| Ic | A numeric value. Ic score of identified CTS (gene-set), useful when plot is TRUE. Default is NULL. |
| genes | A character vector of identified CTS gene unique ids. |
| B | An integer indicating number of times to run this simulation, default 1000. |
| ylim | An integer vector of length 2. Default is NULL. |
| main | A character vector. The title of the plot. Default is NULL. |
| fun | An optional character string indicating the R function to calculate correlations for all possible pairs of columns of a matrix. When using "BioTIP", The method is modified to ignore missing values, analogous to how <code>cor(X, use = "pairwise.complete.obs")</code> works. Note that the "BioTIP" option only function together with <code>shrink = TRUE</code> . |
| shrink | A flag specifying whether to shrink the correlation or not. This approach uses the method outlined by Schafer and Strimmer in "A Shrinkage Approach to Large-Scale Covariance Matrix Estimation and Implications for Functional Genomics" (2005) Comparing to <code>fun='cor'</code> , the 'BioTIP' method without shrinkage is modified to ignore missing values, analogous to how <code>cor(X, use = "pairwise.complete.obs")</code> works. |
| use | An optional character string, when <code>fun=="cor"</code> , it gives a method for computing covariances in the presence of missing values. This must be (an abbreviation of) one of the strings "everything", "all.obs", "complete.obs", "na.or.complete", or "pairwise.complete.obs". |
| output | A string. Please select from 'Ic', 'PCCg', or 'PCCs'. Uses 'Ic' by default. 'PCCg' is the PCC between genes (numerator) and 'PCCs' is PCC between samples (denominator) |

Index

* datasets

- cod, [4](#)
- gencode, [6](#)
- GSE6136_cli, [21](#)
- GSE6136_matrix, [22](#)
- ILEF, [22](#)
- intron, [23](#)

avg.cor.shrink, [2](#)

cluster_walktrap, [8](#)

cod, [4](#)

communities, [9](#)

cor.shrink, [4](#)

dist, [8](#)

gencode, [6](#)

getBiotypes, [6](#)

getCluster_methods, [8](#), [13](#), [16](#)

getCTS, [9](#)

getIc, [10](#)

getIc.new, [12](#)

getMaxMCImember, [10](#), [11](#), [13](#), [15](#), [28](#)

getMaxStats, [10](#), [15](#), [30](#)

getMCI, [15](#), [16](#), [16](#), [28](#)

getMCI_inner, [18](#)

getNetwork, [8](#), [19](#)

getReadthrough, [20](#)

GSE6136_cli, [21](#)

GSE6136_matrix, [22](#)

hclust, [8](#)

igraph, [17](#)

ILEF, [22](#)

intron, [23](#)

KMedoids, [8](#)

matplotlib, [27](#)

optimize.sd_selection, [23](#), [33](#)

plot.igraph, [9](#)

plot_Ic_Simulation, [28](#)

plot_MCI_Simulation, [30](#)

plot_SS_Simulation, [31](#)

plotBar_MCI, [25](#)

plotIc, [26](#)

plotMaxMCI, [27](#)

sd_selection, [9](#), [17](#), [24](#), [32](#)

simulation_Ic, [29](#), [31](#), [35](#)

simulation_Ic_sample, [36](#)

simulationMCI, [30](#), [34](#)