

Package ‘CancerSubtypes’

October 17, 2017

Type Package

Title Cancer subtypes identification, validation and visualization
based on genomic data

Version 1.2.5

Date 2016-06-30

Author Taosheng Xu<taosheng.x@gmail.com>, Thuc
Le<Thuc.Le@unisa.edu.au>

Maintainer Taosheng Xu<taosheng.x@gmail.com>

Depends R (>= 3.4), sigclust, NMF

Imports SNFtool, iCluster, survival, cluster, impute, limma,
ConsensusClusterPlus, grDevices

Description CancerSubtypes integrates the current common computational
biology methods for cancer subtypes identification and provides
a standardized framework for cancer subtype analysis based on
the genomic datasets.

License GPL (>= 2)

Suggests BiocGenerics, knitr, RTCGA.mRNA, RTCGA.clinical

VignetteBuilder knitr

biocViews Clustering, Software, Visualization, GeneExpression

RoxygenNote 6.0.1

NeedsCompilation no

R topics documented:

data.checkDistribution	2
data.imputation	3
data.normalization	4
DiffExp.limma	4
drawHeatmap	5
ExecuteCC	7
ExecuteCNMF	9
ExecuteiCluster	10
ExecuteSNF	12
ExecuteSNF.CC	13
FSbyCox	14

FSbyMAD	15
FSbyPCA	16
FSbyVar	17
GeneExp	18
miRNAExp	18
saveFigure	19
sigclustTest	20
silhouette_SimilarityMatrix	21
spectralAlg	22
status	23
survAnalysis	23
time	24

Index **25**

`data.checkDistribution`
Data check distribution

Description

Data check distribution

Usage

`data.checkDistribution(Data)`

Arguments

Data A matrix representing the genomic data such as gene expression data, miRNA expression data. For the matrix, the rows represent the genomic features, and the columns represent the samples.

Value

A plot describes the mean, variance and Median Absolute Deviation (MAD) distribution of features.

Examples

```
data(GeneExp)
data.checkDistribution(GeneExp)
```

data.imputation	<i>Data imputation</i>
-----------------	------------------------

Description

Data imputation for features with missing values

Usage

```
data.imputation(Data, fun = "median")
```

Arguments

- | | |
|------|--|
| Data | A matrix representing the genomic data such as gene expression data, miRNA expression data.
For the matrix, the rows represent the genomic features, and the columns represent the samples. |
| fun | <p>A character value representing the imputation type. The optional values are shown below:</p> <ul style="list-style-type: none"> • "median". The NAs will be replaced by the median of the existing values of this feature in all samples. • "mean". The NAs will be replaced by the mean of the existing values of this feature in all samples. • "microarray". It will apply the "impute" package to impute the missing values. This is a common way to process the missing observation for MicroArray dataset. |

Value

The data matrix after imputation (without NAs).

Examples

```
Data=matrix(runif(1000),nrow = 50,ncol = 20)
geneName=paste("Gene", 1:50, sep = " ")
sampleName=paste("Sample", 1:20, sep = " ")
rownames(Data)=geneName
colnames(Data)=sampleName
index=sample(c(1:1000),60)
Data[index]=NA
result=data.imputation(Data,fun="median")
```

`data.normalization` *Data normalization*

Description

Conduct normalization for dataset.

Usage

```
data.normalization(Data, type = "feature_Median", log2 = FALSE)
```

Arguments

Data	A matrix representing the genomic data such as gene expression data, miRNA expression data. For the matrix, the rows represent the genomic features, and the columns represent the samples.
type	A character value representing the normalization type. The optional values are shown below: <ul style="list-style-type: none"> • "feature_Median". The default value. Normalize dataset by sweeping the median values of each feature. • "feature_Mean". Normalize dataset by sweeping the mean values of each feature. • "feature_zscore". Conduct z_score normalization for each feature. • "sample_zscore". Conduct z_score normalization for each samples.
log2	A logical value. If TRUE, the data is transform as $\log_2(x+1)$. This is commonly used for RNAseq data.

Value

The normalized data matrix.

Examples

```
data(GeneExp)
result=data.normalization(GeneExp,type="feature_Median",log2=FALSE)
```

DiffExp.limma *DiffExp.limma*

Description

Differently Expression Analysis for genomic data. We apply limma package to conduct the analysis.

Usage

```
DiffExp.limma(Tumor_Data, Normal_Data, group = NULL, topk = NULL,
              RNAseq = FALSE)
```

Arguments

Tumor_Data	A matrix representing the genomic data of cancer samples such as gene expression data, miRNA expression data. For the matrix, the rows represent the genomic features, and the columns represent the cancer samples.
Normal_Data	A matrix representing the genomic data of Normal samples. For the matrix, the rows represent the genomic features corresponding to the Tumor_Data, and the columns represent the normal samples.
group	A vector representing the subtype of each tumor sample in the Tumor_Data. The length of group is equal to the column number of Tumor_Data.
topk	The top number of different expression features that we want to extract in the return result.
RNAseq	A bool type representing the input datatype is a RNASeq or not. Default is FALSE for microarray data.

Value

A list representing the differently expression for each subtype comparing to the Normal group.

Author(s)

Xu,Taosheng <taosheng.x@gmail.com>,Thuc Le <Thuc.Le@unisa.edu.au>

References

Smyth, Gordon K. "Limma: linear models for microarray data." *Bioinformatics and computational biology solutions using R and Bioconductor*. Springer New York, 2005. 397-420.

Examples

```
data(GeneExp)
data(miRNAExp)
GBM=list(GeneExp=GeneExp,miRNAExp=miRNAExp)
result=ExecuteSNF(GBM, clusterNum=3, K=20, alpha=0.5, t=20)
group=result$group
#####Fabricate a normal group by extracting some samples from the cancer dataset
#####for demonstrating the examples.
Normal_Data=GeneExp[,sample(1:100,20)]
result=DiffExp.limma(Tumor_Data=GeneExp,Normal_Data=Normal_Data,group=group,topk=NULL,RNAseq=FALSE)
```

drawHeatmap

Generate heatmaps

Description

Generate heatmap for datasets.

Usage

```
drawHeatmap(data, group = NULL, silhouette = NULL, scale = "no",
  labRow = NULL, labCol = NULL, color = colorRampPalette(c("green",
  "black", "red"))(300), Title = NA)
```

Arguments

data	A matrix representing the genomic data such as gene expression data, miRNA expression data. For the matrix, the rows represent the genomic features, and the columns represent the samples.
group	A vector representing the subtype on each sample. The default is NULL. If it is not NULL, the samples will be rearrangement according to the subtypes in the heatmap.
silhouette	An object of class silhouette. It is a result from function silhouette() or silhouette_SimilarityMatrix(). The default is NULL. If it is not NULL, an annotation will be drawn to show the silhouette width for each sample.
scale	A string for data normalization type before heatmap drawing. The optional values are shown below: <ul style="list-style-type: none"> • "no". No normalization. This is default. • "z_score". Normalize data by z_score of features. • "max_min". Normalize each feature by $(\text{value}-\text{min})/(\text{max}-\text{min})$.
labRow	labels for the rows. Possible values are: <ul style="list-style-type: none"> • NULL. The default value. It will use the row names of the matrix for the heatmap labels. • NA. No row label will be shown. • A list of labels.
labCol	labels for the columns. See labRow.
color	color specification for the heatmap.
Title	A string for the Main title of the heatmap.

Details

We applied the R package "NMF" function "aheatmap()" as the heatmap drawer.

Value

A heatmap

Author(s)

Xu, Taosheng <taosheng.x@gmail.com>, Thuc Le <Thuc.Le@unisa.edu.au>

References

Gaujoux, Renaud, and Cathal Seoighe. "A flexible R package for nonnegative matrix factorization." *BMC bioinformatics* 11.1 (2010): 1.

Examples

```
### SNF result analysis
data(GeneExp)
data(miRNAExp)
data(time)
data(status)
GBM=list(GeneExp=GeneExp, miRNAExp=miRNAExp)
```

```
result=ExecuteSNF(GBM, clusterNum=3, K=20, alpha=0.5, t=20)
group=result$group
distanceMatrix=result$distanceMatrix
silhouette=silhouette_SimilarityMatrix(group, distanceMatrix)
drawHeatmap(GeneExp,group,silhouette=silhouette,scale="max_min",Title="GBM Gene Expression")
drawHeatmap(GeneExp,group,silhouette=silhouette,scale="max_min",
            color="RdYlBu",Title="GBM Gene Expression")
```

ExecuteCC*Execute Consensus Clustering*

Description

This function is based on the R package "ConsensusClusterPlus". We write a shell to unify the input and output format. It is helpful for the standardized flow of cancer subtypes analysis and validation. The parameters are compatible to the original R package "ConsensusClusterPlus" function "ConsensusClusterPlus()".

Please note: we add a new parameter "clusterNum" which represents the result with cancer subtypes group we want to return.

Usage

```
ExecuteCC(clusterNum, d, maxK = 10, clusterAlg = "hc",
         distance = "pearson", title = "ConsensusClusterResult", reps = 500,
         pItem = 0.8, pFeature = 1, plot = "png", innerLinkage = "average",
         finallinkage = "average", writeTable = FALSE, weightsItem = NULL,
         weightsFeature = NULL, verbose = FALSE, corUse = "everything")
```

Arguments

- | | |
|-------------------------|---|
| <code>clusterNum</code> | A integer representing the return cluster number, this value should be less than <code>maxClusterNum(maxK)</code> . This is the only additional parameter in our function compared to the original R package "ConsensusClusterPlus". All the other parameters are compatible to the function "ConsensusClusterPlus()". |
| <code>d</code> | data to be clustered; either a data matrix where columns=items/samples and rows are features. For example, a gene expression matrix of genes in rows and microarrays in columns, or <code>ExpressionSet</code> object, or a distance object (only for cases of no feature resampling)
Please Note: We add a new data type (<code>list</code>) for this parameter. Please see details and examples. |
| <code>maxK</code> | integer value. maximum cluster number for Consensus Clustering Algorithm to evaluate. |
| <code>clusterAlg</code> | character value. cluster algorithm. 'hc' heirarchical (<code>hclust</code>), 'pam' for partitioning around medoids, 'km' for k-means upon data matrix, 'kmdist' for k-means upon distance matrices (former km option), or a function that returns a clustering. |
| <code>distance</code> | character value. 'pearson': (1 - Pearson correlation), 'spearman' (1 - Spearman correlation), 'euclidean', 'binary', 'maximum', 'canberra', 'minkowski' or custom distance function. |
| <code>title</code> | character value for output directory. This title can be an absolute or relative path |

reps	integer value. number of subsamples(in other words, The iteration number of each cluster number)
pItem	Please refer to the "ConsensusClusterPlus" package for detailed information.
pFeature	Please refer to the "ConsensusClusterPlus" package for detailed information.
plot	Please refer to the "ConsensusClusterPlus" package for detailed information.
innerLinkage	Please refer to the "ConsensusClusterPlus" package for detailed information.
finalLinkage	Please refer to the "ConsensusClusterPlus" package for detailed information.
writeTable	Please refer to the "ConsensusClusterPlus" package for detailed information.
weightsItem	Please refer to the "ConsensusClusterPlus" package for detailed information.
weightsFeature	Please refer to the "ConsensusClusterPlus" package for detailed information.
verbose	Please refer to the "ConsensusClusterPlus" package for detailed information.
corUse	Please refer to the "ConsensusClusterPlus" package for detailed information.

Details

If the data is a list containing the matched mutli-genomics data matrices like the input as "ExecuteCluster()" and "ExecuteSNF()", we use "z-score" to normalize features for each data matrix first. Then all the normalized data matrices from the data list are concatenated according to samples. The concatenated data matrix is the samples with a long features (all features in the data list). Our purpose is to make convenient comparing the different method with same dataset format. See examples.

Value

A list with the following elements.

- `group`
A vector represent the group of cancer subtypes. The order is corresponding to the the samples in the data matrix.
- `distanceMatrix`
It is a sample similarity matrix. The more large value between samples in the matrix, the more similarity the samples are.
- `originalResult`
The clustering result of the original function "ConsensusClusterPlus()"

References

Monti, S., Tamayo, P., Mesirov, J., Golub, T. (2003) Consensus Clustering: A Resampling-Based Method for Class Discovery and Visualization of Gene Expression Microarray Data. Machine Learning, 52, 91-118.

See Also

`ConsensusClusterPlus`

Examples

```
### The input dataset is a single gene expression matrix.
data(GeneExp)
data(miRNAExp)
result1=ExecuteCC(clusterNum=3,d=GeneExp,maxK=10,clusterAlg="hc",distance="pearson",title="GBM")
result1$group

### The input dataset is multi-genomics data as a list
GBM=list(GeneExp=GeneExp,miRNAExp=miRNAExp)
result2=ExecuteCC(clusterNum=3,d=GBM,maxK=5,clusterAlg="hc",distance="pearson",title="GBM")
result2$group
```

ExecuteCNMF

Execute Consensus NMF (Nonnegative matrix factorization)

Description

Brunet applied nonnegative matrix factorization (NMF) to analyze the Gene MicroArray dataset in 2004. In the original paper, the author proved that NMF is an efficient method for distinct molecular patterns identification and provides a powerful method for class discovery. This method was implemented in an R package "NMF". Here we applied the "NMF" package to conduct the cancer subtypes identification. We write a shell to unify the input and output format. It is helpful for the standardized flow of cancer subtypes analysis and validation. The R package "NMF" should be installed.

Usage

```
ExecuteCNMF(datasets, clusterNum, nrun = 30)
```

Arguments

datasets	A data matrix or a list containing data matrices. For each data matrix, the rows represent genomic features, and the columns represent samples. If the matrices have negative values, first the negative values will be set to zero to get a matrix 1; all the positive values will be set to zero to get the matrix 2; then a new matrix with all positive values will be get by concatenating matrix1 and -matrix2.
clusterNum	Number of subtypes for the samples
nrun	Number of runs to perform NMF. A default of 30 runs are performed, allowing the computation of a consensus matrix that is used in selecting the best result for cancer subtypes identification as Consensus Clustering method.

Details

If the data is a list containing the matched mutli-genomics data matrices like the input as "ExecuteCluster()" and "ExecuteSNF()", The data matrices in the list are concatenated according to samples. The concatenated data matrix is the samples with a long features (all features in the data list). Our purpose is to make convenient comparing the different method with same dataset format. See examples.

Value

A list with the following elements.

- `group`
A vector represent the group of cancer subtypes. The order is corresponding to the samples in the data matrix.
- `distanceMatrix`
It is a sample similarity matrix. The more large value between samples in the matrix, the more similarity the samples are.
- `originalResult` A NMFfitX class from the result of function "nmf()".

References

[1] Brunet, Jean-Philippe, Pablo Tamayo, Todd R Golub, and Jill P Mesirov. "Metagenes and Molecular Pattern Discovery Using Matrix Factorization." *Proceedings of the National Academy of Sciences* 101, no. 12 (2004):4164-69.

[2] Gaujoux, Renaud, and Cathal Seoighe. "A Flexible R Package for Nonnegative Matrix Factorization." *BMC Bioinformatics* 11 (2010): 367. doi:10.1186/1471-2105-11-367.

See Also

[nmf](#)

Examples

```
data(GeneExp)
#To save the execution time, the nrun is set to 5, but the recommended value is 30.
result=ExecuteCNMF(GeneExp,clusterNum=3,nrun=5)
result$group
```

ExecuteiCluster

Execute iCluster (Integrative clustering of multiple genomic data)

Description

Shen (2009) proposed a latent variable regression with a lasso constraint for joint modeling of multiple omics data types to identify common latent variables that can be used to cluster patient samples into biologically and clinically relevant disease subtypes.

This function is based on the R package "iCluster". The R package "iCluster" should be installed. We write a shell to unify the input and output format. It is helpful for the standardized flow of cancer subtypes analysis and validation. The parameters is compatible to the original R package "iCluster" function "iCluster2()".

Please note: The data matrices are transposed in our function comparing to the original R package "iCluster" on the behalf of the unified input format with other functions. We try to build a standardized flow for cancer subtypes analysis and validation.

Usage

```
ExecuteiCluster(datasets, k, lambda = NULL, scale = TRUE, scalar = FALSE,
  max.iter = 10)
```

Arguments

<code>datasets</code>	A list containing data matrices. For each data matrix, the rows represent genomic features, and the columns represent samples. In order to unify the input parameter with other clustering methods, the data matrices are transposed comparing to the definition in the original "iCluster" package.
<code>k</code>	Number of subtypes for the samples
<code>lambda</code>	Penalty term for the coefficient matrix of the iCluster model
<code>scale</code>	Logical value. If true, the genomic features in the matrix is centered.
<code>scalar</code>	Logical value. If true, a degenerate version assuming scalar covariance matrix is used.
<code>max.iter</code>	maximum iteration for the EM algorithm

Details

For iCluster algorithm, it cannot process high-dimensional data, otherwise it is very very time-consuming or reports a mistake. Based on test, it could smoothly run for the matrix with around 1500 features. Normally it need feature selection step first to reduce feature number.

Value

A list with the following elements.

- `group`
A vector represent the group of cancer subtypes. The order is corresponding to the samples in the data matrix.
- `originalResult`
The clustering result of the original function "iCluster2()"

References

Ronglai Shen, Adam Olshen, Marc Ladanyi. (2009). Integrative clustering of multiple genomic data types using a joint latent variable model with application to breast and lung cancer subtype analysis. *Bioinformatics* 25, 2906-2912.

Ronglai Shen, Qianxing Mo, Nikolaus Schultz, Venkatraman E. Seshan, Adam B. Olshen, Jason Huse, Marc Ladanyi, Chris Sander. (2012). Integrative Subtype Discovery in Glioblastoma Using iCluster. *PLoS ONE* 7, e35236

See Also

[iCluster2](#)

Examples

```
data(GeneExp)
data(miRNAExp)
data1=FSbyVar(GeneExp, cut.type="topk", value=500)
data2=FSbyVar(miRNAExp, cut.type="topk", value=100)
GBM=list(GeneExp=data1, miRNAExp=data2)
result=ExecuteiCluster(datasets=GBM, k=3, lambda=list(0.44,0.33,0.28))
result$group
```

ExecuteSNF

*Execute SNF(Similarity Network Fusion)***Description**

SNF is a multi-omics data processing method that constructs a fusion patient similarity network by integrating the patient similarity obtained from each of the genomic data types. SNF calculates the similarity between patients using each single data type separately. The similarities between patients from different data types are then integrated by a cross-network diffusion process to construct the fusion patient similarity matrix. Finally, a clustering method is applied to the fusion patient similarity matrix to cluster patients into different groups, which imply different cancer subtypes. This function is based on the R package "SNFtool". The R package "SNFtool" should be installed. We write a function to integrate the clustering process and unify the input and output format. It is helpful for the standardized flow of cancer subtypes analysis and validation.

Please note: The data matrices are transposed in our function comparing to the original R package "SNFtools". We try to build a standardized flow for cancer subtypes analysis and validation.

Usage

```
ExecuteSNF(datasets, clusterNum, K = 20, alpha = 0.5, t = 20,
           plot = TRUE)
```

Arguments

datasets	A list containing data matrices. For each data matrix, the rows represent genomic features, and the columns represent samples.
clusterNum	A integer representing the return cluster number
K	Number of nearest neighbors
alpha	Variance for local model
t	Number of iterations for the diffusion process
plot	Logical value. If true, draw the heatmap for the distance matrix with samples ordered to form clusters.

Value

A list with the following elements.

- group
A vector represents the group of cancer subtypes. The order is corresponding to the samples in the data matrix.
- distanceMatrix
It is a samplesimilarity matrix. The more large value between samples in the matrix, the more similarity the samples are.

References

B Wang, A Mezlini, F Demir, M Fiume, T Zu, M Brudno, B Haibe-Kains, A Goldenberg (2014) Similarity Network Fusion: a fast and effective method to aggregate multiple data types on a genome wide scale. Nature Methods. Online. Jan 26, 2014

See Also

[affinityMatrix SNF](#)

Examples

```
data(GeneExp)
data(miRNAExp)
GBM=list(GeneExp=GeneExp,miRNAExp=miRNAExp)
result=ExecuteSNF(GBM, clusterNum=3, K=20, alpha=0.5, t=20)
result$group
```

ExecuteSNF.CC	<i>Execute the combined SNF (Similarity Network Fusion) and Consensus clustering</i>
---------------	--

Description

This function is a combined process of SNF and Consensus Clustering for cancer subtypes identification. First it applied SNF to get the fusion patients similarity matrix. Then use this fusion patients similarity matrix as the sample distance for Consensus Clustering.

Usage

```
ExecuteSNF.CC(datasets, clusterNum, K = 20, alpha = 0.5, t = 20,
  maxK = 10, pItem = 0.8, reps = 500, title = "ConsensusClusterResult",
  plot = "png", finalLinkage = "average")
```

Arguments

datasets	Same as ExecuteSNF
clusterNum	Same as ExecuteSNF
K	Same as ExecuteSNF
alpha	Same as ExecuteSNF
t	Same as ExecuteSNF
maxK	Same as ExecuteCC
pItem	Same as ExecuteCC
reps	Same as ExecuteCC
title	Same as ExecuteCC
plot	Same as ExecuteCC
finalLinkage	Same as ExecuteCC

Value

Same as the ExecuteCC(). A list with the following elements.

- **group**
A vector represent the group of cancer subtypes. The order is corresponding to the samples in the data matrix.
- **distanceMatrix**
It is a sample dissimilarity matrix. The more large value between samples in the matrix, the more dissimilarity the samples are.
- **originalResult**
The clustering result of the original function "ConsensusClusterPlus()"

See Also

[ExecuteSNF](#) [ExecuteCC](#)

Examples

```
data(GeneExp)
data(miRNAExp)
GBM=list(GeneExp,miRNAExp)
result=ExecuteSNF.CC(GBM, clusterNum=3, K=20, alpha=0.5, t=20,
                    maxK = 5, pItem = 0.8, reps=500,
                    title = "GBM", plot = "png",
                    finalLinkage = "average")
result$group
```

FSbyCox

Biological feature (such as gene) selection based on Cox regression model.

Description

Cox model (Proportional hazard model) is a statistical approach for survival risk analysis. We applied the univariate Cox model for feature selection. The proportional hazard assumption test is used to evaluate the significant level of each biological feature related to the survival result for samples. Eventually, the most significant genes are selected for clustering analysis.

Usage

```
FSbyCox(Data, time, status, cutoff = 0.05)
```

Arguments

Data	A data matrix representing the genomic data measured in a set of samples. For the matrix, the rows represent the genomic features, and the columns represent the samples.
time	A numeric vector representing the survival time (days) of a set of samples. Note that the order of the time should map the samples in the Data matrix.

status	A numeric vector representing the survival status of a set of samples. 0=alive/censored, 1=dead. Note that the order of the time should map the samples in the Data matrix.
cutoff	A numeric value in (0,1) representing whether the significant feature X_i is selected according to the Proportional Hazards Assumption p-value of the feature X_i . If $p\text{-value}(X_i) < \text{cutoff}$, the features X_i will be selected for downstream analysis. Normally the significant level is set to 0.05.

Value

A data matrix, extracted a subset with significant features from the input data matrix. The rows represent the significant features, and the columns represents the samples.

Author(s)

Xu, Taosheng <taosheng.x@gmail.com>, Thuc Le <Thuc.Le@unisa.edu.au>

References

Andersen, P. and Gill, R. (1982). Cox's regression model for counting processes, a large sample study. *Annals of Statistics* 10, 1100-1120.
 Therneau, T., Grambsch, P., *Modeling Survival Data: Extending the Cox Model*. Springer-Verlag, 2000.

Examples

```
data(GeneExp)
data(time)
data(status)
data1=FSbyCox(GeneExp, time, status, cutoff=0.05)
```

FSbyMAD	<i>Biological feature (such as gene) selection based on the most variant Median Absolute Deviation (MAD).</i>
---------	---

Description

Biological feature (such as gene) selection based on the most variant Median Absolute Deviation (MAD).

Usage

```
FSbyMAD(Data, cut.type = "topk", value)
```

Arguments

Data	A data matrix representing the genomic data measured in a set of samples. For the matrix, the rows represent the genomic features, and the columns represents the samples.
cut.type	A character value representing the selection type. The optional values are shown below:

- "topk"
 - "cutoff"
- value A numeric value.
 If the cut.type="topk", the top number of value features are selected.
 If the cut.type="cutoff", the features with (MAD>value) are selected.

Value

An extracted subset data matrix with the most variant MAD features from the input data matrix.

Author(s)

Xu,Taosheng <taosheng.x@gmail.com>, Thuc Le <Thuc.Le@unisa.edu.au>

Examples

```
data(GeneExp)
data1=FSbyMAD(GeneExp, cut.type="topk",value=1000)
```

FSbyPCA

Biological feature (such as gene) dimension reduction and extraction based on Principal Component Analysis.

Description

This function is based on the prcomp(), we write a shell for it and make it easy to use on genomic data.

Usage

```
FSbyPCA(Data, PC_percent = 1, scale = TRUE)
```

Arguments

- Data A data matrix representing the genomic data measured in a set of samples. For the matrix, the rows represent the genomic features, and the columns represents the samples.
- PC_percent A numeric values in [0,1] representing the ratio of principal component is selected.
- scale A bool variable, If true, the Data is normalized before PCA.

Value

A new matrix with full or part Principal Component in new projection space.

Author(s)

Xu,Taosheng <taosheng.x@gmail.com>,Thuc Le <Thuc.Le@unisa.edu.au>

Examples

```
data(GeneExp)
data1=FSbyPCA(GeneExp, PC_percent=0.9, scale = TRUE)
```

 FSbyVar

Biological feature (such as gene) selection based on the most variance.

Description

Biological feature (such as gene) selection based on the most variance.

Usage

```
FSbyVar(Data, cut.type = "topk", value)
```

Arguments

Data	A data matrix representing the genomic data measured in a set of samples. For the matrix, the rows represent the genomic features, and the columns represents the samples.
cut.type	A character value representing the selection type. The optional values are shown below: <ul style="list-style-type: none"> • "topk" • "cutoff"
value	A numeric value. If the cut.type="topk", the top number of value features are selected. If the cut.type="cutoff", the features with (var>value) are selected.

Value

An extracted subset data matrix with most variance features from the input data matrix.

Author(s)

Xu, Taosheng <taosheng.x@gmail.com>, Thuc Le <Thuc.Le@unisa.edu.au>

Examples

```
data(GeneExp)
data1=FSbyVar(GeneExp, cut.type="topk", value=1000)
```

GeneExp

Dataset: Gene expression

Description

A glioblastoma (GBM) gene expression dataset downloaded from TCGA. This is a small dataset with 1500 genes and 100 cancer samples extracted from gene expression data for examples.

Format

A data matrix

Details

- Rows are genes
- Columns are cancer samples

Examples

```
data(GeneExp)
```

miRNAExp

Dataset: miRNA expression

Description

A glioblastoma (GBM) miRNA expression dataset downloaded from TCGA. This is a small miRNA expression dataset with 470 miRNAs and 100 cancer samples extracted from miRNA expression data for examples.

Format

A data matrix

Details

- Rows are miRNAs
- Columns are cancer samples

Examples

```
data(miRNAExp)
```

saveFigure	<i>This function save the figure in the current plot.</i>
------------	---

Description

This function save the figure in the current plot.

Usage

```
saveFigure(foldername = NULL, filename = "saveFig", image_width = 10,  
           image_height = 10, image_res = 300)
```

Arguments

foldername	Character values. It specifies the folder name which will be created in the present working path.
filename	Character values. It specifies the saved file name.
image_width	the figure width
image_height	the figure height
image_res	the figure resolution

Value

A * .png file in the specified folder.

Author(s)

Xu,Taosheng <taosheng.x@gmail.com>,Thuc Le <Thuc.Le@unisa.edu.au>

Examples

```
data(GeneExp)  
data(miRNAExp)  
data(time)  
data(status)  
GBM=list(GeneExp=GeneExp,miRNAExp=miRNAExp)  
result=ExecuteSNF(GBM, clusterNum=3, K=20, alpha=0.5, t=20)  
group=result$group  
distanceMatrix=result$distanceMatrix  
p_value=survAnalysis(mainTitle="GBM",time,status,group,  
                    distanceMatrix=distanceMatrix,similarity=TRUE)  
saveFigure(foldername="GBM",filename="GBM",image_width=10,image_height=10,image_res=300)
```

sigclustTest

*A statistical method for testing the significance of clustering results.***Description**

SigClust (Statistical significance of clustering) is a statistical method for testing the significance of clustering results. SigClust can be applied to assess the statistical significance of splitting a data set into two clusters. SigClust studies whether clusters are really there, using the 2-means ($k = 2$) clustering index as a statistic. It assesses the significance of clustering by simulation from a single null Gaussian distribution. Null Gaussian parameters are estimated from the data. Here we apply the SigClust to assess the statistical significance of pairwise subtypes. "sigclust" package should be installed.

Usage

```
sigclustTest(Data, group, nsim = 1000, nrep = 1, icovest = 1)
```

Arguments

Data	A data matrix representing the genomic data measured in a set of samples. For the matrix, the rows represent the genomic features, and the columns represents the samples.
group	The subtypes label of each sample
nsim	This is a parameter inherited from sigclust() in "sigclust" Package. Number of simulated Gaussian samples to estimate the distribution of the clustering index for the main p-value computation.
nrep	This is a parameter inherited from sigclust() in "sigclust" Package. Number of steps to use in 2-means clustering computations (default=1, chosen to optimize speed).
icovest	This is a parameter inherited from sigclust() in "sigclust" Package. Covariance estimation type: 1. Use a soft threshold method as constrained MLE (default); 2. Use sample covariance estimate (recommended when diagnostics fail); 3. Use original background noise threshold estimate (from Liu, et al, (2008)) ("hard thresholding").

Value

A matrix indicates the p-value between pairwise subtypes.

Author(s)

Xu, Taosheng <taosheng.x@gmail.com>, Thuc Le <Thuc.Le@unisa.edu.au>

References

- Liu, Yufeng, Hayes, David Neil, Nobel, Andrew and Marron, J. S, 2008, Statistical Significance of Clustering for High-Dimension, Low-Sample Size Data, Journal of the American Statistical Association 103(483) 1281-1293.
- Huang, Hanwen, Yufeng Liu, Ming Yuan, and J. S. Marron. "Statistical Significance of Clustering Using Soft Thresholding." Journal of Computational and Graphical Statistics, no. just-accepted (2014): 00-00.

See Also[sigclust](#)**Examples**

```
data(GeneExp)
data(miRNAExp)
data(time)
data(status)
GBM=list(GeneExp=GeneExp,miRNAExp=miRNAExp)
result=ExecuteSNF(GBM, clusterNum=3, K=20, alpha=0.5, t=20)
group=result$group
sigclust1=sigclustTest(miRNAExp,group, nsim=500, nrep=1, icovest=3)
sigclust2=sigclustTest(miRNAExp,group, nsim=1000, nrep=1, icovest=1)
```

`silhouette_SimilarityMatrix`

Compute or Extract Silhouette Information from Clustering based on similarity matrix.

Description

Silhouette refers to a method of interpretation and validation of consistency within clusters of data. The technique provides a succinct graphical representation of how well each object lies within its cluster (From Wiki).

Note that: This function is a rewriting version of the function "silhouette()" in R package cluster. The original function "silhouette()" is to compute the silhouette information based on a dissimilarity matrix. Here the silhouette_SimilarityMatrix() is to solve the computation based on the similarity matrix. The result of the silhouette_SimilarityMatrix() is compatible to the function "Silhouette()".

Usage

```
silhouette_SimilarityMatrix(group, similarity_matrix)
```

Arguments

`group` A vector represent the cluster label for a set of samples.
`similarity_matrix` A similarity matrix between samples

Details

For each observation i , the return `sil[i,]` contains the cluster to which i belongs as well as the neighbor cluster of i (the cluster, not containing i , for which the average dissimilarity between its observations and i is minimal), and the silhouette width $s(i)$ of the observation.

Value

An object, `sil`, of class `silhouette` which is an $[n \times 3]$ matrix with attributes. The colnames correspondingly are `c("cluster", "neighbor", "sil_width")`.

Author(s)

Xu,Taosheng <taosheng.x@gmail.com>,Thuc Le <Thuc.Le@unisa.edu.au>

References

Rousseeuw, P.J. (1987) Silhouettes: A graphical aid to the interpretation and validation of cluster analysis. J. Comput. Appl. Math., 20, 53-65.

See Also

[silhouette](#)

Examples

```
data(GeneExp)
data(miRNAExp)
GBM=list(GeneExp=GeneExp,miRNAExp=miRNAExp)
result=ExecuteSNF(GBM, clusterNum=3, K=20, alpha=0.5, t=20)
sil=silhouette_SimilarityMatrix(result$group, result$distanceMatrix)
plot(sil)
###If use the silhouette(), the result is wrong because the input is a similarity matrix.
sil1=silhouette(result$group, result$distanceMatrix)
plot(sil1) ##wrong result
```

spectralAlg

This is an internal function but need to be exported for the function ExecuteSNF.CC() call.

Description

This is Spectral Clustering Algorithm extracted from SNFtools package spectralClustering() with a tiny modification.

Usage

```
spectralAlg(affinity, K, type = 3)
```

Arguments

affinity	Similarity matrix
K	Number of clusters
type	The variants of spectral clustering to use.

Value

A vector consisting of cluster labels of each sample.

Examples

```
####see the spectralClustering() in SNFtool package for the detail example.
data(miRNAExp)
Dist1=SNFtool::dist2(t(miRNAExp),t(miRNAExp))
W1 = SNFtool::affinityMatrix(Dist1, 20, 0.5)
group=spectralAlg(W1,3, type = 3)
```

status

Dataset: Survival status

Description

- A vector representing the survival status for GBM cancer patients matched with the "Gene-Exp" and "miRNAExp" . 0=alive or censored, 1=dead

Format

A numeric vector

Examples

```
data(status)
```

survAnalysis

Survival analysis(Survival curves, Log-rank test) and compute Silhouette information for cancer subtypes

Description

Survival analysis is a very common tool to explain and validate the cancer subtype identification result. It provides the significance testing and graphical display for the verification of the survival patterns between the identified cancer subtypes.

Usage

```
survAnalysis(mainTitle = "Survival Analysis", time, status, group,
  distanceMatrix = NULL, similarity = TRUE)
```

Arguments

mainTitle	A character will display in the result plot.
time	A numeric vector representing the survival time (days) of a set of samples.
status	A numeric vector representing the survival status of a set of samples. 0=alive/censored, 1=dead.
group	A vector represent the cluster label for a set of samples.
distanceMatrix	A data matrix represents the similarity matrix or dissimilarity matrix between samples. If NULL, it will not compute silhouette width and draw the plot.
similarity	A logical value. If TRUE, the distanceMatrix is a similarity distance matrix between samples. Otherwise a dissimilarity distance matrix between samples

Value

The log-rank test p-value

Author(s)

Xu,Taosheng <taosheng.x@gmail.com>,Thuc Le <Thuc.Le@unisa.edu.au>

Examples

```
### SNF result analysis
data(GeneExp)
data(miRNAExp)
data(time)
data(status)
data1=FSbyCox(GeneExp,time,status,cutoff=0.05)
data2=FSbyCox(miRNAExp,time,status,cutoff=0.05)
GBM=list(GeneExp=data1,miRNAExp=data2)
result=ExecuteSNF(GBM, clusterNum=3, K=20, alpha=0.5, t=20)
group=result$group
distanceMatrix=result$distanceMatrix
p_value=survAnalysis(mainTitle="GBM",time,status,group,
                    distanceMatrix=distanceMatrix,similarity=TRUE)
```

time

Dataset: Survival time

Description

- A vector representing the right censored survival time (days) for GBM cancer patients matched with the "GeneExp" and "miRNAExp" datasets.

Format

A numeric vector

Examples

```
data(time)
```


Index

*Topic **datasets**

- GeneExp, [18](#)
- miRNAExp, [18](#)
- status, [23](#)
- time, [24](#)

affinityMatrix, [13](#)

data.checkDistribution, [2](#)
data.imputation, [3](#)
data.normalization, [4](#)
DiffExp.limma, [4](#)
drawHeatmap, [5](#)

ExecuteCC, [7](#), [14](#)
ExecuteCNMF, [9](#)
ExecuteiCluster, [10](#)
ExecuteSNF, [12](#), [14](#)
ExecuteSNF.CC, [13](#)

FSbyCox, [14](#)
FSbyMAD, [15](#)
FSbyPCA, [16](#)
FSbyVar, [17](#)

GeneExp, [18](#)

iCluster2, [11](#)

miRNAExp, [18](#)

nmf, [10](#)

saveFigure, [19](#)
sigclust, [21](#)
sigclustTest, [20](#)
silhouette, [22](#)
silhouette_SimilarityMatrix, [21](#)
SNF, [13](#)
spectralAlg, [22](#)
status, [23](#)
survAnalysis, [23](#)

time, [24](#)