

Package ‘CAMTHC’

April 15, 2019

Type Package

Title Convex Analysis of Mixtures for Tissue Heterogeneity
Characterization

Version 1.0.0

Author Lulu Chen <luluchen@vt.edu>

Maintainer Lulu Chen <luluchen@vt.edu>

biocViews Software, CellBiology, GeneExpression

Description An R package for tissue heterogeneity characterization by convex analysis of mixtures (CAM). It provides basic functions to perform unsupervised deconvolution on mixture expression profiles by CAM and some auxiliary functions to help understand the subpopulation-specific results. It also implements functions to perform supervised deconvolution based on prior knowledge of molecular markers, S matrix or A matrix. Combining molecular markers from CAM and from prior knowledge can achieve semi-supervised deconvolution of mixtures.

License GPL-2

Encoding UTF-8

RoxygenNote 6.0.1

Depends R (>= 3.5)

Suggests knitr, rmarkdown, BiocStyle, testthat, GEOquery, rgl

VignetteBuilder knitr

Imports methods, rJava, BiocParallel, stats, Biobase,
SummarizedExperiment, corpcor, geometry, NMF, DMwR, pcaPP,
apcluster, graphics

SystemRequirements Java (>= 1.8)

BugReports <https://github.com/Lululuella/CAMTHC/issues>

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

git_branch RELEASE_3_8

git_last_commit 381fbfa

git_last_commit_date 2018-10-30

Date/Publication 2019-04-15

R topics documented:

CAMTHC-package	2
AfromMarkers	3
AS-accessor	4
CAM	5
CAMASest	6
CAMASObj-class	8
CAMMGCluster	9
CAMMGObj-class	10
CAMObj-class	10
CAMPrep	10
CAMPrepObj-class	12
cornerSort	12
MDL	13
MDLObj-class	14
MGsforA	15
MGstatistic	16
PCAmat	17
ratMix3	18
simplexplot	18
XWProj	19
Index	21

CAMTHC-package

*CAMTHC: A package for Tissue Heterogeneity Characterization.***Description**

The core function in this package is `CAM` which achieves fully unsupervised deconvolution on mixture expression profiles. Each step in `CAM` can also be performed separately by `CAMPrep`, `CAMMGCluster` and `CAMASest` in a more flexible workflow. `MGstatistic` can help extract a complete marker list from CAM results. `MDL` can help decide the underlying subpopulation number. With other functions, e.g. `AfromMarkers` and `MGstatistic`, this package can also perform supervised deconvolution based on prior knowledge of molecular markers, subpopulation-specific expression matrix (S) or proportion matrix (A). Semi-supervised deconvolution can be achieved by combining molecular markers from CAM and from prior knowledge to analyze mixture expressions.

References

Wang, N., Hoffman, E. P., Chen, L., Chen, L., Zhang, Z., Liu, C., ... Wang, Y. (2016). Mathematical modelling of transcriptional heterogeneity identifies novel markers and subpopulations in complex tissues. *Scientific Reports*, 6, 18909. <http://doi.org/10.1038/srep18909>

AfromMarkers

Proportion matrix estimation from marker genes

Description

This function estimates proportion matrix (A matrix) from observed mixture expression data based on marker genes.

Usage

```
AfromMarkers(data, MGlist, scaleRecover = TRUE)
```

Arguments

data	A data set that will be internally coerced into a matrix. Each row is a gene and each column is a sample. data should be in non-log linear space with non-negative numerical values (i.e. ≥ 0). Missing values are not supported. All-zero rows will be removed internally.
MGlist	A list of vectors, each of which contains known markers and/or CAM-detected markers for one subpopulation.
scaleRecover	If TRUE, scale ambiguity of each column vector in A matrix is removed based on sum-to-one constraint on each row.

Details

With the expression levels of subpopulation-specific marker genes, the relative proportions of constituent subpopulations are estimated by spatial median using [l1median](#). Marker genes could be from unsupervised/supervised detection or from literatures. Scale ambiguity is optionally removed based on sum-to-one constraint of rows.

Value

Return the estimated proportion matrix (A matrix).

Examples

```
#obtain data and marker genes
data(ratMix3)
S <- ratMix3$S
pMGstat <- MGstatistic(S, c("Liver", "Brain", "Lung"))
pMGlist.FC <- lapply(c("Liver", "Brain", "Lung"), function(x)
  rownames(pMGstat)[pMGstat$idx == x & pMGstat$OVE.FC > 10])

#estimate A matrix from markers
Aest <- AfromMarkers(ratMix3$X, pMGlist.FC)
```

Description

Accessors to proportion matrix and subpopulation-specific expression matrix estimated by CAM.

Usage

```
Amat(x, ...)

Smat(x, ...)

## S4 method for signature 'CAMObj'
Amat(x, k, usingPCA = TRUE)

## S4 method for signature 'CAMASObj'
Amat(x, usingPCA = TRUE)

## S4 method for signature 'CAMObj'
Smat(x, k, usingPCA = TRUE)

## S4 method for signature 'CAMASObj'
Smat(x, usingPCA = TRUE)
```

Arguments

<code>x</code>	a <code>CAMObj</code> object or a <code>CAMASObj</code> object.
<code>...</code>	additional argument list.
<code>k</code>	subpopulation number
<code>usingPCA</code>	If TRUE, A matrix is estimated by transforming dimension-reduced A matrix back to original space. Otherwise, A matrix is directly estimated in original data space. The default is TRUE.

Value

Estimated A matrix or S matrix.

Examples

```
#obtain data
data(ratMix3)
data <- ratMix3$X

rCAM <- CAM(data, K = 3, dim.rdc = 3, thres.low = 0.30, thres.high = 0.95)
Aest <- Amat(rCAM, 3)
Sest <- Smat(rCAM, 3)

Aest <- Amat(slot(rCAM, "ASestResult")[[1]])
Sest <- Smat(slot(rCAM, "ASestResult")[[1]])
```

Description

This function performs a fully unsupervised computational deconvolution to identify marker genes that define each of the multiple subpopulations, and estimate the proportions of these subpopulations in the mixture tissues as well as their respective expression profiles.

Usage

```
CAM(data, K = NULL, corner.strategy = 2, dim.rdc = 10, thres.low = 0.05,
     thres.high = 0.95, cluster.method = c("K-Means", "apcluster"),
     cluster.num = 50, MG.num.thres = 20, lof.thres = 0.02, cores = NULL)
```

Arguments

<code>data</code>	Matrix of mixture expression profiles. Data frame, SummarizedExperiment or ExpressionSet object will be internally coerced into a matrix. Each row is a gene and each column is a sample. Data should be in non-log linear space with non-negative numerical values (i.e. ≥ 0). Missing values are not supported. All-zero rows will be removed internally.
<code>K</code>	The candidate subpopulation number(s), e.g. <code>K = 2:8</code> .
<code>corner.strategy</code>	The method to find corners of convex hull. 1: minimum sum of margin-of-errors; 2: minimum sum of reconstruction errors. The default is 2.
<code>dim.rdc</code>	Reduced data dimension; should be not less than maximum candidate <code>K</code> .
<code>thres.low</code>	The lower bound of percentage of genes to keep for CAM with ranked norm. The value should be between 0 and 1. The default is 0.05.
<code>thres.high</code>	The higher bound of percentage of genes to keep for CAM with ranked norm. The value should be between 0 and 1. The default is 0.95.
<code>cluster.method</code>	The method to do clustering. The default "K-Means" will use kmeans . The alternative "apcluster" will use apclusterK-methods .
<code>cluster.num</code>	The number of clusters; should be much larger than <code>K</code> . The default is 50.
<code>MG.num.thres</code>	The clusters with the gene number smaller than <code>MG.num.thres</code> will be treated as outliers. The default is 20.
<code>lof.thres</code>	Remove local outlier using lofactor . <code>MG.num.thres</code> is used as the number of neighbors in the calculation of the local outlier factors. The default value 0.02 will remove top 2% local outliers. Zero value will disable lof.
<code>cores</code>	The number of system cores for parallel computing. If not provided, one core for each element in <code>K</code> will be invoked. Zero value will disable parallel computing.

Details

This function includes three necessary steps to decompose a matrix of mixture expression profiles: data preprocessing, marker gene cluster search, and matrix decomposition. They are implemented in [CAMPrep](#), [CAMMGCluster](#) and [CAMASest](#), separately. More details can be found in the help document of each function.

For this function, you need to specify the range of possible subpopulation numbers and the percentage of low/high-expressed genes to be removed. Typically, 30% ~ 50% low-expressed genes can be removed from gene expression data. The removal of high-expressed genes has much less impact on results, and usually set to be 0% ~ 10%.

This function can also analyze other molecular expression data, such as proteomics data. Much less low-expressed proteins need to be removed, e.g. 0% ~ 10%, due to a limited number of proteins without missing values.

Value

An object of class "CAMObj" containing the following components:

PrepResult	An object of class "CAMPrepObj" containing data preprocessing results from CAMPrep function.
MGResult	A list of "CAMMGObj" objects containing marker gene detection results from CAMMGCluster function for each K value.
ASestResult	A list of "CAMASObj" objects containing estimated proportions, subpopulation-specific expressions and mdl values from CAMASest function for each K value.

Examples

```
#obtain data
data(ratMix3)
data <- ratMix3$X

#set seed to generate reproducible results
set.seed(111)

#CAM with known subpopulation number
rCAM <- CAM(data, K = 3, dim.rdc = 3, thres.low = 0.30, thres.high = 0.95)
#A larger dim.rdc can improve performance but increase time complexity

#CAM with a range of subpopulation number
rCAM <- CAM(data, K = 2:5, dim.rdc = 10, thres.low = 0.30, thres.high = 0.95)

#Use "apcluster" to aggregate gene vectors in CAM
rCAM <- CAM(data, K = 2:5, dim.rdc = 10, thres.low = 0.30, thres.high = 0.95,
            cluster.method = 'apcluster')
```

CAMASest

A and S matrix estimation by CAM

Description

This function estimates A and S matrix based on marker gene clusters detected by CAM.

Usage

```
CAMASest(MGResult, PrepResult, data, corner.strategy = 2)
```

Arguments

<code>MGResult</code>	An object of class " <code>CAMMGObj</code> " obtained from <code>CAMMGCluster</code> function.
<code>PrepResult</code>	An object of class " <code>CAMPRepObj</code> " obtained from <code>CAMPRep</code> function.
<code>data</code>	Matrix of mixture expression profiles which need to be the same as the input of <code>CAMPRep</code> . Data frame, <code>SummarizedExperiment</code> or <code>ExpressionSet</code> object will be internally coerced into a matrix. Each row is a gene and each column is a sample. Data should be in non-log linear space with non-negative numerical values (i.e. ≥ 0). Missing values are not supported. All-zero rows will be removed internally.
<code>corner.strategy</code>	The method to detect corner clusters. 1: minimum sum of margin-of-errors; 2: minimum sum of reconstruction errors. The default is 2.

Details

This function is used internally by `CAM` function to estimate proportion matrix (A), subpopulation-specific expression matrix (S) and mdl values. It can also be used when you want to perform CAM step by step.

The mdl values are calculated in three approaches: (1) based on data and A matrix in dimension-reduced space; (2) based on original data with A matrix estimated by transforming dimension-reduced A matrix back to original space; (3) based on original data with A directly estimated in original space. A and S matrix in original space estimated from the latter two approaches are returned. mdl is the sum of two terms: code length of data under the model and code length of model. Both mdl value and the first term (code length of data) will be returned.

Value

An object of class "`CAMASObj`" containing the following components:

<code>Aest</code>	Estimated proportion matrix from Approach 2.
<code>Sest</code>	Estimated subpopulation-specific expression matrix from Approach 2.
<code>Aest.proj</code>	Estimated proportion matrix from Approach 2, before removing scale ambiguity.
<code>Ascale</code>	The estimated scales to remove scale ambiguity of each column vector in <code>Aest</code> . Sum-to-one constraint on each row of <code>Aest</code> is used for scale estimation.
<code>Aest0</code>	Estimated proportion matrix from Approach 3.
<code>Sest0</code>	Estimated subpopulation-specific expression matrix from Approach 3.
<code>Aest0.proj</code>	Estimated proportion matrix from Approach 3, before removing scale ambiguity.
<code>Ascale0</code>	The estimated scales to remove scale ambiguity of each column vector in <code>Aest0</code> . Sum-to-one constraint on each row of <code>Aest0</code> is used for scale estimation.
<code>datalength</code>	Three values for code length of data. The first is calculated based on dimension-reduced data. The second and third are based on the original data.
<code>mdl</code>	Three mdl values. The first is calculated based on dimension-reduced data. The second and third are based on the original data.

Examples

```

#obtain data
data(ratMix3)
data <- ratMix3$X

#preprocess data
rPrep <- CAMPrep(data, dim.rdc = 3, thres.low = 0.30, thres.high = 0.95)

#Marker gene cluster detection with a fixed K
rMGC <- CAMMGCluster(3, rPrep)

#A and S matrix estimation
rASest <- CAMASest(rMGC, rPrep, data)

```

CAMASObj-class

*Class "CAMASObj"***Description**

An S4 class for storing estimated proportions, subpopulation-specific expressions and mdl values. The mdl values are calculated in three approaches: (1) based on data and A matrix in dimension-reduced space; (2) based on original data with A matrix estimated by transforming dimension-reduced A matrix back to original space; (3) based on original data with A directly estimated in original space. A and S matrix in original space estimated from the latter two approaches are returned. mdl is the sum of two terms: code length of data under the model and code length of model. Both mdl value and the first term (code length of data) will be returned.

Slots

Aest Estimated proportion matrix from Approach 2.

Sest Estimated subpopulation-specific expression matrix from Approach 2.

Aest.proj Estimated proportion matrix from Approach 2, before removing scale ambiguity.

Ascale The estimated scales to remove scale ambiguity of each column vector in Aest. Sum-to-one constraint on each row of Aest is used for scale estimation.

Aest0 Estimated proportion matrix from Approach 3.

Sest0 Estimated subpopulation-specific expression matrix from Approach 3.

Aest0.proj Estimated proportion matrix from Approach 3, before removing scale ambiguity.

Ascale0 The estimated scales to remove scale ambiguity of each column vector in AestO. Sum-to-one constraint on each row of AestO is used for scale estimation.

datalength Three values for code length of data. The first is calculated based on dimension-reduced data. The second and third are based on the original data.

mdl Three mdl values. The first is calculated based on dimension-reduced data. The second and third are based on the original data.

CAMMGCluster	<i>MG cluster detection for CAM</i>
--------------	-------------------------------------

Description

This function finds corner clusters as MG clusters (clusters containing marker genes).

Usage

```
CAMMGCluster(K, PrepResult, nComb = 200)
```

Arguments

K	The candidate subpopulation number.
PrepResult	An object of class "CAMPrepObj" obtained from CAMPrep function.
nComb	The number of possible combinations of clusters as corner clusters. Within these possible combinations ranked by margin errors, we can further select the best one based on reconstruction errors. The default is 200.

Details

This function is used internally by [CAM](#) function to detect clusters containing marker genes, or used when you want to perform CAM step by step.

This function provides two solutions. The first is the combination of clusters yielding the minimum sum of margin-of-errors for cluster centers. In the second, nComb possible combinations are selected by ranking sum of margin-of-errors for cluster centers. Then the best one is selected based on reconstruction errors of all data points in original space.

Value

An object of class "CAMMGObj" containing the following components:

idx	Two numbers which are two solutions' ranks by sum of margin-of-error.
corner	The indexes of clusters as detected corners. Each row is a solution.
error	Two rows. The first row is sum of margin-of-errors for nComb possible combinations. The second row is reconstruction errors for nComb possible combinations.

Examples

```
#obtain data
data(ratMix3)
data <- ratMix3$X

#preprocess data
rPrep <- CAMPrep(data, dim.rdc = 3, thres.low = 0.30, thres.high = 0.95)

#Marker gene cluster detection with a fixed K = 3
rMGC <- CAMMGCluster(3, rPrep)
```

CAMMGObj-class	<i>Class "CAMMGObj"</i>
----------------	-------------------------

Description

An S4 class for storing marker gene detection results.

Slots

`idx` Two numbers which are two solutions' ranks by sum of margin-of-error.

`corner` The indexes of clusters as detected corners. Each row is a solution.

`error` Two rows. The first row is sum of margin-of-errors for nComb possible combinations. The second row is reconstruction errors for nComb possible combinations.

CAMObj-class	<i>Class "CAMObj"</i>
--------------	-----------------------

Description

An S4 class for storing results of CAM.

Slots

`PrepResult` An object of class "[CAMPrepObj](#)" storing data preprocessing results from [CAMPrep](#) function.

`MGResult` A list of "[CAMMGObj](#)" objects storing marker gene detection results from [CAMMGCluster](#) function for each candidate subpopulation number.

`ASestResult` A list of "[CAMASObj](#)" objects storing estimated proportions, subpopulation-specific expressions and mdl values from [CAMASest](#) function for each candidate subpopulation number.

CAMPrep	<i>Data preprocessing for CAM</i>
---------	-----------------------------------

Description

This function perform preprocessing for CAM, including norm-based filtering, dimension deduction, perspective projection, local outlier removal and aggregation of gene expression vectors by clustering.

Usage

```
CAMPrep(data, dim.rdc = 10, thres.low = 0.05, thres.high = 0.95,
  cluster.method = c("K-Means", "apcluster"), cluster.num = 50,
  MG.num.thres = 20, lof.thres = 0)
```

Arguments

<code>data</code>	Matrix of mixture expression profiles. Data frame, SummarizedExperiment or ExpressionSet object will be internally coerced into a matrix. Each row is a gene and each column is a sample. Data should be in non-log linear space with non-negative numerical values (i.e. ≥ 0). Missing values are not supported. All-zero rows will be removed internally.
<code>dim.rdc</code>	Reduced data dimension; should be not less than maximum candidate K.
<code>thres.low</code>	The lower bound of percentage of genes to keep for CAM with ranked norm. The value should be between 0 and 1. The default is 0.05.
<code>thres.high</code>	The higher bound of percentage of genes to keep for CAM with ranked norm. The value should be between 0 and 1. The default is 0.95.
<code>cluster.method</code>	The method to do clustering. The default "K-Means" will use <code>kmeans</code> function. The alternative "apcluster" will use <code>apclusterK-methods</code> .
<code>cluster.num</code>	The number of clusters; should be much larger than K. The default is 50.
<code>MG.num.thres</code>	The clusters with the gene number smaller than <code>MG.num.thres</code> will be treated as outliers. The default is 20.
<code>lof.thres</code>	Remove local outlier using <code>lofactor</code> function. <code>MG.num.thres</code> is used as the number of neighbors in the calculation of the local outlier factors. The default value 0.02 will remove top 2% local outliers. Zero value will disable lof.

Details

This function is used internally by `CAM` function to preprocess data, or used when you want to perform CAM step by step.

Low/high-expressed genes are filtered by their L2-norm ranks. Dimension reduction is slightly different from PCA. The first loading vector is forced to be $c(1,1,\dots,1)$ with unit norm normalization. The remaining are eigenvectors from PCA in the space orthogonal to the first vector. Perspective projection is to project dimension-reduced gene expression vectors to the hyperplane orthogonal to $c(1,0,\dots,0)$, i.e., the first axis in the new coordinate system. local outlier removal is optional to exclude outliers in simplex formed after perspective projection. Finally, gene expression vectors are aggregated by clustering to further reduce the impact of noise/outlier and help improve the efficiency of simplex corner detection.

Value

An object of class "`CAMPrepObj`" containing the following components:

<code>Valid</code>	logical vector to indicate the genes left after filtering.
<code>Xprep</code>	Preprocessed data matrix.
<code>Xproj</code>	Preprocessed data matrix after perspective projection.
<code>W</code>	The matrix whose rows are loading vectors.
<code>cluster</code>	cluster results including two vectors. The first indicates the cluster to which each gene is allocated. The second is the number of genes in each cluster.
<code>c.outlier</code>	The clusters with the gene number smaller than <code>MG.num.thres</code> .
<code>centers</code>	The centers of candidate corner clusters (candidate clusters containing marker genes).

Examples

```
#obtain data
data(ratMix3)
data <- ratMix3$X

#set seed to generate reproducible results
set.seed(111)

#preprocess data
rPrep <- CAMPrep(data, dim.rdc = 3, thres.low = 0.30, thres.high = 0.95)
```

CAMPrepObj-class	<i>Class "CAMPrepObj"</i>
------------------	---------------------------

Description

An S4 class for storing data preprocessing results.

Slots

Valid logical vector to indicate the genes left after filtering.
 Xprep Preprocessed data matrix.
 Xproj Preprocessed data matrix after perspective projection.
 W The matrix whose rows are loading vectors.
 cluster cluster results including two vectors. The first indicates the cluster to which each gene is allocated. The second is the number of genes in each cluster.
 c.outlier The clusters with the gene number smaller than MG.num.thres.
 centers The centers of candidate corner clusters (candidate clusters containing marker genes).

cornerSort	<i>Candidate combinations as corners</i>
------------	--

Description

Given a set of data points, return possible combinations of data points as corners. These combinations are selected by ranking the sum of margin-of-errors.

Usage

```
cornerSort(X, K, nComb)
```

Arguments

X	A matrix of data. Each column is a data point.
K	The number of corner points.
nComb	The number of returned combinations of data points as corners. All combinations will be returned if the number of all combinations is less than nComb.

Details

This function is to detect K corner points from M data points by conducting an exhaustive combinatorial search (with total C_M^K combinations), based on a convex-hull-to-data fitting criterion: sum of margin-of-errors. `nComb` combinations are returned for further selection based on reconstruction errors of all data points in original space.

The function is implemented in Java with R-to-Java interface provided by `rJava` package. It relies on `NonNegativeLeastSquares` class in `Parallel Java Library` (<https://www.cs.rit.edu/~ark/pj.shtml>).

Value

A list containing the following components:

<code>idx</code>	A matrix to show the indexes of data points in combinations to construct a convex hull. Each column is one combination.
<code>error</code>	A vector of margin-of-error sums for each combination.

Examples

```
data <- matrix(c(0.1,0.2,1.0,0.0,0.0,0.5,0.3,
               0.1,0.7,0.0,1.0,0.0,0.5,0.3,
               0.8,0.1,0.0,0.0,1.0,0.0,0.4), nrow =3, byrow = TRUE)
topconv <- cornerSort(data, 3, 10)
```

 MDL

Minimum Description Length

Description

This function obtains minimum description length (mdl) values for each candidate subpopulation number.

Usage

```
MDL(CAMResult, mdl.method = 3)

## S4 method for signature 'MDLObj,missing'
plot(x, data.term = FALSE, ...)
```

Arguments

<code>CAMResult</code>	Result from <code>CAM</code> function.
<code>mdl.method</code>	Approach to calculate mdl values; should be 1, 2, or 3. The default is 3.
<code>x</code>	An object of class " <code>MDLObj</code> " from <code>MDL</code> .
<code>data.term</code>	If true, plot data term (code length of data under model).
<code>...</code>	All other arguments are passed to the plotting command.

Details

This function extracts minimum description length (mdl) values from the result of `CAM` function, which contains mdl values from three approaches for each candidate subpopulation number. For more details about three approaches, refer to [CAMASest](#).

mdl is code length of data under the model plus code length of model. Both mdl value and the first term about data are returned.

Value

An object of class "`MDLObj`" containing the following components:

<code>K</code>	The candidate subpopulation numbers.
<code>datalengths</code>	For each model with a certain subpopulation number, code length of data under the model.
<code>mdls</code>	mdl value for each model with a certain subpopulation number.

Examples

```
#obtain data
data(ratMix3)
data <- ratMix3$X

#Analysis by CAM
rCAM <- CAM(data, K = 2:5, thres.low = 0.30, thres.high = 0.95)

#extract mdl values
MDL(rCAM)
MDL(rCAM, 1)
MDL(rCAM, 2)

#plot MDL curves
plot(MDL(rCAM))
plot(MDL(rCAM), data.term = TRUE) #with data length curve
```

MDLObj-class

Class "MDLObj"

Description

An S4 class for storing mdl values.

Slots

<code>K</code>	The candidate subpopulation numbers.
<code>datalengths</code>	For each model with a certain subpopulation number, code length of data under the model.
<code>mdls</code>	mdl value for each model with a certain subpopulation number.

MGsforA

Marker genes detected by CAM for estimating A

Description

This function returns marker genes detected by CAM for estimating A.

Usage

```
MGsforA(CAMResult = NULL, K = NULL, PrepResult = NULL, MGResult = NULL,
        corner.strategy = 2)
```

Arguments

CAMResult	Result from CAM .
K	The candidate subpopulation number.
PrepResult	An object of class "CAMPrepObj" from CAMPrep .
MGResult	An object of class "CAMMGObj" from CAMMGCluster .
corner.strategy	The method to detect corner clusters. 1: minimum sum of margin-of-errors; 2: minimum sum of reconstruction errors. The default is 2.

Details

This function needs to specify CAMResult and K, or PrepResult and MGResult. The returned marker genes are those used by CAM for estimating A. To obtain a more complete marker gene list, please refer to [MGstatistic](#).

Value

A list of vectors, each of which contains marker genes for one subpopulation.

Examples

```
#obtain data and run CAM
data(ratMix3)
data <- ratMix3$X
rCAM <- CAM(data, K = 3, dim.rdc= 3, thres.low = 0.30, thres.high = 0.95)
#obtain marker genes detected by CAM for estimating A
MGlist <- MGsforA(rCAM, K = 3)

#obtain data and run CAM step by step
rPrep <- CAMPrep(data, dim.rdc= 3, thres.low = 0.30, thres.high = 0.95)
rMGC <- CAMMGCluster(3, rPrep)
#obtain marker genes detected by CAM for estimating A
MGlist <- MGsforA(PrepResult = rPrep, MGResult = rMGC)
```

Description

This function computes One-Versus-Everyone Fold Change (OVE-FC) from subpopulation-specific expression profiles. Bootstrapping is optional.

Usage

```
MGstatistic(data, A = NULL, boot.alpha = NULL, nboot = 1000,
            cores = NULL)
```

Arguments

<code>data</code>	A data set that will be internally coerced into a matrix. Each row is a gene and each column is a sample. Data should be in non-log linear space with non-negative numerical values (i.e. ≥ 0). Missing values are not supported. All-zero rows will be removed internally.
<code>A</code>	When data are mixture expression profiles, A is estimated proportion matrix or prior proportion matrix. When data are pure expression profiles, A is a phenotype vector to indicate which subpopulation each sample belongs to.
<code>boot.alpha</code>	Alpha for bootstrapped OVE-FC confidence interval. The default is 0.05.
<code>nboot</code>	The number of boots.
<code>cores</code>	The number of system cores for parallel computing. If not provided, the default back-end is used.

Details

This function calculates OVE-FC and bootstrapped OVE-FC which can be used to identify markers from all genes.

Value

A data frame containing the following components:

<code>idx</code>	Numbers or phenotypes indicating which subpopulation each gene could be a marker for. If A is a proportion matrix without column name, numbers are returned. Otherwise, phenotypes.
<code>OVE.FC</code>	One-versus-Everyone fold change (OVE-FC)
<code>OVE.FC.alpha</code>	lower confidence bound of bootstrapped OVE-FC at alpha level.

Examples

```
#data are mixture expression profiles, A is proportion matrix
data(ratMix3)
MGstat <- MGstatistic(ratMix3$X, ratMix3$A)

MGstat <- MGstatistic(ratMix3$X, ratMix3$A, boot.alpha = 0.05) #enable boot
```



```

#data are pure expression profiles without replicates
MGstat <- MGstatistic(ratMix3$S) #boot is not applicable

#data are pure expression profiles with phenotypes
S <- matrix(rgamma(3000,0.1,0.1), 1000, 3)
S <- S[, c(1,1,1,2,2,2,3,3,3,3)] + rnorm(1000*10, 0, 0.5)
MGstat <- MGstatistic(S, c(1,1,1,2,2,2,3,3,3,3), boot.alpha = 0.05)

```

PCAmat

Dimension-reduction loading matrix accessor

Description

Accessor to Dimension-reduction loading matrix.

Usage

```

PCAmat(x, ...)

## S4 method for signature 'CAMObj'
PCAmat(x)

## S4 method for signature 'CAMPRepObj'
PCAmat(x)

```

Arguments

x a [CAMObj](#) object or a [CAMPRepObj](#) object
... additional argument list.

Value

The matrix whose rows are loading vectors for dimension reduction.

Examples

```

#obtain data
data(ratMix3)
data <- ratMix3$X

rCAM <- CAM(data, K = 3, dim.rdc = 3, thres.low = 0.30, thres.high = 0.95)
W <- PCAmat(rCAM)
W <- PCAmat(slot(rCAM, "PrepResult"))

```

 ratMix3

Gene expression data downsampled from GSE19380

Description

Rat brain, liver and lung biospecimens derived from one animal at the cRNA homogenate level in different proportions. 3 technical replicates each. We downsample the original data to 10000 probes/probesets and 7 mixtures. Proportions used in experiments and pure expression profiles are also included.

Usage

```
data(ratMix3)
```

Format

A list with three matrices: mixture profiles (X), mixing proportions (A) and pure profiles (S).

References

Shen-Orr et al. (2010) Nat Methods 2010 Apr;7(4):287-9. PMID: 20208531

 simplexplot

The plot of scatter simplex

Description

This function shows scatter simplex of mixture expressions.

Usage

```
simplexplot(data, A, MGlister = NULL, corner.order = NULL,
  data.col = "gray", corner.col = "red", ...)
```

Arguments

data	A data set that will be internally coerced into a matrix. Each row is a gene and each column is a sample. Data should be in non-log linear space with non-negative numerical values (i.e. ≥ 0). Missing values are not supported. All-zero rows will be removed internally.
A	Prior/Estimated proportion matrix.
MGlister	A list of vectors, each of which contains known markers and/or CAM-detected markers for one subpopulation.
corner.order	The order to show simplex corners counterclockwise.
data.col	The color for data points. The default is "gray".
corner.col	The color for corner points. The default is "red".
...	All other arguments are passed to the plotting command.

Details

This function can show the scatter simplex and detected marker genes in a 2D plot. The corners in the high-dimensional simplex will still locate at extreme points of low-dimensional simplex. These corners will follow the order set by `corner.order` to display in the plot counterclockwise.

Value

A plot to the current device.

Examples

```
#obtain data, A matrix, marker genes
data(ratMix3)
data <- ratMix3$X
A <- ratMix3$A
pMGstat <- MGstastic(ratMix3$S, c("Liver","Brain","Lung"))
pMGlist.FC <- lapply(c("Liver","Brain","Lung"), function(x)
  rownames(pMGstat)[pMGstat$idx == x & pMGstat$OVE.FC > 10])

#plot simplex for data
simplexplot(data, A)
simplexplot(data, A, MGlist = pMGlist.FC) #Color marker genes in simplex plot

#set differnt corner order and colors
simplexplot(data, A, MGlist = pMGlist.FC, corner.order = c(2,1,3),
  data.col = "blue", corner.col = c("red","orange","green"))
```

XWProj

Perspective projection to obtain simplex

Description

This function reduces data dimension by loading matrix and then project dimension-reduced data to the hyperplane orthogonal to $c(1,0,\dots,0)$, i.e., the first axis in the new coordinate system..

Usage

```
XWProj(data, W)
```

Arguments

<code>data</code>	A data set that will be internally coerced into a matrix. Each row is a gene and each column is a sample. Missing values are not supported. All-zero rows will be removed internally.
<code>W</code>	The matrix whose rows are loading vectors; should be obtained from CAM/CAMPRep function with accessor method PCAmat .

Details

This function can project gene expression vectors to simplex plot generated by [CAM/CAMPRep](#). Using slot `Xproj` in "[CAMPRepObj](#)" can only show the simplex of genes after filtering. This function helps observe all genes in simplex plot.

Value

The data after perspective projection.

Examples

```
#obtain data
data(ratMix3)
data <- ratMix3$X

#preprocess data
rPrep <- CAMPrep(data, dim.rdc = 3, thres.low = 0.50, thres.high = 0.90)

#obtain simplex
Xproj <- XWProj(data, PCAmat(rPrep))
#plot simplex in 3d space
#plot3d(Xproj[,-1]) #The first dimension is constant after projection
```

Index

AfromMarkers, [2, 3](#)
Amat (AS-accessor), [4](#)
Amat, CAMASObj-method (AS-accessor), [4](#)
Amat, CAMObj-method (AS-accessor), [4](#)
AS-accessor, [4](#)

CAM, [2, 5, 7, 9, 11, 13–15, 19](#)
CAMASest, [2, 5, 6, 6, 10, 14](#)
CAMASObj, [4, 6, 7, 10](#)
CAMASObj (CAMASObj-class), [8](#)
CAMASObj-class, [8](#)
CAMMGCluster, [2, 5–7, 9, 10, 15](#)
CAMMGObj, [6, 7, 9, 10](#)
CAMMGObj (CAMMGObj-class), [10](#)
CAMMGObj-class, [10](#)
CAMObj, [4, 6, 17](#)
CAMObj (CAMObj-class), [10](#)
CAMObj-class, [10](#)
CAMPRep, [2, 5–7, 9, 10, 10, 15, 19](#)
CAMPRepObj, [6, 7, 9–11, 17, 19](#)
CAMPRepObj (CAMPRepObj-class), [12](#)
CAMPRepObj-class, [12](#)
CAMTHC (CAMTHC-package), [2](#)
CAMTHC-package, [2](#)
cornerSort, [12](#)

kmeans, [5, 11](#)

l1median, [3](#)
lofactor, [5, 11](#)

MDL, [2, 13, 13](#)
MDLObj, [13, 14](#)
MDLObj (MDLObj-class), [14](#)
MDLObj-class, [14](#)
MGsforA, [15](#)
MGstatistic, [2, 15, 16](#)

PCAmat, [17, 19](#)
PCAmat, CAMObj-method (PCAmat), [17](#)
PCAmat, CAMPRepObj-method (PCAmat), [17](#)
plot, MDLObj, missing-method (MDL), [13](#)

ratMix3, [18](#)

simplexplot, [18](#)
Smat (AS-accessor), [4](#)
Smat, CAMASObj-method (AS-accessor), [4](#)
Smat, CAMObj-method (AS-accessor), [4](#)

XWProj, [19](#)