

Package ‘PMixClus’

October 23, 2014

Type Package
Title Penalized MODOEL-Based Clustering for RNA-Seq count data
Version 1.0
Date 2014-10-21
Author Ye Tian
Maintainer Ye Tian <tian0049@e.ntu.edu.sg>
Description Penalized NB or Poisson mixture model with application to cluster RNA-Seq count data
Suggests MBCluster.Seq
License GPL- (>=3)
LazyData true

R topics documented:

Cutree	1
HH.Tree	2
plotHH.tree	3
PM.pretreat	4
PMixClus	4
seq_count	6
Index	7

Cutree	<i>Cut the Hybrid-Hierarchical tree.</i>
--------	--

Description

Cut the HH tree at specific level to get the relative cluster labels.

Usage

```
Cutree(tree, group, level)
```

Arguments

tree	The HH tree structure. It should be the value from HH.Tree .
group	The cluster labels for bottom children. It should be the same with arguments group in HH.Tree .
level	The targetted level. It cannot be larger than the height of the tree.

Value

The cluster labels at the targetted level of HH tree.

Examples

```
# data("seq_count")
# seq_data = PM.pretreat(count = seq_count)
# out.NB = PMixClus(x, 1:3, exp(seq(0, 3, length = 5)))
# tree = HH.Tree(seq_data, out.NB$gt.hat[[3]], out.NB$phi.hat[3, ], out.NB$group[3, ])
# Cutree(tree, out.NB$group[3, ], 2)
```

HH.Tree	<i>Construct the Hybrid-Hierarchical tree.</i>
---------	--

Description

Construct the HH tree with the bottom children from penalized model-based method.

Usage

```
HH.Tree(data, gt, phi, group, model = "NB", disp.domain = c(1e-06, 4))
```

Arguments

data	The well prepared data from PM.pretreat
gt	The estimated parameters from PMixClus , e.g. <code>gt[[1]]</code> is the parameters from PMixClus when <code>k = k.init[1]</code> .
phi	The estimated dispersion parameters.
group	The bottom group labels for HH.
model	Select the model between negative binomial and Poisson. If choose Poisson, <code>model = "Pois"</code> .
disp.domain	Same with <code>disp.domail</code> in PMixClus .

Value

A $k \times 3$ matrix is returned, where k is the number of clusters in the bottom. The first two columns in each row show the two clusters merged in that step and the third column shows the distance of these merged two clusters.

References

Si, Y. and Liu, P. and Li, P. and Thomas, P. B. (2014) Model-based clustering for RNA-seq data. *Bioinformatics*, 30, 197-205

Examples

```
# data("seq_count")
# seq_data = PM.pretreat(count = seq_count)
# out.NB = PMixClus(x, 1:3, exp(seq(0, 3, length = 5)))
# op = which.min(out.NB$BIC)      #Find the minimal BIC or EBIC.
# tree = HH.Tree(seq_data, out.NB$gt.hat[[op]], out.NB$phi.hat[op, ], out.NB$group[op, ])
```

plotHH.tree

Plot the Hybrid-Hierarchical tree.

Description

Plot the Hybrid-Hierarchical tree to visualize the clustering structure. The bottom children are samples that are clustered by penalized model-based method. The samples with same color are in the the same cluster in reality if the true clusters are know in advance. The sample names are shown under the bottom bars.

Usage

```
plotHH.tree(sample.names, tree, group, group.true = NULL, cex.leaf = 0.8,
  tree.title = "Hybrid-Hierarchical Tree")
```

Arguments

sample.names	The sample names. It should be consistent with sample names outputted by PM.pretreat .
tree	The HH tree stucture. It should be the value from HH.Tree .
group	The group labels for bottom children. It should be the same with arguments group in HH.Tree .
group.true	The true group labels if they are known in advance.
cex.leaf	The font size of sample names.
tree.title	The title for this figure.

References

Si, Y. and Liu, P. and Li, P. and Thomas, P. B. (2014) Model-based clustering for RNA-seq data. *Bioinformatics*, 30, 197-205

Examples

```
# data("seq_count")
# seq_data = PM.pretreat(count = seq_count)
# out.NB = PMixClus(x, 1:3, exp(seq(0, 3, length = 5)))
# op = which.min(out.NB$BIC)      #Find the minimal BIC or EBIC.
# tree = HH.Tree(seq_data, out.NB$gt.hat[[op]], out.NB$phi.hat[op, ], out.NB$group[op, ])
# group.true = c(rep(1, 5), rep(2, 5))
# plotHH.tree(rownames(data$count), tree, out.NB$group[op, ], group.true = group.true)
```

PM.pretreat	<i>Filter the data and compute the size factor for each sample by median ratio.</i>
-------------	---

Description

This function can filter out genes that have very low expression and estimate the size factor for each sample by median ratio method. The output data is used in most of other functions in this package.

Usage

```
PM.pretreat(count, csum = 5, sizefactor = NULL)
```

Arguments

count	The $N \times P$ RNA-Seq count matrix with N samples and P genes.
csum	The threshold value to filter out the low expressed genes. The function will remove the genes, the sum of whose read counts over all samples are less than csum.
sizefactor	The size factor for each sample. When sizefactor = NULL, the function will estimate the size factors by median ratio method.

Value

count	The filtered data with row names and column names.
sizefactor	The size factors estimated by this function or inputted by users.

References

Anders, S. and Huber, W. (2010). Differential expression analysis for sequence count data. *Genome Biology*, 11, R106

Examples

```
# data("seq_count")
# PM.pretreat(count = seq_count)
```

PMixClus	<i>Penalized Poisson or NB model based clustering for count data.</i>
----------	---

Description

Given the prepared data from [PM.pretreat](#), this function applied penalized Poisson or NB mixture model to do the clustering. The model selection criterions modified BIC and modified EBIC are provided.

Usage

```
PMixClus(data, k.init, lambda, model = "NB", s.update = FALSE,
  MAX_iter = 100, threshold = 1e-07, disp.domain = c(1e-06, 4),
  is.BIC = "BIC")
```

Arguments

<code>data</code>	The well prepared data from PM.pretreat .
<code>k.init</code>	The vector of the number of clusters that the user wants to select from, e.g. <code>k.init = 1:3</code> .
<code>lambda</code>	The vector of tuning parameter that the user wants to select from, e.g. <code>lambda = exp(seq(0, 3, 1e5))</code> .
<code>model</code>	The distribution used inside the function. By default, the function applies the NB distribution. If <code>model = "Pois"</code> , then the Poisson distribution will be applied.
<code>s.update</code>	Whether to update size factors inside the EM for Poisson model.
<code>MAX_iter</code>	The maximal number of iteration for EM.
<code>threshold</code>	The stop threshold for EM algorithm.
<code>disp.domain</code>	The domain in which the function searches for MLE of dispersion parameters. The default value is suitable for most cases in application.
<code>is.BIC</code>	The model selection criterion applied in the function. By default, <code>is.BIC = "BIC"</code> , modified BIC is applied; if <code>is.BIC = "EBIC"</code> , the modified EBIC is applied.

Value

<code>BIC</code>	The value of modified BIC or modified EBIC corresponding to <code>k.init</code> .
<code>lambda.sel</code>	The selected tuning parameters by modified BIC or modified EBIC.
<code>theta.cdifff</code>	The matrix to store fluctuation of genes over different clusters. The rows of this matrix correspond to <code>k.init</code> vector respectively and the columns represent genes. In each row the elements are ordered from higher values to lower values.
<code>gnames.st</code>	The sorted gene names according to <code>theta.cdifff</code> . This matrix store the gene names for every element in the <code>theta.cdifff</code> matrix.
<code>group</code>	The group labels for samples. Each row includes the group labels when input corresponding <code>k</code> in <code>k.init</code> .
<code>sizefactor</code>	The size factors. The rows of this matrix represent results from inputting the corresponding elements in <code>k.init</code> . By default, this matrix is fixed on <code>data\$sizefactor</code> from PM.pretreat ; if <code>model = "Pois"</code> and <code>s.update = T</code> , the size factors are estimated in the EM.
<code>gt.hat</code>	The list contains estimated mean/size factor. This is the necessary argument in HH.Tree .
<code>phi.hat</code>	The estimates of dispersion parameters if <code>model = "NB"</code> . This is the necessary argument in HH.Tree .

References

Pan, W. and Shen, X. (2007). Penalized model-based clustering with application to variable selection. *Journal of Machine Learning Research*, 8, 1145-1164.

Examples

```
# data("seq_count")
# seq_data = PM.pretreat(count = seq_count)
# out.NB = PMixClus(seq_data, 1:3, exp(seq(0, 3, length = 5)))
# out.pois = PMixClus(seq_data, 1:3, exp(seq(0, 3, length = 5)), model = "Pois", s.update=TRUE)
# op = which.min(out.NB$BIC) #Find the minimal BIC or EBIC.
# k.init[op] #The selected number of clusters by BIC or EBIC.
# sum(out.NB$theta.cdifff[op, ] == 0) #The number of genes excluded when \code{k = k.init[op]}.
# out.NB$group[op, ] #The group labels when \code{k = k.init[op]}.
```

seq_count

Simulated RNA-Seq count data

Description

The data set consist of 1000 genes with 5 biological replicates in each of 2 treatments.

Usage

seq_count

Format

This data is 10×1000 matrix, in which first 5 samples in the same treatment and the other 5 samples in the other treatment.

Index

*Topic **datasets**

seq_count, [6](#)

Cutree, [1](#)

HH.Tree, [2](#), [2](#), [3](#), [5](#)

plotHH.tree, [3](#)

PM.pretreat, [2-4](#), [4](#), [5](#)

PMixClus, [2](#), [4](#)

seq_count, [6](#)