

Package ‘oncomix’

April 12, 2018

Title Identifying Genes Overexpressed in Subsets of Tumors from
Tumor-Normal mRNA Expression Data

Version 1.0.0

Author Daniel Pique, John Grealley, Jessica Mar

Maintainer Daniel Pique <daniel.pique@med.einstein.yu.edu>

Description This package helps identify mRNAs that are overexpressed in subsets of tumors relative to normal tissue. Ideal inputs would be paired tumor-normal data from the same tissue from many patients (>15 pairs). This unsupervised approach relies on the observation that oncogenes are characteristically overexpressed in only a subset of tumors in the population, and may help identify oncogene candidates purely based on differences in mRNA expression between previously unknown subtypes.

Depends R (>= 3.4.0)

License GPL-3

Encoding UTF-8

LazyData true

RoxygenNote 6.0.1

Suggests knitr, rmarkdown, testthat, RMySQL

Imports ggplot2, ggrepel, RColorBrewer, mclust, stats,
SummarizedExperiment

VignetteBuilder knitr

biocViews GeneExpression, Sequencing

NeedsCompilation no

R topics documented:

exprNmIIsf	2
exprTumIIsf	2
mixModelParams	3
oncoMixBimodal	3
oncoMixIdeal	4
oncoMixTraditionalDE	5
plotGeneHist	5
queryRes	6
scatterMixPlot	6
topGeneQuants	7

Index**9**

exprNmIIsof	<i>Human Breast Cancer RNA-sequencing data from TCGA - Adj. Normal Tissue</i>
-------------	---

Description

These RNA-sequencing expression data were obtained from the Cancer Genome Atlas project (now the Genomic Data Commons (GDC)). The sequencing data were generated from adjacent normal breast tissue data from 113 patients with breast cancer. Quantification of RNA expression values was performed using standard GDC pipelines. The expression values are reported in transcripts per million reads. Out of an initial 73,599 RNA transcripts, 700 are included as part of this dataset. These 700 transcripts represent a random subset of the transcripts with at least 20 samples. Rows contain anonymized patient identifiers, while columns contain UCSC gene symbols.

Author(s)

Daniel Pique <daniel.pique@med.einstein.yu.edu>

References

<https://gdc.cancer.gov/>

exprTumIsof	<i>Human Breast Cancer RNA-sequencing data from TCGA - Tumor Tissue</i>
-------------	---

Description

These RNA-sequencing expression data were obtained from the Cancer Genome Atlas project (now the Genomic Data Commons (GDC)). The sequencing data were generated from breast carcinoma samples from 113 patients. Quantification of RNA expression values was performed using standard GDC pipelines. The expression values are reported in transcripts per million reads. Out of an initial 73,599 RNA transcripts, 700 are included as part of this dataset. These 700 transcripts represent a random subset of the transcripts with at least 20 anonymized patient identifiers, while columns contain UCSC gene symbols.

Author(s)

Daniel Pique <daniel.pique@med.einstein.yu.edu>

References

<https://gdc.cancer.gov/>

mixModelParams	<i>Generate the parameters for two 2-component Gaussian mixture models with equal variances</i>
----------------	---

Description

This function allows you to generate the parameters for two 2-component Gaussian mixture model with equal variances from 2 matrices of data with a priori labels (eg tumor vs normal.) This application was originally intended for matrices of gene expression data treated with 2 conditions.

Usage

```
mixModelParams(exprNm1, exprTum)
```

Arguments

exprNm1	A dataframe (S3 or S4), matrix, or SummarizedExperiment object containing normal data with patients as columns and genes as rows.
exprTum	A dataframe (S3 or S4), matrix, or SummarizedExperiment object containing tumor data with patients as columns and genes as rows.

Value

Returns a dataframe, each element of which contains the 12 mixture model parameters for each gene in an n x 12 matrix, where n is the number of genes.

Examples

```
exprNm1 <- as.data.frame(matrix(data=rgamma(n=150, shape=2, rate=2),
  nrow=10, ncol=15))
colnames(exprNm1) <- paste0("patientN", seq_len(ncol(exprNm1)))
rownames(exprNm1) <- paste0("gene", seq_len(nrow(exprNm1)))

exprTum <- as.data.frame(matrix(data=rgamma(n=150, shape=4, rate=3),
  nrow=10, ncol=15))
colnames(exprTum) <- paste0("patientT", seq_len(ncol(exprTum)))
rownames(exprTum) <- paste0("gene", seq_len(nrow(exprTum)))

mmParams <- mixModelParams(exprNm1, exprTum)
```

oncoMixBimodal	<i>Creating a schematic of a 2-component mixture model</i>
----------------	--

Description

This function allows you to generate a plot

Usage

```
oncoMixBimodal(means = c(3, 7))
```

Arguments

means Set the values for the difference between parameter means

Value

Returns a ggplot object that shows a 2-component Gaussian mixture model

Examples

```
oncoMixBimodal(means=c(3,7))  
oncoMixBimodal(means=c(3,10))
```

oncoMixIdeal *Creating an ideal oncomix gene*

Description

This function allows you to generate a plot

Usage

```
oncoMixIdeal(means = c(3, 7))
```

Arguments

means Set the difference between parameter means for the overexpressed (oe) group.
Defaults to c(3,7)

Value

Returns a ggplot object that shows the statistical model for an idealized/theoretical oncogene candidate mRNA that is overexpressed in a subset of tumors

Examples

```
oncoMixIdeal(means=c(3,10))  
oncoMixIdeal(means=c(2,18.5))
```

oncoMixTraditionalDE *Creating a schematic of a traditional diff. expression experiment*

Description

This function allows you to generate a schematic of the assumptions of a traditional DE experiment between two known groups.

Usage

```
oncoMixTraditionalDE(means = c(3, 7))
```

Arguments

means Set the values for the difference between parameter means

Value

Returns a ggplot object that shows the traditional method (2 sample t-test) for mRNA differential expression.

Examples

```
oncoMixTraditionalDE(means=c(3,7))
oncoMixTraditionalDE(means=c(3,10))
```

plotGeneHist *Plot a histogram of gene expression values from tumor and adjacent normal tissue.*

Description

This function allows you to plot a histogram of gene expression values from tumor and adjacent normal tissue with the option of including the best fitting Gaussian curve.

Usage

```
plotGeneHist(mmParams, exprNml, exprTum, isof)
```

Arguments

mmParams The output from the getMixModelParams function.

exprNml A dataframe (S3 or S4), matrix, or SummarizedExperiment object containing normal data with patients as columns and genes as rows.

exprTum A dataframe (S3 or S4), matrix, or SummarizedExperiment object containing tumor data with patients as columns and genes as rows.

isof The gene isoform to visualize

Value

Returns a histogram of the gene expression values from the two groups.

See Also

[mixModelParams](#)

Examples

```
exprNm1 <- as.data.frame(matrix(data=rgamma(n=150, shape=2, rate=2),
  nrow=10, ncol=15))
colnames(exprNm1) <- paste0("patientN", seq_len(ncol(exprNm1)))
rownames(exprNm1) <- paste0("gene", seq_len(nrow(exprNm1)))

exprTum <- as.data.frame(matrix(data=rgamma(n=150, shape=4, rate=3),
  nrow=10, ncol=15))
colnames(exprTum) <- paste0("patientT", seq_len(ncol(exprTum)))
rownames(exprTum) <- paste0("gene", seq_len(nrow(exprTum)))
mmParams <- mixModelParams(exprNm1, exprTum)
isof <- rownames(mmParams)[1]
plotGeneHist(mmParams, exprNm1, exprTum, isof)
```

queryRes

Oncogene Database Mapping Gene Symbol to UCSC ID (kgID)

Description

These data were downloaded in September 2017 from the url listed below and represent a mapping between gene symbols in ongene, an oncogene database curated from the scientific literature, and gene identifiers from the University of California, Santa Cruz's genomic database.

Author(s)

Min Zhao <mzhao@usc.edu.au>

References

http://ongene.bioinfo-minzhao.org/ongene_human.txt

scatterMixPlot

Generate a scatter plot with the output from mixModelParams

Description

This function allows you to generate the parameters for two 2-component mixture models with equal variances

Usage

```
scatterMixPlot(mmParams, selIndThresh = 1, geneLabels = NULL)
```

Arguments

mmParams	The output from the mixModelParams function. Will utilize the deltaMu2 and deltaMu1 rows.
selIndThresh	This is the selectivity index threshold to use. All genes with SI values above this threshold will be highlighted in purple. Specify either selIndThresh or geneLabels (not both simultaneously).
geneLabels	A character vector of gene names used to label the genes with that name on the scatter plot. Specify either selIndThresh or geneLabels (not both simultaneously).

Value

Returns a ggplot scatter object that can be plotted

See Also

[mixModelParams](#)

Examples

```
exprNm1 <- as.data.frame(matrix(data=rgamma(n=150, shape=2, rate=2),
  nrow=10, ncol=15))
colnames(exprNm1) <- paste0("patientN", seq_len(ncol(exprNm1)))
rownames(exprNm1) <- paste0("gene", seq_len(nrow(exprNm1)))

exprTum <- as.data.frame(matrix(data=rgamma(n=150, shape=4, rate=3),
  nrow=10, ncol=15))
colnames(exprTum) <- paste0("patientT", seq_len(ncol(exprTum)))
rownames(exprTum) <- paste0("gene", seq_len(nrow(exprTum)))

mmParams <- mixModelParams(exprNm1, exprTum)
scatterMixPlot(mmParams)
```

topGeneQuants

Identify genes that meet pre-specified quantiles

Description

This function allows you to subset genes that are above pre-specified quantiles and that most closely resemble the distribution of oncogenes.

Usage

```
topGeneQuants(mmParams, delMu2Thr = 90, delMu1Thr = 10, siThr = 0.99)
```

Arguments

mmParams	The output from the mixModelParams function.
delMu2Thr	The percentile threshold for the deltaMu2 statistic. All genes exceeding this percentile threshold will be selected.

deltMu1Thr	The percentile threshold for the deltaMu1 statistic. All genes exceeding this percentile threshold will be selected.
siThr	The threshold for the selectivity index statistic (between 0-1). All genes exceeding this threshold will be selected.

Value

Returns a dataframe containing all genes meeting the prespecified thresholds.

See Also

[mixModelParams](#)

Examples

```
exprNm1 <- as.data.frame(matrix(data=rgamma(n=150, shape=2, rate=2),
  nrow=10, ncol=15))
colnames(exprNm1) <- paste0("patientN", seq_len(ncol(exprNm1)))
rownames(exprNm1) <- paste0("gene", seq_len(nrow(exprNm1)))

exprTum <- as.data.frame(matrix(data=rgamma(n=150, shape=4, rate=3),
  nrow=10, ncol=15))
colnames(exprTum) <- paste0("patientT", seq_len(ncol(exprTum)))
rownames(exprTum) <- paste0("gene", seq_len(nrow(exprTum)))

mmParams <- mixModelParams(exprNm1, exprTum)
topGeneQuants(mmParams)
```


Index

- *Topic **Breast**
 - exprNmIsof, 2
 - exprTumIsof, 2
- *Topic **Database**
 - queryRes, 6
- *Topic **Gaussian,**
 - plotGeneHist, 5
- *Topic **Oncogenes,**
 - queryRes, 6
- *Topic **RNA**
 - exprNmIsof, 2
 - exprTumIsof, 2
- *Topic **Tissue**
 - exprNmIsof, 2
- *Topic **Tumor**
 - exprTumIsof, 2
- *Topic **differential**
 - oncoMixTraditionalDE, 5
- *Topic **expression,**
 - exprNmIsof, 2
 - exprTumIsof, 2
- *Topic **expression**
 - oncoMixTraditionalDE, 5
- *Topic **idealized,**
 - oncoMixBimodal, 3
 - oncoMixIdeal, 4
 - oncoMixTraditionalDE, 5
- *Topic **mixture-model,**
 - mixModelParams, 3
- *Topic **oncoMix,**
 - oncoMixBimodal, 3
 - oncoMixIdeal, 4
 - oncoMixTraditionalDE, 5
 - plotGeneHist, 5
 - scatterMixPlot, 6
- *Topic **oncomix,**
 - mixModelParams, 3
- *Topic **subsetting**
 - topGeneQuants, 7
- *Topic **theoretical,**
 - oncoMixTraditionalDE, 5
- *Topic **theoretical**
 - oncoMixBimodal, 3
 - oncoMixIdeal, 4
- *Topic **two-component**
 - mixModelParams, 3
 - plotGeneHist, 5
 - scatterMixPlot, 6
- *Topic **visualization,**
 - plotGeneHist, 5
 - scatterMixPlot, 6
- exprNmIsof, 2
- exprTumIsof, 2
- mixModelParams, 3, 6–8
- oncoMixBimodal, 3
- oncoMixIdeal, 4
- oncoMixTraditionalDE, 5
- plotGeneHist, 5
- queryRes, 6
- scatterMixPlot, 6
- topGeneQuants, 7