

Package ‘RTCGAToolbox’

April 12, 2018

Type Package

Title A new tool for exporting TCGA Firehose data

Version 2.8.0

Author Mehmet Kemal Samur

Maintainer Marcel Ramos <marcel.ramos@roswellpark.org>

Description Managing data from large scale projects such as The Cancer Genome Atlas (TCGA) for further analysis is an important and time consuming step for research projects. Several efforts, such as Firehose project, make TCGA pre-processed data publicly available via web services and data portals but it requires managing, downloading and preparing the data for following steps. We developed an open source and extensible R based data client for Firehose pre-processed data and demonstrated its use with sample case studies. Results showed that RTCGAToolbox could improve data management for researchers who are interested with TCGA data. In addition, it can be integrated with other analysis pipelines for following data analysis.

License file LICENSE

LazyData true

LazyLoad yes

Depends R (>= 3.4.0)

Imports Biobase, BiocGenerics, data.table (>= 1.9.4), GenomicRanges, GenomeInfoDb, httr, IRanges, limma (>= 3.18), methods, plyr, RaggedExperiment, RCircos, RCurl, RJSONIO, S4Vectors, SummarizedExperiment, survival, XML

Suggests BiocStyle, Homo.sapiens, knitr, MultiAssayExperiment, readr, rmarkdown

biocViews DifferentialExpression, GeneExpression, Sequencing

URL <http://mksamur.github.io/RTCGAToolbox/>

BugReports <https://github.com/mksamur/RTCGAToolbox/issues>

VignetteBuilder knitr

RoxygenNote 6.0.1

Collate 'RTCGAToolbox-Class.R' 'RTCGAToolbox.R' 'utils.R'
'biocExtract.R' 'getCNGECorrelation.R'
'getDiffExpressedGenes.R' 'getFirehoseAnalyzeDates.R'
'getFirehoseData.R' 'getFirehoseDatasets.R'

'getFirehoseRunningDates.R' 'getMutationRate.R' 'getReport.R'
 'getSurvival.R' 'selectType.R'

NeedsCompilation no

R topics documented:

biocExtract	2
CorResult-class	4
DGEResult-class	4
FirehoseCGHArray-class	4
FirehoseData-class	5
FirehoseGISTIC-class	6
FirehoseMethylationArray-class	7
FirehosemRNAArray-class	7
getCNGECorrelation	7
getData	8
getDiffExpressedGenes	9
getFirehoseAnalyzeDates	10
getFirehoseData	10
getFirehoseDatasets	12
getFirehoseRunningDates	12
getMutationRate	13
getReport	13
getSurvival	14
hg19.ucsc.gene.locations	15
RTCGASample	15
RTCGAToolbox	16
selectType	16
showResults	17
showResults,CorResult-method	18
showResults,DGEResult-method	18

Index **20**

biocExtract	<i>Extract and convert data from a FirehoseData object to a Bioconductor object</i>
-------------	---

Description

This function processes data from a [FirehoseData](#) object. Raw data is converted to a conventional Bioconductor object. The function returns either a [SummarizedExperiment](#) or a [RaggedExperiment](#) class object. In cases where there are multiple platforms in a data type, an attempt to consolidate datasets will be made based on matching dimension names. For ranged data, this functionality is provided with more control as part of the [RaggedExperiment](#) features. See [RaggedExperiment](#) for more details.

Usage

```
biocExtract(object, type = c("clinical", "RNASeqGene", "miRNASeqGene",  
  "RNASeq2GeneNorm", "CNASNP", "CNVSNP", "CNASeq", "CNACGH", "Methylation",  
  "Mutation", "mRNAArray", "miRNAArray", "RPPAArray", "GISTIC", "GISTICA",  
  "GISTICT"))
```

Arguments

object	A FirehoseData object from which to extract data.
type	The type of data to extract from the "FirehoseData" object, see type section.

Value

Either an [SummarizedExperiment](#) object or a [RaggedExperiment](#) object.

type

Choices include:

- clinical - Get the clinical data slot
- RNASeqGene - RNASeqGene
- RNASeq2GeneNorm - Normalized
- miRNASeqGene - micro RNA SeqGene
- CNASNP - Copy Number Alteration
- CNVSNP - Copy Number Variation
- CNASeq - Copy Number Alteration
- CNACGH - Copy Number Alteration
- Methylation - Methylation
- mRNAArray - Messenger RNA
- miRNAArray - micro RNA
- RPPAArray - Reverse Phase Protein Array
- Mutation - Mutations
- GISTICA - GISTIC v2 ('AllByGene' only)
- GISTICT - GISTIC v2 ('ThresholdedByGene' only)
- GISTIC - GISTIC v2 scores and probabilities (both)

Author(s)

Marcel Ramos <marcel.ramos@roswellpark.org>

Examples

```
## Not run:  
coadmut <- getFirehoseData("COAD", runDate = "20151101", Mutation = TRUE)  
biocExtract(coadmut, "Mutation")  
  
## End(Not run)
```

CorResult-class	<i>An S4 class to store correlations between gene expression level and copy number data</i>
-----------------	---

Description

An S4 class to store correlations between gene expression level and copy number data

Slots

Dataset A cohort name

Correlations Results data frame

DGEResult-class	<i>An S4 class to store differential gene expression results</i>
-----------------	--

Description

An S4 class to store differential gene expression results

Slots

Dataset Dataset name

Toptable Results data frame

FirehoseCGHArray-class	<i>An S4 class to store data from CGA platforms</i>
------------------------	---

Description

An S4 class to store data from CGA platforms

Slots

Filename Platform name

DataMatrix A data frame that stores the CGH data.

FirehoseData-class *An S4 class to store main data object from client function.*

Description

An S4 class to store main data object from client function.

Usage

```
## S4 method for signature 'FirehoseData'
show(object)

## S4 method for signature 'FirehoseData'
getData(object, type, platform)

## S4 method for signature 'FirehoseData'
updateObject(object, ..., verbose = FALSE)

## S4 method for signature 'FirehoseData'
selectType(object, dataType)
```

Arguments

object	A FirehoseData object
type	A data type to be extracted
platform	An index for data types that may come from multiple platforms (such as mRNAArray), for GISTIC data, one of the options: 'AllByGene' or 'Thresholded-ByGene'
...	additional arguments for updateObject
verbose	logical (default FALSE) whether to print extra messages
dataType	An available data type, see object show method

Methods (by generic)

- show: show method
- getData: Get a matrix or data.frame from FirehoseData
- updateObject: Update an old RCGAToolbox FirehoseData object to the most recent API
- selectType: Extract data type

Slots

Dataset A cohort name
runDate Standard data run date from [getFirehoseRunningDates](#)
gistic2Date Analyze running date from [getFirehoseAnalyzeDates](#)
clinical clinical data frame
RNASeqGene Gene level expression data matrix from RNAseq
RNASeq2GeneNorm Gene level expression data matrix from RNAseq (RSEM)

miRNASeqGene miRNA expression data from matrix smallRNAseq
 CNASNP A data frame to store somatic copy number alterations from SNP array platform
 CNVSNP A data frame to store germline copy number variants from SNP array platform
 CNASeq A data frame to store somatic copy number alterations from sequencing platform
 CNACGH A list that stores FirehoseCGHArray object for somatic copy number alterations from CGH platform
 Methylation A list that stores FirehoseMethylationArray object for methylation data
 mRNAArray A list that stores FirehosemRNAArray object for gene expression data from microarray
 miRNAArray A list that stores FirehosemiRNAArray object for miRNA expression data from microarray
 RPPAArray A list that stores FirehosemRNAArray object for RPPA data
 Mutation A data frame for mutation information from sequencing data
 GISTIC A FirehoseGISTIC object to store processed copy number data
 BarcodeUUID A data frame that stores the Barcodes, UUIDs and Short sample identifiers

FirehoseGISTIC-class *An S4 class to store processed copy number data. (Data processed by using GISTIC2 algorithm)*

Description

An S4 class to store processed copy number data. (Data processed by using GISTIC2 algorithm)

Usage

```
## S4 method for signature 'FirehoseGISTIC'
updateObject(object, ..., verbose = FALSE)
```

Arguments

object	A FirehoseGISTIC object
...	additional arguments for updateObject
verbose	logical (default FALSE) whether to print extra messages

Methods (by generic)

- updateObject: Update an old FirehoseGISTIC object to the most recent API

Slots

Dataset Cohort name
 AllByGene A data frame that stores continuous copy number
 ThresholdedByGene A data frame for discrete copy number data

 FirehoseMethylationArray-class

An S4 class to store data from methylation platforms

Description

An S4 class to store data from methylation platforms

Slots

Filename Platform name

DataMatrix A data frame that stores the methylation data.

FirehosemRNAArray-class

An S4 class to store data from array (mRNA, miRNA etc.) platforms

Description

An S4 class to store data from array (mRNA, miRNA etc.) platforms

Slots

Filename Platform name

DataMatrix A data matrix that stores the expression data.

getCNGECorrelation

Perform correlation analysis between gene expression and copy number data

Description

getCNGECorrelation returns a list that stores the results correlation between gene expression and copy number data.

Usage

```
getCNGECorrelation(dataObject, adj.method = "BH", adj.pval = 0.05,
  raw.pval = 0.05)
```

Arguments

dataObject This must be FirehoseData object.

adj.method Raw p value adjustment methods (Default "BH")

adj.pval Adjusted p value cut off for results table (Default 0.05)

raw.pval raw p value cut off for results table (Default 0.05)

Value

Returns a list that stores results for each dataset

Examples

```
data(RTCGASample)
corRes <- getCNGECorrelation(RTCGASample)
corRes
showResults(corRes[[1]])
```

getData

Extract data from FirehoseData object

Description

A go-to function for getting top level information from a [FirehoseData](#) object. Available datatypes for a particular object can be seen by entering the object name in the console ('show' method).

Usage

```
getData(object, type, platform)
```

Arguments

object	A FirehoseData object
type	A data type to be extracted
platform	An index for data types that may come from multiple platforms (such as mRNAArray), for GISTIC data, one of the options: 'AllByGene' or 'Thresholded-ByGene'

Value

Returns matrix or data.frame depending on data type

Examples

```
data(RTCGASample)
getData(RTCGASample, "clinical")
getData(RTCGASample, "RNASeqGene")
```

getDiffExpressedGenes *Perform differential gene expression analysis for mRNA expression data.*

Description

getDiffExpressedGenes returns a list that stores the results for each dataset.

Usage

```
getDiffExpressedGenes(dataObject, DrawPlots = TRUE, adj.method = "BH",  
  adj.pval = 0.05, raw.pval = 0.05, logFC = 2, hmTopUpN = 100,  
  hmTopDownN = 100, meanFilter = 10)
```

Arguments

dataObject	This must be FirehoseData object.
DrawPlots	A logical parameter to draw heatmaps and volcano plots.
adj.method	Raw p value adjustment methods (Default "BH")
adj.pval	Adjusted p value cut off for results table (Default 0.05)
raw.pval	raw p value cut off for results table (Default 0.05)
logFC	log fold change cut off for results table (Default 2)
hmTopUpN	Max number of up regulated genes in heatmap (Default 100)
hmTopDownN	Max number of down regulated genes in heatmap (Default 100)
meanFilter	Mean read counts for each gene to filter not expressed genes (Default 10)

Value

Returns a list that stores results for each dataset

Examples

```
data(RTCGASample)  
dgegenes <- getDiffExpressedGenes(RTCGASample)  
dgegenes  
showResults(dgegenes[[1]])  
dgegenes <- showResults(dgegenes[[1]])  
head(dgegenes)
```

```
getFirehoseAnalyzeDates
```

Get data analyze dates.

Description

getFirehoseAnalyzeDates returns the character vector for analyze release dates.

Usage

```
getFirehoseAnalyzeDates(last = NULL)
```

Arguments

last To list last n dates. (Default NULL)

Value

A character vector for dates.

Examples

```
getFirehoseAnalyzeDates(last=2)
```

```
getFirehoseData            Get data from Firehose portal.
```

Description

getFirehoseData returns FirehoseData object that stores TCGA data.

Usage

```
getFirehoseData(dataset, runDate = "20160128", gistic2Date = "20160128",
  RNASeqGene = FALSE, clinical = TRUE, miRNASeqGene = FALSE,
  RNASeq2GeneNorm = FALSE, CNASNP = FALSE, CNVSNP = FALSE,
  CNASeq = FALSE, CNACGH = FALSE, Methylation = FALSE, Mutation = FALSE,
  mRNAArray = FALSE, miRNAArray = FALSE, RPPAArray = FALSE,
  RNAseqNorm = "raw_counts", RNAseq2Norm = "normalized_count",
  forceDownload = FALSE, destdir = ".", fileSizeLimit = 500,
  getUUIDs = FALSE)
```

Arguments

dataset	A cohort name. All dataset names can be accessible via getFirehoseDatasets
runDate	Standard data run dates. Date list can be accessible via getFirehoseRunningDates
gistic2Date	Analyze running dates for GISTIC processed copy number data. Date list can be accessible via getFirehoseAnalyzeDates
RNASeqGene	Logical (default FALSE) parameter for RNAseq data.

clinical	Logical (default TRUE) parameter for clinical data.
miRNASeqGene	Logical (default FALSE) parameter for smallRNAseq data.
RNASeq2GeneNorm	Logical (default FALSE) parameter for RNAseq v2 (RSEM processed) data.
CNASNP	Logical (default FALSE) parameter for somatic copy number alterations data from SNP array.
CNVSNP	Logical (default FALSE) parameter for germline copy number variants data from SNP array.
CNASeq	Logical (default FALSE) parameter for somatic copy number alterations data from sequencing.
CNACGH	Logical (default FALSE) parameter for somatic copy number alterations data from CGH.
Methylation	Logical (default FALSE) parameter for methylation data.
Mutation	Logical (default FALSE) parameter for mutation data from sequencing.
mRNAArray	Logical (default FALSE) parameter for mRNA expression data from microarray.
miRNAArray	Logical (default FALSE) parameter for miRNA expression data from microarray.
RPPAArray	Logical (default FALSE) parameter for RPPA data
RNAseqNorm	RNAseq data normalization method. (Default raw_counts)
RNAseq2Norm	RNAseq v2 data normalization method. (Default normalized_count)
forceDownload	A logic (Default FALSE) key to force download RTCGAToolbox every time. By default if you download files into your working directory once than RTCGAToolbox using local files next time.
destdir	Directory in which to store the resulting downloaded file. Defaults to current working directory.
fileSizeLimit	Files that are larger than set value (megabyte) won't be downloaded (Default: 500)
getUUIDs	Logical key to get UUIDs from barcode (Default: FALSE)

Details

This is a main client function to download data from Firehose TCGA portal.

Value

A FirehoseData data object that stores data for selected data types.

Examples

```
# Sample Dataset
data(RTCGASample)
RTCGASample
## Not run:
BRCAdata <- getFirehoseData(dataset="BRCA",
runDate="20140416",gistic2Date="20140115",
RNASeqGene=TRUE,clinical=TRUE,mRNAArray=TRUE,Mutation=TRUE)

## End(Not run)
```

getFirehoseDatasets *Get list of TCGA cohorts.*

Description

getFirehoseDatasets returns a character array for cohorts.

Usage

```
getFirehoseDatasets()
```

Value

A character string

Examples

```
getFirehoseDatasets()
```

getFirehoseRunningDates
 Get standard data running dates.

Description

getFirehoseRunningDates returns the character vector for standard data release dates.

Usage

```
getFirehoseRunningDates(last = NULL)
```

Arguments

last To list last n dates. (Default NULL)

Value

A character vector for dates.

Examples

```
getFirehoseRunningDates()  
getFirehoseRunningDates(last=2)
```

getMutationRate	<i>Make a table for mutation rate of each gene in the cohort</i>
-----------------	--

Description

Make a table for mutation rate of each gene in the cohort

Usage

```
getMutationRate(dataObject)
```

Arguments

dataObject This must be FirehoseData object.

Value

Returns a data table

Examples

```
data(RTCGASample)
mutRate <- getMutationRate(dataObject=RTCGASample)
mutRate <- mutRate[order(mutRate[,2],decreasing = TRUE),]
head(mutRate)
## Not run:
```

getReport	<i>Draws a circle plot into working directory</i>
-----------	---

Description

getReport draws a circle plot into your workin director to show log fold changes for differentially expressed genes, copy number alterations and mutations.

Usage

```
getReport(dataObject, DGEResult1 = NULL, DGEResult2 = NULL, geneLocations)
```

Arguments

dataObject This must be FirehoseData object.
DGEResult1 Differential gene expression results object (Optional)
DGEResult2 Differential gene expression results object (Optional)
geneLocations Gene coordinates.

Value

Draws a circle plot

Examples

```

data(RTCGASample)
require("Homo.sapiens")
locations <- genes(Homo.sapiens,columns="SYMBOL")
locations <- as.data.frame(locations)
locations <- locations[,c(6,1,5,2:3)]
locations <- locations[!is.na(locations[,1]),]
locations <- locations[!duplicated(locations[,1]),]
rownames(locations) <- locations[,1]
t1 <- getDiffExpressedGenes(RTCGASample)
## Not run:
getReport(dataObject=RTCASample,DGEResult1=t1[[1]],geneLocations=locations)

## End(Not run)

```

getSurvival

Perform survival analysis based on gene expression data

Description

getSurvival draws a KM plot and show survival analysis results between groups that are defined by gene expression data

Usage

```
getSurvival(dataObject, numberOfGroups = 2, geneSymbols, sampleTimeCensor)
```

Arguments

dataObject This must be FirehoseData object.

numberOfGroups Can be set as 2 or 3. (Default 2) Order and divide samples into n groups by using gene expression data.

geneSymbols Gene symbol that is going to be tested

sampleTimeCensor a data frame that stores clinical data. First column should store sample IDs, second column should have time and third column should have event information. For more information please see vignette.

Value

Draws a KM plot

Examples

```

## get data with getFirehoseData function and call survival analysis
## Always check clinical data file for structural changes

data(RTCGASample)
clinicData <- getData(RTCGASample,"clinical")
clinicData = clinicData[,3:5]
clinicData[is.na(clinicData[,3]),3] = clinicData[is.na(clinicData[,3]),2]
survData <- data.frame(Samples=rownames(clinicData),Time=as.numeric(clinicData[,3]),
Censor=as.numeric(clinicData[,1]))
getSurvival(dataObject=RTCASample,geneSymbols=c("FCGBP"),sampleTimeCensor=survData)

```

hg19.ucsc.gene.locations

Gene coordinates for circle plot.

Description

A dataset containing the gene coordinates The variables are as follows:

Format

A data frame with 28454 rows and 5 variables

Details

- GeneSymbol. Gene symbols
- Chromosome. Chromosome name
- Strand. Gene strand on chromosome
- Start. Gene location on chromosome
- End. Gene location on chromosome

RTCGASample

A sample data object for sample codes.

Description

A FirehoseData object for running sample codes. The variables are as follows:

Example dataset not biologically meaningful

Usage

RTCGASample

Format

A FirehoseData data object

Details

- a2. A sample data object

 RTCGAToolbox

RTCGAToolbox: A New Tool for Exporting TCGA Firehose Data

Description

Managing data from large-scale projects (such as The Cancer Genome Atlas (TCGA) for further analysis is an important and time consuming step for research projects. Several efforts, such as the Firehose project, make TCGA pre-processed data publicly available via web services and data portals, but this information must be managed, downloaded and prepared for subsequent steps. We have developed an open source and extensible R based data client for pre-processed data from the Firehose, and demonstrate its use with sample case studies. Results show that our RTCGAToolbox can facilitate data management for researchers interested in working with TCGA data. The RTCGAToolbox can also be integrated with other analysis pipelines for further data processing.

Details

The main function you're likely to need from **RTCGAToolbox** is `getFirehoseData`. Otherwise refer to the vignettes to see how to use the **RTCGAToolbox**

Author(s)

Mehmet Kemal Samur

 selectType

Accessor function for the FirehoseData object

Description

An accessor function for the `FirehoseData` class. An argument will specify the data type to return. See `FirehoseData-class` for more details.

Usage

```
selectType(object, dataType)
```

Arguments

object	A <code>FirehoseData</code> class object
dataType	A data type, see details.

Details

- clinical - Get the clinical data slot
- RNASeqGene - RNASeqGene
- RNASeq2GeneNorm - Normalized
- miRNASeqGene - micro RNA SeqGene
- CNASNP - Copy Number Alteration
- CNVSNP - Copy Number Variation

- CNASeq - Copy Number Alteration
- CNACGH - Copy Number Alteration
- Methylation - Methylation
- mRNAArray - Messenger RNA
- miRNAArray - micro RNA
- RPPAArray - Reverse Phase Protein Array
- Mutation - Mutations
- GISTIC - GISTIC v2 scores and probabilities

Value

The data type element of the FirehoseData object

showResults	<i>Export toptable or correlation data frame</i>
-------------	--

Description

Export toptable or correlation data frame

Usage

```
showResults(object)
```

Arguments

object A [DGEResult](#) or [CorResult](#) object

Value

Returns toptable or correlation data frame

Examples

```
data(RTCGASample)
dgeRes = getDiffExpressedGenes(RTCGASample)
dgeRes
showResults(dgeRes[[1]])
```

showResults,CorResult-method

Export toptable or correlation data frame

Description

Export toptable or correlation data frame

Usage

```
## S4 method for signature 'CorResult'  
showResults(object)
```

Arguments

object A [DGEResult](#) or [CorResult](#) object

Value

Returns correlation results data frame

Examples

```
data(RTCGASample)  
corRes = getCNGECorrelation(RTCGASample,adj.pval = 1,raw.pval = 1)  
corRes  
showResults(corRes[[1]])
```

showResults,DGEResult-method

Export toptable or correlation data frame

Description

Export toptable or correlation data frame

Usage

```
## S4 method for signature 'DGEResult'  
showResults(object)
```

Arguments

object A [DGEResult](#) or [CorResult](#) object

Value

Returns toptable for DGE results

Examples

```
data(RTCGASample)
dgeRes = getDiffExpressedGenes(RTCGASample)
dgeRes
showResults(dgeRes[[1]])
```

Index

- *Topic **data**
 - RTCGASample, 15
- biocExtract, 2
- CorResult, 17, 18
- CorResult-class, 4
- DGEResult, 17, 18
- DGEResult-class, 4
- FirehoseCGHArray-class, 4
- FirehoseData, 2, 8, 16
- FirehoseData-class, 5, 16
- FirehoseGISTIC-class, 6
- FirehoseMethylationArray-class, 7
- FirehosemRNAArray-class, 7
- getCNGECorrelation, 7
- getData, 8
- getData, FirehoseData-method (FirehoseData-class), 5
- getDiffExpressedGenes, 9
- getFirehoseAnalyzeDates, 5, 10, 10
- getFirehoseData, 10, 16
- getFirehoseDatasets, 10, 12
- getFirehoseRunningDates, 5, 10, 12
- getMutationRate, 13
- getReport, 13
- getSurvival, 14
- hg19.ucsc.gene.locations, 15
- RaggedExperiment, 2, 3
- RTCGASample, 15
- RTCGAToolbox, 16
- RTCGAToolbox-package (RTCGAToolbox), 16
- selectType, 16
- selectType, FirehoseData-method (FirehoseData-class), 5
- show, FirehoseData-method (FirehoseData-class), 5
- showResults, 17
- showResults, CorResult, CorResult-method (showResults, CorResult-method), 18
- showResults, CorResult-method, 18
- showResults, DGEResult, DGEResult-method (showResults, DGEResult-method), 18
- showResults, DGEResult-method, 18
- SummarizedExperiment, 2, 3
- updateObject, FirehoseData-method (FirehoseData-class), 5
- updateObject, FirehoseGISTIC-method (FirehoseGISTIC-class), 6