

# Package ‘NetworkDataCompanion’

November 13, 2024

**Title** Tools for Analyzing TCGA and GTEx Data

**Version** 0.0.0.9000

**Maintainer** Kate Hoff Shutta <kshutta@hsph.harvard.edu>

**Description** An R library of utilities for performing analyses on TCGA and GTEx data using the Network Zoo (<https://netzoo.github.io>).

**License** GPL-3

**biocViews**

**Depends** AnnotationDbi,  
data.table,  
dplyr,  
edgeR,  
EpiSCORE,  
GenomicDataCommons,  
huge,  
magrittr,  
methods,  
org.Hs.eg.db,  
presto,  
recount,  
recount3,  
stringr,  
TCGAPurityFiltering,  
TCGAutils,  
testthat,  
tidyr

**Remotes** pmandros/TCGAPurityFiltering,  
immunogenomics/presto,  
aet21/EpiSCORE

**Encoding** UTF-8

**Roxygen** list(markdown = TRUE)

**RoxygenNote** 7.3.2

**Suggests** knitr,  
rmarkdown

**VignetteBuilder** knitr

## R topics documented:

convertBetaToM . . . . .	2
CreateNetworkDataCompanionObject . . . . .	3
estimateCellCountsEpiSCORE . . . . .	3
extractSampleAndGeneInfo . . . . .	4
filterBarcodesIntersection . . . . .	5
filterChromosome . . . . .	5
filterDuplicatesSeqDepth . . . . .	6
filterDuplicatesSeqDepthOther . . . . .	7
filterGenesByTPM . . . . .	7
filterGenesProteins . . . . .	8
filterPurity . . . . .	9
filterSampleType . . . . .	9
filterTumorType . . . . .	10
geneNameToENSG . . . . .	11
getGeneInfo . . . . .	11
logCPMNormalization . . . . .	12
logTPMNormalization . . . . .	13
mapBarcodeToBarcode . . . . .	13
mapProbesToGenes . . . . .	14
NetworkDataCompanion-class . . . . .	15
probeToMeanPromoterMethylation . . . . .	16
<b>Index</b>	<b>17</b>

---

convertBetaToM	<i>Convert methylation beta values to M-values.</i>
----------------	---

---

### Description

This function uses the typical logit base 2 transformation to convert from methylation beta values (in the [0,1] range) to m-values (on the real line). The formula is  $m = \log_2(\text{beta}/(1-\text{beta}))$ .

### Usage

```
## S4 method for signature 'NetworkDataCompanion'
convertBetaToM(methylation_betas)
```

### Arguments

methylation\_betas  
A numeric vector of values in the range [0,1].

### Value

A numeric vector of m-values corresponding to the converted values of methylation\_betas.

### Author(s)

Kate Hoff Shutta (kshutta@hsph.harvard.edu)

**References**

<https://bmcbioinformatics.biomedcentral.com/articles/10.1186/1471-2105-11-587>

---

CreateNetworkDataCompanionObject

*Constructor for the NetworkDataCompanion object.*

---

**Description**

This function is used to construct a NetworkDataCompanion object. The member functions of this object are the functions of this package.

**Usage**

```
CreateNetworkDataCompanionObject(clinical_patient_file, project_name)
```

**Arguments**

clinical\_patient\_file

Path to a comma-separated file containing clinical data for the samples of interest.

project\_name     A character string that identifies the project.

**Value**

A NetworkDataCompanion object.

**Author(s)**

Panagiotis Mandros (pmandros@hsph.harvard.edu)

---

estimateCellCountsEpiSCORE

*Run the EpiSCORE algorithm to estimate cell type proportions.*

---

**Description**

This function applies the 'constAvBetaTSS' and 'wRPC' functions from the EpiSCORE R package within the TCGA data structure. The 'wRPC' parameters used are the defaults: 'useW=TRUE', 'wth=0.4', and 'maxit=100'.

**Usage**

```
## S4 method for signature 'NetworkDataCompanion'
estimateCellCountsEpiSCORE(methylation_betas, tissue, array = "450k")
```

**Arguments**

methylation_betas	A data frame of methylation beta values, with CGs in rows and samples in columns. The first column must be "probeID" and contain the Illumina probeIDs matching the specified array or a subset thereof.
tissue	Tissue type. Must be one of the tissues with a reference available in EpiSCORE. Acceptable values are "Bladder", "Brain", "Breast", "Colon", "Heart", "Kidney", "Liver", "Lung", "OE", "Pancreas_6ct", "Pancreas_9ct", "Prostate", "Skin".
array	Methylation array identifier. Acceptable values are "450k" or "850k" (EPIC).

**Value**

A data frame containing samples in rows and estimated cell type proportions in columns. The first two columns are the TCGA barcode and the TCGA UUID.

**Author(s)**

Kate Hoff Shutta (kshutta@hsph.harvard.edu)

**References**

Teschendorff, A.E., Zhu, T., Breeze, C.E. et al. EPIScore: cell type deconvolution of bulk tissue DNA methylomes from single-cell RNA-Seq data. *Genome Biol* 21, 221 (2020). <https://doi.org/10.1186/s13059-020-02126-9>

---

extractSampleAndGeneInfo

*Extracts experiment-specific information and metadata from ranged summarized experiment object.*

---

**Description**

Extracts experiment-specific information and metadata from ranged summarized experiment object.

**Usage**

```
## S4 method for signature 'NetworkDataCompanion'
extractSampleAndGeneInfo(expression_rds_obj)
```

**Arguments**

expression_rds_obj	A ranged summarized experiment object
--------------------	---------------------------------------

**Value**

rds_sample_info	metadata about the samples (columns)
rds_gene_info	metadata about the genes (rows)

**Author(s)**

Jonas Fischer (jfischer@hsph.harvard.edu)

---

`filterBarcodesIntersection`

*Convenience wrapper function for `mapBarcodeToBarcode` that applies the function directly to two data frames.*

---

### Description

This function returns a list of the two argument data frames, intersected, and the second frame ordered to match the first. NOTE: Ordering is done based on columns, which are expected to be named by TCGA barcodes.

### Usage

```
## S4 method for signature 'NetworkDataCompanion'
filterBarcodesIntersection(exp1, exp2)
```

### Arguments

<code>exp1</code>	A matrix or dataframe with TCGA barcodes as the <code>colnames</code> attribute.
<code>exp2</code>	A matrix or dataframe with TCGA barcodes as the <code>colnames</code> attribute.

### Details

No additional details at this time.

### Value

<code>mappedExp1</code>	A data frame filtered to include only columns with TCGA barcodes that are in both <code>colnames(exp1)</code> and <code>colnames(exp2)</code> .
<code>mappedExp2</code>	A data frame filtered to include only columns with TCGA barcodes that are in both <code>colnames(exp1)</code> and <code>colnames(exp2)</code> . IMPORTANT: The columns of <code>mappedExp2</code> are ordered to match the column ordering of <code>mappedExp1</code> .

### Author(s)

Jonas Fischer (jfisher@hsph.harvard.edu)

---

`filterChromosome`

*Filter for genes in a particular chromosome or chromosomes.*

---

### Description

This function filters for genes in a particular chromosome or chromosomes.

### Usage

```
## S4 method for signature 'NetworkDataCompanion'
filterChromosome(rds_gene_info, chroms)
```

**Arguments**

- `rds_gene_info` A data frame extracted from a `RangedSummarizedExperiment` object containing expression data using the `extractSampleAndGeneInfo` function.
- `chroms` A character vector of chromosomes. Must exactly match chromosomes in the `seqname` attribute of `rds_gene_info`.

**Value**

Integer vector of the row indices (genes) in `rds_gene_info` to keep.

**Author(s)**

Jonas Fischer (jfisher@hsph.harvard.edu)

---

`filterDuplicatesSeqDepth`

*A function to filter duplicates based on RNA sequencing depth.*

---

**Description**

This function filters out duplicates based on RNA-seq, keeping the samples with maximum read depth. Returns indices of samples to KEEP.

**Usage**

```
## S4 method for signature 'NetworkDataCompanion'  
filterDuplicatesSeqDepth(expression_count_matrix)
```

**Arguments**

- `expression_count_matrix`  
Matrix of count data from RNA-seq experiment, with genes in rows and samples in columns.

**Value**

Integer vector of indices to keep, corresponding to columns of `expression_count_matrix`.

**Author(s)**

Jonas Fischer(jfisher@hsph.harvard.edu)

---

`filterDuplicatesSeqDepthOther`

*A version of filterDuplicatesSeqDepth to handle the case when sequencing depth is not available.*

---

### Description

This function takes a random duplicate if no info is available on sequencing depth for all vials.

### Usage

```
## S4 method for signature 'NetworkDataCompanion'
filterDuplicatesSeqDepthOther(expression_count_matrix, tcga_barcodes)
```

### Arguments

`expression_count_matrix` Matrix of count data from RNA-seq experiment, with genes in rows and samples in columns.

`tcga_barcodes` List of TCGA barcodes for filtering.

### Value

Integer vector of indices to KEEP in `tcga_barcodes`.

### Author(s)

Jonas Fischer (jfischer@hsph.harvard.edu)

---

`filterGenesByTPM`

*Filter genes based on minimum expression level (TPM) across samples.*

---

### Description

Filter all genes which have at least `tpm_threshold` TPM scores in at least `sample_fraction` of samples.

### Usage

```
## S4 method for signature 'NetworkDataCompanion'
filterGenesByTPM(expression_tpm_matrix, tpm_threshold, sample_fraction)
```

**Arguments**

- `expression_tpm_matrix` A data frame extracted from a `RangedSummarizedExperiment` object containing expression data using the `extractSampleAndGeneInfo` function.
- `tpm_threshold` Numeric  $> 0$ . Genes with TPM below this values in more than `sample_fraction` of the data will be excluded from the analysis.
- `sample_fraction` Numeric in  $[0,1]$ . Genes with TPM below `tpm_threshold` in more than this fraction of the data will be excluded from the analysis.

**Value**

Integer vector indexing the rows of `expression_tpm_matrix` that correspond to genes that should be kept.

**Author(s)**

Jonas Fischer (jfischer@hsph.harvard.edu)

---

`filterGenesProteins`     *Filtering protein coding genes through an rds object.*

---

**Description**

Filtering protein coding genes through an rds object.

**Usage**

```
## S4 method for signature 'NetworkDataCompanion'
filterGenesProteins(rds_gene_info)
```

**Arguments**

- `rds_gene_info` rds info object of genes, usually extracted from row information from recount retrieved rds expression objects.

**Value**

Array of indices that correspond to the protein coding genes in the rds gene info table.

**Author(s)**

Jonas Fischer (jfischer@hsph.harvard.edu)



---

filterPurity	<i>Filter samples based on tumor purity.</i>
--------------	--

---

**Description**

This function filters a character vector of TCGA barcodes for tumor purity based on a particular method and threshold.

**Usage**

```
## S4 method for signature 'NetworkDataCompanion'  
filterPurity(TCGA_barcodes, method="ESTIMATE", threshold=.6)
```

**Arguments**

TCGA_barcodes	Character vector of TCGA barcodes that the user wishes to filter based on tumor purity.
method	One of "ESTIMATE", "ABSOLUTE", "LUMP", "IHC", or "CPE". Default is "ESTIMATE".
threshold	Threshold for purity-based filtering. Samples with a purity below threshold will be filtered out.

**Details**

Describe the method options.

**Value**

Integer vector of indices indicating which samples in TCGA\_bar codes should be kept.

**Author(s)**

Panagiotis Mandros (pmandros@hsph.harvard.edu)

**References**

This code is based on the TCGAPurityFiltering package found at <https://github.com/pmandros/TCGAPurityFiltering>.

---

filterSampleType	<i>Filter samples based on sample type.</i>
------------------	---

---

**Description**

This function filters samples based on TCGA sample types.

**Usage**

```
## S4 method for signature 'NetworkDataCompanion'  
filterSampleType(TCGA_barcodes, types_of_samples)
```

**Arguments**

TCGA\_barcodes    A character vector of TCGA barcodes.  
 types\_of\_samples    A character vector representing the types of samples to select.

**Details**

Candidate values for types\_of\_samples are characters of the form "01", "02", etc. that correspond to TCGA sample type codes. Valid arguments for types\_of\_samples can be found on the GDC website: <https://gdc.cancer.gov/resources-tcga-users/tcga-code-tables/sample-type-codes> and a table mapping sample types to values can be found by using `NetworkDataCompanion::getSampleTypeMap()`.

**Value**

Named list of containing "index", an integer vector of indices in TCGA\_barcodes to keep, and "type", a character vector of sample type corresponding to each index.

**Author(s)**

Jonas Fischer (jfisher@hsph.harvard.edu)

---

filterTumorType	<i>Filter samples based on tumor type.</i>
-----------------	--

---

**Description**

This function filters samples based on tumor type. Some examples are: "Primary Tumor", "Solid Tissue Normal", "Primary Blood Derived Cancer - Peripheral Blood". This function is particularly useful for excluding normal samples from analyses.

**Usage**

```
## S4 method for signature 'NetworkDataCompanion'
filterTumorType(TCGA_barcodes, type_of_tumor, rds_info)
```

**Arguments**

TCGA\_barcodes    A character vector of TCGA barcodes.  
 type\_of\_tumor    A string representing the type of tumor to select. Currently, only a single tumor type is supported.  
 rds\_info    A data frame extracted from a `RangedSummarizedExperiment` object containing expression data using the `extractSampleAndGeneInfo` function.

**Details**

Candidate values for type\_of\_tumor: "Primary Tumor", "Solid Tissue Normal", "Primary Blood Derived Cancer - Peripheral Blood", etc. There are other options that show up in TCGA that are not listed here. Make sure it is an exact match - check spaces, case, etc.

**Value**

Integer vector of indices in TCGA\_barcodes to keep.

**Author(s)**

Jonas Fischer (jfischer@hsph.harvard.edu)

---

geneNameToENSG	<i>Convert from gene name to Ensembl stable id.</i>
----------------	---

---

**Description**

Given an input character vector of gene names, this function converts them to Ensembl stable IDs. Note from <https://useast.ensembl.org/Help/Faq?id=488>: "An Ensembl stable ID consists of five parts: ENS(species)(object type)(identifier).(version)."

**Usage**

```
## S4 method for signature 'NetworkDataCompanion'
geneNameToENSG(gene_names, version = FALSE)
```

**Arguments**

gene_names	Character vector of gene names.
version	Boolean; retrieve Ensembl version along with Ensembl identifier.

**Value**

Character vector of Ensembl stable IDs.

**Author(s)**

Panagiotis Mandros (pmandros@hsph.harvard.edu)

**References**

<https://useast.ensembl.org/index.html>

---

getGeneInfo	<i>Retrieve a variety of gene information based on gene name or Ensembl stable ID.</i>
-------------	--

---

**Description**

This function uses the gene\_mapping attribute of the NetworkDataCompanion object to provide information on seqid, source, start, end, strand, gene\_id, gene\_name, gene\_type, and gene\_id\_no\_ver.

**Usage**

```
## S4 method for signature 'NetworkDataCompanion'
getGeneInfo(gene_names_or_ids)
```

**Arguments**

gene\_names\_or\_ids

A character vector of gene names or Ensembl stable IDs.

**Details**

This function will determine the input type based on the presence of the string "ENSG".

**Value**

A data frame with rows representing genes and columns representing gene attributes (e.g., source, start, end.)

**Author(s)**

Panagiotis Mandros (pmandros@hsph.harvard.edu)

---

logCPMNormalization     *Function to compute CPM values from raw counts.*

---

**Description**

This function computes CPM values from raw expression counts using the edgeR package as a backend.

**Usage**

```
## S4 method for signature 'NetworkDataCompanion'  
logCPMNormalization(exp_count_mat)
```

**Arguments**

exp\_count\_mat     Matrix or data.frame of raw expression counts.

**Value**

counts	The original count matrix passed as argument to this function.
CPM	CPM transformed values of the same shape as the count matrix.
logCPM	log(CPM + 1) transformed values of the same shape as the count matrix.

**Author(s)**

Jonas Fischer (jfischer@hsph.harvard.edu)

**References**

<https://bioconductor.org/packages/release/bioc/html/edgeR.html>

**See Also**

[edgeR](#).

---

logTPMNormalization	<i>Function for within-array normalization of a RangedSummarizedExperiment object with log transcripts per million (TPM) normalization.</i>
---------------------	---

---

### Description

Returns a named list with raw counts (useful for duplicate filtering based on sequencing depth, see `?filterDuplicatesSeqDepth`), TPM, (useful for TPM-based filtering, see `?filterGenesByTPM`), and the actual log TPM. A pseudocount of 1 is added to each TPM value for this function, so returned "log TPM" values actually correspond to  $\log(\text{TPM} + 1)$ .

### Usage

```
## S4 method for signature 'NetworkDataCompanion'
logTPMNormalization(expression_rds_obj)
```

### Arguments

`expression_rds_obj`  
A `RangedSummarizedExperiment` object.

### Value

<code>counts</code>	A data frame of RNA sequencing counts matching the row and column ordering of <code>expression_rds_obj</code> .
<code>TPM</code>	A data frame of TPM matching the row and column ordering of <code>expression_rds_obj</code> .
<code>logTPM</code>	A data frame of log-transformed TPM with pseudocounts (i.e., $\log(\text{TPM} + 1)$ ) matching the row and column ordering of <code>expression_rds_obj</code> .

### Author(s)

Jonas Fischer (jfisher@hsph.harvard.edu)

### References

`RangedSummarizedExperiment` documentation

---

mapBarcodeToBarcode	<i>Helper function for mapping two sets of TCGA barcodes to each other.</i>
---------------------	---

---

### Description

There are 4 different pieces of information returned in a named list that are all useful depending on the context in which they are used.

`is_inter1` is an indicator (boolean) vector of the same length as `bc1` that indicates which elements of `bc1` are present in `bc2`.

`ids1` indicates where to find each barcode of `bc1` in `bc2`, returning NA if there is no match. That is, `ids1[i] != NA`, then `bc1[i] == bc2[ids1[i]]`.

The same information is provided for `bc2`.

**Usage**

```
## S4 method for signature 'NetworkDataCompanion'
mapBarcodeToBarcode(bc1,bc2)
```

**Arguments**

bc1	Character vector of barcodes in the first set.
bc2	Character vector of barcodes in the second set.

**Value**

is_inter1	Boolean vector of the same length as bc1 that indicates which elements of bc1 are present in bc2.
idcs1	Integer vector of the same length as bc1 that indicates where to find each barcode of bc1 in bc2, returning NA if there is no match. That is, <code>idcs1[i] != NA</code> , then <code>bc1[i] == bc2[idcs1[i]]</code>
is_inter2	Boolean vector of the same length as bc2 that indicates which elements of bc2 are present in bc1.
idcs2	Integer vector of the same length as bc2 that indicates where to find each barcode of bc2 in bc1, returning NA if there is no match. That is, <code>idcs2[i] != NA</code> , then <code>bc2[i] == bc1[idcs2[i]]</code> .

**Note**

For example, if you want to map experiment 1 onto experiment 2, keeping only the information for samples that are present in both, and reordering the first experiment to match the samples of the second, you can do:

```
exp1[,is_inter1] # this will remove samples that are not in experiment 2)
exp2[,idcs1[is_inter1]] # this will remove samples that are not in exp1 and reorder to match exp1
```

**Author(s)**

Jonas Fischer (jfisher@hsph.harvard.edu)

---

mapProbesToGenes	<i>Maps input probe IDs to gene TSS within a certain range upstream and downstream.</i>
------------------	---

---

**Description**

Maps input probe IDs to gene TSS within a certain range upstream and downstream.

**Usage**

```
## S4 method for signature 'NetworkDataCompanion'
mapProbesToGenes(probelist, rangeUp, rangeDown, localManifestPath=NA)
```

**Arguments**

probelist	A character vector of Illumina array probes, e.g., c("cg03636183", "cg19859270").
rangeUp	The number of base pairs upstream to search for a TSS. Must be a non-negative number.
rangeDown	The number of base pairs downstream to search for a TSS. Must be a non-negative number.
localManifestPath	If you wish to use a manifest file other than the Illumina manifest found at <a href="https://zhouserver.research.chop.edu/InfiniumAnnotation/20210615/HM450/HM450.hg38.manifest.gencode.v36.tsv.gz">https://zhouserver.research.chop.edu/InfiniumAnnotation/20210615/HM450/HM450.hg38.manifest.gencode.v36.tsv.gz</a> , you can pass a path to that file here. It should be formatted in the same way as the Illumina manifest.

**Value**

A matrix with four columns: probeID, geneName, ensemblID, distToTSS. When a probe maps to more than one TSS within the upstream and downstream parameters provided, the geneName, ensemblID, and distToTSS columns will contain lists of genes separated by a semicolon(";"). Ordering of the lists matches between the three columns.

**Author(s)**

Kate Hoff Shutta (kshutta@hsph.harvard.edu)

**References**

<https://zwdzwd.github.io/InfiniumAnnotation>

---

NetworkDataCompanion-class

*NetworkDataCompanion Reference Class*

---

**Description**

NetworkDataCompanion is a reference class that provides fields for handling data related to TCGA and GTEx projects.

**Details**

A package for easy and reproducible wrangling of TCGA and GTEx data.

**Fields**

TCGA\_purities A data.frame containing TCGA sample purity information.  
 clinical\_patient\_data A data.frame with clinical data for each patient.  
 project\_name A character vector with the name of the project.  
 gene\_mapping A data.frame that maps gene identifiers between datasets.  
 sample\_type\_mapping A data.frame for sample type classification and mapping.

Methods

Placeholder for methods documentation. Define each method here once methods are added.

Author(s)

Viola Fanfani (<vfanfani@hsph.harvard.edu>), Jonas Fischer (<jfischer@hsph.harvard.edu>),  
Panagiotis Mandros (<pmandros@hsph.harvard.edu>), Soel Micheletti (<smicheletti@hsph.harvard.edu>),  
Kate Hoff Shutta (<kshutta@hsph.harvard.edu>)

References

<https://www.biorxiv.org/content/10.1101/2024.11.05.622163v1.abstract>

---

probeToMeanPromoterMethylation  
*probeToMeanPromoterMethylation*

---

Description

Calculates the average promoter methylation within a certain window around the transcription start site (TSS), as defined by the input probe\_gene\_map.

Usage

```
## S4 method for signature 'NetworkDataCompanion'  
probeToMeanPromoterMethylation(methylation_betas, probe_gene_map, genesOfInterest)
```

Arguments

methylation\_betas  
A data frame of methylation beta values, with CGs in rows and samples in columns. The first column must be "probeID" and contain the Illumina probeIDs matching the probe\_gene\_map argument or a subset thereof.

probe\_gene\_map  
Output from mapProbesToGenes, or otherwise a bespoke matrix with four columns: probeID, geneName, ensemblID, distToTSS.

genesOfInterest  
Character vector of gene names for which mean promoter methylation should be calculated.

Value

Matrix of samples in rows and genes in columns. row.names stores sample names and colnames stores gene names.

Author(s)

Kate Hoff Shutta (kshutta@hsph.harvard.edu)



# Index

`convertBetaToM`, [2](#)  
`CreateNetworkDataCompanionObject`, [3](#)  
  
`edgeR`, [12](#)  
`estimateCellCountsEpiSCORE`, [3](#)  
`extractSampleAndGeneInfo`, [4](#)  
  
`filterBarcodesIntersection`, [5](#)  
`filterChromosome`, [5](#)  
`filterDuplicatesSeqDepth`, [6](#)  
`filterDuplicatesSeqDepthOther`, [7](#)  
`filterGenesByTPM`, [7](#)  
`filterGenesProteins`, [8](#)  
`filterPurity`, [9](#)  
`filterSampleType`, [9](#)  
`filterTumorType`, [10](#)  
  
`geneNameToENSG`, [11](#)  
`getGeneInfo`, [11](#)  
  
`logCPMNormalization`, [12](#)  
`logTPMNormalization`, [13](#)  
  
`mapBarcodeToBarcode`, [13](#)  
`mapProbesToGenes`, [14](#)  
  
`NetworkDataCompanion`  
    (`NetworkDataCompanion-class`),  
    [15](#)  
`NetworkDataCompanion-class`, [15](#)  
  
`probeToMeanPromoterMethylation`, [16](#)