

# Package ‘mia’

November 25, 2022

**Type** Package

**Version** 1.6.0

**Title** Microbiome analysis

**Description** mia implements tools for microbiome analysis based on the SummarizedExperiment, SingleCellExperiment and TreeSummarizedExperiment infrastructure. Data wrangling and analysis in the context of taxonomic data is the main scope. Additional functions for common task are implemented such as community indices calculation and summarization.

**biocViews** Microbiome, Software, DataImport

**License** Artistic-2.0 | file LICENSE

**Encoding** UTF-8

**LazyData** false

**Depends** R (>= 4.0), SummarizedExperiment, SingleCellExperiment, TreeSummarizedExperiment (>= 1.99.3), MultiAssayExperiment

**Imports** methods, stats, utils, MASS, ape, decontam, vegan, BiocGenerics, S4Vectors, IRanges, Biostrings, DECIPHER, BiocParallel, DelayedArray, DelayedMatrixStats, scuttle, scatter, DirichletMultinomial, rlang, dplyr, tibble, tidyr

**Suggests** testthat, knitr, patchwork, BiocStyle, yaml, phyloseq, dada2, stringr, biomformat, reldist, ade4, microbiomeDataSets, rmarkdown

**URL** <https://github.com/microbiome/mia>

**BugReports** <https://github.com/microbiome/mia/issues>

**Roxygen** list(markdown = TRUE)

**RoxygenNote** 7.2.1

**VignetteBuilder** knitr

**git\_url** <https://git.bioconductor.org/packages/mia>

**git\_branch** RELEASE\_3\_16

**git\_last\_commit** 41e31ad

**git\_last\_commit\_date** 2022-11-01

**Date/Publication** 2022-11-25

**Author** Felix G.M. Ernst [aut] (<<https://orcid.org/0000-0001-5064-0928>>),  
 Sudarshan A. Shetty [aut] (<<https://orcid.org/0000-0001-7280-9915>>),  
 Tuomas Borman [aut, cre] (<<https://orcid.org/0000-0002-8563-8884>>),  
 Leo Lahti [aut] (<<https://orcid.org/0000-0001-5537-637X>>),  
 Yang Cao [ctb],  
 Nathan D. Olson [ctb],  
 Levi Waldron [ctb],  
 Marcel Ramos [ctb],  
 Héctor Corrada Bravo [ctb],  
 Jayaram Kancherla [ctb],  
 Domenick Braccia [ctb]

**Maintainer** Tuomas Borman <tuomas.v.borman@utu.fi>

## R topics documented:

mia-package . . . . .	3
agglomerate-methods . . . . .	3
calculateDMN . . . . .	5
calculateJSD . . . . .	9
calculateOverlap . . . . .	11
calculateUnifrac . . . . .	12
estimateDivergence . . . . .	15
estimateDiversity . . . . .	17
estimateDominance . . . . .	22
estimateEvenness . . . . .	26
estimateRichness . . . . .	29
getAbundance . . . . .	32
getExperimentCrossAssociation . . . . .	34
getPrevalence . . . . .	39
isContaminant . . . . .	44
loadFromMetaphlan . . . . .	46
loadFromMothur . . . . .	48
loadFromQIIME2 . . . . .	49
makePhyloseqFromTreeSE . . . . .	51
makeTreeSEFromBiom . . . . .	53
makeTreeSEFromDADA2 . . . . .	54
makeTreeSEFromPhyloseq . . . . .	55
meltAssay . . . . .	56
merge-methods . . . . .	57
mergeSEs . . . . .	59
mia-datasets . . . . .	62
perSampleDominantTaxa . . . . .	65
relabundance . . . . .	67
runCCA . . . . .	68
runDPCoA . . . . .	70
runNMDS . . . . .	72

splitByRanks . . . . .	75
splitOn . . . . .	77
subsampleCounts . . . . .	79
subsetSamples . . . . .	81
summaries . . . . .	82
taxonomy-methods . . . . .	85
transformCounts . . . . .	88

## Index 94

---

mia-package	mia <i>Package</i> .
-------------	----------------------

---

### Description

mia implements tools for microbiome analysis based on the SummarizedExperiment, SingleCellExperiment and TreeSummarizedExperiment infrastructure. Data wrangling and analysis in the context of taxonomic data is the main scope. Additional functions for common task are implemented such as community indices calculation and summarization.

### See Also

[TreeSummarizedExperiment](#) class

---

agglomerate-methods	Agglomerate data using taxonomic information
---------------------	--

---

### Description

agglomerateByRank can be used to sum up data based on the association to certain taxonomic ranks given as rowData. Only available [taxonomyRanks](#) can be used.

### Usage

```
## S4 method for signature 'SummarizedExperiment'
agglomerateByRank(
  x,
  rank = taxonomyRanks(x)[1],
  onRankOnly = FALSE,
  na.rm = FALSE,
  empty.fields = c(NA, "", " ", "\t", "-", "_"),
  ...
)

## S4 method for signature 'SingleCellExperiment'
agglomerateByRank(x, ..., altexp = NULL, strip_altexp = TRUE)

## S4 method for signature 'TreeSummarizedExperiment'
agglomerateByRank(x, ..., agglomerateTree = FALSE)
```

**Arguments**

<code>x</code>	a <a href="#">SummarizedExperiment</a> object
<code>rank</code>	a single character defining a taxonomic rank. Must be a value of <code>taxonomyRanks()</code> function.
<code>onRankOnly</code>	TRUE or FALSE: Should information only from the specified rank be used or from ranks equal and above? See details. (default: <code>onRankOnly = FALSE</code> )
<code>na.rm</code>	TRUE or FALSE: Should taxa with an empty rank be removed? Use it with caution, since empty entries on the selected rank will be dropped. This setting can be tweaked by defining <code>empty.fields</code> to your needs. (default: <code>na.rm = TRUE</code> )
<code>empty.fields</code>	a character value defining, which values should be regarded as empty. (Default: <code>c(NA, "", " ", "\t")</code> ). They will be removed if <code>na.rm = TRUE</code> before agglomeration.
<code>...</code>	arguments passed to <code>agglomerateByRank</code> function for <code>SummarizedExperiment</code> objects, <a href="#">mergeRows</a> and <a href="#">sumCountsAcrossFeatures</a> . <ul style="list-style-type: none"> <li>• <code>remove_empty_ranks</code>A single boolean value for selecting whether to remove those columns of <code>rowData</code> that include only NAs after agglomeration. (By default: <code>remove_empty_ranks = FALSE</code>)</li> <li>• <code>make_unique</code>A single boolean value for selecting whether to make row-names unique. (By default: <code>make_unique = TRUE</code>)</li> </ul>
<code>altexp</code>	String or integer scalar specifying an alternative experiment containing the input data.
<code>strip_altexp</code>	TRUE or FALSE: Should alternative experiments be removed prior to agglomeration? This prevents to many nested alternative experiments by default (default: <code>strip_altexp = TRUE</code> )
<code>agglomerateTree</code>	TRUE or FALSE: should <code>rowTree()</code> also be agglomerated? (Default: <code>agglomerateTree = FALSE</code> )

**Details**

Based on the available taxonomic data and its structure setting `onRankOnly = TRUE` has certain implications on the interpretability of your results. If no loops exist (loops meaning two higher ranks containing the same lower rank), the results should be comparable. you can check for loops using [detectLoop](#).

Agglomeration sum up values of assays at specified taxonomic level. Certain assays, e.g. those that include binary or negative values, can lead to meaningless values, when values are summed. In those cases, consider doing agglomeration first and then transformation.

**Value**

A taxonomically-agglomerated, optionally-pruned object of the same class as `x`.

**See Also**

[mergeRows](#), [sumCountsAcrossFeatures](#)

**Examples**

```
data(GlobalPatterns)
# print the available taxonomic ranks
colnames(rowData(GlobalPatterns))
taxonomyRanks(GlobalPatterns)

# agglomerate at the Family taxonomic rank
x1 <- agglomerateByRank(GlobalPatterns, rank="Family")
## How many taxa before/after agglomeration?
nrow(GlobalPatterns)
nrow(x1)

# with agglomeration of the tree
x2 <- agglomerateByRank(GlobalPatterns, rank="Family",
                        agglomerateTree = TRUE)
nrow(x2) # same number of rows, but
rowTree(x1) # ... different
rowTree(x2) # ... tree

# If assay contains binary or negative values, summing might lead to meaningless
# values, and you will get a warning. In these cases, you might want to do
# agglomeration again at chosen taxonomic level.
tse <- transformSamples(GlobalPatterns, method = "pa")
tse <- agglomerateByRank(tse, rank = "Genus")
tse <- transformSamples(tse, method = "pa")

# removing empty labels by setting na.rm = TRUE
sum(is.na(rowData(GlobalPatterns)$Family))
x3 <- agglomerateByRank(GlobalPatterns, rank="Family", na.rm = TRUE)
nrow(x3) # different from x2

# Because all the rownames are from the same rank, rownames do not include
# prefixes, in this case "Family:".
print(rownames(x3[1:3,]))

# To add them, use getTaxonomyLabels function.
rownames(x3) <- getTaxonomyLabels(x3, with_rank = TRUE)
print(rownames(x3[1:3,]))

# use 'remove_empty_ranks' to remove columns that include only NAs
x4 <- agglomerateByRank(GlobalPatterns, rank="Phylum", remove_empty_ranks = TRUE)
head(rowData(x4))

## Look at enterotype dataset...
data(enterotype)
## print the available taxonomic ranks. Shows only 1 rank available
## not useful for agglomerateByRank
taxonomyRanks(enterotype)
```

---

calculatedMNM                      *Dirichlet-Multinomial Mixture Model: Machine Learning for Microbiome Data*

---

## Description

These functions are accessors for functions implemented in the [DirichletMultinomial](#) package

## Usage

```
calculateDMN(x, ...)

## S4 method for signature 'ANY'
calculateDMN(
  x,
  k = 1,
  BPPARAM = SerialParam(),
  seed = runif(1, 0, .Machine$integer.max),
  ...
)

## S4 method for signature 'SummarizedExperiment'
calculateDMN(
  x,
  assay_name = abund_values,
  abund_values = exprs_values,
  exprs_values = "counts",
  transposed = FALSE,
  ...
)

runDMN(x, name = "DMN", ...)

getDMN(x, name = "DMN", ...)

## S4 method for signature 'SummarizedExperiment'
getDMN(x, name = "DMN")

bestDMNFit(x, name = "DMN", type = c("laplace", "AIC", "BIC"), ...)

## S4 method for signature 'SummarizedExperiment'
bestDMNFit(x, name = "DMN", type = c("laplace", "AIC", "BIC"))

getBestDMNFit(x, name = "DMN", type = c("laplace", "AIC", "BIC"), ...)

## S4 method for signature 'SummarizedExperiment'
getBestDMNFit(x, name = "DMN", type = c("laplace", "AIC", "BIC"))
```

```
calculateDMNgroup(x, ...)  
  
## S4 method for signature 'ANY'  
calculateDMNgroup(  
  x,  
  variable,  
  k = 1,  
  seed = runif(1, 0, .Machine$integer.max),  
  ...  
)  
  
## S4 method for signature 'SummarizedExperiment'  
calculateDMNgroup(  
  x,  
  variable,  
  assay_name = abund_values,  
  abund_values = exprs_values,  
  exprs_values = "counts",  
  transposed = FALSE,  
  ...  
)  
  
performDMNgroupCV(x, ...)  
  
## S4 method for signature 'ANY'  
performDMNgroupCV(  
  x,  
  variable,  
  k = 1,  
  seed = runif(1, 0, .Machine$integer.max),  
  ...  
)  
  
## S4 method for signature 'SummarizedExperiment'  
performDMNgroupCV(  
  x,  
  variable,  
  assay_name = abund_values,  
  abund_values = exprs_values,  
  exprs_values = "counts",  
  transposed = FALSE,  
  ...  
)
```

### Arguments

x                    a numeric matrix with samples as rows or a [SummarizedExperiment](#) object.  
...                   optional arguments not used.

k	the number of Dirichlet components to fit. See <a href="#">dmn</a>
BPPARAM	A <a href="#">BiocParallelParam</a> object specifying whether the UniFrac calculation should be parallelized.
seed	random number seed. See <a href="#">dmn</a>
assay_name	a single character value for specifying which assay to use for calculation.
abund_values	a single character value for specifying which assay to use for calculation. (Please use <code>assay_name</code> instead. At some point <code>abund_values</code> will be disabled.)
exprs_values	a single character value for specifying which assay to use for calculation. (Please use <code>assay_name</code> instead.)
transposed	Logical scalar, is x transposed with samples in rows?
name	the name to store the result in <a href="#">metadata</a>
type	the type of measure used for the goodness of fit. One of 'laplace', 'AIC' or 'BIC'.
variable	a variable from <code>colData</code> to use as a grouping variable. Must be a character of factor.

### Value

`calculateDMN` and `getDMN` return a list of DMN objects, one element for each value of `k` provided.

`bestDMNFit` returns the index for the best fit and `getBestDMNFit` returns a single DMN object.

`calculateDMNgroup` returns a [DMNGroup](#) object

`performDMNgroupCV` returns a `data.frame`

### See Also

[DMN-class](#), [DMNGroup-class](#), [dmn](#), [dmngroup](#), [cvdmn](#), [accessors for DMN objects](#)

### Examples

```
f1 <- system.file(package="DirichletMultinomial", "extdata", "Twins.csv")
counts <- as.matrix(read.csv(f1, row.names=1))
f1 <- system.file(package="DirichletMultinomial", "extdata", "TwinStudy.t")
pheno0 <- scan(f1)
lvls <- c("Lean", "Obese", "Overwt")
pheno <- factor(lvls[pheno0 + 1], levels=lvls)
colData <- DataFrame(pheno = pheno)

tse <- TreeSummarizedExperiment(assays = list(counts = counts),
                               colData = colData)

#
dmn <- calculateDMN(tse)
dmn[[1L]]

# since this take a bit of resources to calculate for k > 1, the data is
```



```

# loaded
## Not run:
tse <- runDMN(tse, name = "DMN", k = 1:7)

## End(Not run)
data(dmn_se)
# Convert SE to TreeSE to enable all features of mia and TreeSE object
# (Optional step at this point)
dmn_tse <- as(dmn_se, "TreeSummarizedExperiment")
names(metadata(dmn_tse))

# return a list of DMN objects
getDMN(dmn_tse)
# return, which objects fits best
bestDMNFit(dmn_tse, type = "laplace")
# return the model, which fits best
getBestDMNFit(dmn_tse, type = "laplace")

```

---

calculateJSD

*Calculate the Jensen-Shannon Divergence*


---

## Description

This function calculates the Jensen-Shannon Divergence (JSD) in a [SummarizedExperiment](#) object.

## Usage

```

## S4 method for signature 'ANY'
calculateJSD(x, ...)

## S4 method for signature 'SummarizedExperiment'
calculateJSD(
  x,
  assay_name = abund_values,
  abund_values = exprs_values,
  exprs_values = "counts",
  transposed = FALSE,
  ...
)

runJSD(x, BPPARAM = SerialParam(), chunkSize = nrow(x))

```

## Arguments

x	a numeric matrix or a <a href="#">SummarizedExperiment</a> .
...	optional arguments not used.
assay_name	a single character value for specifying which assay to use for calculation.

abund_values	a single character value for specifying which assay to use for calculation. (Please use assay_name instead. At some point abund_values will be disabled.)
exprs_values	a single character value for specifying which assay to use for calculation. (Please use assay_name instead.)
transposed	Logical scalar, is x transposed with cells in rows?
BPPARAM	A <code>BiocParallelParam</code> object specifying whether the JSD calculation should be parallelized.
chunkSize	an integer scalar, defining the size of data send to the individual worker. Only has an effect, if BPPARAM defines more than one worker. (default: chunkSize = nrow(x))

**Value**

a sample-by-sample distance matrix, suitable for NMDS, etc.

**Author(s)**

Susan Holmes <susan@stat.stanford.edu>. Adapted for phyloseq by Paul J. McMurdie. Adapted for mia by Felix G.M. Ernst

**References**

Jensen-Shannon Divergence and Hilbert space embedding. Bent Fuglede and Flemming Topsøe University of Copenhagen, Department of Mathematics <http://www.math.ku.dk/~topsøe/ISIT2004JSD.pdf>

**See Also**

[http://en.wikipedia.org/wiki/Jensen-Shannon\\_divergence](http://en.wikipedia.org/wiki/Jensen-Shannon_divergence)

**Examples**

```
data(enterotype)
library(scater)

jsd <- calculateJSD(enterotype)
class(jsd)
head(jsd)

enterotype <- runMDS(enterotype, FUN = calculateJSD, name = "JSD",
                    exprs_values = "counts")
head(reducedDim(enterotype))
head(attr(reducedDim(enterotype), "eig"))
attr(reducedDim(enterotype), "GOF")
```

---

calculateOverlap	<i>Estimate overlap</i>
------------------	-------------------------

---

### Description

This function calculates overlap for all sample-pairs in a [SummarizedExperiment](#) object.

### Usage

```
calculateOverlap(  
  x,  
  assay_name = abund_values,  
  abund_values = "counts",  
  detection = 0,  
  ...  
)  
  
## S4 method for signature 'SummarizedExperiment'  
calculateOverlap(  
  x,  
  assay_name = abund_values,  
  abund_values = "counts",  
  detection = 0,  
  ...  
)  
  
runOverlap(x, ...)  
  
## S4 method for signature 'SummarizedExperiment'  
runOverlap(x, name = "overlap", ...)
```

### Arguments

x	a <a href="#">SummarizedExperiment</a> object containing a tree.
assay_name	A single character value for selecting the <a href="#">assay</a> to calculate the overlap.
abund_values	a single character value for specifying which assay to use for calculation. (Please use <code>assay_name</code> instead. At some point <code>abund_values</code> will be disabled.)
detection	A single numeric value for selecting detection threshold for absence/presence of features. Feature that has abundance under threshold in either of samples, will be discarded when evaluating overlap between samples.
...	Optional arguments not used.
name	A single character value specifying the name of overlap matrix that is stored in <code>reducedDim(x)</code> .

### Details

This function calculates overlap between all the sample-pairs. Overlap reflects similarity between sample-pairs.

When overlap is calculated using relative abundances, the higher the value the higher the similarity is. When using relative abundances, overlap value 1 means that all the abundances of features are equal between two samples, and 0 means that samples have completely different relative abundances.

### Value

calculateOverlap returns sample-by-sample distance matrix. runOverlap returns x that includes overlap matrix in its reducedDim.

### Author(s)

Leo Lahti and Tuomas Borman. Contact: [microbiome.github.io](https://microbiome.github.io)

### See Also

[calculateJSD](#) [calculateUnifrac](#)

### Examples

```
data(esophagus)
tse <- esophagus
tse <- transformSamples(tse, method = "relabundance")
overlap <- calculateOverlap(tse, assay_name = "relabundance")
overlap

# Store result to reducedDim
tse <- runOverlap(tse, assay_name = "relabundance", name = "overlap_between_samples")
head(reducedDims(tse)$overlap_between_samples)
```

---

calculateUnifrac

*Calculate weighted or unweighted (Fast) Unifrac distance*

---

### Description

This function calculates the (Fast) Unifrac distance for all sample-pairs in a [TreeSummarizedExperiment](#) object.

**Usage**

```

calculateUnifrac(x, tree, ...)

## S4 method for signature 'ANY,phylo'
calculateUnifrac(
  x,
  tree,
  weighted = FALSE,
  normalized = TRUE,
  BPPARAM = SerialParam()
)

## S4 method for signature 'TreeSummarizedExperiment,missing'
calculateUnifrac(
  x,
  assay_name = abund_values,
  abund_values = exprs_values,
  exprs_values = "counts",
  tree_name = "phylo",
  transposed = FALSE,
  ...
)

runUnifrac(
  x,
  tree,
  weighted = FALSE,
  normalized = TRUE,
  BPPARAM = SerialParam()
)

```

**Arguments**

x	a numeric matrix or a <a href="#">TreeSummarizedExperiment</a> object containing a tree. Please note that <code>runUnifrac</code> expects a matrix with samples per row and not per column. This is implemented to be compatible with other distance calculations such as <code>dist</code> as much as possible.
tree	if x is a matrix, a <a href="#">phylo</a> object matching the matrix. This means that the phylo object and the columns should relate to the same type of features (aka. microorganisms).
...	optional arguments not used.
weighted	TRUE or FALSE: Should use weighted-Unifrac calculation? Weighted-Unifrac takes into account the relative abundance of species/taxa shared between samples, whereas unweighted-Unifrac only considers presence/absence. Default is FALSE, meaning the unweighted-Unifrac distance is calculated for all pairs of samples.

normalized	TRUE or FALSE: Should the output be normalized such that values range from 0 to 1 independent of branch length values? Default is TRUE. Note that (unweighted) Unifrac is always normalized by total branch-length, and so this value is ignored when weighted == FALSE.
BPPARAM	A <a href="#">BiocParallelParam</a> object specifying whether the Unifrac calculation should be parallelized.
assay_name	a single character value for specifying which assay to use for calculation.
abund_values	a single character value for specifying which assay to use for calculation. (Please use assay_name instead. At some point abund_values will be disabled.)
exprs_values	a single character value for specifying which assay to use for calculation. (Please use assay_name instead.)
tree_name	a single character value for specifying which tree will be used in calculation. (By default: tree_name = "phylo")
transposed	Logical scalar, is x transposed with cells in rows, i.e., is Unifrac distance calculated based on rows (FALSE) or columns (TRUE). (By default: transposed = FALSE)

### Details

Please note that if calculateUnifrac is used as a FUN for runMDS, the argument ntop has to be set to nrow(x).

### Value

a sample-by-sample distance matrix, suitable for NMDS, etc.

### Author(s)

Paul J. McMurdie. Adapted for mia by Felix G.M. Ernst

### References

<http://bmf.colorado.edu/unifrac/>

The main implementation (Fast Unifrac) is adapted from the algorithm's description in:

Hamady, Lozupone, and Knight, "Fast UniFrac: facilitating high-throughput phylogenetic analyses of microbial communities including analysis of pyrosequencing and PhyloChip data." *The ISME Journal* (2010) 4, 17–27.

See also additional descriptions of Unifrac in the following articles:

Lozupone, Hamady and Knight, "Unifrac - An Online Tool for Comparing Microbial Community Diversity in a Phylogenetic Context.", *BMC Bioinformatics* 2006, 7:371

Lozupone, Hamady, Kelley and Knight, "Quantitative and qualitative (beta) diversity measures lead to different insights into factors that structure microbial communities." *Appl Environ Microbiol.* 2007

Lozupone C, Knight R. "Unifrac: a new phylogenetic method for comparing microbial communities." *Appl Environ Microbiol.* 2005 71 (12):8228-35.

## Examples

```
data(esophagus)
library(scater)
calculateUnifrac(esophagus, weighted = FALSE)
calculateUnifrac(esophagus, weighted = TRUE)
calculateUnifrac(esophagus, weighted = TRUE, normalized = FALSE)
# for using calculateUnifrac in conjunction with runMDS the tree argument
# has to be given separately. In addition, subsetting using ntop must
# be disabled
esophagus <- runMDS(esophagus, FUN = calculateUnifrac, name = "Unifrac",
                   tree = rowTree(esophagus),
                   exprs_values = "counts",
                   ntop = nrow(esophagus))
reducedDim(esophagus)
```

---

estimateDivergence	<i>Estimate divergence</i>
--------------------	----------------------------

---

## Description

Estimate divergence against a given reference sample.

## Usage

```
estimateDivergence(
  x,
  assay_name = abund_values,
  abund_values = "counts",
  name = "divergence",
  reference = "median",
  FUN = vegan::vegdist,
  method = "bray",
  ...
)

## S4 method for signature 'SummarizedExperiment'
estimateDivergence(
  x,
  assay_name = abund_values,
  abund_values = "counts",
  name = "divergence",
  reference = "median",
  FUN = vegan::vegdist,
  method = "bray",
  ...
)
```

## Arguments

x	a <a href="#">SummarizedExperiment</a> object.
assay_name	the name of the assay used for calculation of the sample-wise estimates.
abund_values	a single character value for specifying which assay to use for calculation. (Please use assay_name instead. At some point abund_values will be disabled.)
name	a name for the column of the colData the results should be stored in. By default, name is "divergence".
reference	a numeric vector that has length equal to number of features, or a non-empty character value; either 'median' or 'mean'. reference specifies the reference that is used to calculate divergence. by default, reference is "median".
FUN	a function for distance calculation. The function must expect the input matrix as its first argument. With rows as samples and columns as features. By default, FUN is <code>vegan::vegdist</code> .
method	a method that is used to calculate the distance. Method is passed to the function that is specified by FUN. By default, method is "bray".
...	optional arguments

## Details

Microbiota divergence (heterogeneity / spread) within a given sample set can be quantified by the average sample dissimilarity or beta diversity with respect to a given reference sample.

This measure is sensitive to sample size. Subsampling or bootstrapping can be applied to equalize sample sizes between comparisons.

## Value

x with additional `colData` named `*name*`

## Author(s)

Leo Lahti and Tuomas Borman. Contact: [microbiome.github.io](https://microbiome.github.io)

## See Also

[plotColData](#)

- [estimateRichness](#)
- [estimateEvenness](#)
- [estimateDominance](#)



**Examples**

```

data(GlobalPatterns)
tse <- GlobalPatterns

# By default, reference is median of all samples. The name of column where results
# is "divergence" by default, but it can be specified.
tse <- estimateDivergence(tse)

# The method that are used to calculate distance in divergence and
# reference can be specified. Here, euclidean distance and dist function from
# stats package are used. Reference is the first sample.
tse <- estimateDivergence(tse, name = "divergence_first_sample",
                        reference = assays(tse)$counts[,1],
                        FUN = stats::dist, method = "euclidean")

# Reference can also be median or mean of all samples.
# By default, divergence is calculated by using median. Here, mean is used.
tse <- estimateDivergence(tse, name = "divergence_average", reference = "mean")

# All three divergence results are stored in colData.
colData(tse)

```

---

estimateDiversity      *Estimate (alpha) diversity measures*

---

**Description**

Several functions for calculating (alpha) diversity indices, including the vegan package options and some others.

**Usage**

```

estimateDiversity(
  x,
  assay_name = abund_values,
  abund_values = "counts",
  index = c("coverage", "fisher", "gini_simpson", "inverse_simpson",
            "log_modulo_skewness", "shannon"),
  name = index,
  ...
)

## S4 method for signature 'SummarizedExperiment'
estimateDiversity(
  x,
  assay_name = abund_values,
  abund_values = "counts",

```

```

    index = c("coverage", "fisher", "gini_simpson", "inverse_simpson",
              "log_modulo_skewness", "shannon"),
    name = index,
    ...,
    BPPARAM = SerialParam()
)

## S4 method for signature 'TreeSummarizedExperiment'
estimateDiversity(
  x,
  assay_name = abund_values,
  abund_values = "counts",
  index = c("coverage", "faith", "fisher", "gini_simpson", "inverse_simpson",
            "log_modulo_skewness", "shannon"),
  name = index,
  tree_name = "phylo",
  ...,
  BPPARAM = SerialParam()
)

estimateFaith(
  x,
  tree = "missing",
  assay_name = abund_values,
  abund_values = "counts",
  name = "faith",
  ...
)

## S4 method for signature 'SummarizedExperiment,phylo'
estimateFaith(
  x,
  tree,
  assay_name = abund_values,
  abund_values = "counts",
  name = "faith",
  node_lab = NULL,
  ...
)

## S4 method for signature 'TreeSummarizedExperiment,missing'
estimateFaith(
  x,
  assay_name = abund_values,
  abund_values = "counts",
  name = "faith",
  tree_name = "phylo",
  ...
)

```

)

**Arguments**

x	a <a href="#">SummarizedExperiment</a> object or <a href="#">TreeSummarizedExperiment</a> . The latter is recommended for microbiome data sets and tree-based alpha diversity indices.
assay_name	the name of the assay used for calculation of the sample-wise estimates.
abund_values	a single character value for specifying which assay to use for calculation. (Please use assay_name instead. At some point abund_values will be disabled.)
index	a character vector, specifying the diversity measures to be calculated.
name	a name for the column(s) of the colData the results should be stored in. By default this will use the original names of the calculated indices.
...	optional arguments: <ul style="list-style-type: none"> <li>• threshold A numeric value in the unit interval, determining the threshold for coverage index. By default, threshold is 0.9.</li> <li>• quantile Arithmetic abundance classes are evenly cut up to to this quantile of the data. The assumption is that abundances higher than this are not common, and they are classified in their own group. By default, quantile is 0.5.</li> <li>• num_of_classes The number of arithmetic abundance classes from zero to the quantile cutoff indicated by quantile. By default, num_of_classes is 50.</li> </ul>
BPPARAM	A <a href="#">BiocParallelParam</a> object specifying whether calculation of estimates should be parallelized.
tree_name	a single character value for specifying which rowTree will be used to calculate faith index. (By default: tree_name = "phylo")
tree	A phylogenetic tree that is used to calculate 'faith' index. If x is a <a href="#">TreeSummarizedExperiment</a> , rowTree(x) is used by default.
node_lab	NULL or a character vector specifying the links between rows and node labels of tree. If a certain row is not linked with the tree, missing instance should be noted as NA. When NULL, all the rownames should be found from the tree. (By default: node_lab = NULL)

**Details**

The available indices include the 'Coverage', 'Faith's phylogenetic diversity', 'Fisher alpha', 'Gini-Simpson', 'Inverse Simpson', 'log-modulo skewness', and 'Shannon' indices. See details for more information and references.

Alpha diversity is a joint quantity that combines elements or community richness and evenness. Diversity increases, in general, when species richness or evenness increase.

By default, this function returns all indices.

- 'coverage' Number of species needed to cover a given fraction of the ecosystem (50 percent by default). Tune this with the threshold argument.

- 'faith' Faith's phylogenetic alpha diversity index measures how long the taxonomic distance is between taxa that are present in the sample. Larger values represent higher diversity. Using this index requires rowTree. (Faith 1992)
- 'fisher' Fisher's alpha; as implemented in `vegan::fisher.alpha`. (Fisher et al. 1943)
- 'gini\_simpson' Gini-Simpson diversity i.e.  $1 - \lambda$ , where  $\lambda$  is the Simpson index, calculated as the sum of squared relative abundances. This corresponds to the diversity index 'simpson' in `vegan::diversity`. This is also called Gibbs–Martin, or Blau index in sociology, psychology and management studies. The Gini-Simpson index (1-lambda) should not be confused with Simpson's dominance (lambda), Gini index, or inverse Simpson index (1/lambda).
- 'inverse\_simpson' Inverse Simpson diversity:  $1/\lambda$  where  $\lambda = \sum(p^2)$  and p refers to relative abundances. This corresponds to the diversity index 'invsimpson' in `vegan::diversity`. Don't confuse this with the closely related Gini-Simpson index
- 'log\_modulo\_skewness' The rarity index characterizes the concentration of species at low abundance. Here, we use the skewness of the frequency distribution of arithmetic abundance classes (see Magurran & McGill 2011). These are typically right-skewed; to avoid taking log of occasional negative skews, we follow Locey & Lennon (2016) and use the log-modulo transformation that adds a value of one to each measure of skewness to allow logarithmization.
- 'shannon' Shannon diversity (entropy).

### Value

x with additional `colData` named \*name\*

### Author(s)

Leo Lahti and Tuomas Borman. Contact: [microbiome.github.io](https://microbiome.github.io)

### References

- Beisel J-N. et al. (2003) A Comparative Analysis of Diversity Index Sensitivity. *Internal Rev. Hydrobiol.* 88(1):3-15. [https://portais.ufg.br/up/202/o/2003-comparative\\_ennes\\_index.pdf](https://portais.ufg.br/up/202/o/2003-comparative_ennes_index.pdf)
- Bulla L. (1994) An index of diversity and its associated diversity measure. *Oikos* 70:167–171
- Faith D.P. (1992) Conservation evaluation and phylogenetic diversity. *Biological Conservation* 61(1):1-10.
- Fisher R.A., Corbet, A.S. & Williams, C.B. (1943) The relation between the number of species and the number of individuals in a random sample of animal population. *Journal of Animal Ecology* 12, 42-58.
- Locey K.J. & Lennon J.T. (2016) Scaling laws predict global microbial diversity. *PNAS* 113(21):5970-5975.
- Magurran A.E., McGill BJ, eds (2011) Biological Diversity: Frontiers in Measurement and Assessment. (Oxford Univ Press, Oxford), Vol 12.
- Smith B. & Wilson JB. (1996) A Consumer's Guide to Diversity Indices. *Oikos* 76(1):70-82.

**See Also**[plotColData](#)

- [estimateRichness](#)
- [estimateEvenness](#)
- [estimateDominance](#)
- [diversity](#)
- [estimateR](#)

**Examples**

```

data(GlobalPatterns)
tse <- GlobalPatterns

# All index names as known by the function
index <- c("shannon","gini_simpson","inverse_simpson", "coverage", "fisher",
"faith", "log_modulo_skewness")

# Corresponding polished names
name <- c("Shannon","GiniSimpson","InverseSimpson", "Coverage", "Fisher",
"Faith", "LogModSkewness")

# Calculate diversities
tse <- estimateDiversity(tse, index = index)

# The colData contains the indices with their code names by default
colData(tse)[, index]

# Removing indices
colData(tse)[, index] <- NULL

# 'threshold' can be used to determine threshold for 'coverage' index
tse <- estimateDiversity(tse, index = "coverage", threshold = 0.75)
# 'quantile' and 'num_of_classes' can be used when
# 'log_modulo_skewness' is calculated
tse <- estimateDiversity(tse, index = "log_modulo_skewness",
    quantile = 0.75, num_of_classes = 100)

# It is recommended to specify also the final names used in the output.
tse <- estimateDiversity(tse,
    index = c("shannon", "gini_simpson", "inverse_simpson", "coverage",
    "fisher", "faith", "log_modulo_skewness"),
    name = c("Shannon", "GiniSimpson", "InverseSimpson", "Coverage",
    "Fisher", "Faith", "LogModSkewness"))

# The colData contains the indices by their new names provided by the user
colData(tse)[, name]

# Compare the indices visually
pairs(colData(tse)[, name])

```

```

# Plotting the diversities - use the selected names
library(scater)
plotColData(tse, "Shannon")
# ... by sample type
plotColData(tse, "Shannon", "SampleType")
## Not run:
# combining different plots
library(patchwork)
plot_index <- c("Shannon","GiniSimpson")
plots <- lapply(plot_index,
                plotColData,
                object = tse,
                x = "SampleType",
                colour_by = "SampleType")
plots <- lapply(plots,"+",
               theme(axis.text.x = element_text(angle=45,hjust=1)))
names(plots) <- plot_index
plots$Shannon + plots$GiniSimpson + plot_layout(guides = "collect")

## End(Not run)

```

---

estimateDominance      *Estimate dominance measures*

---

### Description

This function calculates community dominance indices. This includes the ‘Absolute’, ‘Berger-Parker’, ‘Core abundance’, ‘Gini’, ‘McNaughton’s’, ‘Relative’, and ‘Simpson’s’ indices.

### Usage

```

estimateDominance(
  x,
  assay_name = abund_values,
  abund_values = "counts",
  index = c("absolute", "dbp", "core_abundance", "gini", "dmn", "relative",
            "simpson_lambda"),
  ntaxa = 1,
  aggregate = TRUE,
  name = index,
  ...,
  BPPARAM = SerialParam()
)

## S4 method for signature 'SummarizedExperiment'
estimateDominance(
  x,
  assay_name = abund_values,
  abund_values = "counts",

```

```

index = c("absolute", "dbp", "core_abundance", "gini", "dmn", "relative",
  "simpson_lambda"),
ntaxa = 1,
aggregate = TRUE,
name = index,
...,
BPPARAM = SerialParam()
)

```

## Arguments

x	a <a href="#">SummarizedExperiment</a> object
assay_name	A single character value for selecting the <a href="#">assay</a> to calculate the sample-wise estimates.
abund_values	a single character value for specifying which assay to use for calculation. (Please use assay_name instead. At some point abund_values will be disabled.)
index	a character vector, specifying the indices to be calculated.
ntaxa	Optional and only used for the Absolute and Relative dominance indices: The n-th position of the dominant taxa to consider (default: ntaxa = 1). Disregarded for the indices “dbp”, “core_abundance”, “Gini”, “dmn”, and “Simpson”.
aggregate	Optional and only used for the Absolute, dbp, Relative, and dmn dominance indices: Aggregate the values for top members selected by ntaxa or not. If TRUE, then the sum of relative abundances is returned. Otherwise the relative abundance is returned for the single taxa with the indicated rank (default: aggregate = TRUE). Disregarded for the indices “core_abundance”, “gini”, “dmn”, and “simpson”.
name	A name for the column(s) of the colData where the calculated Dominance indices should be stored in.
...	additional arguments currently not used.
BPPARAM	A <a href="#">BiocParallelParam</a> object specifying whether calculation of estimates should be parallelized. (Currently not used)

## Details

A dominance index quantifies the dominance of one or few species in a community. Greater values indicate higher dominance.

Dominance indices are in general negatively correlated with alpha diversity indices (species richness, evenness, diversity, rarity). More dominant communities are less diverse.

estimateDominance calculates the following community dominance indices:

- ‘absolute’ Absolute index equals to the absolute abundance of the most dominant n species of the sample (specify the number with the argument ntaxa). Index gives positive integer values.
- ‘dbp’ Berger-Parker index (See Berger & Parker 1970) calculation is a special case of the ‘relative’ index. dbp is the relative abundance of the most abundant species of the sample. Index gives values in interval 0 to 1, where bigger value represent greater dominance.

$$dbp = \frac{N_1}{N_{tot}}$$

where  $N_1$  is the absolute abundance of the most dominant species and  $N_{tot}$  is the sum of absolute abundances of all species.

- 'core\_abundance' Core abundance index is related to core species. Core species are species that are most abundant in all samples, i.e., in whole data set. Core species are defined as those species that have prevalence over 50\ species must be prevalent in 50\ calculate the core abundance index. Core abundance index is sum of relative abundances of core species in the sample. Index gives values in interval 0 to 1, where bigger value represent greater dominance.

$$core\_abundance = \frac{N_{core}}{N_{tot}}$$

where  $N_{core}$  is the sum of absolute abundance of the core species and  $N_{tot}$  is the sum of absolute abundances of all species.

- 'gini' Gini index is probably best-known from socio-economic contexts (Gini 1921). In economics, it is used to measure, for example, how unevenly income is distributed among population. Here, Gini index is used similarly, but income is replaced with abundance.

If there is small group of species that represent large portion of total abundance of microbes, the inequality is large and Gini index closer to 1. If all species has equally large abundances, the equality is perfect and Gini index equals 0. This index should not be confused with Gini-Simpson index, which quantifies diversity.

- 'dmn' McNaughton's index is the sum of relative abundances of the two most abundant species of the sample (McNaughton & Wolf, 1970). Index gives values in the unit interval:

$$dmn = (N_1 + N_2)/N_{tot}$$

where  $N_1$  and  $N_2$  are the absolute abundances of the two most dominant species and  $N_{tot}$  is the sum of absolute abundances of all species.

- 'relative' Relative index equals to the relative abundance of the most dominant n species of the sample (specify the number with the argument ntaxa). This index gives values in interval 0 to 1.

$$relative = N_1/N_{tot}$$

where  $N_1$  is the absolute abundance of the most dominant species and  $N_{tot}$  is the sum of absolute abundances of all species.

- 'simpson\_lambda' Simpson's (dominance) index or Simpson's lambda is the sum of squared relative abundances. This index gives values in the unit interval. This value equals the probability that two randomly chosen individuals belongs to the same species. The higher the probability, the greater the dominance (See e.g. Simpson 1949).

$$lambda = \sum(p^2)$$

where p refers to relative abundances.

There is also a more advanced Simpson dominance index (Simpson 1949). However, this is not provided and the simpler squared sum of relative abundances is used instead as the alternative index is not in the unit interval and it is highly correlated with the simpler variant implemented here.



**Value**

x with additional `colData` named `*name*`

**Author(s)**

Leo Lahti and Tuomas Borman. Contact: [microbiome.github.io](https://microbiome.github.io)

**References**

Berger WH & Parker FL (1970) Diversity of Planktonic Foraminifera in Deep-Sea Sediments. *Science* 168(3937):1345-1347. doi: 10.1126/science.168.3937.1345

Gini C (1921) Measurement of Inequality of Incomes. *The Economic Journal* 31(121): 124-126. doi: 10.2307/2223319

McNaughton, SJ and Wolf LL. (1970). Dominance and the niche in ecological systems. *Science* 167:13, 1–139

Simpson EH (1949) Measurement of Diversity. *Nature* 163(688). doi: 10.1038/163688a0

**See Also**

- [estimateRichness](#)
- [estimateEvenness](#)
- [estimateDiversity](#)

**Examples**

```
data(esophagus)

# Calculates Simpson's lambda (can be used as a dominance index)
esophagus <- estimateDominance(esophagus, index="simpson_lambda")

# Shows all indices
colData(esophagus)

# Indices must be written correctly (e.g. dbp, not dbp), otherwise an error
# gets thrown
## Not run: esophagus <- estimateDominance(esophagus, index="dbp")
# Calculates dbp and Core Abundance indices
esophagus <- estimateDominance(esophagus, index=c("dbp", "core_abundance"))
# Shows all indices
colData(esophagus)
# Shows dbp index
colData(esophagus)$dbp
# Deletes dbp index
colData(esophagus)$dbp <- NULL
# Shows all indices, dbp is deleted
colData(esophagus)
# Deletes all indices
colData(esophagus) <- NULL
```

```

# Calculates all indices
esophagus <- estimateDominance(esophagus)
# Shows all indices
colData(esophagus)
# Deletes all indices
colData(esophagus) <- NULL

# Calculates all indices with explicitly specified names
esophagus <- estimateDominance(esophagus,
  index = c("dbp", "dmn", "absolute", "relative",
            "simpson_lambda", "core_abundance", "gini"),
  name = c("BergerParker", "McNaughton", "Absolute", "Relative",
           "SimpsonLambda", "CoreAbundance", "Gini")
)
# Shows all indices
colData(esophagus)

```

---

estimateEvenness	<i>Estimate Evenness measures</i>
------------------	-----------------------------------

---

## Description

This function calculates community evenness indices. These include the ‘Camargo’, ‘Pielou’, ‘Simpson’, ‘Evar’ and ‘Bulla’ evenness measures. See details for more information and references.

## Usage

```

estimateEvenness(
  x,
  assay_name = abund_values,
  abund_values = "counts",
  index = c("pielou", "camargo", "simpson_evenness", "evar", "bulla"),
  name = index,
  ...
)

## S4 method for signature 'SummarizedExperiment'
estimateEvenness(
  x,
  assay_name = abund_values,
  abund_values = "counts",
  index = c("camargo", "pielou", "simpson_evenness", "evar", "bulla"),
  name = index,
  ...,
  BPPARAM = SerialParam()
)

```

**Arguments**

x	a <code>SummarizedExperiment</code> object
assay_name	A single character value for selecting the <code>assay</code> used for calculation of the sample-wise estimates.
abund_values	a single character value for specifying which assay to use for calculation. (Please use <code>assay_name</code> instead. At some point <code>abund_values</code> will be disabled.)
index	a character vector, specifying the evenness measures to be calculated.
name	a name for the column(s) of the <code>colData</code> the results should be stored in.
...	optional arguments: <ul style="list-style-type: none"> <li>• <code>threshold</code> a numeric threshold. assay values below or equal to this threshold will be set to zero.</li> </ul>
BPPARAM	A <code>BiocParallelParam</code> object specifying whether calculation of estimates should be parallelized.

**Details**

Evenness is a standard index in community ecology, and it quantifies how evenly the abundances of different species are distributed. The following evenness indices are provided:

By default, this function returns all indices.

The available evenness indices include the following (all in lowercase):

- `'camargo'` Camargo's evenness (Camargo 1992)
- `'simpson_evenness'` Simpson's evenness is calculated as inverse Simpson diversity ( $1/\lambda$ ) divided by observed species richness  $S$ :  $(1/\lambda)/S$ .
- `'pielou'` Pielou's evenness (Pielou, 1966), also known as Shannon or Shannon-Weaver/Wiener/Weiner evenness;  $H/\ln(S)$ . The Shannon-Weaver is the preferred term; see Spellerberg and Fedor (2003).
- `'evar'` Smith and Wilson's Evar index (Smith & Wilson 1996).
- `'bulla'` Bulla's index (O) (Bulla 1994).

Desirable statistical evenness metrics avoid strong bias towards very large or very small abundances; are independent of richness; and range within the unit interval with increasing evenness (Smith & Wilson 1996). Evenness metrics that fulfill these criteria include at least `camargo`, `simpson`, `smith-wilson`, and `bulla`. Also see Magurran & McGill (2011) and Beisel et al. (2003) for further details.

**Value**

x with additional `colData` named `*name*`

## References

- Beisel J-N. et al. (2003) A Comparative Analysis of Evenness Index Sensitivity. *Internal Rev. Hydrobiol.* 88(1):3-15. URL: [https://portais.ufg.br/up/202/o/2003-comparative\\_evenness\\_index.pdf](https://portais.ufg.br/up/202/o/2003-comparative_evenness_index.pdf)
- Bulla L. (1994) An index of evenness and its associated diversity measure. *Oikos* 70:167–171.
- Camargo, JA. (1992) New diversity index for assessing structural alterations in aquatic communities. *Bull. Environ. Contam. Toxicol.* 48:428–434.
- Locey KJ and Lennon JT. (2016) Scaling laws predict global microbial diversity. *PNAS* 113(21):5970-5975; doi:10.1073/pnas.1521291113.
- Magurran AE, McGill BJ, eds (2011) *Biological Diversity: Frontiers in Measurement and Assessment* (Oxford Univ Press, Oxford), Vol 12.
- Pielou, EC. (1966) The measurement of diversity in different types of biological collections. *J Theoretical Biology* 13:131–144.
- Smith B and Wilson JB. (1996) A Consumer's Guide to Evenness Indices. *Oikos* 76(1):70-82.
- Spellerberg and Fedor (2003). A tribute to Claude Shannon (1916 –2001) and a plea for more rigorous use of species richness, species diversity and the 'Shannon–Wiener' Index. *Alpha Ecology & Biogeography* 12, 177–197.

## See Also

[plotColData](#)

- [estimateRichness](#)
- [estimateDominance](#)
- [estimateDiversity](#)

## Examples

```
data(esophagus)
tse <- esophagus

# Specify index and their output names
index <- c("pielou", "camargo", "simpson_evenness", "evar", "bulla")
name <- c("Pielou", "Camargo", "SimpsonEvenness", "Evar", "Bulla")

# Estimate evenness and give polished names to be used in the output
tse <- estimateEvenness(tse, index = index, name = name)

# Check the output
head(colData(tse))
```

---

estimateRichness	<i>Estimate richness measures</i>
------------------	-----------------------------------

---

## Description

Several functions for calculation of community richness indices available via wrapper functions. They are implemented via the vegan package.

## Usage

```
estimateRichness(
  x,
  assay_name = abund_values,
  abund_values = "counts",
  index = c("ace", "chao1", "hill", "observed"),
  name = index,
  detection = 0,
  ...,
  BPPARAM = SerialParam()
)

## S4 method for signature 'SummarizedExperiment'
estimateRichness(
  x,
  assay_name = abund_values,
  abund_values = "counts",
  index = c("ace", "chao1", "hill", "observed"),
  name = index,
  detection = 0,
  ...,
  BPPARAM = SerialParam()
)
```

## Arguments

x	a <a href="#">SummarizedExperiment</a> object.
assay_name	the name of the assay used for calculation of the sample-wise estimates.
abund_values	a single character value for specifying which assay to use for calculation. (Please use assay_name instead. At some point abund_values will be disabled.)
index	a character vector, specifying the richness measures to be calculated.
name	a name for the column(s) of the colData the results should be stored in.
detection	a numeric value for selecting detection threshold for the abundances. The default detection threshold is 0.
...	additional parameters passed to estimateRichness
BPPARAM	A <a href="#">BiocParallelParam</a> object specifying whether calculation of estimates should be parallelized.

## Details

These include the 'ace', 'Chao1', 'Hill', and 'Observed' richness measures. See details for more information and references.

The richness is calculated per sample. This is a standard index in community ecology, and it provides an estimate of the number of unique species in the community. This is often not directly observed for the whole community but only for a limited sample from the community. This has led to alternative richness indices that provide different ways to estimate the species richness.

Richness index differs from the concept of species diversity or evenness in that it ignores species abundance, and focuses on the binary presence/absence values that indicate simply whether the species was detected.

The function takes all index names in full lowercase. The user can provide the desired spelling through the argument `name` (see examples).

The following richness indices are provided.

- 'ace' Abundance-based coverage estimator (ACE) is another nonparametric richness index that uses sample coverage, defined based on the sum of the probabilities of the observed species. This method divides the species into abundant (more than 10 reads or observations) and rare groups in a sample and tends to underestimate the real number of species. The ACE index ignores the abundance information for the abundant species, based on the assumption that the abundant species are observed regardless of their exact abundance. We use here the bias-corrected version (O'Hara 2005, Chiu et al. 2014) implemented in `estimateR`. For an exact formulation, see `estimateR`. Note that this index comes with an additional column with standard error information.
- 'chao1' This is a nonparametric estimator of species richness. It assumes that rare species carry information about the (unknown) number of unobserved species. We use here the bias-corrected version (O'Hara 2005, Chiu et al. 2014) implemented in `estimateR`. This index implicitly assumes that every taxa has equal probability of being observed. Note that it gives a lower bound to species richness. The bias-corrected for an exact formulation, see `estimateR`. This estimator uses only the singleton and doubleton counts, and hence it gives more weight to the low abundance species. Note that this index comes with an additional column with standard error information.
- 'hill' Effective species richness aka Hill index (see e.g. Chao et al. 2016). Currently only the case 1D is implemented. This corresponds to the exponent of Shannon diversity. Intuitively, the effective richness indicates the number of species whose even distribution would lead to the same diversity than the observed community, where the species abundances are unevenly distributed.
- 'observed' The *observed richness* gives the number of species that is detected above a given detection threshold in the observed sample (default 0). This is conceptually the simplest richness index. The corresponding index in the `vegan` package is "richness".

## Value

x with additional `colData` named `*name*`

## Author(s)

Leo Lahti. Contact: [microbiome.github.io](https://microbiome.github.io)

## References

- Chao A. (1984) Non-parametric estimation of the number of classes in a population. *Scand J Stat.* 11:265–270.
- Chao A, Chun-Huo C, Jost L (2016). Phylogenetic Diversity Measures and Their Decomposition: A Framework Based on Hill Numbers. *Biodiversity Conservation and Phylogenetic Systematics*, Springer International Publishing, pp. 141–172, doi:10.1007/978-3-319-22461-9\_8.
- Chiu, C.H., Wang, Y.T., Walther, B.A. & Chao, A. (2014). Improved nonparametric lower bound of species richness via a modified Good-Turing frequency formula. *Biometrics* 70, 671-682.
- O’Hara, R.B. (2005). Species richness estimators: how many species can dance on the head of a pin? *J. Anim. Ecol.* 74, 375-386.

## See Also

[plotColData](#)

- [estimateR](#)

## Examples

```
data(esophagus)

# Calculates all richness indices by default
esophagus <- estimateRichness(esophagus)

# Shows all indices
colData(esophagus)

# Shows Hill index
colData(esophagus)$hill

# Deletes hill index
colData(esophagus)$hill <- NULL

# Shows all indices, hill is deleted
colData(esophagus)

# Delete the remaining indices
colData(esophagus)[, c("observed", "chao1", "ace")] <- NULL

# Calculates observed richness index and saves them with specific names
esophagus <- estimateRichness(esophagus,
  index = c("observed", "chao1", "ace", "hill"),
  name = c("Observed", "Chao1", "ACE", "Hill"))

# Show the new indices
colData(esophagus)

# Deletes all colData (including the indices)
colData(esophagus) <- NULL
```

```

# Calculate observed richness excluding singletons (detection limit 1)
esophagus <- estimateRichness(esophagus, index="observed", detection = 1)

# Deletes all colData (including the indices)
colData(esophagus) <- NULL

# Indices must be written correctly (all lowercase), otherwise an error
# gets thrown
## Not run: esophagus <- estimateRichness(esophagus, index="ace")

# Calculates Chao1 and ACE indices only
esophagus <- estimateRichness(esophagus, index=c("chao1", "ace"),
                             name=c("Chao1", "ACE"))

# Deletes all colData (including the indices)
colData(esophagus) <- NULL

# Names of columns can be chosen arbitrarily, but the length of arguments
# must match.
esophagus <- estimateRichness(esophagus,
                             index = c("ace", "chao1"),
                             name = c("index1", "index2"))

# Shows all indices
colData(esophagus)

```

---

getAbundance

*Get abundance values by "SampleID" or "FeatureID"*


---

## Description

These are functions for extracting abundances present in `assay(x)`. These functions are convenience wrapper around subsetting columns or rows from `assay(x, name)`.

## Usage

```

getAbundanceSample(
  x,
  sample_id,
  assay_name = abund_values,
  abund_values = "counts"
)

## S4 method for signature 'SummarizedExperiment'
getAbundanceSample(
  x,
  sample_id = NULL,
  assay_name = abund_values,
  abund_values = "counts"
)

```



```

)

getAbundanceFeature(
  x,
  feature_id,
  assay_name = abund_values,
  abund_values = "counts"
)

## S4 method for signature 'SummarizedExperiment'
getAbundanceFeature(
  x,
  feature_id = NULL,
  assay_name = abund_values,
  abund_values = "counts"
)

```

### Arguments

x	A <a href="#">SummarizedExperiment</a> object.
sample_id	A “SampleID” from which user wants to extract the abundances of “FeatureID”. This is essentially a column name in <code>assay(x)</code> .
assay_name	a character value to select an <a href="#">assayNames</a>
abund_values	a single character value for specifying which assay to use for calculation. (Please use <code>assay_name</code> instead. At some point <code>abund_values</code> will be disabled.)
feature_id	A “FeatureID” for which user wants to extract the abundances from all of “SampleID” in <a href="#">assayNames</a> . This is essentially a rowname in <code>assay(x)</code> .

### Details

`getAbundanceSample` returns abundance values for all “FeatureIDs” in a user specified “SampleID”.

`getAbundanceFeature` returns abundance values in all “SampleIDs” for user specified “FeatureID”.

### Value

`getAbundanceSample` and `getAbundanceFeature` return a numeric matrix of the abundance values for all “SampleIDs”/“FeatureIDs”

### Author(s)

Sudarshan A. Shetty

### Examples

```

# getAbundanceSample
data(GlobalPatterns)

```

```

getAbundanceSample(GlobalPatterns,
                    sample_id = 'CC1',
                    assay_name = 'counts')
# getAbundanceFeature
getAbundanceFeature(GlobalPatterns,
                    feature_id = '522457',
                    assay_name = 'counts')

```

---

```
getExperimentCrossAssociation
```

*Calculate correlations between features of two experiments.*

---

### **Description**

Calculate correlations between features of two experiments.

### **Usage**

```

getExperimentCrossAssociation(x, ...)

## S4 method for signature 'MultiAssayExperiment'
getExperimentCrossAssociation(
  x,
  experiment1 = 1,
  experiment2 = 2,
  assay_name1 = abund_values1,
  abund_values1 = "counts",
  assay_name2 = abund_values2,
  abund_values2 = "counts",
  altexp1 = NULL,
  altexp2 = NULL,
  colData_variable1 = NULL,
  colData_variable2 = NULL,
  MARGIN = 1,
  method = c("kendall", "spearman", "categorical", "pearson"),
  mode = "table",
  p_adj_method = c("fdr", "BH", "bonferroni", "BY", "hochberg", "holm", "hommel", "none"),
  p_adj_threshold = NULL,
  cor_threshold = NULL,
  sort = FALSE,
  filter_self_correlations = FALSE,
  verbose = TRUE,
  test_significance = FALSE,
  show_warnings = TRUE,
  paired = FALSE,
  ...
)

```

```

## S4 method for signature 'SummarizedExperiment'
getExperimentCrossAssociation(x, experiment2 = x, ...)

testExperimentCrossAssociation(x, ...)

## S4 method for signature 'ANY'
testExperimentCrossAssociation(x, ...)

testExperimentCrossCorrelation(x, ...)

## S4 method for signature 'ANY'
testExperimentCrossCorrelation(x, ...)

getExperimentCrossCorrelation(x, ...)

## S4 method for signature 'ANY'
getExperimentCrossCorrelation(x, ...)

```

## Arguments

x	A <a href="#">MultiAssayExperiment</a> or <a href="#">SummarizedExperiment</a> object.
...	Additional arguments: <ul style="list-style-type: none"> <li>• symmetric A single boolean value for specifying if measure is symmetric or not. When symmetric = TRUE, associations are calculated only for unique variable-pairs, and they are assigned to corresponding variable-pair. This decreases the number of calculations in 2-fold meaning faster execution. (By default: symmetric = FALSE)</li> <li>• association_FUN A function that is used to calculate (dis-)similarity between features. Function must take matrix as an input and give numeric values as an output. Adjust method and other parameters correspondingly. Supported functions are, for example, stats::dist and vegan::vegdist.</li> </ul>
experiment1	A single character or numeric value for selecting the experiment 1 from experiments(x) of MultiassayExperiment object. (By default: experiment1 = 1)
experiment2	A single character or numeric value for selecting the experiment 2 from experiments(x) of MultiAssayExperiment object or altExp(x) of TreeSummarizedExperiment object. Alternatively, experiment2 can also be TreeSE object when x is TreeSE object. (By default: experiment2 = 2 when x is MAE and experiment2 = x when x is TreeSE)
assay_name1	A single character value for selecting the <a href="#">assay</a> of experiment 1 to be transformed. (By default: assay_name1 = "counts")
abund_values1	a single character value for specifying which assay of experiment 1 to use for calculation. (Please use assay_name1 instead. At some point abund_values1 will be disabled.)
assay_name2	A single character value for selecting the <a href="#">assay</a> of experiment 2 to be transformed. (By default: assay_name2 = "counts")

<code>abund_values2</code>	a single character value for specifying which assay of experiment 2 to use for calculation. (Please use <code>assay_name2</code> instead. At some point <code>abund_values2</code> will be disabled.)
<code>altexp1</code>	A single numeric or character value specifying alternative experiment from the <code>altExp</code> of experiment 1. If NULL, then the experiment is itself and <code>altExp</code> option is disabled. (By default: <code>altexp1 = NULL</code> )
<code>altexp2</code>	A single numeric or character value specifying alternative experiment from the <code>altExp</code> of experiment 2. If NULL, then the experiment is itself and <code>altExp</code> option is disabled. (By default: <code>altexp2 = NULL</code> )
<code>colData_variable1</code>	A character value specifying column(s) from <code>colData</code> of experiment 1. If <code>colData_variable1</code> is used, <code>assay_name1</code> is disabled. (By default: <code>colData_variable1 = NULL</code> )
<code>colData_variable2</code>	A character value specifying column(s) from <code>colData</code> of experiment 2. If <code>colData_variable2</code> is used, <code>assay_name2</code> is disabled. (By default: <code>colData_variable2 = NULL</code> )
<code>MARGIN</code>	A single numeric value for selecting if association are calculated row-wise / for features (1) or column-wise / for samples (2). Must be 1 or 2. (By default: <code>MARGIN = 1</code> )
<code>method</code>	A single character value for selecting association method ('kendall', 'pearson', or 'spearman' for continuous/numeric; 'categorical' for discrete) (By default: <code>method = "kendall"</code> )
<code>mode</code>	A single character value for selecting output format Available formats are 'table' and 'matrix'. (By default: <code>mode = "table"</code> )
<code>p_adj_method</code>	A single character value for selecting adjustment method of p-values. Passed to <code>p.adjust</code> function. (By default: <code>p_adj_method = "fdr"</code> )
<code>p_adj_threshold</code>	A single numeric value (from 0 to 1) for selecting adjusted p-value threshold for filtering. (By default: <code>p_adj_threshold = NULL</code> )
<code>cor_threshold</code>	A single numeric absolute value (from 0 to 1) for selecting correlation threshold for filtering. (By default: <code>cor_threshold = NULL</code> )
<code>sort</code>	A single boolean value for selecting whether to sort features or not in result matrices. Used method is hierarchical clustering. (By default: <code>sort = FALSE</code> )
<code>filter_self_correlations</code>	A single boolean value for selecting whether to filter out correlations between identical items. Applies only when correlation between experiment itself is tested, i.e., when assays are identical. (By default: <code>filter_self_correlations = FALSE</code> )
<code>verbose</code>	A single boolean value for selecting whether to get messages about progress of calculation.
<code>test_significance</code>	A single boolean value for selecting whether to test statistical significance of associations.



```

# Show first 5 entries
head(result, 5)

# testExperimentCorrelation returns also significances
# filter_self_correlations = TRUE filters self correlations
# With p_adj_threshold it is possible to filter those features that do not have
# any correlations that have p-value under threshold
result <- testExperimentCrossAssociation(mae[[1]], experiment2 = mae[[1]], method = "pearson",
                                       filter_self_correlations = TRUE,
                                       p_adj_threshold = 0.05)

# Show first 5 entries
head(result, 5)

# Also getExperimentCrossAssociation returns significances when
# test_significance = TRUE
# Warnings can be suppressed by using show_warnings = FALSE
result <- getExperimentCrossAssociation(mae[[1]], experiment2 = mae[[2]], method = "pearson",
                                       assay_name2 = "nmr",
                                       mode = "matrix", test_significance = TRUE,
                                       show_warnings = FALSE)

# Returned value is a list of matrices
names(result)

# Calculate Bray-Curtis dissimilarity between samples. If dataset includes
# paired samples, you can use paired = TRUE.
result <- getExperimentCrossAssociation(mae[[1]], mae[[1]], MARGIN = 2, paired = FALSE,
                                       association_FUN = vegan::vegdist, method = "bray")

# If experiments are equal and measure is symmetric (e.g., taxa1 vs taxa2 == taxa2 vs taxa1),
# it is possible to speed-up calculations by calculating association only for unique
# variable-pairs. Use "symmetric" to choose whether to measure association for only
# other half of of variable-pairs.
result <- getExperimentCrossAssociation(mae, experiment1 = "microbiota", experiment2 = "microbiota",
                                       assay_name1 = "counts", assay_name2 = "counts",
                                       symmetric = TRUE)

# For big data sets, calculation might take long. To make calculations quicker, you can take
# a random sample from data. In a complex biological problems, random sample
# can describe the data enough. Here our random sample is 30 % of whole data.
sample_size <- 0.3
tse <- mae[[1]]
tse_sub <- tse[ sample( seq_len( nrow(tse) ), sample_size * nrow(tse) ), ]
result <- testExperimentCrossAssociation(tse_sub)

# It is also possible to choose variables from colData and calculate association
# between assay and sample metadata or between variables of sample metadata
mae[[1]] <- estimateDiversity(mae[[1]])
# colData_variable works similarly to assay_name. Instead of fetching an assay
# named assay_name from assay slot, it fetches a column named colData_variable
# from colData.
result <- getExperimentCrossAssociation(mae[[1]], assay_name1 = "counts",

```

```
colData_variable2 = c("shannon", "coverage"))
```

---

```
getPrevalence           Calculation prevalence information for features across samples
```

---

### Description

These functions calculate the population prevalence for taxonomic ranks in a [SummarizedExperiment-class](#) object.

### Usage

```
getPrevalence(x, ...)

## S4 method for signature 'ANY'
getPrevalence(x, detection = 0, include_lowest = FALSE, sort = FALSE, ...)

## S4 method for signature 'SummarizedExperiment'
getPrevalence(
  x,
  assay_name = abund_values,
  abund_values = "counts",
  as_relative = TRUE,
  rank = NULL,
  ...
)

getPrevalentTaxa(x, ...)

## S4 method for signature 'ANY'
getPrevalentTaxa(x, prevalence = 50/100, include_lowest = FALSE, ...)

## S4 method for signature 'SummarizedExperiment'
getPrevalentTaxa(
  x,
  rank = NULL,
  prevalence = 50/100,
  include_lowest = FALSE,
  ...
)

getPrevalentFeatures(x, ...)

## S4 method for signature 'ANY'
getPrevalentFeatures(x, ...)
```

```
getRareTaxa(x, ...)  
  
## S4 method for signature 'ANY'  
getRareTaxa(x, prevalence = 50/100, include_lowest = FALSE, ...)  
  
## S4 method for signature 'SummarizedExperiment'  
getRareTaxa(x, rank = NULL, prevalence = 50/100, include_lowest = FALSE, ...)  
  
getRareFeatures(x, ...)  
  
## S4 method for signature 'ANY'  
getRareFeatures(x, ...)  
  
subsetByPrevalentTaxa(x, ...)  
  
## S4 method for signature 'SummarizedExperiment'  
subsetByPrevalentTaxa(x, rank = NULL, ...)  
  
subsetByPrevalentFeatures(x, ...)  
  
## S4 method for signature 'ANY'  
subsetByPrevalentFeatures(x, ...)  
  
subsetByRareTaxa(x, ...)  
  
## S4 method for signature 'SummarizedExperiment'  
subsetByRareTaxa(x, rank = NULL, ...)  
  
subsetByRareFeatures(x, ...)  
  
## S4 method for signature 'ANY'  
subsetByRareFeatures(x, ...)  
  
getPrevalentAbundance(  
  x,  
  assay_name = abund_values,  
  abund_values = "relabundance",  
  ...  
)  
  
## S4 method for signature 'ANY'  
getPrevalentAbundance(  
  x,  
  assay_name = abund_values,  
  abund_values = "relabundance",  
  ...  
)
```



```

## S4 method for signature 'SummarizedExperiment'
getPrevalentAbundance(
  x,
  assay_name = abund_values,
  abund_values = "counts",
  ...
)

agglomerateByPrevalence(x, ...)

## S4 method for signature 'SummarizedExperiment'
agglomerateByPrevalence(
  x,
  rank = taxonomyRanks(x)[1L],
  other_label = "Other",
  ...
)

```

### Arguments

x	a <a href="#">SummarizedExperiment</a> object
detection	Detection threshold for absence/presence. Either an absolute value compared directly to the values of x or a relative value between 0 and 1, if <code>as_relative = TRUE</code> .
include_lowest	logical scalar: Should the lower boundary of the detection and prevalence cutoffs be included? (default: FALSE)
sort	logical scalar: Should the result be sorted by prevalence? (default: FALSE)
assay_name	A single character value for selecting the <a href="#">assay</a> to use for prevalence calculation.
abund_values	a single character value for specifying which assay to use for calculation. (Please use <code>assay_name</code> instead. At some point <code>abund_values</code> will be disabled.)
as_relative	logical scalar: Should the detection threshold be applied on compositional (relative) abundances? (default: TRUE)
rank, ...	additional arguments <ul style="list-style-type: none"> <li>• If <code>!is.null(rank)</code> arguments are passed on to <a href="#">agglomerateByRank</a>. See <a href="#">?agglomerateByRank</a> for more details.</li> <li>• for <code>getPrevalentTaxa</code>, <code>getRareTaxa</code>, <code>subsetByPrevalentTaxa</code> and <code>subsetByRareTaxa</code> additional parameters passed to <code>getPrevalence</code></li> <li>• for <code>getPrevalentAbundance</code> additional parameters passed to <code>getPrevalentTaxa</code></li> </ul>
prevalence	Prevalence threshold (in 0 to 1). The required prevalence is strictly greater by default. To include the limit, set <code>include_lowest</code> to TRUE.
other_label	A single character valued used as the label for the summary of non-prevalent taxa. (default: <code>other_label = "Other"</code> )

## Details

getPrevalence calculates the relative frequency of samples that exceed the detection threshold. For SummarizedExperiment objects, the prevalence is calculated for the selected taxonomic rank, otherwise for the rows. The absolute population prevalence can be obtained by multiplying the prevalence by the number of samples (ncol(x)). If as\_relative = TRUE the relative frequency (between 0 and 1) is used to check against the detection threshold.

The core abundance index from getPrevalentAbundance gives the relative proportion of the core species (in between 0 and 1). The core taxa are defined as those that exceed the given population prevalence threshold at the given detection level as set for getPrevalentTaxa.

subsetPrevalentTaxa and subsetRareTaxa return a subset of x. The subset includes the most prevalent or rare taxa that are calculated with getPrevalentTaxa or getRareTaxa respectively.

getPrevalentTaxa returns taxa that are more prevalent with the given detection threshold for the selected taxonomic rank.

getRareTaxa returns complement of getPrevalentTaxa.

## Value

subsetPrevalentTaxa and subsetRareTaxa return subset of x.

All other functions return a named vectors:

- getPrevalence returns a numeric vector with the names being set to either the row names of x or the names after agglomeration.
- getPrevalentAbundance returns a numeric vector with the names corresponding to the column name of x and include the joint abundance of prevalent taxa.
- getPrevalentTaxa and getRareTaxa return a character vector with only the names exceeding the threshold set by prevalence, if the rownames of x is set. Otherwise an integer vector is returned matching the rows in x.

## Author(s)

Leo Lahti For getPrevalentAbundance: Leo Lahti and Tuomas Borman. Contact: [microbiome.github.io](https://microbiome.github.io)

## References

A Salonen et al. The adult intestinal core microbiota is determined by analysis depth and health status. Clinical Microbiology and Infection 18(S4):16 20, 2012. To cite the R package, see citation('mia')

## See Also

[agglomerateByRank](#), [getTopTaxa](#)

**Examples**

```
data(GlobalPatterns)
tse <- GlobalPatterns
# Get prevalence estimates for individual ASV/OTU
prevalence.frequency <- getPrevalence(tse,
                                     detection = 0,
                                     sort = TRUE,
                                     as_relative = TRUE)

head(prevalence.frequency)

# Get prevalence estimates for phylums
# - the getPrevalence function itself always returns population frequencies
prevalence.frequency <- getPrevalence(tse,
                                     rank = "Phylum",
                                     detection = 0,
                                     sort = TRUE,
                                     as_relative = TRUE)

head(prevalence.frequency)

# - to obtain population counts, multiply frequencies with the sample size,
# which answers the question "In how many samples is this phylum detectable"
prevalence.count <- prevalence.frequency * ncol(tse)
head(prevalence.count)

# Detection threshold 1 (strictly greater by default);
# Note that the data (GlobalPatterns) is here in absolute counts
# (and not compositional, relative abundances)
# Prevalence threshold 50 percent (strictly greater by default)
prevalent <- getPrevalentTaxa(tse,
                             rank = "Phylum",
                             detection = 10,
                             prevalence = 50/100,
                             as_relative = FALSE)

head(prevalent)

# Gets a subset of object that includes prevalent taxa
altExp(tse, "prevalent") <- subsetByPrevalentTaxa(tse,
                                                  rank = "Family",
                                                  detection = 0.001,
                                                  prevalence = 0.55,
                                                  as_relative = TRUE)

altExp(tse, "prevalent")

# getRareTaxa returns the inverse
rare <- getRareTaxa(tse,
                  rank = "Phylum",
                  detection = 1/100,
                  prevalence = 50/100,
                  as_relative = TRUE)

head(rare)

# Gets a subset of object that includes rare taxa
```

```

altExp(tse, "rare") <- subsetByRareTaxa(tse,
                                     rank = "Class",
                                     detection = 0.001,
                                     prevalence = 0.001,
                                     as_relative = TRUE)

altExp(tse, "rare")

# Names of both experiments, prevalent and rare, can be found from slot altExpNames
tse

data(esophagus)
getPrevalentAbundance(esophagus, assay_name = "counts")

# data can be aggregated based on prevalent taxonomic results
agglomerateByPrevalence(tse,
                        rank = "Phylum",
                        detection = 1/100,
                        prevalence = 50/100,
                        as_relative = TRUE)

```

---

isContaminant

*decontam functions*


---

## Description

The decontam functions `isContaminant` and `isNotContaminant` are made available for [SummarizedExperiment](#) objects.

## Usage

```

## S4 method for signature 'SummarizedExperiment'
isContaminant(
  seqtab,
  assay_name = abund_values,
  abund_values = "counts",
  name = "isContaminant",
  concentration = NULL,
  control = NULL,
  batch = NULL,
  threshold = 0.1,
  normalize = TRUE,
  detailed = TRUE,
  ...
)

## S4 method for signature 'SummarizedExperiment'
isNotContaminant(
  seqtab,

```

```

    assay_name = abund_values,
    abund_values = "counts",
    name = "isNotContaminant",
    control = NULL,
    threshold = 0.5,
    normalize = TRUE,
    detailed = FALSE,
    ...
)

addContaminantQC(x, name = "isContaminant", ...)

## S4 method for signature 'SummarizedExperiment'
addContaminantQC(x, name = "isContaminant", ...)

addNotContaminantQC(x, name = "isNotContaminant", ...)

## S4 method for signature 'SummarizedExperiment'
addNotContaminantQC(x, name = "isNotContaminant", ...)

```

## Arguments

seqtab, x	a <a href="#">SummarizedExperiment</a>
assay_name	A single character value for selecting the <a href="#">assay</a> to use.
abund_values	a single character value for specifying which assay to use for calculation. (Please use <code>assay_name</code> instead. At some point <code>abund_values</code> will be disabled.)
name	A name for the column of the <code>colData</code> in which the contaminant information should be stored.
concentration	NULL or a single character value. Defining a column with numeric values from the <code>colData</code> to use as concentration information. (default: <code>concentration = NULL</code> )
control	NULL or a single character value. Defining a column with logical values from the <code>colData</code> to define control and non-control samples. (default: <code>control = NULL</code> )
batch	NULL or a single character value. Defining a column with values interpretable as a factor from the <code>colData</code> to use as batch information. (default: <code>batch = NULL</code> )
threshold	numeric scalar. See <a href="#">decontam:isContaminant</a> or <a href="#">decontam:isNotContaminant</a>
normalize, detailed	logical scalar. See <a href="#">decontam:isContaminant</a> or <a href="#">decontam:isNotContaminant</a>
...	<ul style="list-style-type: none"> <li>for <code>isContaminant/ isNotContaminant</code>: arguments passed on to <a href="#">decontam:isContaminant</a> or <a href="#">decontam:isNotContaminant</a></li> <li>for <code>addContaminantQC/addNotContaminantQC</code>: arguments passed on to <code>isContaminant/ isNotContaminant</code></li> </ul>

**Value**

for isContaminant/ isNotContaminant a DataFrame or for addContaminantQC/addNotContaminantQC a modified object of class(x)

**See Also**

[decontam:isContaminant](#), [decontam:isNotContaminant](#)

**Examples**

```
data(esophagus)
# setup of some mock data
colData(esophagus)$concentration <- c(1,2,3)
colData(esophagus)$control <- c(FALSE,FALSE,TRUE)

isContaminant(esophagus,
              method = "frequency",
              concentration = "concentration")
esophagus <- addContaminantQC(esophagus,
                             method = "frequency",
                             concentration = "concentration")

colData(esophagus)

isNotContaminant(esophagus, control = "control")
esophagus <- addNotContaminantQC(esophagus, control = "control")
colData(esophagus)
```

---

loadFromMetaphlan	<i>Import Metaphlan results to</i> TreeSummarizedExperiment
-------------------	---

---

**Description**

Import Metaphlan results to TreeSummarizedExperiment

**Arguments**

file	a single character value defining the file path of the Metaphlan file. The file must be in merged Metaphlan format.
sample_meta	a single character value defining the file path of the sample metadata file. The file must be in tsv format (default: sample_meta = NULL).
phy_tree	a single character value defining the file path of the phylogenetic tree. (default: phy_tree = NULL).
...	additional arguments: <ul style="list-style-type: none"> <li>• assay_name: A single character value for naming <a href="#">assay</a> (default: assay_name = "counts")</li> </ul>

- `abund_values`: A single character value for specifying which assay to use for calculation. (Please use `assay_name` instead. At some point `abund_values` will be disabled.)
- `removeTaxaPrefixes`: TRUE or FALSE: Should taxonomic prefixes be removed? (default: `removeTaxaPrefixes = FALSE`)

## Details

Import Metaphlan results. Input must be in merged Metaphlan format. Data is imported so that data at the lowest rank is imported as a `TreeSummarizedExperiment` object. Data at higher rank is imported as a `SummarizedExperiment` objects which are stored to `altExp` of `TreeSummarizedExperiment` object.

## Value

A `TreeSummarizedExperiment` object

## Author(s)

Leo Lahti and Tuomas Borman. Contact: [microbiome.github.io](https://microbiome.github.io)

## References

Beghini F, McIver LJ, Blanco-Míguez A, Dubois L, Asnicar F, Maharjan S, Mailyan A, Manghi P, Scholz M, Thomas AM, Valles-Colomer M, Weingart G, Zhang Y, Zolfo M, Huttenhower C, Franzosa EA, & Segata N (2021) Integrating taxonomic, functional, and strain-level profiling of diverse microbial communities with bioBakery 3. *Elife* 10:e65088. doi: 10.7554/eLife.65088

## See Also

[makeTreeSEFromPhyloseq](#) [makeTreeSEFromBiom](#) [makeTreeSEFromDADA2](#) [loadFromQIIME2](#) [loadFromMothur](#)

## Examples

```
# (Data is from tutorial
# https://github.com/biobakery/biobakery/wiki/metaphlan3#merge-outputs)

# File path
file_path <- system.file("extdata", "merged_abundance_table.txt", package = "mia")
# Import data
tse <- loadFromMetaphlan(file_path)
# Data at the lowest rank
tse
# Data at higher rank is stored in altExp
altExps(tse)
# Higher rank data is in SE format, for example, Phylum rank
altExp(tse, "Phylum")
```

---

loadFromMothur	<i>Import Mothur results as a TreeSummarizedExperiment</i>
----------------	--

---

### Description

This method creates a `TreeSummarizedExperiment` object from Mothur files provided as input.

### Usage

```
loadFromMothur(sharedFile, taxonomyFile = NULL, designFile = NULL)
```

### Arguments

<code>sharedFile</code>	a single character value defining the file path of the feature table to be imported. The File has to be in shared file format as defined in Mothur documentation.
<code>taxonomyFile</code>	a single character value defining the file path of the taxonomy table to be imported. The File has to be in taxonomy file or constaxonomy file format as defined in Mothur documentation. (default: <code>taxonomyFile = NULL</code> ).
<code>designFile</code>	a single character value defining the file path of the sample metadata to be imported. The File has to be in desing file format as defined in Mothur documentation. (default: <code>designFile = NULL</code> ).

### Details

Results exported from Mothur can be imported as a `SummarizedExperiment` using `loadFromMothur`. Except for the `sharedFile`, the other data types, `taxonomyFile`, and `designFile`, are optional, but are highly encouraged to be provided.

### Value

A `TreeSummarizedExperiment` object

### Author(s)

Leo Lahti and Tuomas Borman. Contact: [microbiome.github.io](https://microbiome.github.io)

### References

<https://mothur.org/> [https://mothur.org/wiki/shared\\_file/](https://mothur.org/wiki/shared_file/) [https://mothur.org/wiki/taxonomy\\_file/](https://mothur.org/wiki/taxonomy_file/) [https://mothur.org/wiki/constaxonomy\\_file/](https://mothur.org/wiki/constaxonomy_file/) [https://mothur.org/wiki/design\\_file/](https://mothur.org/wiki/design_file/)

### See Also

[makeTreeSEFromPhyloseq](#) [makeTreeSEFromBiom](#) [makeTreeSEFromDADA2](#) [loadFromQIIME2](#)



**Examples**

```

# Abundance table
counts <- system.file("extdata", "mothur_example.shared", package = "mia")
# Taxa table (in "cons.taxonomy" or "taxonomy" format)
taxa <- system.file("extdata", "mothur_example.cons.taxonomy", package = "mia")
#taxa <- system.file("extdata", "mothur_example.taxonomy", package = "mia")
# Sample meta data
meta <- system.file("extdata", "mothur_example.design", package = "mia")

# Creates se object from files
se <- loadFromMothur(counts, taxa, meta)
# Convert SE to TreeSE
tse <- as(se, "TreeSummarizedExperiment")
tse

```

---

loadFromQIIME2	<i>Import QIIME2 results to TreeSummarizedExperiment</i>
----------------	--

---

**Description**

Results exported from QIIME2 can be imported as a `TreeSummarizedExperiment` using `loadFromQIIME2`. Except for the `featureTableFile`, the other data types, `taxonomyTableFile`, `refSeqFile` and `phyTreeFile`, are optional, but are highly encouraged to be provided.

Import the QIIME2 artifacts to R.

**Usage**

```

loadFromQIIME2(
  featureTableFile,
  taxonomyTableFile = NULL,
  sampleMetaFile = NULL,
  featureNamesAsRefSeq = TRUE,
  refSeqFile = NULL,
  phyTreeFile = NULL,
  ...
)

readQZA(file, temp = tempdir(), ...)

```

**Arguments**

`featureTableFile`  
a single character value defining the file path of the feature table to be imported.

`taxonomyTableFile`  
a single character value defining the file path of the taxonomy table to be imported. (default: `taxonomyTableFile = NULL`).

sampleMetaFile	a single character value defining the file path of the sample metadata to be imported. The file has to be in tsv format. (default: sampleMetaFile = NULL).
featureNamesAsRefSeq	TRUE or FALSE: Should the feature names of the feature table be regarded as reference sequences? This setting will be disregarded, if refSeqFile is not NULL. If the feature names do not contain valid DNA characters only, the reference sequences will not be set.
refSeqFile	a single character value defining the file path of the reference sequences for each feature. (default: refSeqFile = NULL).
phyTreeFile	a single character value defining the file path of the phylogenetic tree. (default: phyTreeFile = NULL).
...	additional arguments: <ul style="list-style-type: none"> <li>temp: the temporary directory used for decompressing the data. (default: tempdir())</li> <li>removeTaxaPrefixes: TRUE or FALSE: Should taxonomic prefixes be removed? (default: removeTaxaPrefixes = FALSE)</li> </ul>
file	character, path of the input qza file. Only files in format of BIOMV210DirFmt (feature table), TSVTaxonomyDirectoryFormat (taxonomic table), NewickDirectoryFormat (phylogenetic tree ) and DNASequencesDirectoryFormat (representative sequences) are supported right now.
temp	character, a temporary directory in which the qza file will be decompressed to, default tempdir().

### Details

Both arguments featureNamesAsRefSeq and refSeqFile can be used to define reference sequences of features. featureNamesAsRefSeq is only taken into account, if refSeqFile is NULL. No reference sequences are tried to be created, if featureNameAsRefSeq is FALSE and refSeqFile is NULL.

### Value

A [TreeSummarizedExperiment](#) object

matrix object for feature table, `DataFrame` for taxonomic table, `ape::phylo` object for phylogenetic tree, `Biostrings::DNAStrngSet` for representative sequences of taxa.

### Author(s)

Yang Cao

### References

Bolyen E et al. 2019: Reproducible, interactive, scalable and extensible microbiome data science using QIIME 2. *Nature Biotechnology* 37: 852–857. <https://doi.org/10.1038/s41587-019-0209-9>  
<https://qiime2.org>

**See Also**

[makeTreeSEFromPhyloseq](#) [makeTreeSEFromBiom](#) [makeTreeSEFromDADA2](#) [loadFromMothur](#)

**Examples**

```
featureTableFile <- system.file("extdata", "table.qza", package = "mia")
taxonomyTableFile <- system.file("extdata", "taxonomy.qza", package = "mia")
sampleMetaFile <- system.file("extdata", "sample-metadata.tsv", package = "mia")
phyTreeFile <- system.file("extdata", "tree.qza", package = "mia")
refSeqFile <- system.file("extdata", "refseq.qza", package = "mia")
tse <- loadFromQIIME2(
  featureTableFile = featureTableFile,
  taxonomyTableFile = taxonomyTableFile,
  sampleMetaFile = sampleMetaFile,
  refSeqFile = refSeqFile,
  phyTreeFile = phyTreeFile
)

tse
# Read individual files
featureTableFile <- system.file("extdata", "table.qza", package = "mia")
taxonomyTableFile <- system.file("extdata", "taxonomy.qza", package = "mia")
sampleMetaFile <- system.file("extdata", "sample-metadata.tsv", package = "mia")

assay <- readQZA(featureTableFile)
rowdata <- readQZA(taxonomyTableFile, removeTaxaPrefixes = TRUE)
coldata <- read.table(sampleMetaFile, header = TRUE, sep = "\t", comment.char = "")

# Assign rownames
rownames(coldata) <- coldata[, 1]
coldata[, 1] <- NULL

# Order coldata based on assay
coldata <- coldata[match(colnames(assay), rownames(coldata)), ]

# Create SE from individual files
se <- SummarizedExperiment(assays = list(assay), rowData = rowdata, colData = coldata)
se
```

---

makePhyloseqFromTreeSE

*Create a phyloseq object from a TreeSummarizedExperiment object*

---

**Description**

This function creates a phyloseq object from a TreeSummarizedExperiment object. By using assay\_name, it is possible to specify which table from assay is added to the phyloseq object.

## Usage

```
makePhyloseqFromTreeSE(x, ...)  
  
## S4 method for signature 'SummarizedExperiment'  
makePhyloseqFromTreeSE(x, assay_name = abund_values, abund_values = "counts")  
  
## S4 method for signature 'TreeSummarizedExperiment'  
makePhyloseqFromTreeSE(x, ...)  
  
makePhyloseqFromTreeSummarizedExperiment(x, ...)  
  
## S4 method for signature 'ANY'  
makePhyloseqFromTreeSummarizedExperiment(x, ...)
```

## Arguments

x	a TreeSummarizedExperiment object
...	additional arguments
assay_name	A single character value for selecting the <a href="#">assay</a> to be included in the phyloseq object that is created. By default, it is counts table.
abund_values	a single character value for specifying which assay to use for calculation. (Please use assay_name instead. At some point abund_values will be disabled.)

## Details

makePhyloseqFromTreeSE is used for creating a phyloseq object from TreeSummarizedExperiment object.

## Value

An object of class Phyloseq object.

## Author(s)

Leo Lahti and Tuomas Borman. Contact: [microbiome.github.io](https://microbiome.github.io)

## Examples

```
# Get tse object  
data(GlobalPatterns)  
tse <- GlobalPatterns  
  
# Create a phyloseq object from it  
phy <- makePhyloseqFromTreeSE(tse)  
phy  
  
# By default the chosen table is counts, but if there are other tables,  
# they can be chosen with assay_name.
```

```
# Counts relative abundances table
tse <- transformCounts(tse, method = "relabundance")
phy2 <- makePhyloseqFromTreeSE(tse, assay_name = "relabundance")
phy2
```

---

makeTreeSEFromBiom      *Loading a biom file*

---

## Description

For convenience a few functions are available to convert data from a 'biom' file or object into a [TreeSummarizedExperiment](#)

## Usage

```
loadFromBiom(file, removeTaxaPrefixes = FALSE)

makeTreeSEFromBiom(obj, removeTaxaPrefixes = FALSE, ...)

makeTreeSummarizedExperimentFromBiom(obj, ...)
```

## Arguments

file	biom file location
removeTaxaPrefixes	TRUE or FALSE: Should taxonomic prefixes be removed? (default removeTaxaPrefixes = FALSE)
obj	object of type <a href="#">biom</a>
...	optional arguments (not used).

## Value

An object of class [TreeSummarizedExperiment](#)

## See Also

[makeTreeSEFromPhyloseq](#) [makeTreeSEFromDADA2](#) [loadFromQIIME2](#) [loadFromMothur](#)

## Examples

```
if(requireNamespace("biomformat")) {
  library(biomformat)
  # load from file
  rich_dense_file = system.file("extdata", "rich_dense_otu_table.biom",
                                package = "biomformat")
  se <- loadFromBiom(rich_dense_file, removeTaxaPrefixes = TRUE)
```

```

# load from object
x1 <- biomformat::read_biom(rich_dense_file)
se <- makeTreeSEFromBiom(x1)
# Convert SE to TreeSE
tse <- as(se, "TreeSummarizedExperiment")
tse
}

```

---

makeTreeSEFromDADA2    *Coerce 'DADA2' results to TreeSummarizedExperiment*

---

## Description

makeTreeSEFromDADA2 is a wrapper for the mergePairs function from the dada2 package.

## Usage

```

makeTreeSEFromDADA2(...)

makeTreeSummarizedExperimentFromDADA2(...)

```

## Arguments

...                    See mergePairs function for more details.

## Details

A count matrix is constructed via makeSequenceTable(mergePairs(...)) and rownames are dynamically created as ASV(N) with N from 1 to nrow of the count tables. The colnames and rownames from the output of makeSequenceTable are stored as colnames and in the referenceSeq slot of the TreeSummarizedExperiment, respectively.

## Value

An object of class TreeSummarizedExperiment

## See Also

[makeTreeSEFromPhyloseq](#) [makeTreeSEFromBiom](#) [loadFromQIIME2](#) [loadFromMothur](#)

## Examples

```

if(requireNamespace("dada2")) {
  fnF <- system.file("extdata", "sam1F.fastq.gz", package="dada2")
  fnR = system.file("extdata", "sam1R.fastq.gz", package="dada2")
  dadaF <- dada2::dada(fnF, selfConsist=TRUE)
  dadaR <- dada2::dada(fnR, selfConsist=TRUE)

  tse <- makeTreeSEFromDADA2(dadaF, fnF, dadaR, fnR)
  tse
}

```

---

`makeTreeSEFromPhyloseq`*Coerce a phyloseq object to a TreeSummarizedExperiment*

---

**Description**

`makeTreeSEFromPhyloseq` converts phyloseq objects into `TreeSummarizedExperiment` objects.

**Usage**

```
makeTreeSEFromPhyloseq(obj)
```

```
makeTreeSummarizedExperimentFromPhyloseq(obj)
```

```
## S4 method for signature 'ANY'
```

```
makeTreeSummarizedExperimentFromPhyloseq(obj)
```

**Arguments**

`obj` a phyloseq object

**Details**

All data stored in a phyloseq object is transferred.

**Value**

An object of class `TreeSummarizedExperiment`

**See Also**

[makeTreeSEFromBiom](#) [makeTreeSEFromDADA2](#) [loadFromQIIME2](#) [loadFromMothur](#)

**Examples**

```
if (requireNamespace("phyloseq")) {  
  data(GlobalPatterns, package="phyloseq")  
  makeTreeSEFromPhyloseq(GlobalPatterns)  
  data(enterotype, package="phyloseq")  
  makeTreeSEFromPhyloseq(enterotype)  
  data(esophagus, package="phyloseq")  
  makeTreeSEFromPhyloseq(esophagus)  
}
```

---

meltAssay

---

*Converting a [SummarizedExperiment](#) object into a long data.frame*


---

## Description

meltAssay Converts a [SummarizedExperiment](#) object into a long data.frame which can be used for tidyverse-tools.

## Usage

```
meltAssay(
  x,
  assay_name = abund_values,
  abund_values = "counts",
  add_row_data = NULL,
  add_col_data = NULL,
  feature_name = "FeatureID",
  sample_name = "SampleID",
  ...
)

## S4 method for signature 'SummarizedExperiment'
meltAssay(
  x,
  assay_name = abund_values,
  abund_values = "counts",
  add_row_data = NULL,
  add_col_data = NULL,
  feature_name = "FeatureID",
  sample_name = "SampleID",
  ...
)
```

## Arguments

x	A numeric matrix or a <a href="#">SummarizedExperiment</a>
assay_name	a character value to select an <a href="#">assayNames</a>
abund_values	a single character value for specifying which assay to use for calculation. (Please use assay_name instead. At some point abund_values will be disabled.)
add_row_data	NULL, TRUE or a character vector to select information from the rowData to add to the molten assay data. If add_row_data = NULL no data will be added, if add_row_data = TRUE all data will be added and if add_row_data is a character vector, it will be used to subset to given column names in rowData. (default: add_row_data = NULL)



add_col_data	NULL, TRUE or a character vector to select information from the colData to add to the molten assay data. If add_col_data = NULL no data will be added, if add_col_data = TRUE all data will be added and if add_col_data is a character vector, it will be used to subset to given column names in colData. (default: add_col_data = NULL)
feature_name	a character scalar to use as the output's name for the feature identifier. (default: feature_name = "FeatureID")
sample_name	a character scalar to use as the output's name for the sample identifier. (default: sample_name = "SampleID")
...	optional arguments: <ul style="list-style-type: none"> <li>• check_names A boolean value passed to data.frame function's check.name argument. Determines if sample names are checked that they are syntactically valid variable names and are not duplicated. If they are not, sample names are modified. (default: check_names = TRUE)</li> </ul>

### Details

If the colData contains a column "SampleID" or the rowData contains a column "FeatureID", they will be renamed to "SampleID\_col" and "FeatureID\_row", if row names or column names are set.

### Value

A tibble with the molten data. The assay values are given in a column named like the selected assay assay\_name. In addition, a column "FeatureID" will contain the rownames, if set, and analogously a column "SampleID" with the colnames, if set

### Author(s)

Sudarshan A. Shetty

### Examples

```
data(GlobalPatterns)
molten_tse <- meltAssay(GlobalPatterns,
                        assay_name = "counts",
                        add_row_data = TRUE,
                        add_col_data = TRUE
                        )

molten_tse
```

---

merge-methods

*Merge a subset of the rows or columns of a SummarizedExperiment*

---

### Description

mergeRows/mergeCols merge data on rows or columns of a SummarizedExperiment as defined by a factor alongside the chosen dimension. Metadata from the rowData or colData are retained as defined by archetype.

**Usage**

```

mergeRows(x, f, archetype = 1L, ...)

mergeCols(x, f, archetype = 1L, ...)

## S4 method for signature 'SummarizedExperiment'
mergeRows(x, f, archetype = 1L, ...)

## S4 method for signature 'SummarizedExperiment'
mergeCols(x, f, archetype = 1L, ...)

## S4 method for signature 'TreeSummarizedExperiment'
mergeRows(x, f, archetype = 1L, mergeTree = FALSE, mergeRefSeq = FALSE, ...)

## S4 method for signature 'TreeSummarizedExperiment'
mergeCols(x, f, archetype = 1L, mergeTree = FALSE, ...)

```

**Arguments**

x	a <a href="#">SummarizedExperiment</a> or a <a href="#">TreeSummarizedExperiment</a>
f	A factor for merging. Must be the same length as <code>nrow(x)/ncol(x)</code> . Rows/Cols corresponding to the same level will be merged. If <code>length(levels(f)) == nrow(x)/ncol(x)</code> , x will be returned unchanged.
archetype	Of each level of f, which element should be regarded as the archetype and metadata in the columns or rows kept, while merging? This can be single integer value or an integer vector of the same length as <code>levels(f)</code> . (Default: <code>archetype = 1L</code> , which means the first element encountered per factor level will be kept)
...	optional arguments: <ul style="list-style-type: none"> <li>passed onto <a href="#">sumCountsAcrossFeatures</a>, except <code>subset_row</code>, <code>subset_col</code></li> </ul>
mergeTree	TRUE or FALSE: should to <code>rowTree()</code> also be merged? (Default: <code>mergeTree = FALSE</code> )
mergeRefSeq	TRUE or FALSE: should a consensus sequence calculate? If set to FALSE, the result from archetype is returned; If set to TRUE the result from <a href="#">DECIPHER::ConsensusSequence</a> is returned. (Default: <code>mergeRefSeq = FALSE</code> )

**Details**

[assay](#) are agglomerated, i.e.. summed up. Other than counts / absolute values might lead to meaningless values.

These functions are similar to [sumCountsAcrossFeatures](#). However, additional support for `TreeSummarizedExperiment` was added and science field agnostic names were used. In addition the archetype argument lets the user select how to preserve row or column data.

For merge data of assays the function from `scuttle` are used.

**Value**

an object with the same class `x` with the specified entries merged into one entry in all relevant components.

**See Also**

[sumCountsAcrossFeatures](#)

**Examples**

```
data(esophagus)
esophagus
plot(rowTree(esophagus))
# get a factor for merging
f <- factor(regmatches(rownames(esophagus),
                      regexpr("[0-9]*_[0-9]*", rownames(esophagus))))
merged <- mergeRows(esophagus, f, mergeTree = TRUE)
plot(rowTree(merged))
#
data(GlobalPatterns)
GlobalPatterns
merged <- mergeCols(GlobalPatterns, colData(GlobalPatterns)$SampleType)
merged
```

---

mergeSEs

---

*Merge SE objects into single SE object.*


---

**Description**

Merge SE objects into single SE object.

**Usage**

```
mergeSEs(x, ...)
```

```
## S4 method for signature 'SimpleList'
mergeSEs(
  x,
  assay_name = "counts",
  join = "full",
  missing_values = NA,
  collapse_samples = FALSE,
  verbose = TRUE,
  ...
)
```

```
## S4 method for signature 'SummarizedExperiment'
mergeSEs(x, y = NULL, ...)
```

```

## S4 method for signature 'list'
mergeSEs(x, ...)

full_join(x, ...)

## S4 method for signature 'ANY'
full_join(x, ...)

inner_join(x, ...)

## S4 method for signature 'ANY'
inner_join(x, ...)

left_join(x, ...)

## S4 method for signature 'ANY'
left_join(x, ...)

right_join(x, ...)

## S4 method for signature 'ANY'
right_join(x, ...)

```

### Arguments

x	a <a href="#">SummarizedExperiment</a> object or a list of <a href="#">SummarizedExperiment</a> objects.
...	optional arguments (not used).
assay_name	A single character value for selecting the <a href="#">assay</a> to be merged. (By default: assay_name = "counts")
join	A single character value for selecting the joining method. Must be 'full', 'inner', 'left', or 'right'. 'left' and 'right' are disabled when more than two objects are being merged. (By default: join = "full")
missing_values	NA, 0, or a single character values specifying the notation of missing values. (By default: missing_values = NA)
collapse_samples	A boolean value for selecting whether to collapse identically named samples to one. (By default: collapse_samples = FALSE)
verbose	A single boolean value to choose whether to show messages. (By default: verbose = TRUE)
y	a <a href="#">SummarizedExperiment</a> object when x is a <a href="#">SummarizedExperiment</a> object. Disabled when x is a list.

### Details

This function merges multiple [SummarizedExperiment](#) objects. It combines `rowData`, `assays`, and `colData` so that the output includes each unique row and column ones. The merging is done based

on rownames and colnames. rowTree and colTree are preserved if linkage between rows/cols and the tree is found.

Equally named rows are interpreted as equal. Further matching based on rowData is not done. For samples, collapsing is disabled by default meaning that equally named samples that are stored in different objects are interpreted as unique. Collapsing can be enabled with `collapse_samples = TRUE` when equally named samples describe the same sample.

If, for example, all rows are not shared with individual objects, there are missing values in assays. The notation of missing can be specified with the `missing_values` argument. If input consists of `TreeSummarizedExperiment` objects, also `rowTree`, `colTree`, and `referenceSeq` are preserved if possible. The data is preserved if all the rows or columns can be found from it.

Compared to `cbind` and `rbind` `mergeSEs` allows more freely merging since `cbind` and `rbind` expect that rows and columns are matching, respectively.

You can choose joining methods from `'full'`, `'inner'`, `'left'`, and `'right'`. In all the methods, all the samples are included in the result object. However, with different methods, it is possible to choose which rows are included.

- `full` – all unique features
- `inner` – all shared features
- `left` – all the features of the first object
- `right` – all the features of the second object

You can also do e.g., a full join by using a function `full_join` which is an alias for `mergeSEs`. Also other joining methods have dplyr-like aliases.

The output depends on the input. If the input contains `SummarizedExperiment` object, then the output will be `SummarizedExperiment`. When all the input objects belong to `TreeSummarizedExperiment`, the output will be `TreeSummarizedExperiment`.

## Value

A single `SummarizedExperiment` object.

## Author(s)

Leo Lahti and Tuomas Borman. Contact: [microbiome.github.io](https://microbiome.github.io)

## See Also

- `TreeSummarizedExperiment::cbind`
- `TreeSummarizedExperiment::rbind`
- [full\\_join](#)
- [inner\\_join](#)
- [left\\_join](#)
- [right\\_join](#)

## Examples

```

data(GlobalPatterns)
data(esophagus)
data(enterotype)

# Take only subsets so that it wont take so long
tse1 <- GlobalPatterns[1:100, ]
tse2 <- esophagus
tse3 <- enterotype[1:100, ]

# Merge two TreeSEs
tse <- mergeSEs(tse1, tse2)

# Merge a list of TreeSEs
list <- SimpleList(tse1, tse2, tse3)
tse <- mergeSEs(list, assay_name = "counts", missing_values = 0)
tse

# With 'join', it is possible to specify the merging method. Subsets are used
# here just to show the functionality
tse_temp <- mergeSEs(tse[1:10, 1:10], tse[5:100, 11:20], join = "left")
tse_temp

# You can also do a left_join by using alias "left_join"
tse_temp <- left_join(tse[1:10, 1:10], tse[5:100, 11:20])

# If your objects contain samples that describe one and same sample,
# you can collapse equally named samples to one by specifying 'collapse_samples'
tse_temp <- inner_join(list(tse[1:10, 1], tse[1:20, 1], tse[1:5, 1]),
                      collapse_samples = TRUE)
tse_temp

```

---

mia-datasets

*mia datasets*

---

## Description

These datasets are conversions of the **phyloseq** datasets `GlobalPatterns`, `enterotype`, and `esophagus` into the [TreeSummarizedExperiment](#) data container.

`dmn_se` contains an example `SummarizedExperiment` derived from data in the **DirichletMultinomial** package. See `?calculatedDMN` for more details.

Global patterns of 16S rRNA diversity at a depth of millions of sequences per sample (2011)

This work compared the microbial communities from 25 environmental samples and three known “mock communities” at a an averag depth of 3.1 million reads per sample. Authors reproduced diversity patterns seen in many other published studies, while investigating technical issues/bias by applying the same techniques to simulated microbial communities of known composition. Many

thanks to J. Gregory Caporaso for providing the OTU-clustered data files for inclusion in the **phyloseq** package, from which this data `TreeSummarizedExperiment` version was then converted.

The enterotype data of the human gut microbiome (Arumugam et al. 2011) includes taxonomic profiling for 280 fecal samples from 22 subjects based on shotgun DNA sequencing. The authors claimed that the data naturally clumps into three community-level clusters, or “enterotypes”, that are not immediately explained by sequencing technology or demographic features of the subjects. A later addendum (2014) the authors stated that enterotypes "should not be seen as discrete clusters, but as a way of stratifying samples to reduce complexity."

The esophagus data set from Pei et al. (2004) includes 3 samples from 3 human adults based on biopsies analysed with 16S rDNA PCR. The 16S rRNA sequence processing has been provided in the `mothur` wiki at the link below.

This data set from Turnbaugh et al. (2009) was used to introduce Dirichlet Multinomial Mixtures (DMM) for microbiota stratification by Holmes et al. (2012).

PeerJ data by Potbhare et al. (2022) includes skin microbial profiles of 58 volunteers with multiple factors. 16S r-RNA sequencing of V3-V4 regions was done to generate millions of read using illumina platform. A standard bioinformatic and statistical analysis done to explore skin bacterial diversity and its association with age, diet, geographical locations. The authors investigated significant association of skin microbiota with individual's geographical location.

The HintikkaXO dataset contains high-throughput profiling data from 40 rat samples, including 39 biomarkers, 38 metabolites (NMR), and 12706 OTUs from 318 species, measured from Cecum. This is diet comparison study with High/Low fat diet and xylo-oligosaccharide supplementation. Column metadata is common for all experiments (microbiota, metabolites, biomarkers) and includes the following fields:

## Usage

```
data(GlobalPatterns)
```

```
data(enterotype)
```

```
data(esophagus)
```

```
data(dmn_se)
```

```
data(peerj13075)
```

```
data(HintikkaXOData)
```

## Format

An object of class `TreeSummarizedExperiment` with 19216 rows and 26 columns.

An object of class `TreeSummarizedExperiment` with 553 rows and 280 columns.

An object of class `TreeSummarizedExperiment` with 58 rows and 3 columns.

An object of class `SummarizedExperiment` with 130 rows and 278 columns.

An object of class `TreeSummarizedExperiment` with 674 rows and 58 columns.

An object of class `MultiAssayExperiment` of length 3.

## Details

- Sample: Sample ID (character)
- Rat: Rat ID (factor)
- Site: Site of measurement ("Cecum"); single value
- Diet: Diet group (factor; combination of the Fat and XOS fields)
- Fat: Fat in Diet (factor; Low/High)
- XOS: XOS Diet Supplement (numeric; 0/1)

Row metadata of the microbiota data contains taxonomic information on the Phylum, Class, Order, Family, Genus, Species, and OTU levels.

Biomarker data contains 39 biomarkers.

Metabolite data contains 38 NMR metabolites.

## Author(s)

Caporaso, J. G., et al.

Arumugam, M., Raes, J., et al.

Pei et al. <zhiheng.pei@med.nyu.edu>.

Turnbaugh, PJ et al.

Potbhare, R., et al.

Leo Lahti et al.

## References

Caporaso, J. G., et al. (2011). Global patterns of 16S rRNA diversity at a depth of millions of sequences per sample. *PNAS*, 108, 4516-4522. <http://www.pnas.org/content/108/suppl.1/4516.short>

Arumugam, M., et al. (2011). Enterotypes of the human gut microbiome. *Nature*, 473(7346), 174-180. <http://www.nature.com/doi/10.1038/nature09944> Supplemental information includes subject data. OTU-clustered data was initially downloaded from the publicly-accessible: [http://www.bork.embl.de/Docu/Arumugam\\_et\\_al\\_2011/downloads.html](http://www.bork.embl.de/Docu/Arumugam_et_al_2011/downloads.html)

Arumugam, M., et al. (2014). Addendum: Enterotypes of the human gut microbiome. *Nature* 506, 516 (2014). <https://doi.org/10.1038/nature13075>

Pei, Z., Bini, E. J., Yang, L., Zhou, M., Francois, F., & Blaser, M. J. (2004). Bacterial biota in the human distal esophagus. *Proceedings of the National Academy of Sciences of the United States of America*, 101(12), 4250-4255. <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC384727>

McMurdie, J. & Holmes, S. (2013) *phyloseq*: An R Package for reproducible interactive analysis and graphics of microbiome census data. *PLoS ONE*. 8(4):e61217. <https://doi.org/10.1371/journal.pone.0061217>

Mothur-processed files and the sequence data can be downloaded at: [http://www.mothur.org/wiki/Esophageal\\_community\\_analysis](http://www.mothur.org/wiki/Esophageal_community_analysis)

Holmes I, Harris K, Quince C (2012). Dirichlet Multinomial Mixtures: Generative Models for Microbial Metagenomics. *PLoS ONE* 7(2): e30126. <https://doi.org/10.1371/journal.pone.0030126>



Turnbaugh PJ, Hamady M, Yatsunencko T, Cantarel BL, Duncan A, et al. (2009). A core gut microbiome in obese and lean twins. *Nature* 457: 480–484. <https://doi.org/10.1038/nature07540>

Potbhare, R., RaviKumar, A., Munukka, E., Lahti, L., & Ashma, R. (2022). Skin microbiota diversity among genetically unrelated individuals of Indian origin. *PeerJ*, 10, e13075. <https://peerj.com/articles/13075/> Supplemental information includes OTU table and taxonomy table and publicly-accessible from: DOI:10.7717/peerj.13075/supp-1 DOI:10.7717/peerj.13075/supp-2

Hintikka L et al. (2021): Xylo-oligosaccharides in prevention of hepatic steatosis and adipose tissue inflammation: associating taxonomic and metabolomic patterns in fecal microbiotas with bi-clustering. *International Journal of Environmental Research and Public Health* 18(8):4049 <https://doi.org/10.3390/ijerph18084049>

---

perSampleDominantTaxa *Get dominant taxa*

---

## Description

These functions return information about the most dominant taxa in a [SummarizedExperiment](#) object.

## Usage

```
perSampleDominantTaxa(  
  x,  
  assay_name = abund_values,  
  abund_values = "counts",  
  rank = NULL,  
  ...  
)  
  
## S4 method for signature 'SummarizedExperiment'  
perSampleDominantTaxa(  
  x,  
  assay_name = abund_values,  
  abund_values = "counts",  
  rank = NULL,  
  ...  
)  
  
perSampleDominantFeatures(x, ...)  
  
## S4 method for signature 'SummarizedExperiment'  
perSampleDominantFeatures(x, ...)  
  
addPerSampleDominantTaxa(x, name = "dominant_taxa", ...)
```

```
## S4 method for signature 'SummarizedExperiment'
addPerSampleDominantTaxa(x, name = "dominant_taxa", ...)

addPerSampleDominantFeatures(x, ...)

## S4 method for signature 'SummarizedExperiment'
addPerSampleDominantFeatures(x, ...)
```

### Arguments

x	A <a href="#">SummarizedExperiment</a> object.
assay_name	A single character value for selecting the <a href="#">assay</a> to use for identifying dominant taxa.
abund_values	a single character value for specifying which assay to use for calculation. (Please use <code>assay_name</code> instead. At some point <code>abund_values</code> will be disabled.)
rank	A single character defining a taxonomic rank. Must be a value of the output of <code>taxonomyRanks()</code> .
...	Additional arguments passed on to <code>agglomerateByRank()</code> when rank is specified.
name	A name for the column of the <code>colData</code> where the dominant taxa will be stored in when using <code>addPerSampleDominantTaxa</code> .

### Details

`addPerSampleDominantTaxa` extracts the most abundant taxa in a [SummarizedExperiment](#) object, and stores the information in the `colData`. It is a wrapper for `perSampleDominantTaxa`.

With `rank` parameter, it is possible to agglomerate taxa based on taxonomic ranks. E.g. if 'Genus' rank is used, all abundances of same Genus are added together, and those families are returned. See `agglomerateByRank()` for additional arguments to deal with missing values or special characters.

### Value

`perSampleDominantTaxa` returns a named character vector `x` while `addPerSampleDominantTaxa` returns [SummarizedExperiment](#) with additional column in `colData` named `*name*`.

### Author(s)

Leo Lahti, Tuomas Borman and Sudarshan A. Shetty.

### Examples

```
data(GlobalPatterns)
x <- GlobalPatterns

# Finds the dominant taxa.
sim.dom <- perSampleDominantTaxa(x, rank="Genus")
```

```
# Add information to colData
x <- addPerSampleDominantTaxa(x, rank = "Genus", name="dominant_genera")
colData(x)
```

---

relabundance	<i>Getter / setter for relative abundance data</i>
--------------	--

---

### Description

relabundance is a getter/setter for relative abundance stored in the assay slot 'relabundance' of a [TreeSummarizedExperiment](#) object. This is a shortcut function for `assay(x, "relabundance")`.

### Usage

```
relabundance(x, ...)
```

```
relabundance(x) <- value
```

```
## S4 method for signature 'SummarizedExperiment'
```

```
relabundance(x)
```

```
## S4 replacement method for signature 'SummarizedExperiment'
```

```
relabundance(x) <- value
```

### Arguments

x	a <a href="#">TreeSummarizedExperiment</a> object
...	optional arguments not used currently.
value	a matrix to store as the the 'relabundance' assay

### Value

For relabundance the matrix stored with the name "relabundance".

### Examples

```
data(GlobalPatterns)
# Calculates relative abundances
GlobalPatterns <- relAbundanceCounts(GlobalPatterns)
# Fetches calculated relative abundances
head(relabundance(GlobalPatterns))
```

---

runCCA	<i>Canonical Correspondence Analysis</i>
--------	--

---

**Description**

These functions perform Canonical Correspondence Analysis on data stored in a SummarizedExperiment.

**Usage**

```
calculateCCA(x, ...)  
  
runCCA(x, ...)  
  
calculateRDA(x, ...)  
  
runRDA(x, ...)  
  
## S4 method for signature 'ANY'  
calculateCCA(x, formula, variables, scale = TRUE, ...)  
  
## S4 method for signature 'SummarizedExperiment'  
calculateCCA(  
  x,  
  formula,  
  variables,  
  ...,  
  assay_name = abund_values,  
  abund_values = exprs_values,  
  exprs_values = "counts"  
)  
  
## S4 method for signature 'SingleCellExperiment'  
runCCA(x, ..., altexp = NULL, name = "CCA")  
  
## S4 method for signature 'ANY'  
calculateRDA(x, formula, variables, ...)  
  
## S4 method for signature 'SummarizedExperiment'  
calculateRDA(  
  x,  
  formula,  
  variables,  
  ...,  
  assay_name = abund_values,  
  abund_values = exprs_values,  
  exprs_values = "counts"  
)
```

```
## S4 method for signature 'SingleCellExperiment'
runRDA(x, ..., altexp = NULL, name = "RDA")
```

### Arguments

x	For calculate* a <a href="#">SummarizedExperiment</a> or a numeric matrix with columns as samples For run* a <a href="#">SingleCellExperiment</a> or a derived object.
...	additional arguments passed to <code>vegan::cca</code> or <code>vegan::dbrda</code>
formula	If x is a <a href="#">SummarizedExperiment</a> a formula can be supplied. Based on the right-hand side of the given formula <code>colData</code> is subset to variables. variables and formula can be missing, which turns the CCA analysis into a CA analysis and <code>dbrDA</code> into <code>PCoA/MDS</code> .
variables	When x is a <a href="#">SummarizedExperiment</a> , variables can be used to specify variables from <code>colData</code> . When x is a matrix, variables is a data.frame or an object coercible to one containing the variables to use. All variables are used. Please subset, if you want to consider only some of them. variables and formula can be missing, which turns the CCA analysis into a CA analysis and <code>dbrDA</code> into <code>PCoA/MDS</code> .
scale	a logical scalar, should the expression values be standardized? scale is disabled when using *RDA functions. Please scale before performing RDA (Check examples.)
assay_name	a single character value for specifying which assay to use for calculation.
abund_values	a single character value for specifying which assay to use for calculation. (Please use <code>assay_name</code> instead. At some point <code>abund_values</code> will be disabled.)
exprs_values	a single character value for specifying which assay to use for calculation. (Please use <code>assay_name</code> instead.)
altexp	String or integer scalar specifying an alternative experiment containing the input data.
name	String specifying the name to be used to store the result in the <code>reducedDims</code> of the output.

### Details

\*CCA functions utilize `vegan::cca` and \*RDA functions `vegan::dbrDA`. By default `dbrDA` is done with euclidean distances which equals to RDA.

### Value

For `calculateCCA` a matrix with samples as rows and CCA dimensions as columns  
For `runCCA` a modified x with the results stored in `reducedDim` as the given name

**See Also**

For more details on the actual implementation see [cca](#) and [dbrda](#)

**Examples**

```
library(scater)
data(GlobalPatterns)
GlobalPatterns <- runCCA(GlobalPatterns, data ~ SampleType)
plotReducedDim(GlobalPatterns,"CCA", colour_by = "SampleType")

GlobalPatterns <- runRDA(GlobalPatterns, data ~ SampleType)
plotReducedDim(GlobalPatterns,"CCA", colour_by = "SampleType")

# To scale values when using *RDA functions, use transformFeatures
tse <- GlobalPatterns
tse <- transformFeatures(tse, method = "z")
# Data might include taxa that do not vary. Remove those because after z-transform
# their value is NA
tse <- tse[ rowSums( is.na( assay(tse, "z") ) ) == 0, ]
# Calculate RDA
tse <- runRDA(tse, formula = data ~ SampleType,
              assay_name = "z", name = "rda_scaled", na.action = na.omit)
# Plot
plotReducedDim(tse,"rda_scaled", colour_by = "SampleType")
```

---

runDPCoA

---

*Calculation of Double Principal Correspondance analysis*


---

**Description**

Double Principal Correspondance analysis is made available via the ade4 package in typical fashion. Results are stored in the reducedDims and are available for all the expected functions.

**Usage**

```
calculateDPCoA(x, y, ...)

## S4 method for signature 'ANY,ANY'
calculateDPCoA(
  x,
  y,
  ncomponents = 2,
  ntop = NULL,
  subset_row = NULL,
  scale = FALSE,
  transposed = FALSE,
  ...
)
```

```
## S4 method for signature 'TreeSummarizedExperiment,missing'
calculateDPCoA(
  x,
  ...,
  assay_name = abund_values,
  abund_values = exprs_values,
  exprs_values = "counts",
  tree_name = "phylo"
)

runDPCoA(x, ..., altexp = NULL, name = "DPCoA")
```

### Arguments

x	For calculateDPCoA, a numeric matrix of expression values where rows are features and columns are cells. Alternatively, a TreeSummarizedExperiment containing such a matrix. For runDPCoA a <a href="#">TreeSummarizedExperiment</a> containing the expression values as well as a rowTree to calculate y using <a href="#">cophenetic.phylo</a> .
y	a dist or a symmetric matrix compatible with ade4:dpcoa
...	Currently not used.
ncomponents	Numeric scalar indicating the number of DPCoA dimensions to obtain.
ntop	Numeric scalar specifying the number of features with the highest variances to use for dimensionality reduction. Alternatively NULL, if all features should be used. (default: ntop = NULL)
subset_row	Vector specifying the subset of features to use for dimensionality reduction. This can be a character vector of row names, an integer vector of row indices or a logical vector.
scale	Logical scalar, should the expression values be standardized?
transposed	Logical scalar, is x transposed with cells in rows?
assay_name	a single character value for specifying which assay to use for calculation.
abund_values	a single character value for specifying which assay to use for calculation. (Please use assay_name instead. At some point abund_values will be disabled.)
exprs_values	a single character value for specifying which assay to use for calculation. (Please use assay_name instead.)
tree_name	a single character value for specifying which rowTree will be used in calculation. (By default: tree_name = "phylo")
altexp	String or integer scalar specifying an alternative experiment containing the input data.
name	String specifying the name to be used to store the result in the reducedDims of the output.

**Details**

In addition to the reduced dimension on the features, the reduced dimension for samples are returned as well as `sample_red` attribute. `eig`, `feature_weights` and `sample_weights` are returned as attributes as well.

**Value**

For `calculateDPCoA` a matrix with samples as rows and CCA dimensions as columns

For `runDPCoA` a modified `x` with the results stored in `reducedDim` as the given name

**See Also**

[plotReducedDim](#) [reducedDims](#)

**Examples**

```
data(esophagus)
dpcoa <- calculateDPCoA(esophagus)
head(dpcoa)

esophagus <- runDPCoA(esophagus)
reducedDims(esophagus)

library(scater)
plotReducedDim(esophagus, "DPCoA")
```

---

runNMDS

*Perform non-metric MDS on sample-level data*

---

**Description**

Perform non-metric multi-dimensional scaling (nMDS) on samples, based on the data in a `SingleCellExperiment` object.

**Usage**

```
calculateNMDS(x, ...)

## S4 method for signature 'ANY'
calculateNMDS(
  x,
  FUN = vegdist,
  nmdsFUN = c("isoMDS", "monoMDS"),
  ncomponents = 2,
  ntop = 500,
  subset_row = NULL,
  scale = FALSE,
```



```

    transposed = FALSE,
    keep_dist = FALSE,
    ...
)

## S4 method for signature 'SummarizedExperiment'
calculateNMDS(
  x,
  ...,
  assay_name = abund_values,
  abund_values = exprs_values,
  exprs_values = "counts",
  FUN = vegdist
)

## S4 method for signature 'SingleCellExperiment'
calculateNMDS(
  x,
  ...,
  assay_name = abund_values,
  abund_values = exprs_values,
  exprs_values = "counts",
  dimred = NULL,
  n_dimred = NULL,
  FUN = vegdist
)

runNMDS(x, ..., altexp = NULL, name = "NMDS")

plotNMDS(x, ..., ncomponents = 2)

```

## Arguments

x	For calculateNMDS, a numeric matrix of expression values where rows are features and columns are cells. Alternatively, a <code>TreeSummarizedExperiment</code> containing such a matrix. For runNMDS a <a href="#">SingleCellExperiment</a>
...	additional arguments to pass to FUN and nmdsFUN.
FUN	a function or character value with a function name returning a <code>dist</code> object
nmdsFUN	a character value to choose the scaling implementation, either “isoMDS” for <code>MASS::isoMDS</code> or “monoMDS” for <code>vegan::monoMDS</code>
ncomponents	Numeric scalar indicating the number of NMDS dimensions to obtain.
ntop	Numeric scalar specifying the number of features with the highest variances to use for dimensionality reduction.
subset_row	Vector specifying the subset of features to use for dimensionality reduction. This can be a character vector of row names, an integer vector of row indices or a logical vector.



```
#
data(esophagus)
esophagus <- runNMDS(esophagus, FUN = vegan::vegdist, name = "BC")
esophagus <- runNMDS(esophagus, FUN = vegan::vegdist, name = "euclidean",
                      method = "euclidean")
reducedDims(esophagus)
```

---

splitByRanks

*Split/Unsplit a SingleCellExperiment by taxonomic ranks*


---

## Description

splitByRanks takes a SummarizedExperiment, splits it along the taxonomic ranks, aggregates the data per rank, converts the input to a SingleCellExperiment objects and stores the aggregated data as alternative experiments.

## Usage

```
splitByRanks(x, ...)

## S4 method for signature 'SummarizedExperiment'
splitByRanks(x, ranks = taxonomyRanks(x), na.rm = TRUE, ...)

## S4 method for signature 'SingleCellExperiment'
splitByRanks(x, ranks = taxonomyRanks(x), na.rm = TRUE, ...)

## S4 method for signature 'TreeSummarizedExperiment'
splitByRanks(x, ranks = taxonomyRanks(x), na.rm = TRUE, ...)

unsplitByRanks(x, ...)

## S4 method for signature 'SingleCellExperiment'
unsplitByRanks(x, ranks = taxonomyRanks(x), keep_reducedDims = FALSE, ...)

## S4 method for signature 'TreeSummarizedExperiment'
unsplitByRanks(x, ranks = taxonomyRanks(x), keep_reducedDims = FALSE, ...)
```

## Arguments

x	a <a href="#">SummarizedExperiment</a> object
...	arguments passed to <a href="#">agglomerateByRank</a> function for SummarizedExperiment objects and other functions. See <a href="#">agglomerateByRank</a> for more details.
ranks	a character vector defining taxonomic ranks. Must all be values of <a href="#">taxonomyRanks()</a> function.
na.rm	TRUE or FALSE: Should taxa with an empty rank be removed? Use it with caution, since results with NA on the selected rank will be dropped. This setting can be tweaked by defining <code>empty.fields</code> to your needs. (default: <code>na.rm = TRUE</code> )

keep\_reducedDims

TRUE or FALSE: Should the reducedDims(x) be transferred to the result? Please note, that this breaks the link between the data used to calculate the reduced dims. (default: keep\_reducedDims = FALSE)

## Details

unsplitByRanks takes these alternative experiments and flattens them again into a single SummarizedExperiment.

splitByRanks will use by default all available taxonomic ranks, but this can be controlled by setting ranks manually. NA values are removed by default, since they would not make sense, if the result should be used for unsplitByRanks at some point. The input data remains unchanged in the returned SingleCellExperiment objects.

unsplitByRanks will remove any NA value on each taxonomic rank so that no ambiguous data is created. In addition, a column taxonomicLevel is created or overwritten in the rowData to specify from which alternative experiment this originates from. This can also be used for [splitAltExps](#) to split the result along the same factor again. The input data from the base objects is not returned, only the data from the altExp(). Be aware that changes to rowData of the base object are not returned, whereas only the colData of the base object is kept.

## Value

For splitByRanks: SummarizedExperiment objects in a SimpleList.

For unsplitByRanks: x, with rowData and assay data replaced by the unsplit data. colData of x is kept as well and any existing rowTree is dropped as well, since existing rowLinks are not valid anymore.

## See Also

[splitOn](#) [unsplitOn](#) [mergeRows](#), [sumCountsAcrossFeatures](#), [agglomerateByRank](#), [altExps](#), [splitAltExps](#)

## Examples

```
data(GlobalPatterns)
# print the available taxonomic ranks
taxonomyRanks(GlobalPatterns)

# splitByRanks
altExps(GlobalPatterns) <- splitByRanks(GlobalPatterns)
altExps(GlobalPatterns)
altExp(GlobalPatterns, "Kingdom")
altExp(GlobalPatterns, "Species")

# unsplitByRanks
x <- unsplitByRanks(GlobalPatterns)
x
```

---

splitOn	<i>Split TreeSummarizedExperiment column-wise or row-wise based on grouping variable</i>
---------	--

---

## Description

Split TreeSummarizedExperiment column-wise or row-wise based on grouping variable

## Usage

```
splitOn(x, ...)

## S4 method for signature 'SummarizedExperiment'
splitOn(x, f = NULL, ...)

## S4 method for signature 'SingleCellExperiment'
splitOn(x, f = NULL, ...)

## S4 method for signature 'TreeSummarizedExperiment'
splitOn(x, f = NULL, update_rowTree = FALSE, ...)

unsplitOn(x, ...)

## S4 method for signature 'list'
unsplitOn(x, update_rowTree = FALSE, ...)

## S4 method for signature 'SimpleList'
unsplitOn(x, update_rowTree = FALSE, ...)

## S4 method for signature 'SingleCellExperiment'
unsplitOn(x, altExpNames = names(altExps(x)), keep_reducedDims = FALSE, ...)
```

## Arguments

x	A <a href="#">SummarizedExperiment</a> object or a list of <a href="#">SummarizedExperiment</a> objects.
...	Arguments passed to mergeRows/mergeCols function for SummarizedExperiment objects and other functions. See <a href="#">mergeRows</a> for more details. <ul style="list-style-type: none"> <li>• use_names A single boolean value to select whether to name elements of list by their group names.</li> </ul>
f	A single character value for selecting the grouping variable from rowData or colData or a factor or vector with the same length as one of the dimensions. If f matches with both dimensions, MARGIN must be specified. Split by cols is not encouraged, since this is not compatible with storing the results in altExps.
update_rowTree	TRUE or FALSE: Should the rowTree be updated based on splitted data? Option is enabled when x is a TreeSummarizedExperiment object or a list of such objects. (By default: update_rowTree = FALSE)

`altExpNames` a character vector specifying the alternative experiments to be unsplit. (By default: `altExpNames = names(altExps(x))`)

`keep_reducedDims` TRUE or FALSE: Should the `reducedDims(x)` be transferred to the result? Please note, that this breaks the link between the data used to calculate the reduced dims. (By default: `keep_reducedDims = FALSE`)

### Details

`splitOn` split data based on grouping variable. Splitting can be done column-wise or row-wise. The returned value is a list of `SummarizedExperiment` objects; each element containing members of each group.

### Value

For `splitOn`: `SummarizedExperiment` objects in a `SimpleList`.

For `unsplitOn`: `x`, with `rowData` and `assay` data replaced by the unsplit data. `colData` of `x` is kept as well and any existing `rowTree` is dropped as well, since existing `rowLinks` are not valid anymore.

### Author(s)

Leo Lahti and Tuomas Borman. Contact: [microbiome.github.io](https://microbiome.github.io)

### See Also

[splitByRanks](#) [unsplitByRanks](#) [mergeRows](#), [sumCountsAcrossFeatures](#), [agglomerateByRank](#), [altExps](#), [splitAltExps](#)

### Examples

```
data(GlobalPatterns)
tse <- GlobalPatterns
# Split data based on SampleType.
se_list <- splitOn(tse, f = "SampleType")

# List of SE objects is returned.
se_list

# Create arbitrary groups
rowData(tse)$group <- sample(1:10, nrow(tse), replace = TRUE)
colData(tse)$group <- sample(1:10, ncol(tse), replace = TRUE)

# Split based on rows
# Each element is named based on their group name. If you don't want to name
# elements, use use_name = FALSE. Since "group" can be found from rowData and colData
# you must use MARGIN.
se_list <- splitOn(tse, f = "group", use_names = FALSE, MARGIN = 1)

# When column names are shared between elements, you can store the list to altExps
altExps(tse) <- se_list
```

```

altExps(tse)

# If you want to split on columns and update rowTree, you can do
se_list <- splitOn(tse, f = colData(tse)$group, update_rowTree = TRUE)

# If you want to combine groups back together, you can use unsplitBy
unsplitOn(se_list)

```

---

subsampleCounts	<i>Subsample Counts</i>
-----------------	-------------------------

---

### Description

subsampleCounts will randomly subsample counts in SummarizedExperiment and return the a modified object in which each sample has same number of total observations/counts/reads.

### Usage

```

subsampleCounts(
  x,
  assay_name = abund_values,
  abund_values = "counts",
  min_size = min(colSums2(assay(x))),
  seed = runif(1, 0, .Machine$integer.max),
  replace = TRUE,
  name = "subsampling",
  verbose = TRUE,
  ...
)

## S4 method for signature 'SummarizedExperiment'
subsampleCounts(
  x,
  assay_name = abund_values,
  abund_values = "counts",
  min_size = min(colSums2(assay(x))),
  seed = runif(1, 0, .Machine$integer.max),
  replace = TRUE,
  name = "subsampling",
  verbose = TRUE,
  ...
)

```

**Arguments**

x	A SummarizedExperiment object.
assay_name	A single character value for selecting the SummarizedExperiment assay used for random subsampling. Only counts are useful and other transformed data as input will give meaningless output.
abund_values	a single character value for specifying which assay to use for calculation. (Please use assay_name instead. At some point abund_values will be disabled.)
min_size	A single integer value equal to the number of counts being simulated this can equal to lowest number of total counts found in a sample or a user specified number.
seed	A random number seed for reproducibility of sampling.
replace	Logical Default is TRUE. The default is with replacement (replace=TRUE). See <a href="#">phyloseq::rarefy_even_depth</a> for details on implications of this parameter.
name	A single character value specifying the name of transformed abundance table.
verbose	Logical Default is TRUE. When TRUE an additional message about the random number used is printed.
...	additional arguments not used

**Details**

Although the subsampling approach is highly debated in microbiome research, we include the `subsampleCounts` function because there may be some instances where it can be useful. Note that the output of `subsampleCounts` is not the equivalent as the input and any result have to be verified with the original dataset.

**Value**

`subsampleCounts` return x with subsampled data.

**Author(s)**

Sudarshan A. Shetty and Felix G.M. Ernst

**References**

- McMurdie PJ, Holmes S. Waste not, want not: why rarefying microbiome data is inadmissible. *PLoS computational biology*. 2014 Apr 3;10(4):e1003531.
- Gloor GB, Macklaim JM, Pawlowsky-Glahn V & Egozcue JJ (2017) Microbiome Datasets Are Compositional: And This Is Not Optional. *Frontiers in Microbiology* 8: 2224. doi: 10.3389/fmicb.2017.02224
- Weiss S, Xu ZZ, Peddada S, Amir A, Bittinger K, Gonzalez A, Lozupone C, Zaneveld JR, Vázquez-Baeza Y, Birmingham A, Hyde ER. Normalization and microbial differential abundance strategies depend upon data characteristics. *Microbiome*. 2017 Dec;5(1):1-8.



**Examples**

```

# When samples in TreeSE are less than specified min_size, they will be removed.
# If after subsampling features are not present in any of the samples,
# they will be removed.
data("GlobalPatterns")
tse <- GlobalPatterns
tse.subsampled <- subsampleCounts(tse,
                                min_size = 60000,
                                name = "subsampled",
                                seed = 123)

tse.subsampled
dim(tse)
dim(tse.subsampled)

```

---

subsetSamples

*Subset functions*


---

**Description**

To make a transition from phyloseq easier, the subsetSamples and subsetFeatures functions are implemented. To avoid name clashes they are named differently.

**Usage**

```

subsetSamples(x, ...)

subsetFeatures(x, ...)

subsetTaxa(x, ...)

## S4 method for signature 'SummarizedExperiment'
subsetSamples(x, ...)

## S4 method for signature 'SummarizedExperiment'
subsetFeatures(x, ...)

## S4 method for signature 'SummarizedExperiment'
subsetTaxa(x, ...)

```

**Arguments**

```

x          a SummarizedExperiment object
...       See subset. drop is not supported.

```

**Details**

However, the use of these functions is discouraged since subsetting using `[]` works on both dimension at the same time, is more flexible and is used throughout R to subset data with two or more dimension. Therefore, these functions will be removed in Bioconductor release 3.15 (April, 2022).

**Value**

A subset of `x`

**Examples**

```
data(GlobalPatterns)
subsetSamples(GlobalPatterns, colData(GlobalPatterns)$SampleType == "Soil")
# Vector that is used to specify subset must not include NAs
subsetFeatures(GlobalPatterns, rowData(GlobalPatterns)$Phylum == "Actinobacteria" &
               !is.na(rowData(GlobalPatterns)$Phylum))
```

---

summaries

*Summarizing microbiome data*


---

**Description**

To query a SummarizedExperiment for interesting features, several functions are available.

**Usage**

```
getTopTaxa(
  x,
  top = 5L,
  method = c("mean", "sum", "median"),
  assay_name = abund_values,
  abund_values = "counts",
  na.rm = TRUE,
  ...
)

## S4 method for signature 'SummarizedExperiment'
getTopTaxa(
  x,
  top = 5L,
  method = c("mean", "sum", "median", "prevalence"),
  assay_name = abund_values,
  abund_values = "counts",
  na.rm = TRUE,
  ...
)
```

```

getTopFeatures(x, ...)

## S4 method for signature 'SummarizedExperiment'
getTopFeatures(x, ...)

getUniqueTaxa(x, ...)

## S4 method for signature 'SummarizedExperiment'
getUniqueTaxa(x, rank = NULL, ...)

getUniqueFeatures(x, ...)

## S4 method for signature 'SummarizedExperiment'
getUniqueFeatures(x, ...)

countDominantTaxa(x, group = NULL, ...)

## S4 method for signature 'SummarizedExperiment'
countDominantTaxa(x, group = NULL, ...)

countDominantFeatures(x, ...)

## S4 method for signature 'SummarizedExperiment'
countDominantFeatures(x, ...)

## S4 method for signature 'SummarizedExperiment'
summary(object, assay_name = abund_values, abund_values = "counts")

```

### Arguments

x	A <a href="#">SummarizedExperiment</a> object.
top	Numeric value, how many top taxa to return. Default return top five taxa.
method	Specify the method to determine top taxa. Either sum, mean, median or prevalence. Default is 'mean'.
assay_name	a character value to select an <a href="#">assayNames</a> By default it expects count data.
abund_values	a single character value for specifying which assay to use for calculation. (Please use assay_name instead. At some point abund_values will be disabled.)
na.rm	For getTopTaxa logical argument for calculation method specified to argument method. Default is TRUE.
...	Additional arguments passed on to <a href="#">agglomerateByRank()</a> when rank is specified for countDominantTaxa.
rank	A single character defining a taxonomic rank. Must be a value of the output of <a href="#">taxonomyRanks()</a> .
group	With group, it is possible to group the observations in an overview. Must be one of the column names of colData.

object            A [SummarizedExperiment](#) object.

### Details

The `getTopTaxa` extracts the most top abundant “FeatureID”s in a [SummarizedExperiment](#) object.

The `getUniqueTaxa` is a basic function to access different taxa at a particular taxonomic rank.

`countDominantTaxa` returns information about most dominant taxa in a tibble. Information includes their absolute and relative abundances in whole data set.

The `summary` will return a summary of counts for all samples and features in [SummarizedExperiment](#) object.

### Value

The `getTopTaxa` returns a vector of the most top abundant “FeatureID”s

The `getUniqueTaxa` returns a vector of unique taxa present at a particular rank

The `countDominantTaxa` returns an overview in a tibble. It contains dominant taxa in a column named `*name*` and its abundance in the data set.

The `summary` returns a list with two tibbles

### Author(s)

Leo Lahti, Tuomas Borman and Sudarshan A. Shetty

### See Also

[getPrevalentTaxa](#)

[perCellQCMetrics](#), [perFeatureQCMetrics](#), [addPerCellQC](#), [addPerFeatureQC](#), [quickPerCellQC](#)

### Examples

```
data(GlobalPatterns)
top_taxa <- getTopTaxa(GlobalPatterns,
                      method = "mean",
                      top = 5,
                      assay_name = "counts")

top_taxa

# Use 'detection' to select detection threshold when using prevalence method
top_taxa <- getTopTaxa(GlobalPatterns,
                      method = "prevalence",
                      top = 5,
                      abund_values = "counts",
                      detection = 100)

top_taxa

# Top taxa os specific rank
getTopTaxa(agglomerateByRank(GlobalPatterns,
                             rank = "Genus",
                             na.rm = TRUE))
```

```

# Gets the overview of dominant taxa
dominant_taxa <- countDominantTaxa(GlobalPatterns,
                                   rank = "Genus")

dominant_taxa

# With group, it is possible to group observations based on specified groups
# Gets the overview of dominant taxa
dominant_taxa <- countDominantTaxa(GlobalPatterns,
                                   rank = "Genus",
                                   group = "SampleType",
                                   na.rm = TRUE)

dominant_taxa

# Get an overview of sample and taxa counts
summary(GlobalPatterns)

# Get unique taxa at a particular taxonomic rank
# sort = TRUE means that output is sorted in alphabetical order
# With na.rm = TRUE, it is possible to remove NAs
# sort and na.rm can also be used in function getTopTaxa
getUniqueTaxa(GlobalPatterns, "Phylum", sort = TRUE)

```

---

taxonomy-methods

*Functions for accessing taxonomic data stored in rowData.*


---

## Description

These function work on data present in `rowData` and define a way to represent taxonomic data alongside the features of a `SummarizedExperiment`.

## Usage

```

TAXONOMY_RANKS

taxonomyRanks(x)

## S4 method for signature 'SummarizedExperiment'
taxonomyRanks(x)

taxonomyRankEmpty(
  x,
  rank = taxonomyRanks(x)[1L],
  empty.fields = c(NA, "", " ", "\t", "-", "_")
)

## S4 method for signature 'SummarizedExperiment'

```

```

taxonomyRankEmpty(
  x,
  rank = taxonomyRanks(x)[1],
  empty.fields = c(NA, "", " ", "\t", "-", "_")
)

checkTaxonomy(x, ...)

## S4 method for signature 'SummarizedExperiment'
checkTaxonomy(x)

getTaxonomyLabels(x, ...)

## S4 method for signature 'SummarizedExperiment'
getTaxonomyLabels(
  x,
  empty.fields = c(NA, "", " ", "\t", "-", "_"),
  with_rank = FALSE,
  make_unique = TRUE,
  resolve_loops = FALSE,
  ...
)

taxonomyTree(x, ...)

## S4 method for signature 'SummarizedExperiment'
taxonomyTree(x)

addTaxonomyTree(x, ...)

## S4 method for signature 'SummarizedExperiment'
addTaxonomyTree(x)

mapTaxonomy(x, ...)

## S4 method for signature 'SummarizedExperiment'
mapTaxonomy(x, taxa = NULL, from = NULL, to = NULL, use_grep1 = FALSE)

IdTaxaToDataFrame(from)

```

### Arguments

<code>x</code>	a <code>SummarizedExperiment</code> object
<code>rank</code>	a single character defining a taxonomic rank. Must be a value of <code>taxonomyRanks()</code> function.
<code>empty.fields</code>	a character value defining, which values should be regarded as empty. (Default: <code>c(NA, "", " ", "\t")</code> ). They will be removed if <code>na.rm = TRUE</code> before agglomeration.

...	optional arguments not used currently.
with_rank	TRUE or FALSE: Should the level be add as a suffix? For example: "Phylum:Crenarchaeota" (default: with_rank = FALSE)
make_unique	TRUE or FALSE: Should the labels be made unique, if there are any duplicates? (default: make_unique = TRUE)
resolve_loops	TRUE or FALSE: Should resolveLoops be applied to the taxonomic data? Please note that has only an effect, if the data is unique. (default: resolve_loops = TRUE)
taxa	a character vector, which is used for subsetting the taxonomic information. If no information is found, NULL is returned for the individual element. (default: NULL)
from	<ul style="list-style-type: none"> <li>• For mapTaxonomy: a scalar character value, which must be a valid taxonomic rank. (default: NULL)</li> <li>• otherwise a Taxa object as returned by <a href="#">IdTaxa</a></li> </ul>
to	a scalar character value, which must be a valid taxonomic rank. (default: NULL)
use_grepl	TRUE or FALSE: should pattern matching via grepl be used? Otherwise literal matching is used. (default: FALSE)

### Format

a character vector of length 8 containing the taxonomy ranks recognized. In functions this is used as case insensitive.

### Details

taxonomyRanks returns, which columns of rowData(x) are regarded as columns containing taxonomic information.

taxonomyRankEmpty checks, if a selected rank is empty of information.

checkTaxonomy checks, if taxonomy information is valid and whether it contains any problems. This is a soft test, which reports some diagnostic and might mature into a data validator used upon object creation.

getTaxonomyLabels generates a character vector per row consisting of the lowest taxonomic information possible. If data from different levels, is to be mixed, the taxonomic level is prepended by default.

taxonomyTree generates a phylo tree object from the available taxonomic information. Internally it uses [toTree](#) and [resolveLoop](#) to sanitize data if needed.

IdTaxaToDataFrame extracts taxonomic results from results of [IdTaxa](#).

mapTaxonomy maps the given features (taxonomic groups; taxa) to the specified taxonomic level (to argument) in rowData of the SummarizedExperiment data object (i.e. rowData(x)[, taxonomyRanks(x)]). If the argument to is not provided, then all matching taxonomy rows in rowData will be returned. This function allows handy conversions between different

Taxonomic information from the IdTaxa function of DECIPHER package are returned as a special class. With as(taxa, "DataFrame") the information can be easily converted to a DataFrame compatible with storing the taxonomic information a rowData. Please note that the assigned confidence information are returned as metadata and can be accessed using metadata(df)\$confidence.

**Value**

- taxonomyRanks: a character vector with all the taxonomic ranks found in colnames(rowData(x))
- taxonomyRankEmpty: a logical value
- mapTaxonomy: a list per element of taxa. Each element is either a DataFrame, a character or NULL. If all character results have the length of one, a single character vector is returned.

**See Also**

[agglomerateByRank](#), [toTree](#), [resolveLoop](#)

**Examples**

```
data(GlobalPatterns)
GlobalPatterns
taxonomyRanks(GlobalPatterns)

checkTaxonomy(GlobalPatterns)

table(taxonomyRankEmpty(GlobalPatterns,"Kingdom"))
table(taxonomyRankEmpty(GlobalPatterns,"Species"))

getTaxonomyLabels(GlobalPatterns[1:20,])

# mapTaxonomy
## returns the unique taxonomic information
mapTaxonomy(GlobalPatterns)
# returns specific unique taxonomic information
mapTaxonomy(GlobalPatterns, taxa = "Escherichia")
# returns information on a single output
mapTaxonomy(GlobalPatterns, taxa = "Escherichia",to="Family")

# adding a rowTree() based on the available taxonomic information. Please
# note that any tree already stored in rowTree() will be overwritten.
x <- GlobalPatterns
x <- addTaxonomyTree(x)
x
```

---

transformCounts

*Transform Counts*


---

**Description**

Variety of transformations for abundance data, stored in assay. See details for options.



**Usage**

```
transformSamples(  
  x,  
  assay_name = abund_values,  
  abund_values = "counts",  
  method = c("clr", "rclr", "hellinger", "log10", "pa", "rank", "relabundance"),  
  name = method,  
  pseudocount = NULL,  
  threshold = 0  
)  
  
## S4 method for signature 'SummarizedExperiment'  
transformSamples(  
  x,  
  assay_name = abund_values,  
  abund_values = "counts",  
  method = c("clr", "rclr", "hellinger", "log10", "pa", "rank", "relabundance"),  
  name = method,  
  pseudocount = NULL,  
  threshold = 0  
)  
  
transformCounts(  
  x,  
  assay_name = abund_values,  
  abund_values = "counts",  
  method = c("clr", "rclr", "hellinger", "log10", "pa", "rank", "relabundance"),  
  name = method,  
  pseudocount = NULL,  
  threshold = 0  
)  
  
## S4 method for signature 'SummarizedExperiment'  
transformCounts(  
  x,  
  assay_name = abund_values,  
  abund_values = "counts",  
  method = c("clr", "rclr", "hellinger", "log10", "pa", "rank", "relabundance"),  
  name = method,  
  pseudocount = NULL,  
  threshold = 0  
)  
  
transformFeatures(  
  x,  
  assay_name = abund_values,  
  abund_values = "counts",  
  method = c("log10", "pa", "z"),
```

```

    name = method,
    pseudocount = NULL,
    threshold = 0
)

## S4 method for signature 'SummarizedExperiment'
transformFeatures(
  x,
  assay_name = abund_values,
  abund_values = "counts",
  method = c("log10", "pa", "z"),
  name = method,
  pseudocount = NULL,
  threshold = 0
)

ZTransform(x, ...)

## S4 method for signature 'SummarizedExperiment'
ZTransform(x, ...)

relAbundanceCounts(x, ...)

## S4 method for signature 'SummarizedExperiment'
relAbundanceCounts(x, ...)

```

### Arguments

x	A <a href="#">SummarizedExperiment</a> object.
assay_name	A single character value for selecting the <a href="#">assay</a> to be transformed.
abund_values	a single character value for specifying which assay to use for calculation. (Please use <code>assay_name</code> instead. At some point <code>abund_values</code> will be disabled.)
method	A single character value for selecting the transformation method.
name	A single character value specifying the name of transformed abundance table.
pseudocount	NULL or numeric value deciding whether pseudocount is added. The numeric value specifies the value of pseudocount.
threshold	A numeric value for setting threshold for pa transformation. By default it is 0. (Only used for <code>method = "pa"</code> )
...	additional arguments

### Details

These functions provide a variety of options for transforming abundance data. The transformed data is calculated and stored in a new assay.

Available wrapper functions include:

- transformSamples sample-wise (column-wise) transformation. Alias for transformCounts.
- transformFeatures feature-wise (row-wise) transformation.
- ZTransform Shortcut for Z-transformation.
- relAbundanceCounts Shortcut for fetching relative abundance table.

Altogether, transformCounts or transformSamples and transformFeatures apply transformations to the abundance table (assay). The available transformation methods include:

- 'clr' Centered log ratio (clr) transformation aims to remove compositionality effect; it is also used to skewness and to center data.

If the data contains zeros, pseudocount (commonly the smallest positive value of the data) must be added since clr is a logarithmic transformation that only allows positive values. (See e.g. Gloor et al. 2017.)

$$clr = \log_{10} \frac{xg(x)}{\mu} = \log_{10} x - \log_{10} \mu$$

where  $x$  is a single value,  $g(x)$  is geometric mean of sample-wise values, and  $\mu$  is an arithmetic mean of sample-wise values.

- 'rclr' rclr or robust clr is similar to regular clr. Problem of regular clr is that logarithmic transformations lead to undefined values when zeros are present in the data. In rclr, values are divided by geometric mean of observed taxa and zero values are not taken into account. Zero values will stay as zeroes. Because of high-dimensionality of data, rclr's geometric mean of observed taxa is a good approximation to the true geometric mean. (For details, see Martino et al. 2019.)
- 'hellinger' Hellinger transformation can be used to reduce the impact of extreme data points. It can be utilize for clustering or ordination analysis. (See e.g. Legendre & Gallagher 2001.)

$$hellinger = \sqrt{\frac{x}{x_{tot}}}$$

where  $x$  is a single value and  $x_{tot}$  is the sum of all values

- 'log10' log10 transformation can be used for reducing the skewness of the data.

$$log10 = \log_{10} x$$

where  $x$  is a single value of data.

- 'pa' Transforms table to presence/absence table. All abundances higher than  $\epsilon$  are transformed to 1 (present), otherwise 0 (absent). By default, threshold is 0.
- 'rank' Rank returns ranks of taxa. For each sample, the least abundant taxa get lower value and more abundant taxa bigger value. The implementation is based on the colRanks function with ties.method="first".
- 'relabundance' Transforms abundances to relative. Generally, all microbiome data are compositional. That is, e.g., because all measuring instruments have their capacity limits. To make results comparable with other results, values must be relative. (See e.g. Gloor et al. 2017.)

$$relabundance = \frac{x}{x_{tot}}$$

where  $x$  is a single value and  $x_{tot}$  is the sum of all values.

- 'z' Z-transformation, Z score transformation, or Z-standardization normalizes the data by shifting (to mean  $\mu$ ) and scaling (to standard deviation  $\sigma$ ). Z-transformation can be done with function `ZTransform`. It is done per rows (features / taxa), unlike most other transformations. This is often preceded by `log10p` or `clr` transformation. In other words, single value is standardized with respect of feature's values.

$$z = \frac{x - \mu}{\sigma}$$

where  $x$  is a single value,  $\mu$  is the mean of the feature, and  $\sigma$  is the standard deviation of the feature.

### Value

`transformCounts`, `transformSamples`, `transformFeatures`, `relAbundanceCounts`, and `ZTransform` return `x` with additional, transformed abundance table named `name` in the `assay`.

### Author(s)

Leo Lahti and Tuomas Borman. Contact: [microbiome.github.io](https://microbiome.github.io)

### References

- Gloor GB, Macklaim JM, Pawlowsky-Glahn V & Egozcue JJ (2017) Microbiome Datasets Are Compositional: And This Is Not Optional. *Frontiers in Microbiology* 8: 2224. doi: 10.3389/fmicb.2017.02224
- Legendre P & Gallagher ED (2001) Ecologically meaningful transformations for ordination of species data. *Oecologia* 129: 271-280.
- Martino C, Morton JT, Marotz CA, Thompson LR, Tripathi A, Knight R & Zengler K (2019) A Novel Sparse Compositional Technique Reveals Microbial Perturbations. *mSystems* 4: 1. doi: 10.1128/mSystems.00016-19

### Examples

```
data(esophagus)
x <- esophagus

# By specifying 'method', it is possible to apply different transformations,
# e.g. compositional transformation.
x <- transformSamples(x, method = "relabundance")

# The target of transformation can be specified with "assay_name"
# Pseudocount can be added by specifying 'pseudocount'.

# Get pseudocount; here smallest positive value
mat <- assay(x, "relabundance")
pseudonumber <- min(mat[mat>0])
# Perform CLR
x <- transformSamples(x, assay_name = "relabundance", method = "clr",
  pseudocount = pseudonumber
)
```

```
head(assay(x, "clr"))

# Name of the stored table can be specified.
x <- transformSamples(x, method="hellinger", name="test")
head(assay(x, "test"))

# pa returns presence absence table. With 'threshold', it is possible to set the
# threshold to a desired level. By default, it is 0.
x <- transformSamples(x, method = "pa", threshold = 35)
head(assay(x, "pa"))

# rank returns ranks of taxa. It is calculated column-wise, i.e., per sample
# and using the ties.method="first" from the colRanks function
x <- transformSamples(x, method = "rank")
head(assay(x, "rank"))

# transformCounts is an alias for transformSamples
x <- transformCounts(x, method = "relabundance", name = "test2")
head(assay(x, "test2"))

# In order to use other ranking variants, modify the chosen assay directly:
assay(x, "rank_average", withDimnames = FALSE) <- colRanks(assay(x, "counts"),
                                                           ties.method="average",
                                                           preserveShape = TRUE)

# If you want to do the transformation for features, you can do that by using
x <- transformFeatures(x, method = "log10", name = "log10_features", pseudocount = 1)
head(assay(x, "log10_features"))

# Z-transform can be done for features by using shortcut function
x <- ZTransform(x)
head(assay(x, "z"))

# For visualization purposes it is sometimes done by applying CLR for samples,
# followed by Z transform for taxa
x <- ZTransform(transformCounts(x, method="clr", assay_name = "counts", pseudocount = 1))

# Relative abundances can be also calculated with the dedicated
# relAbundanceCounts function.
x <- relAbundanceCounts(x)
head(assay(x, "relabundance"))
```

# Index

\* **datasets**  
     taxonomy-methods, 85  
 \* **data**  
     mia-datasets, 62  
 ?agglomerateByRank, 41  
 [, 82  
  
 addContaminantQC (isContaminant), 44  
 addContaminantQC, SummarizedExperiment-method  
     (isContaminant), 44  
 addNotContaminantQC (isContaminant), 44  
 addNotContaminantQC, SummarizedExperiment-method  
     (isContaminant), 44  
 addPerCellQC, 84  
 addPerFeatureQC, 84  
 addPerSampleDominantFeatures  
     (perSampleDominantTaxa), 65  
 addPerSampleDominantFeatures, SummarizedExperiment-method  
     (perSampleDominantTaxa), 65  
 addPerSampleDominantTaxa  
     (perSampleDominantTaxa), 65  
 addPerSampleDominantTaxa, SummarizedExperiment-method  
     (perSampleDominantTaxa), 65  
 addTaxonomyTree (taxonomy-methods), 85  
 addTaxonomyTree, SummarizedExperiment-method  
     (taxonomy-methods), 85  
 agglomerate-methods, 3  
 agglomerateByPrevalence  
     (getPrevalence), 39  
 agglomerateByPrevalence, SummarizedExperiment-method  
     (getPrevalence), 39  
 agglomerateByRank, 41, 42, 75, 76, 78, 88  
 agglomerateByRank  
     (agglomerate-methods), 3  
 agglomerateByRank, SingleCellExperiment-method  
     (agglomerate-methods), 3  
 agglomerateByRank, SummarizedExperiment-method  
     (agglomerate-methods), 3  
 agglomerateByRank, TreeSummarizedExperiment-method  
     (agglomerate-methods), 3  
  
 altExps, 76, 78  
 ape::phylo, 50  
 assay, 11, 23, 27, 35, 41, 45, 46, 52, 58, 60,  
     66, 90, 92  
 assayNames, 33, 56, 83  
  
 bestDMNFit (calculateDMN), 6  
 bestDMNFit, SummarizedExperiment-method  
     (calculateDMN), 6  
 BiocParallelParam, 8, 10, 14, 19, 23, 27, 29  
 biom, 53  
 Biostrings::DNAStrSet, 50  
  
 calculateCCA (runCCA), 68  
 calculateCCA, ANY-method (runCCA), 68  
 calculateCCA, SummarizedExperiment-method  
     (runCCA), 68  
 calculateDMN, 5  
 calculateDMN, ANY-method (calculateDMN),  
     6  
 calculateDMN, SummarizedExperiment-method  
     (calculateDMN), 6  
 calculateDMNgroup (calculateDMN), 6  
 calculateDMNgroup, ANY-method  
     (calculateDMN), 6  
 calculateDMNgroup, SummarizedExperiment-method  
     (calculateDMN), 6  
 calculateDPCoA (runDPCoA), 70  
 calculateDPCoA, ANY, ANY-method  
     (runDPCoA), 70  
 calculateDPCoA, TreeSummarizedExperiment, missing-method  
     (runDPCoA), 70  
 calculateJSD, 9, 12  
 calculateJSD, ANY-method (calculateJSD),  
     9  
 calculateJSD, SummarizedExperiment-method  
     (calculateJSD), 9  
 calculateNMDS (runNMDS), 72  
 calculateNMDS, ANY-method (runNMDS), 72

- calculateNMDS, SingleCellExperiment-method (runNMDS), 72
- calculateNMDS, SummarizedExperiment-method (runNMDS), 72
- calculateOverlap, 11
- calculateOverlap, SummarizedExperiment-method (calculateOverlap), 11
- calculateRDA (runCCA), 68
- calculateRDA, ANY-method (runCCA), 68
- calculateRDA, SummarizedExperiment-method (runCCA), 68
- calculateUnifrac, 12, 12
- calculateUnifrac, ANY, phylo-method (calculateUnifrac), 12
- calculateUnifrac, TreeSummarizedExperiment, missing-method (calculateUnifrac), 12
- cca, 70
- checkTaxonomy (taxonomy-methods), 85
- checkTaxonomy, SummarizedExperiment-method (taxonomy-methods), 85
- colData, 16, 20, 25, 27, 30, 66
- cophenetic.phylo, 71
- countDominantFeatures (summaries), 82
- countDominantFeatures, SummarizedExperiment-method (summaries), 82
- countDominantTaxa (summaries), 82
- countDominantTaxa, SummarizedExperiment-method (summaries), 82
  
- dbrda, 70
- DECIPHER: :ConsensusSequence, 58
- decontam: isContaminant, 45, 46
- decontam: isNotContaminant, 45, 46
- detectLoop, 4
- DirichletMultinomial, 6
- dist, 13, 73
- diversity, 21
- dmn, 8
- dmn\_se (mia-datasets), 62
- dmn\_se, (mia-datasets), 62
- DMNGroup, 8
- dmngroup, 8
  
- enterotype (mia-datasets), 62
- esophagus (mia-datasets), 62
- estimateDivergence, 15
- estimateDivergence, SummarizedExperiment-method (estimateDivergence), 15
- estimateDiversity, 17, 25, 28
- estimateDiversity, SummarizedExperiment-method (estimateDiversity), 17
- estimateDiversity, TreeSummarizedExperiment-method (estimateDiversity), 17
- estimateDominance, 16, 21, 22, 28
- estimateDominance, SummarizedExperiment-method (estimateDominance), 22
- estimateEvenness, 16, 21, 25, 26
- estimateEvenness, SummarizedExperiment-method (estimateEvenness), 26
- estimateFaith (estimateDiversity), 17
- estimateFaith, SummarizedExperiment, phylo-method (estimateDiversity), 17
- estimateFaith, TreeSummarizedExperiment, missing-method (estimateDiversity), 17
- estimateR, 21, 30, 31
- estimateRichness, 16, 21, 25, 28, 29
- estimateRichness, SummarizedExperiment-method (estimateRichness), 29
  
- full\_join, 61
- full\_join (mergeSEs), 59
- full\_join, ANY-method (mergeSEs), 59
- getAbundance, 32
- getAbundanceFeature (getAbundance), 32
- getAbundanceFeature, SummarizedExperiment-method (getAbundance), 32
- getAbundanceSample (getAbundance), 32
- getAbundanceSample, SummarizedExperiment-method (getAbundance), 32
- getBestDMNFit (calculateDMN), 6
- getBestDMNFit, SummarizedExperiment-method (calculateDMN), 6
- getDMN (calculateDMN), 6
- getDMN, SummarizedExperiment-method (calculateDMN), 6
- getExperimentCrossAssociation, 34
- getExperimentCrossAssociation, MultiAssayExperiment-method (getExperimentCrossAssociation), 34
- getExperimentCrossAssociation, SummarizedExperiment-method (getExperimentCrossAssociation), 34
- getExperimentCrossCorrelation (getExperimentCrossAssociation), 34
- getExperimentCrossCorrelation, ANY-method (getExperimentCrossAssociation), 34

- 34
- getPrevalence, 39
- getPrevalence, ANY-method
  - (getPrevalence), 39
- getPrevalence, SummarizedExperiment-method
  - (getPrevalence), 39
- getPrevalentAbundance (getPrevalence), 39
- getPrevalentAbundance, ANY-method
  - (getPrevalence), 39
- getPrevalentAbundance, SummarizedExperiment-method
  - (getPrevalence), 39
- getPrevalentFeatures (getPrevalence), 39
- getPrevalentFeatures, ANY-method
  - (getPrevalence), 39
- getPrevalentTaxa, 84
- getPrevalentTaxa (getPrevalence), 39
- getPrevalentTaxa, ANY-method
  - (getPrevalence), 39
- getPrevalentTaxa, SummarizedExperiment-method
  - (getPrevalence), 39
- getRareFeatures (getPrevalence), 39
- getRareFeatures, ANY-method
  - (getPrevalence), 39
- getRarePrevalentFeatures
  - (getPrevalence), 39
- getRareTaxa (getPrevalence), 39
- getRareTaxa, ANY-method (getPrevalence), 39
- getRareTaxa, SummarizedExperiment-method
  - (getPrevalence), 39
- getTaxonomyLabels (taxonomy-methods), 85
- getTaxonomyLabels, SummarizedExperiment-method
  - (taxonomy-methods), 85
- getTopFeatures (summaries), 82
- getTopFeatures, SummarizedExperiment-method
  - (summaries), 82
- getTopTaxa, 42
- getTopTaxa (summaries), 82
- getTopTaxa, SummarizedExperiment-method
  - (summaries), 82
- getUniqueFeatures (summaries), 82
- getUniqueFeatures, SummarizedExperiment-method
  - (summaries), 82
- getUniqueTaxa (summaries), 82
- getUniqueTaxa, SummarizedExperiment-method
  - (summaries), 82
- GlobalPatterns (mia-datasets), 62
- HintikkaXOData (mia-datasets), 62
- IdTaxa, 87
- IdTaxaToDataFrame (taxonomy-methods), 85
- inner\_join, 61
- inner\_join (mergeSEs), 59
- inner\_join, ANY-method (mergeSEs), 59
- isContaminant, 44
- isContaminant, SummarizedExperiment-method
  - (isContaminant), 44
- isNotContaminant, SummarizedExperiment-method
  - (isContaminant), 44
- left\_join, 61
- left\_join (mergeSEs), 59
- left\_join, ANY-method (mergeSEs), 59
- loadFromBiom (makeTreeSEFromBiom), 53
- loadFromMetaphlan, 46
- loadFromMothur, 47, 48, 51, 53–55
- loadFromQIIME2, 47, 48, 49, 53–55
- makePhyloseqFromTreeSE, 51
- makePhyloseqFromTreeSE, SummarizedExperiment-method
  - (makePhyloseqFromTreeSE), 51
- makePhyloseqFromTreeSE, TreeSummarizedExperiment-method
  - (makePhyloseqFromTreeSE), 51
- makePhyloseqFromTreeSummarizedExperiment
  - (makePhyloseqFromTreeSE), 51
- makePhyloseqFromTreeSummarizedExperiment, ANY-method
  - (makePhyloseqFromTreeSE), 51
- makeTreeSEFromBiom, 47, 48, 51, 53, 54, 55
- makeTreeSEFromDADA2, 47, 48, 51, 53, 54, 55
- makeTreeSEFromPhyloseq, 47, 48, 51, 53, 54, 55
- makeTreeSummarizedExperimentFromBiom
  - (makeTreeSEFromBiom), 53
- makeTreeSummarizedExperimentFromDADA2
  - (makeTreeSEFromDADA2), 54
- makeTreeSummarizedExperimentFromPhyloseq
  - (makeTreeSEFromPhyloseq), 55
- makeTreeSummarizedExperimentFromPhyloseq, ANY-method
  - (makeTreeSEFromPhyloseq), 55
- mapTaxonomy (taxonomy-methods), 85
- mapTaxonomy, SummarizedExperiment-method
  - (taxonomy-methods), 85
- MASS::isoMDS, 73, 74
- meltAssay, 56
- meltAssay, SummarizedExperiment-method
  - (meltAssay), 56



- merge-methods, [57](#)
- mergeCols (merge-methods), [57](#)
- mergeCols, SummarizedExperiment-method (merge-methods), [57](#)
- mergeCols, TreeSummarizedExperiment-method (merge-methods), [57](#)
- mergeRows, [4, 76–78](#)
- mergeRows (merge-methods), [57](#)
- mergeRows, SummarizedExperiment-method (merge-methods), [57](#)
- mergeRows, TreeSummarizedExperiment-method (merge-methods), [57](#)
- mergeSEs, [59](#)
- mergeSEs, list-method (mergeSEs), [59](#)
- mergeSEs, SimpleList-method (mergeSEs), [59](#)
- mergeSEs, SummarizedExperiment-method (mergeSEs), [59](#)
- metadata, [8](#)
- mia-datasets, [62](#)
- mia-package, [3](#)
- MultiAssayExperiment, [35](#)
- name, [30](#)
- peerj13075 (mia-datasets), [62](#)
- perCellQCMetrics, [84](#)
- perFeatureQCMetrics, [84](#)
- performDMNgroupCV (calculateDMN), [6](#)
- performDMNgroupCV, ANY-method (calculateDMN), [6](#)
- performDMNgroupCV, SummarizedExperiment-method (calculateDMN), [6](#)
- perSampleDominantFeatures (perSampleDominantTaxa), [65](#)
- perSampleDominantFeatures, SummarizedExperiment-method (perSampleDominantTaxa), [65](#)
- perSampleDominantTaxa, [65](#)
- perSampleDominantTaxa, SummarizedExperiment-method (perSampleDominantTaxa), [65](#)
- phylo, [13](#)
- phyloseq::rarefy\_even\_depth, [80](#)
- plotColData, [16, 21, 28, 31](#)
- plotMDS, [74](#)
- plotNMDS (runNMDS), [72](#)
- plotReducedDim, [72](#)
- quickPerCellQC, [84](#)
- rarifyCounts (subsamplingCounts), [79](#)
- readQZA (loadFromQIIME2), [49](#)
- reducedDims, [72](#)
- relabundance, [67](#)
- relabundance, SummarizedExperiment-method (relabundance), [67](#)
- relabundance<- (relabundance), [67](#)
- relabundance<- , SummarizedExperiment-method (relabundance), [67](#)
- relAbundanceCounts (transformCounts), [88](#)
- relAbundanceCounts, SummarizedExperiment-method (transformCounts), [88](#)
- resolveLoop, [87, 88](#)
- right\_join, [61](#)
- right\_join (mergeSEs), [59](#)
- right\_join, ANY-method (mergeSEs), [59](#)
- runCCA, [68](#)
- runCCA, SingleCellExperiment-method (runCCA), [68](#)
- runDMN (calculateDMN), [6](#)
- runDPCoA, [70](#)
- runJSD (calculateJSD), [9](#)
- runNMDS, [72](#)
- runOverlap (calculateOverlap), [11](#)
- runOverlap, SummarizedExperiment-method (calculateOverlap), [11](#)
- runRDA (runCCA), [68](#)
- runRDA, SingleCellExperiment-method (runCCA), [68](#)
- runUnifrac (calculateUnifrac), [12](#)
- SingleCellExperiment, [69, 73](#)
- splitAltExps, [76, 78](#)
- splitByRanks, [75, 78](#)
- splitByRanks, SingleCellExperiment-method (splitByRanks), [75](#)
- splitByRanks, SummarizedExperiment-method (splitByRanks), [75](#)
- splitByRanks, TreeSummarizedExperiment-method (splitByRanks), [75](#)
- splitOn, [76, 77](#)
- splitOn, SingleCellExperiment-method (splitOn), [77](#)
- splitOn, SummarizedExperiment-method (splitOn), [77](#)
- splitOn, TreeSummarizedExperiment-method (splitOn), [77](#)
- subsamplingCounts, [79](#)
- subsamplingCounts, SummarizedExperiment-method (subsamplingCounts), [79](#)

- subset, [81](#)
- subsetByPrevalentFeatures
  - (getPrevalence), [39](#)
- subsetByPrevalentFeatures, ANY-method
  - (getPrevalence), [39](#)
- subsetByPrevalentTaxa (getPrevalence), [39](#)
- subsetByPrevalentTaxa, SummarizedExperiment-method
  - (getPrevalence), [39](#)
- subsetByRareFeatures (getPrevalence), [39](#)
- subsetByRareFeatures, ANY-method
  - (getPrevalence), [39](#)
- subsetByRareTaxa (getPrevalence), [39](#)
- subsetByRareTaxa, SummarizedExperiment-method
  - (getPrevalence), [39](#)
- subsetFeatures (subsetSamples), [81](#)
- subsetFeatures, SummarizedExperiment-method
  - (subsetSamples), [81](#)
- subsetSamples, [81](#)
- subsetSamples, SummarizedExperiment-method
  - (subsetSamples), [81](#)
- subsetTaxa (subsetSamples), [81](#)
- subsetTaxa, SummarizedExperiment-method
  - (subsetSamples), [81](#)
- sumCountsAcrossFeatures, [4](#), [58](#), [59](#), [76](#), [78](#)
- summaries, [82](#)
- SummarizedExperiment, [4](#), [7](#), [9](#), [11](#), [16](#), [19](#), [23](#), [27](#), [29](#), [33](#), [35](#), [41](#), [44](#), [45](#), [56](#), [58](#), [60](#), [65](#), [66](#), [69](#), [75](#), [77](#), [81](#), [83](#), [84](#), [86](#), [90](#)
- summary, SummarizedExperiment-method
  - (summaries), [82](#)
- taxonomy-methods, [85](#)
- TAXONOMY\_RANKS (taxonomy-methods), [85](#)
- taxonomyRankEmpty (taxonomy-methods), [85](#)
- taxonomyRankEmpty, SummarizedExperiment-method
  - vegan::diversity, [20](#)
  - (taxonomy-methods), [85](#)
- taxonomyRanks, [3](#)
- taxonomyRanks (taxonomy-methods), [85](#)
- taxonomyRanks, SummarizedExperiment-method
  - (taxonomy-methods), [85](#)
- taxonomyTree (taxonomy-methods), [85](#)
- taxonomyTree, SummarizedExperiment-method
  - (taxonomy-methods), [85](#)
- testExperimentCrossAssociation
  - (getExperimentCrossAssociation), [34](#)
- testExperimentCrossAssociation, ANY-method
  - (getExperimentCrossAssociation), [34](#)
- testExperimentCrossCorrelation
  - (getExperimentCrossAssociation), [34](#)
- testExperimentCrossCorrelation, ANY-method
  - (getExperimentCrossAssociation), [34](#)
- toTree, [87](#), [88](#)
- transformCounts, [88](#)
- transformCounts, SummarizedExperiment-method
  - (transformCounts), [88](#)
- transformFeatures (transformCounts), [88](#)
- transformFeatures, SummarizedExperiment-method
  - (transformCounts), [88](#)
- transformSamples (transformCounts), [88](#)
- transformSamples, SummarizedExperiment-method
  - (transformCounts), [88](#)
- TreeSummarizedExperiment, [3](#), [12](#), [13](#), [19](#), [47](#), [48](#), [50](#), [53](#), [58](#), [62](#), [67](#), [71](#)
- twins (mia-datasets), [62](#)
- unsplitByRanks, [78](#)
- unsplitByRanks (splitByRanks), [75](#)
- unsplitByRanks, SingleCellExperiment-method
  - (splitByRanks), [75](#)
- unsplitByRanks, TreeSummarizedExperiment-method
  - (splitByRanks), [75](#)
- unsplitOn, [76](#)
- unsplitOn (splitOn), [77](#)
- unsplitOn, list-method (splitOn), [77](#)
- unsplitOn, SimpleList-method (splitOn), [77](#)
- unsplitOn, SingleCellExperiment-method
  - (splitOn), [77](#)
- ZTransform (transformCounts), [88](#)
- ZTransform, SummarizedExperiment-method
  - (transformCounts), [88](#)
- vegan::fisher.alpha, [20](#)
- vegan::monoMDS, [73](#), [74](#)