

# Package ‘RaceID’

January 22, 2019

**Title** Identification of Cell Types and Inference of Lineage Trees from Single-Cell RNA-Seq Data

**Version** 0.1.3

**Date** 2019-01-22

**Author** Dominic Grün <dominic.gruen@gmail.com>

**Maintainer** Dominic Grün <dominic.gruen@gmail.com>

**Description** Application of 'RaceID' allows inference of cell types and prediction of lineage trees by the StemID2 algorithm. Herman, J.S., Sagar, Grün D. (2018) <DOI:10.1038/nmeth.4662>.

**Depends** R (>= 3.3)

**biocViews**

**Imports** coop, cluster, FateID, fpc, grDevices, ica, igraph, irlba, locfit, methods, MASS, Matrix, pheatmap, quadprog, randomForest, RColorBrewer, Rtsne, vegan

**Suggests** DESeq2, destiny, knitr, rmarkdown, scan

**VignetteBuilder** knitr

**License** GPL-3

**Encoding** UTF-8

**LazyData** true

**RoxygenNote** 6.0.1

**NeedsCompilation** no

**Repository** CRAN

**Date/Publication** 2019-01-22 15:50:03 UTC

## R topics documented:

barplotgene	3
branchcells	3
CCcorrect	4
cellsfromtree	5

clustdiffgenes . . . . .	6
clustexp . . . . .	7
clustheatmap . . . . .	8
compdist . . . . .	9
compentropy . . . . .	10
compfr . . . . .	10
compmedoids . . . . .	11
comppvalue . . . . .	12
compscore . . . . .	13
comptsne . . . . .	13
diffexpnb . . . . .	14
diffgenes . . . . .	16
filterdata . . . . .	17
findoutliers . . . . .	18
getfdata . . . . .	19
getproj . . . . .	19
imputeexp . . . . .	20
intestinalData . . . . .	20
intestinalDataSmall . . . . .	21
lineagegraph . . . . .	21
Ltree-class . . . . .	22
plotbackground . . . . .	24
plotdiffgenes . . . . .	24
plotdiffgenesnb . . . . .	25
plotdimsat . . . . .	26
plotdistanceratio . . . . .	27
plotexpmap . . . . .	27
plotgraph . . . . .	28
plotjaccard . . . . .	29
plotlabelsmap . . . . .	29
plotlinkpv . . . . .	30
plotlinkscore . . . . .	30
plotmap . . . . .	31
plotmarkergenes . . . . .	31
plotoutlierprobs . . . . .	33
plotprojections . . . . .	33
plotsaturation . . . . .	34
plotsensitivity . . . . .	34
plotsilhouette . . . . .	35
plotspantree . . . . .	35
plotsymbolsmap . . . . .	36
projback . . . . .	36
projcells . . . . .	37
projenrichment . . . . .	38
rfcorrect . . . . .	39
SCseq . . . . .	40
varRegression . . . . .	41

---

barplotgene	<i>Gene Expression Barplot</i>
-------------	--------------------------------

---

**Description**

This functions generates a barplot of gene expression across all clusters.

**Usage**

```
barplotgene(object, g, n = NULL, logsc = FALSE)
```

**Arguments**

object	SCseq class object.
g	Individual gene name or vector with a group of gene names corresponding to a subset of valid row names of the ndata slot of the SCseq object.
n	String of characters representing the title of the plot. Default is NULL and the first element of g is chosen.
logsc	logical. If TRUE, then gene expression values are log2-transformed after adding a pseudo-count of 0.1. Default is FALSE and untransformed values are shown.

**Value**

None

---

branchcells	<i>Differential Gene Expression between Links</i>
-------------	---

---

**Description**

This function computes expression z-score between groups of cells from the same cluster residing on different links

**Usage**

```
branchcells(object, br)
```

**Arguments**

object	Ltree class object.
br	List containing two branches, where each component has to be two valid cluster numbers seperated by a . and with one common cluster in the two components. The lower number precedes the larger one, i.e. 1 . 3. For each component, the cluster number need to be ordered in increasing order.

**Value**

A list of four components:

n                    a vector with the number of significant links for each cluster.  
 scl                  a vector with the delta entropy for each cluster.  
 k                    a vector with the StemID score for each cluster.  
 diffgenes          a vector with the StemID score for each cluster.

**Examples**

```
sc <- SCseq(intestinalDataSmall)
sc <- filterdata(sc)
sc <- compdist(sc)
sc <- clustexp(sc)
sc <- findoutliers(sc)
sc <- comptsne(sc)
ltr <- Ltree(sc)
ltr <- compentropy(ltr)
ltr <- projcells(ltr)
ltr <- lineagegraph(ltr)
ltr <- compvalue(ltr)
x <- branchcells(ltr,list("1.3","3.6"))
head(x$diffgenes$z)
plotmap(x$scl)
plotdiffgenes(x$diffgenes,names(x$diffgenes$z)[1])
```

---

 CCcorrect

---

*Dimensional Reduction by PCA or ICA*


---

**Description**

This functions performs dimensional reduction by PCA or ICA and removes components enriched for particular gene sets, e.g. cell cycle related genes genes associated with technical batch effects.

**Usage**

```
CCcorrect(object, vset = NULL, CGenes = NULL, ccor = 0.4, pvalue = 0.01,
  quant = 0.01, nComp = NULL, dimR = FALSE, mode = "pca",
  logscale = FALSE, FSelect = TRUE)
```

**Arguments**

object              SCseq class object.  
 vset                List of vectors with genes sets. The loadings of each component are tested for enrichment in any of these gene sets and if the lower quant or upper 1 - quant fraction of genes ordered by loading is enriched at a p-value < pvalue the component is discarded. Default is NULL.

CGenes	Vector of gene names. If this argument is given, gene sets to be tested for enrichment in PCA- or ICA-components are defined by all genes with a Pearson's correlation of >ccor to a gene in CGenes. The loadings of each component are tested for enrichment in any of these gene sets and if the lower quant or upper 1 - quant fraction of genes ordered by loading is enriched at a p-value < pvalue the component is discarded. Default is NULL.
ccor	Positive number between 0 and 1. Correlation threshold used to determine correlating gene sets for all genes in CGenes. Default is 0.4.
pvalue	Positive number between 0 and 1. P-value cutoff for determining enriched components. See vset or CGenes. Default is 0.01.
quant	Positive number between 0 and 1. Upper and lower fraction of gene loadings used for determining enriched components. See vset or CGenes. Default is 0.01.
nComp	Number of PCA- or ICA-components to use. Default is NULL and the maximal number of components is computed.
dimR	logical. If TRUE, then the number of principal components to use for downstream analysis is derived from a saturation criterion. See function plotdimsat. Default is FALSE and all nComp components are used.
mode	"pca" or "ica" to perform either principal component analysis or independent component analysis. Default is pca.
logscale	logical. If TRUE data are log-transformed prior to PCA or ICA. Default is FALSE.
FSelect	logical. If TRUE, then PCA or ICA is performed on the filtered expression matrix using only the features stored in slotcluster\$features as computed in the function filterdata. See FSelect for function filterdata. Default is TRUE.

### Value

The function returns an updated SCseq object with the principal or independent component matrix written to the slot dimRed\$x of the SCseq object. Additional information on the PCA or ICA is stored in slot dimRed.

### Examples

```
sc <- SCseq(intestinalDataSmall)
sc <- filterdata(sc)
sc <- CCcorrect(sc, dimR=TRUE, nComp=3)
```

### Description

This function extracts a vector of cells on a given differentiation trajectory in pseudo-temporal order determined from the projection coordinates.

**Usage**

```
cellsfromtree(object, z)
```

**Arguments**

object	Ltree class object.
z	Vector of valid cluster numbers ordered along the trajectory.

**Value**

A list of four components:

f	a vector of cells ids ordered along the trajectory defined by z.
g	a vector of integer number. Number i indicates that a cell resides on the link between the i-th and (i+1)-th cluster in z.

**Examples**

```
sc <- SCseq(intestinalDataSmall)
sc <- filterdata(sc)
sc <- compdist(sc)
sc <- clustexp(sc)
sc <- findoutliers(sc)
sc <- comptsne(sc)
ltr <- Ltree(sc)
ltr <- compentropy(ltr)
ltr <- projcells(ltr)
ltr <- lineagegraph(ltr)
ltr <- compvalue(ltr)
x <- cellsfromtree(ltr,c(1,3,6,2))
```

---

clustdiffgenes

*Inference of differentially expressed genes in a cluster*

---

**Description**

This functions computes differentially expressed genes in a cluster by comparing to all remaining cells outside of the cluster based on a negative binomial model of gene expression

**Usage**

```
clustdiffgenes(object, cl, pvalue = 0.01)
```

**Arguments**

object	SCseq class object.
c1	A valid cluster number from the final cluster partition stored in the cpart slot of the SCseq object.
pvalue	Positive real number smaller than one. This is the p-value cutoff for the inference of differential gene expression. Default is 0.01.

**Value**

A data.frame of differentially expressed genes ordered by p-value in increasing order, with four columns:

mean.nc1	mean expression across cells outside of cluster c1.
mean.c1	mean expression across cells within cluster c1.
fc	fold-change of mean expression in cluster c1 versus the remaining cells.
pv	inferred p-value for differential expression.
padj	Benjamini-Hochberg corrected FDR.

**Examples**

```
sc <- SCseq(intestinalDataSmall)
sc <- filterdata(sc)
sc <- compdist(sc)
sc <- clustexp(sc)
sc <- findoutliers(sc)
x <- clustdiffgenes(sc,1)
head(x[x$fc>1,])
```

---

clustexp

*Clustering of single-cell transcriptome data*

---

**Description**

This functions performs the initial clustering of the RaceID3 algorithm.

**Usage**

```
clustexp(object, sat = TRUE, samp = NULL, c1n = NULL, clustnr = 30,
bootnr = 50, rseed = 17000, FUNcluster = "kmedoids")
```

**Arguments**

object	SCseq class object.
sat	logical. If TRUE, then the number of clusters is determined based on finding the saturation point of the mean within-cluster dispersion as a function of the cluster number. Default is TRUE. If FALSE, then cluster number needs to be given as cln.
samp	Number of random sample of cells used for the inference of cluster number and for inferring Jaccard similarities. Default is 1000.
cln	Number of clusters to be used. Default is NULL and the cluster number is inferred by the saturation criterion.
clustnr	Maximum number of clusters for the derivation of the cluster number by the saturation of mean within-cluster-dispersion. Default is 30.
bootnr	Number of bootstrapping runs for clusterboot. Default is 50.
rseed	Integer number. Random seed to enforce reproducible clustering results. Default is 17000.
FUNcluster	Clustering method used by RaceID3. One of "kmedoids", "kmeans", "hclust". Default is "kmedoids".

**Value**

SCseq object with clustering data stored in slot cluster and slot clusterpar. The clustering partition is stored in cluster\$kpart.

**Examples**

```
sc <- SCseq(intestinalDataSmall)
sc <- filterdata(sc)
sc <- compdist(sc)
sc <- clustexp(sc)
```

---

clustheatmap

*Plotting a Heatmap of the Distance Matrix*


---

**Description**

This functions plots a heatmap of the distance matrix grouped by clusters.

**Usage**

```
clustheatmap(object, final = TRUE, hmethod = "single")
```



**Arguments**

object	SCseq class object.
final	logical. If TRUE, then cells are grouped based on final clusters after outlier identification. If FALSE, then initial clusters prior to outlier identification are used for grouping. Default is TRUE.
hmethod	Agglomeration method used for determining the cluster order from hierarchical clustering of the cluster medoids. See hclust function.

**Value**

Returns a vector of cluster numbers ordered as determined by herarchical clustering of cluster the cluster medoids as depicted in the heatmap.

---

compdist	<i>Computing a distance matrix for cell type inference</i>
----------	--

---

**Description**

This functions computes the distance matrix used for cell type inference by RaceID3.

**Usage**

```
compdist(object, metric = "pearson", FSelect = TRUE, knn = NULL)
```

**Arguments**

object	SCseq class object.
metric	Distances are computed from the filtered expression matrix after optional feature selection, dimensional reduction, and/or transformation (batch correction). Possible values for metric are "pearson", "spearman", "logpearson", "euclidean". Default is "pearson". In case of the correlation based methods, the distance is computed as 1 – correlation.
FSelect	Logical parameter. If TRUE, then feature selection is performed prior to RaceID3 analysis. Default is TRUE.
knn	Positive integer number of nearest neighbours used for imputing gene expression values. Default is NULL and no imputing is done.

**Value**

SCseq object with the distance matrix in slot distances. If FSelect=TRUE, the genes used for computing the distance object are stored in slot cluster\$features.

**Examples**

```
sc <- SCseq(intestinalDataSmall)
sc <- filterdata(sc)
sc <- compdist(sc)
```

---

`compentropy`*Compute transcriptome entropy of each cell*

---

**Description**

This function computes the transcriptome entropy for each cell.

**Usage**

```
compentropy(object)
```

**Arguments**

`object`            Ltree class object.

**Value**

An Ltree class object with a vector of entropies for each cell in the same order as column names in slot `sc@ndata`.

**Examples**

```
sc <- SCseq(intestinalDataSmall)
sc <- filterdata(sc)
sc <- compdist(sc)
sc <- clustexp(sc)
sc <- findoutliers(sc)
sc <- comptsne(sc)
ltr <- Ltree(sc)
ltr <- compentropy(ltr)
```

---

`compfr`*Computation of a two dimensional Fruchterman-Rheingold representation*

---

**Description**

This functions performs the computation of a Fruchterman-Rheingold graph layout based on an adjacency matrix derived from the distance object in slot `distances` using the **igraph** package.

**Usage**

```
compfr(object, knn = 10, rseed = 15555)
```

**Arguments**

object	SCseq class object.
knn	Positive integer number of nearest neighbours used for the inference of the Fruchterman-Rheingold layout. Default is 10.
rseed	Integer number. Random seed to enforce reproducible layouts.

**Value**

SCseq object with layout coordinates stored in slot fr.

**Examples**

```
sc <- SCseq(intestinalDataSmall)
sc <- filterdata(sc)
sc <- compdist(sc)
sc <- clustexp(sc)
sc <- findoutliers(sc)
sc <- compfr(sc)
```

---

compmedoids

*Computes Medoids from a Clustering Partition*

---

**Description**

This functions computes cluster medoids given an SCseq object and a clustering partition.

**Usage**

```
compmedoids(object, part)
```

**Arguments**

object	SCseq class object.
part	Clustering partition. A vector of cluster numbers for (a subset of) cells (i.e. column names) of slot ndata from the SCseq object.

**Value**

Returns a list of medoids (column names of slot ndata from the SCseq object) ordered by increasing cluster number.

---

`comppvalue`*Computing P-values for Link Significance*

---

**Description**

This function computes a p-value for the significance (i.e. over-representation of assigned cells) of each inter-cluster link.

**Usage**

```
comppvalue(object, pthr = 0.01, sensitive = FALSE)
```

**Arguments**

<code>object</code>	Ltree class object.
<code>pthr</code>	p-value cutoff for link significance. This threshold is applied for the calculation of link scores reflecting how uniformly a link is occupied by cells.
<code>sensitive</code>	logical. Only relevant when <code>nmode=TRUE</code> in function <code>projcell</code> . If <code>TRUE</code> , then all cells on the most highly significant link are and the link itself are disregarded to test significance of the remaining links with a binomial p-value. Default is <code>FALSE</code> .

**Value**

An Ltree class object with link p-value and occupancy data stored in slot `cdata`.

**Examples**

```
sc <- SCseq(intestinalDataSmall)
sc <- filterdata(sc)
sc <- compdist(sc)
sc <- clustexp(sc)
sc <- findoutliers(sc)
sc <- comptsne(sc)
ltr <- Ltree(sc)
ltr <- compentropy(ltr)
ltr <- projcells(ltr)
ltr <- lineagegraph(ltr)
ltr <- comppvalue(ltr)
```

---

compscore	<i>Compute StemID2 score</i>
-----------	------------------------------

---

### Description

This function extracts the number of links connecting a given cluster to other cluster, the delta median entropy of each cluster (median entropy of a cluster after subtracting the minimum median entropy across all clusters), and the StemID2 score which is the product of both quantities for each cluster.

### Usage

```
compscore(object, nn = 1, scthr = 0, show = TRUE)
```

### Arguments

object	L tree class object.
nn	Positive integer number. Number of higher order neighbors to be included for the determination of links: indirect connections via n-1 intermittant neighbors are allowed. Default is 1.
scthr	Real number between zero and one. Score threshold for links to be included in the calculation. For scthr=0 all significant links are included. The maximum score is one.
show	logical. If TRUE, then plot heatmap of projections. Default is TRUE.

### Value

A list of three components:

links	a vector with the number of significant links for each cluster.
entropy	a vector with the delta entropy for each cluster.
StemIDscore	a vector with the StemID score for each cluster.

---

comptsne	<i>Computation of a two dimensional t-SNE representation</i>
----------	--

---

### Description

This functions performs the computation of a t-SNE map from the distance object in slot distances using the **Rtsne** package.

### Usage

```
comptsne(object, initial_cmd = TRUE, perplexity = 30, rseed = 15555)
```

**Arguments**

object	SCseq class object.
initial_cmd	logical. If TRUE, then the t-SNE map computation is initialized with a configuration obtained by classical multidimensional scaling. Default is TRUE.
perplexity	Positive number. Perplexity of the t-SNE map. Default is 30.
rseed	Integer number. Random seed to enforce reproducible t-SNE map.

**Value**

SCseq object with t-SNE coordinates stored in slot `tsne`.

**Examples**

```
sc <- SCseq(intestinalDataSmall)
sc <- filterdata(sc)
sc <- compdist(sc)
sc <- clustexp(sc)
sc <- findoutliers(sc)
sc <- comptsne(sc)
```

---

diffexpnb

---

*Function for differential expression analysis*


---

**Description**

This function performs differential expression analysis between two sets of single cell transcripts. The inference is based on a noise model or relies on the DESeq2 approach.

**Usage**

```
diffexpnb(x, A, B, DESeq = FALSE, method = "pooled", norm = FALSE,
          vfit = NULL, locreg = FALSE, ...)
```

**Arguments**

x	expression data frame with genes as rows and cells as columns. Gene IDs should be given as row names and cell IDs should be given as column names. This can be a reduced expression table only including the features (genes) to be used in the analysis. This input has to be provided if <code>g</code> (see below) is given and corresponds to a valid gene ID, i. e. one of the rownames of <code>x</code> . The default value is NULL. In this case, cluster identities are highlighted in the plot.
A	vector of cell IDs corresponding column names of <code>x</code> . Differential expression in set A versus set B will be evaluated.
B	vector of cell IDs corresponding column names of <code>x</code> . Differential expression in set A versus set B will be evaluated.

DESeq	logical value. If TRUE, then <b>DESeq2</b> is used for the inference of differentially expressed genes. In this case, it is recommended to provide non-normalized input data x. Default value is FALSE
method	either "per-condition" or "pooled". If DESeq is not used, this parameter determines, if the noise model is fitted for each set separately ("per-condition") or for the pooled set comprising all cells in A and B. Default value is "pooled".
norm	logical value. If TRUE then the total transcript count in each cell is normalized to the minimum number of transcripts across all cells in set A and B. Default value is FALSE.
vfit	function describing the background noise model. Inference of differentially expressed genes can be performed with a user-specified noise model describing the expression variance as a function of the mean expression. Default value is NULL.
locreg	logical value. If FALSE then regression of a second order polynomial is performed to determine the relation of variance and mean. If TRUE a local regression is performed instead. Default value is FALSE.
...	additional arguments to be passed to the low level function <code>DESeqDataSetFromMatrix</code> .

### Value

If DESeq equals TRUE, the function returns the output of **DESeq2**. In this case list of the following two components is returned:

cds	object returned by the <b>DESeq2</b> function <code>DESeqDataSetFromMatrix</code> .
res	data frame containing the results of the <b>DESeq2</b> analysis.

Otherwise, a list of three components is returned:

vf1	a data frame of three columns, indicating the mean $m$ , the variance $v$ and the fitted variance $v_m$ for set A.
vf2	a data frame of three columns, indicating the mean $m$ , the variance $v$ and the fitted variance $v_m$ for set B.
res	a data frame with the results of the differential gene expression analysis with the structure of the DESeq output, displaying mean expression of the two sets, fold change and log2 fold change between the two sets, the p-value for differential expression ( <code>pval</code> ) and the Benjamini-Hochberg corrected false discovery rate ( <code>padj</code> ).

### Examples

```
sc <- SCseq(intestinalDataSmall)
sc <- filterdata(sc)
sc <- compdist(sc)
sc <- clustexp(sc)
sc <- findoutliers(sc)
A <- names(sc@cpart)[sc@cpart %in% c(1,2)]
B <- names(sc@cpart)[sc@cpart %in% c(3)]
y <- diffexpnb(getfdata(sc,n=c(A,B)), A=A, B=B )
```

diffgenes

*Compute Expression Differences between Clusters***Description**

This functions computes expression differences between clusters and ranks genes by z-score differences.

**Usage**

```
diffgenes(object, c11, c12, mincount = 1)
```

**Arguments**

object	SCseq class object.
c11	A vector of valid cluster numbers (contained in the cpart slot of the SCseq object). Represents the first group of the comparison.
c12	A vector of valid cluster numbers (contained in the cpart slot of the SCseq object). Represents the second group of the comparison.
mincount	Minimal normalized expression level of a gene to be included into the analysis. A gene needs to be expressed at this level in at least a single cell.

**Value**

A list with four components:

z	a vector of z-scores in decreasing order with genes up-regulated in c11 appearing at the top of the list.
c11	a data.frame with expression values for cells in c11.
c12	a data.frame with expression values for cells in c12.
c11n	a vector of cluster numbers for cells in c11.
c12n	a vector of cluster numbers for cells in c12.

**Examples**

```
sc <- SCseq(intestinalDataSmall)
sc <- filterdata(sc)
sc <- compdist(sc)
sc <- clustexp(sc)
sc <- findoutliers(sc)
x <- diffgenes(sc,1,2)
head(x$z)
plotdiffgenes(x,names(x$z)[1])
```



---

filterdata	<i>Data filtering</i>
------------	-----------------------

---

### Description

This function allows filtering of genes and cells to be used in the RaceID3 analysis. It also can perform batch effect correction using an internal method or a recently published alternative `mnnCorrect` from the **scran** package.

### Usage

```
filterdata(object, mintotal = 3000, minexpr = 5, minnumber = 5,
           LBatch = NULL, knn = 10, CGenes = NULL, FGenes = NULL, ccor = 0.4,
           bmode = "RaceID")
```

### Arguments

object	SCseq class object.
mintotal	minimum total transcript number required. Cells with less than mintotal transcripts are filtered out. Default is 3000.
minexpr	minimum required transcript count of a gene in at least minnumber cells. All other genes are filtered out. Default is 5.
minnumber	See minexpr. Default is 1.
LBatch	List of experimental batches used for batch effect correction. Each list element contains a vector with cell names (i.e. column names of the input expression data) falling into this batch. Default is NULL, i.e. no batch correction.
knn	Number of nearest neighbors used to infer corresponding cell types in different batches. Default is 10.
CGenes	List of gene names. All genes with correlated expression to any of the genes in CGenes are filtered out for cell type inference. Default is NULL.
FGenes	List of gene names to be filtered out for cell type inference. Default is NULL.
ccor	Correlation coefficient used as a threshold for determining genes correlated to genes in CGenes. Only genes correlating less than ccor to all genes in CGenes are retained for analysis. Default is 0.4.
bmode	Method used for batch effect correction. Any of "RaceID", "scran". Default is "RaceID".

### Value

An SCseq class object with filtered and normalized expression data.

### Examples

```
sc <- SCseq(intestinalDataSmall)
sc <- filterdata(sc)
```

---

`findoutliers`*Inference of outlier cells and final clustering*

---

### Description

This functions performs the outlier identification based on the clusters inferred with the `clustexp` function.

### Usage

```
findoutliers(object, probthr = 0.001, outminc = 5, outlg = 2,
  outdistquant = 0.95)
```

### Arguments

<code>object</code>	SCseq class object.
<code>probthr</code>	outlier probability threshold for a minimum of <code>outlg</code> genes to be an outlier cell. This probability is computed from a negative binomial background model of expression in a cluster. Default is 0.001.
<code>outminc</code>	minimal transcript count of a gene in a clusters to be tested for being an outlier gene. Default is 5.
<code>outlg</code>	Minimum number of outlier genes required for being an outlier cell. Default is 2.
<code>outdistquant</code>	Real number between zero and one. Outlier cells are merged to outlier clusters if their distance smaller than the <code>outdistquant</code> -quantile of the distance distribution of pairs of cells in the original clusters after outlier removal. Default is 0.95.

### Value

SCseq object with outlier data stored in slot `out` and slot `outlierpar`. The final clustering partition is stored in `cpart`.

### Examples

```
sc <- SCseq(intestinalDataSmall)
sc <- filterdata(sc)
sc <- compdist(sc)
sc <- clustexp(sc)
sc <- findoutliers(sc)
```

---

<code>getfdata</code>	<i>Extracting filtered expression data</i>
-----------------------	--

---

**Description**

This functions allows the extraction of a filtered and normalized expression matrix

**Usage**

```
getfdata(object, g = NULL, n = NULL)
```

**Arguments**

<code>object</code>	SCseq class object.
<code>g</code>	Vector of gene names to be included corresponding to a subset of valid row names of the <code>ndata</code> slot of the SCseq object. Default is <code>NULL</code> and data for all genes remaining after filtering by the <code>filterdata</code> function are shown.
<code>n</code>	Vector of valid column names corresponding to a subset of valid column names of the <code>ndata</code> slot of the SCseq object. Default is <code>NULL</code> and data for all cells remaining after filtering by the <code>filterdata</code> function are shown.

**Value**

Matrix of filtered expression data with genes as rows and cells as columns.

---

<code>getproj</code>	<i>Extract Projections of all Cells from a Cluster</i>
----------------------	--

---

**Description**

This function extracts projections of all cells in a cluster and plots a heatmap of these hierarchically clustered projections (rows) to all other clusters (columns). A minimum spanning tree of the cluster centers is overlaid for comparison.

**Usage**

```
getproj(object, i, show = TRUE, zscore = FALSE)
```

**Arguments**

<code>object</code>	Ltree class object.
<code>i</code>	Cluster number. This number has to correspond to one of the RaceID3 clusters included for the StemID2 inference, i.e. to a number present in slot <code>ldata\$lp</code> .
<code>show</code>	logical. If <code>TRUE</code> , then plot heatmap of projections. Default is <code>TRUE</code> .
<code>zscore</code>	logical. If <code>TRUE</code> and <code>show=TRUE</code> , then plot z-score-transformed projections. If <code>TRUE</code> and <code>show=FALSE</code> , then plot untransformed projections. Default is <code>FALSE</code> .

**Value**

A list of two components:

- pr a data.frame of projections for all cells in cluster i (rows) onto all other clusters (columns).
- prz a data.frame of z-transformed projections for all cells in cluster i (rows) onto all other clusters (columns).

---

<code>imputeexp</code>	<i>Imputed expression matrix</i>
------------------------	----------------------------------

---

**Description**

This functions returns an imputed expression matrix based on the imputing computed with `compdist`.

**Usage**

```
imputeexp(object, genes = NULL)
```

**Arguments**

- `object` SCseq class object.
- `genes` vector of valid gene names corresponding to row names of slot `ndata`. Default is `NULL` and imputing is done for all genes.

**Value**

An expression matrix with imputed expression values after size normalization. Genes are in rows and cells in columns.

---

<code>intestinalData</code>	<i>Single-cell transcriptome data of intestinal epithelial cells</i>
-----------------------------	--

---

**Description**

This dataset contains gene expression values, i. e. transcript counts, of 278 intestinal epithelial cells.

**Usage**

```
intestinalData
```

**Format**

A sparse matrix (using the **Matrix**) with cells as columns and genes as rows. Entries are raw transcript counts.

**Value**

None

**References**

Grün et al. (2016) Cell Stem Cell 19(2): 266-77 <DOI:10.1016/j.stem.2016.05.010> ([PubMed](#))

---

intestinalDataSmall     *Single-cell transcriptome data of intestinal epithelial cells*

---

**Description**

This dataset is a smaller subset of the original dataset, which contains gene expression values, i. e. transcript counts, of 278 intestinal epithelial cells. The dataset is included for quick testing and examples. Only cells with >10,000 transcripts per cell and only genes with >20 transcript counts in >10 cells were retained.

**Usage**

```
intestinalDataSmall
```

**Format**

A sparse matrix (using the **Matrix**) with cells as columns and genes as rows. Entries are raw transcript counts.

**Value**

None

**References**

Grün et al. (2016) Cell Stem Cell 19(2): 266-77 <DOI:10.1016/j.stem.2016.05.010> ([PubMed](#))

---

lineagegraph     *Inference of a Lineage Graph*

---

**Description**

This function assembles a lineage graph based on the cell projections onto inter-cluster links.

**Usage**

```
lineagegraph(object)
```

**Arguments**

object            Ltree class object.

**Value**

An Ltree class object with lineage graph-related data stored in slots `ltcoord`, `prttree`, and `cdata`.

**Examples**

```
sc <- SCseq(intestinalDataSmall)
sc <- filterdata(sc)
sc <- compdist(sc)
sc <- clustexp(sc)
sc <- findoutliers(sc)
sc <- comptsne(sc)
ltr <- Ltree(sc)
ltr <- compentropy(ltr)
ltr <- projcells(ltr)
ltr <- lineagegraph(ltr)
```

---

Ltree-class

*The Ltree Class*


---

**Description**

The Ltree class is the central object storing all information generated during lineage tree inference by the StemID algorithm. It comprises a number of slots for a variety of objects.

validity function for Ltree

**Arguments**

object            An Ltree object.

**Slots**

`sc` An SCseq object with the RaceID3 analysis of the single-cell RNA-seq data for which a lineage tree should be derived.

`ldata` List object storing information on the clustering partition, the distance matrix, and the cluster centers in dimensionally-reduced input space and in two-dimensional t-sne space. Elements: `lp`: vector with the filtered partition into clusters after discarding clusters with `cthr` cells or less. `pdi`:matrix with the coordinates of all cells in the embedded space. Clusters with `cthr` transcripts or less were discarded (see function `projcells`). Rows are medoids and columns are coordinates. `cn`: data.frame with the coordinates of the cluster medoids in the embedded space. Clusters with `cthr` transcripts or less were discarded. Rows are medoids and columns are coordinates. `m`: vector with the numbers of the clusters which survived the filtering. `pdil`: data.frame with coordinates of cells in the two-dimensional t-SNE representation computed by RaceID3. Clusters with `cthr` transcripts or less were discarded. Rows are cells and columns

are coordinates. `cn1`: data.frame with the coordinates of the cluster medoids in the two-dimensional t-SNE representation computed by RaceID3. Clusters with `cthr` transcripts or less were discarded. Rows are medoids and columns are coordinates.

`entropy` Vector with transcriptome entropy computed for each cell.

`trproj` List containing two data.frames. Elements: `res`: data.frame with three columns for each cell. The first column `o` shows the cluster of a cell, the second column `l` shows the cluster number for the link the cell is assigned to, and the third column `h` shows the projection as a fraction of the length of the inter-cluster link. Parallel projections are positive, while anti-parallel projections are negative. `rma`: data.frame with all projection coordinates for each cell. Rows are cells and columns are clusters. Projections are given as a fraction of the length of the inter-cluster link. Parallel projections are positive, while anti-parallel projections are negative. The column corresponding to the originating cluster of a cell shows NA.

`par` List of parameters used for the StemID2 analysis.

`prback` data.frame of the same structure as the `trproj$res`. In case randomizations are used to compute significant projections, the projections of all `pdi$shuff` randomizations are appended to this data.frame and therefore the number of rows corresponds to the number of cells multiplied by `pdi$shuf`. See function `projback`.

`prbacka` data.frame reporting the aggregated results of the randomizations with four columns. Column `n` denotes the number of the randomization sample, column `o` and `l` contain the numbers of the originating and the terminal cluster, respectively, for each inter-cluster link and column `count` shows the number of cells assigned to this link in randomization sample `n`. The discrete distribution for the computation of the link p-value is given by the data contained in this object (if `nmode=FALSE`).

`ltcoord` Matrix storing projection coordinates of all cells in the two-dimensional t-SNE space, used for visualization.

`prtree` List with two elements. The first element `l` stores a list with the projection coordinates for each link. The name of each element identifies the link and is composed of two cluster numbers separated by a dot. The second element `n` is a list of the same structure and contains the cell names corresponding to the projection coordinates stored in `l`.

`cdata` list of data.frames, each with cluster ids as rows and columns: `counts` data.frame indicating the number of cells on the links connecting the cluster of origin (rows) to other clusters (columns). `counts.br` data.frame containing the cell counts on cluster connections averaged across the randomized background samples (if `nmode = FALSE`) or as derived from sampling statistics (if `nmode = TRUE`). `pv.e` matrix of enrichment p-values estimated from sampling statistics (if `nmode = TRUE`); entries are 0 if the observed number of cells on the respective link exceeds the  $(1 - \text{pethr})$ -quantile of the randomized background distribution and 0.5 otherwise (if `nmode = FALSE`). `pv.d` matrix of depletion p-values estimated from sampling statistics (if `nmode = TRUE`); entries are 0 if the observed number of cells on the respective link is lower than the `pethr`-quantile of the randomized background distribution and 0.5 otherwise (if `nmode = FALSE`). `pvn.e` matrix of enrichment p-values estimated from sampling statistics (if `nmode = TRUE`); 1-quantile, with the quantile estimated from the number of cells on a link as derived from the randomized background distribution (if `nmode = FALSE`). `pvn.d` matrix of depletion p-values estimated from sampling statistics (if `nmode = TRUE`); quantile estimated from the number of cells on a link as derived from the randomized background distribution (if `nmode = FALSE`).

plotbackground *Plot Background Model*

---

**Description**

This functions produces a scatter plot showing the gene expression variance as a function of the mean and the inferred polynomial fit of the background model computed by RaceID3. It also shows a local regression.

**Usage**

```
plotbackground(object)
```

**Arguments**

object            SCseq class object.

**Value**

None

---

plotdiffgenes *Barplot of differentially expressed genes*

---

**Description**

This functions produces a barplot of differentially expressed genes derived by the function `diffgenes`

**Usage**

```
plotdiffgenes(z, gene)
```

**Arguments**

z                    Output of `diffgenes`  
gene                Valid gene name. Has to correspond to one of the rownames of the `ndata` slot of the SCseq object.

**Value**

None



**Examples**

```

sc <- SCseq(intestinalDataSmall)
sc <- filterdata(sc)
sc <- compdist(sc)
sc <- clustexp(sc)
sc <- findoutliers(sc)
x <- diffgenes(sc,1,2)
head(x$z)
plotdiffgenes(x,names(x$z)[1])

```

---

plotdiffgenesnb      *Function for plotting differentially expressed genes*

---

**Description**

This is a plotting function for visualizing the output of the diffexpnb function.

**Usage**

```

plotdiffgenesnb(x, pthr = 0.05, padj = TRUE, lthr = 0, mthr = -Inf,
  Aname = NULL, Bname = NULL, show_names = TRUE)

```

**Arguments**

x	output of the function diffexpnb.
pthr	real number between 0 and 1. This number represents the p-value cutoff applied for displaying differentially expressed genes. Default value is 0.05. The parameter padj (see below) determines if this cutoff is applied to the uncorrected p-value or to the Benjamini-Hochberg corrected false discovery rate.
padj	logical value. If TRUE, then genes with a Benjamini-Hochberg corrected false discovery rate lower than pthr are displayed. If FALSE, then genes with a p-value lower than pthr are displayed.
lthr	real number between 0 and Inf. Differentially expressed genes are displayed only for log2 fold-changes greater than lthr. Default value is 0.
mthr	real number between -Inf and Inf. Differentially expressed genes are displayed only for log2 mean expression greater than mthr. Default value is -Inf.
Aname	name of expression set A, which was used as input to diffexpnb. If provided, this name is used in the axis labels. Default value is NULL.
Bname	name of expression set B, which was used as input to diffexpnb. If provided, this name is used in the axis labels. Default value is NULL.
show_names	logical value. If TRUE then gene names displayed for differentially expressed genes. Default value is FALSE.

**Value**

None

**Examples**

```

sc <- SCseq(intestinalDataSmall)
sc <- filterdata(sc)
sc <- compdist(sc)
sc <- clustexp(sc)
sc <- findoutliers(sc)
A <- names(sc@cpart)[sc@cpart %in% c(1,2)]
B <- names(sc@cpart)[sc@cpart %in% c(3)]
y <- diffexpnb(getfdata(sc,n=c(A,B)), A=A, B=B )
plotdiffgenesnb(y)

```

---

plotdimsat

*Plotting the Saturation of Explained Variance*


---

**Description**

This functions plots the explained variance as a function of PCA/ICA components computed by the function CCcorrect. The number of components where the change in explained variability upon adding further components approaches linear behaviour demarcates the saturation point and is highlighted in blue.

**Usage**

```
plotdimsat(object, change = TRUE, lim = NULL)
```

**Arguments**

object	SCseq class object.
change	logical. If TRUE then the change in explained variance is plotted. Default is FALSE and the explained variance is shown.
lim	Number of components included for he calculation and shown in the plot. Default is NULL and all components are included.

**Value**

None

---

plotdistanceratio      *Histogram of Cell-to-Cell Distances in Real versus Embedded Space*

---

### Description

This function plots a histogram of the ratios of cell-to-cell distances in the original versus the high-dimensional embedded space used as input for the StemID2 inferences. The embedded space approximates correlation-based distances by Euclidean distances obtained by classical multi-dimensional scaling. A minimum spanning tree of the cluster centers is overlaid for comparison.

### Usage

```
plotdistanceratio(object)
```

### Arguments

object              Ltree class object.

### Value

None.

---

plotexpmap              *Highlighting gene expression in the t-SNE map*

---

### Description

This functions highlights gene expression in a two-dimensional t-SNE map or a Fruchterman-Rheingold graph layout of the single-cell transcriptome data.

### Usage

```
plotexpmap(object, g, n = NULL, logsc = FALSE, imputed = FALSE,
            fr = FALSE, cells = NULL, cex = 1, map = TRUE, leg = TRUE)
```

### Arguments

object              SCseq class object.

g                    Individual gene name or vector with a group of gene names corresponding to a subset of row names of the ndata slot of the SCseq object.

n                    String of characters representing the title of the plot. Default is NULL and the first element of g is chosen.

logsc               logical. If TRUE, then gene expression values are log<sub>2</sub>-transformed after adding a pseudo-count of 0.1. Default is FALSE and untransformed values are shown.

imputed	logical. If TRUE and imputing was done by calling compdist with <code>knn &gt; 0</code> , then imputed expression values are shown. If FALSE, then raw counts are shown. Default is FALSE.
fr	logical. If TRUE then plot t-SNE map, else plot Fruchterman-Rheingold layout.
cells	Vector of valid cell names corresponding to column names of slot <code>nData</code> of the <code>SCseq</code> object. Gene expression is only shown for this subset.
cex	size of data points. Default value is 1.
map	logical. If TRUE then data points are shown. Default value is TRUE.
leg	logical. If TRUE then the legend is shown. Default value is TRUE.

**Value**

None

---

plotgraph	<i>StemID2 Lineage Graph</i>
-----------	------------------------------

---

**Description**

This function plots a graph of lineage trajectories connecting RaceID3 cluster medoids as inferred by StemID2 to approximate the lineage tree. The plot highlights significant links, where colour indicates the level of significance and width indicates the link score. The node colour reflects the level of transcriptome entropy.

**Usage**

```
plotgraph(object, showCells = FALSE, showTsne = TRUE, tp = 0.5,
          scthr = 0)
```

**Arguments**

object	Ltree class object.
showCells	logical. If TRUE, then projections of cells are shown in the plot. Default is FALSE.
showTsne	logical. If TRUE, then show transparent t-SNE map (with transparency <code>tp</code> ) of cells in the background. Default is TRUE.
tp	Real number between zero and one. Level of transparency of the t-SNE map. Default is 0.5. See <code>showTsne</code> .
scthr	Real number between zero and one. Score threshold for links to be shown in the graph. For <code>scthr=0</code> all significant links are shown. The maximum score is one.

**Value**

None.

---

plotjaccard	<i>Plot Jaccard Similarities</i>
-------------	----------------------------------

---

**Description**

This functions plots a barchart of Jaccard similarities for the RaceID3 clusters before outlier identification

**Usage**

```
plotjaccard(object)
```

**Arguments**

object            SCseq class object.

**Value**

None

---

plotlabelsmap	<i>Plot labels in the t-SNE map</i>
---------------	-------------------------------------

---

**Description**

This functions plots cell labels into a two-dimensional t-SNE map or a Fruchterman-Rheingold graph layout of the single-cell transcriptome data.

**Usage**

```
plotlabelsmap(object, labels = NULL, fr = FALSE)
```

**Arguments**

object            SCseq class object.

labels            Vector of labels for all cells to be highlighted in the t-SNE map. The order has to be the same as for the columns in slot `ndata` of the SCseq object. Default is `NULL` and cell names are highlighted.

fr                logical. If `TRUE` then plot t-SNE map, else plot Fruchterman-Rheingold layout.

**Value**

None

---

`plotlinkpv`*Heatmap of Link P-values*

---

**Description**

This function plots a heatmap of link p-values.

**Usage**

```
plotlinkpv(object)
```

**Arguments**

`object`            Ltree class object.

**Value**

None.

---

`plotlinkscore`*Heatmap of Link Scores*

---

**Description**

This function plots a heatmap of link score.

**Usage**

```
plotlinkscore(object)
```

**Arguments**

`object`            Ltree class object.

**Value**

None.

---

plotmap	<i>Plotting a t-SNE map</i>
---------	-----------------------------

---

**Description**

This functions plots a two-dimensional t-SNE map or a Fruchterman-Rheingold graph layout of the single-cell transcriptome data.

**Usage**

```
plotmap(object, final = TRUE, tp = 1, fr = FALSE, cex = 0.5)
```

**Arguments**

object	SCseq class object.
final	logical. If TRUE, then highlight final clusters after outlier identification. If FALSE, then highlight initial clusters prior to outlier identification. Default is TRUE.
tp	Number between 0 and 1 to change transparency of dots in the map. Default is 1.
fr	logical. If TRUE then plot t-SNE map, else plot Fruchterman-Rheingold layout.
cex	size of data points. Default value is 0.5.

**Value**

None

---

plotmarkergenes	<i>Plotting a Heatmap of Marker Gene Expression</i>
-----------------	---

---

**Description**

This functions generates a heatmap of expression for defined group of genes and can highlight the clustering partition and another sample grouping, e.g. origin or cell type.

**Usage**

```
plotmarkergenes(object, genes, imputed = FALSE, cthr = 0, cl = NULL,
  cells = NULL, order.cells = FALSE, aggr = FALSE, norm = FALSE,
  cap = NULL, flo = NULL, samples = NULL, cluster_cols = FALSE,
  cluster_rows = TRUE, cluster_set = FALSE, samples_col = NULL,
  zsc = FALSE, logscale = TRUE)
```

**Arguments**

object	SCseq class object.
genes	A vector with a group of gene names corresponding to a subset of valid row names of the <code>ndata</code> slot of the SCseq object.
imputed	logical. If TRUE and imputing was done by calling <code>compdist</code> with <code>knn &gt; 0</code> , then imputed expression values are shown. If FALSE, then raw counts are shown. Default is FALSE
cthr	Integer number greater or equal zero. Only clusters with <code>&gt;cthr</code> cells are included in the t-SNE map. Default is 0.
c1	Vector of valid cluster numbers contained in slot <code>cpart</code> of the SCseq object. Default is NULL and all clusters with <code>&gt;cthr</code> cells are included.
cells	Vector of valid cell names corresponding to column names of slot <code>ndata</code> of the SCseq object. Gene expression is only shown for this subset. Default is NULL and all cells are included. The set of <code>cells</code> is intersected with the subset of clusters in <code>c1</code> if given.
order.cells	logical. If TRUE, then columns of the heatmap are ordered by cell name and not by cluster number. If <code>cells</code> are given, then columns are ordered as in <code>cells</code> .
aggr	logical. If TRUE, then only average expression is shown for each cluster. Default is FALSE and expression in individual cells is shown.
norm	logical. If TRUE, then expression of each gene across clusters is normalized to 1, in order to depict all genes on the same scale. Default is FALSE.
cap	Numeric. Upper bound for gene expression. All values larger then <code>cap</code> are replaced by <code>cap</code> . Default is NULL and no <code>cap</code> is applied.
flo	Numeric. Lower bound for gene expression. All values smaller then <code>floor</code> are replaced by <code>floor</code> . Default is NULL and no <code>floor</code> is applied.
samples	A vector with a group of sample names for each cell in the same order as the column names of the <code>ndata</code> slot of the SCseq object.
cluster_cols	logical. If TRUE, then columns are clustered. Default is FALSE.
cluster_rows	logical. If TRUE, then rows are clustered. Default is TRUE.
cluster_set	logical. If TRUE then clusters are ordered by hierarchical clustering of the cluster medoids.
samples_col	Vector of colors used for highlighting all samples contained in <code>samples</code> in the heatmap. Default is NULL.
zsc	logical. If TRUE then a z-score transformation is applied. Default is FALSE.
logscale	logical. If TRUE then a log2 transformation is applied. Default is TRUE.

**Value**

None



---

plotoutlierprobs      *Plot Outlier Probabilities*

---

**Description**

This functions plots a barchart of outlier probabilities across all cells in each cluster.

**Usage**

```
plotoutlierprobs(object)
```

**Arguments**

object              SCseq class object.

**Value**

None

---

plotprojections      *Two-Dimensional Map of Cell Projections*

---

**Description**

This function plots the projections of cells on inter-clustering links connecting cluster medoids in a two-dimensional t-SNE representation. A minimum spanning tree of the cluster centers is overlaid for comparison.

**Usage**

```
plotprojections(object)
```

**Arguments**

object              Ltree class object.

**Value**

None.

---

plotsaturation      *Plot Saturation of Within-Cluster Dispersion*

---

**Description**

This functions plots the (change in the) mean within-cluster dispersion as a function of the cluster number and highlights the saturation point inferred based on the saturation criterion applied by RaceID3: The number of clusters where the change in within-cluster dispersion upon adding further clusters approaches linear behaviour demarcates the saturation point and is highlighted in blue.

**Usage**

```
plotsaturation(object, disp = FALSE)
```

**Arguments**

object	SCseq class object.
disp	logical. If FALSE, then the change of the within-cluster dispersion is plotted. if TRUE the actual dispersion is plotted. Default is FALSE

**Value**

None

---

---

plotsensitivity      *Plot Sensitivity*

---

**Description**

This functions plots the number of outliers as a function of the outlier probability.

**Usage**

```
plotsensitivity(object)
```

**Arguments**

object	SCseq class object.
--------	---------------------

**Value**

None

---

plotsilhouette	<i>Plot Cluster Silhouette</i>
----------------	--------------------------------

---

**Description**

This functions produces a silhouette plot for RaceID3 clusters prior or post outlier identification.

**Usage**

```
plotsilhouette(object, final = FALSE)
```

**Arguments**

object	SCseq class object.
final	logical. If TRUE, then plot silhouette coefficients for final clusters after outlier identification. Default is FALSE and silhouette coefficients are plotted for initial clusters.

**Value**

None

---

plotspanntree	<i>Minimum Spanning Tree of RaceID3 clusters</i>
---------------	--

---

**Description**

This function plots a minimum spanning tree of the RaceID3 cluster medoids in a two-dimensional t-SNE representation.

**Usage**

```
plotspanntree(object)
```

**Arguments**

object	Ltree class object.
--------	---------------------

**Value**

None.

---

plotsymbolsmap	<i>Plotting groups as different symbols in the t-SNE map</i>
----------------	--

---

### Description

This functions highlights groups of cells by different symbols in a two-dimensional t-SNE map or a Fruchterman-Rheingold graph layout of the single-cell transcriptome data.

### Usage

```
plotsymbolsmap(object, types, subset = NULL, samples_col = NULL,
               cex = 0.25, fr = FALSE, leg = TRUE, map = TRUE)
```

### Arguments

object	SCseq class object.
types	Vector assigning each cell to a type to be highlighted in the t-SNE map. The order has to be the same as for the columns in slot ndata of the SCseq object. Default is NULL and each cell is highlighted by a different symbol.
subset	Vector containing a subset of types from types to be highlighted in the map. Default is NULL and all types are shown.
samples_col	Vector of colors used for highlighting all samples contained in samples in the map. Default is NULL.
cex	size of data points. Default value is 0.25.
fr	logical. If TRUE then plot t-SNE map, else plot Fruchterman-Rheingold layout.
leg	logical. If TRUE then the legend is shown. Default value is TRUE.
map	logical. If TRUE then data points are shown. Default value is TRUE.

### Value

None

---

projback	<i>Compute Cell Projections for Randomized Background Distribution</i>
----------	--

---

### Description

This function computes the projections of cells onto inter-cluster links for randomized cell positions in a high-dimensional embedded space. Significance of link based on an increased number of cells on a link is inferred based on this background model.

### Usage

```
projback(object, pdishuf = 500, fast = FALSE, rseed = 17000)
```

**Arguments**

object	Ltree class object.
pdishuf	Number of randomizations of cell positions for which to compute projections of cells on inter-cluster links. Default is 2000. No randomizations are needed in this mode and the function will do nothing. Default is TRUE.
fast	logical. If TRUE and nmode=FALSE cells will still be assigned to links based on maximum projections but a fast approximate background model will be used to infer significance. The function will do nothing in this case. Default is FALSE.
rseed	Integer number used as seed to ensure reproducibility of randomizations. Default is 17000.

**Value**

An Ltree class object with all information on randomized cell projections onto links stored in the prbacka slot.

**Examples**

```
sc <- SCseq(intestinalDataSmall)
sc <- filterdata(sc)
sc <- compdist(sc)
sc <- clustexp(sc)
sc <- findoutliers(sc)
sc <- comptsne(sc)
ltr <- Ltree(sc)
ltr <- compentropy(ltr)
ltr <- projcells(ltr, nmode=FALSE)
ltr <- projback(ltr, pdishuf=50)
```

---

projcells

*Compute transcriptome entropy of each cell*

---

**Description**

This function computes the projections of cells onto inter-cluster links in a high-dimensional embedded space.

**Usage**

```
projcells(object, cthr = 5, nmode = TRUE, knn = 3, fr = FALSE)
```

**Arguments**

object	Ltree class object.
cthr	Positive integer number. Clusters to be included into the StemID2 analysis must contain more than cthr cells. Default is 5.

nmode	logical. If TRUE, then a cell of given cluster is assigned to the link to the cluster with the smallest average distance of the knn nearest neighbours within this cluster. Default is TRUE.
knn	Positive integer number. See nmode. Default is 3.
fr	logical. Use Fruchterman-Rheingold layout instead of t-SNE for dimensional-reduction representation of the lineage tree. Default is FALSE.

**Value**

An Ltree class object with all information on cell projections onto links stored in the ldata slot.

**Examples**

```
sc <- SCseq(intestinalDataSmall)
sc <- filterdata(sc)
sc <- compdist(sc)
sc <- clustexp(sc)
sc <- findoutliers(sc)
sc <- comptsne(sc)
ltr <- Ltree(sc)
ltr <- compentropy(ltr)
ltr <- projcells(ltr)
```

---

 projenrichment

*Enrichment of cells on inter-cluster links*


---

**Description**

This function plots a heatmap of the enrichment ratios of cells on significant links.

**Usage**

```
projenrichment(object)
```

**Arguments**

object            Ltree class object.

**Value**

None.

---

`rfcorrect`*Random Forests-based Reclassification*

---

## Description

This functions applies random forests-based reclassification of cell clusters to enhance robustness of the final clusters.

## Usage

```
rfcorrect(object, rfseed = 12345, nbtree = NULL, final = TRUE,  
          nbfactor = 5, ...)
```

## Arguments

<code>object</code>	SCseq class object.
<code>rfseed</code>	Seed for enforcing reproducible results. Default is 12345.
<code>nbtree</code>	Number of trees to be built. Default is NULL and the number of tree is given by the number of cells times <code>nbfactor</code> .
<code>final</code>	logical. If TRUE, then reclassification of cell types using out-of-bag analysis is performed based on the final clusters after outlier identification. If FALSE, then the cluster partition prior to outlier identification is used for reclassification.
<code>nbfactor</code>	Positive integer number. See <code>nbtree</code> .
<code>...</code>	additional input arguments to the <code>randomForest</code> function of the <b>randomForest</b> package.

## Value

The function returns an updated SCseq object with random forests votes written to slot `out$rfvotes`. The clustering partition prior or post outlier identification (slot `cluster$kpart` or `cpart`, if parameter `final` equals FALSE or TRUE, respectively) is overwritten with the partition derived from the reclassification.

## Examples

```
sc <- SCseq(intestinalDataSmall)  
sc <- filterdata(sc)  
sc <- compdist(sc)  
sc <- clustexp(sc)  
sc <- findoutliers(sc)  
sc <- rfcorrect(sc)
```

SCseq

*The SCseq Class***Description**

The SCseq class is the central object storing all information generated during cell type identification with the RaceID3 algorithm. It comprises a number of slots for a variety of objects.

validity function for SCseq

**Arguments**

object            An SCseq object.

**Slots**

expdata The raw expression data matrix with cells as columns and genes as rows in sparse matrix format.

ndata Filtered data with expression normalized to one for each cell.

counts Vector with total transcript counts for each cell in ndata remaining after filtering.

genes Vector with gene names of all genes in ndata remaining after filtering.

dimRed list object object storing information on a feature matrix obtained by dimensional reduction, batch effect correction etc. Component x stores the actual feature matrix.

distances distance (or dis-similarity) matrix computed by RaceID3.

imputed list with two matrices computed for imputing gene expression. The first matrix nn contains the cell indices of the knn nearest neighbours, the second matrix contains the probabilities at which each cell contributes to the imputed gene expression value for the cell corresponding to the columns.

tsne data.frame with coordinates of two-dimensional tsne layout computed by RaceID3.

fr data.frame with coordinates of two-dimensional Fruchterman-Rheingold graphlayout computed by RaceID3.

cluster list storing information on the initial clustering step of the RaceID3 algorithm

background list storing the polynomial fit for the background model of gene expression variability computed by RaceID3, which is used for outlier identification.

out list storing information on outlier cells used for the prediction of rare cell types by RaceID3

cpart vector containing the final clustering (i.e. cell type) partition computed by RaceID3

fc01 vector containing the colour scheme for the RaceID3 clusters

medoids vector containing the cell ids for the cluster medoids

filterpar list containing the parameters used for cell and gene filtering

clusterpar list containing the parameters used for clustering

outlierpar list containing the parameters used for outlier identification



---

varRegression	<i>Linear Regression of Sources of Variability</i>
---------------	--

---

**Description**

This functions regresses out variability associated with particular sources.

**Usage**

```
varRegression(object, vars = NULL, logscale = FALSE, Batch = FALSE)
```

**Arguments**

object	SCseq class object.
vars	data.frame of variables to be regressed out. Each column corresponds to a variable and each variable corresponds to a cell. The object must contain all cells, i.e. column names of the slot <code>ndata</code> from the SCseq object.
logscale	logical. If TRUE data are log-transformed prior to regression. Default is FALSE.
Batch	logical. If TRUE, then the function will regress out batch-associated variability based on genes stored in the <code>filterpar\$BGenes</code> slot of the SCseq object. This requires prior batch correction with the <code>filterdata</code> function using <code>bmode="RaceID"</code> .

**Value**

The function returns an updated SCseq object with the corrected expression matrix written to the slot `dimRed$x` of the SCseq object.

**Examples**

```
sc <- SCseq(intestinalDataSmall)
sc <- filterdata(sc)
b <- sub("(\\_\\d+)$", "", colnames(intestinalData))
vars <- data.frame(row.names=colnames(intestinalData), batch=b)
sc <- varRegression(sc, vars)
```

# Index

## \*Topic **datasets**

- intestinalData, [20](#)
- intestinalDataSmall, [21](#)
  
- barplotgene, [3](#)
- branchcells, [3](#)
  
- CCcorrect, [4](#)
- cellsfromtree, [5](#)
- clustdiffgenes, [6](#)
- clustexp, [7](#)
- clustheatmap, [8](#)
- compdist, [9](#)
- compentropy, [10](#)
- compfr, [10](#)
- compmedoids, [11](#)
- comppvalue, [12](#)
- compscore, [13](#)
- comptsne, [13](#)
  
- diffexpnb, [14](#)
- diffgenes, [16](#)
  
- filterdata, [17](#)
- findoutliers, [18](#)
  
- getfdata, [19](#)
- getproj, [19](#)
  
- imputeexp, [20](#)
- intestinalData, [20](#)
- intestinalDataSmall, [21](#)
  
- lineagegraph, [21](#)
- Ltree (Ltree-class), [22](#)
- Ltree-class, [22](#)
  
- plotbackground, [24](#)
- plotdiffgenes, [24](#)
- plotdiffgenesnb, [25](#)
- plotdimsat, [26](#)
  
- plotdistanceratio, [27](#)
- plotexpmap, [27](#)
- plotgraph, [28](#)
- plotjaccard, [29](#)
- plotlabelsmap, [29](#)
- plotlinkpv, [30](#)
- plotlinkscore, [30](#)
- plotmap, [31](#)
- plotmarkergenes, [31](#)
- plotoutlierprobs, [33](#)
- plotprojections, [33](#)
- plotsaturation, [34](#)
- plotsensitivity, [34](#)
- plotsilhouette, [35](#)
- plotspantree, [35](#)
- plotsymbolsmap, [36](#)
- projback, [36](#)
- projcells, [37](#)
- projenrichment, [38](#)
  
- rfcorrect, [39](#)
  
- SCseq, [40](#)
- SCseq-class (SCseq), [40](#)
  
- varRegression, [41](#)