

# Package ‘scrapper’

January 4, 2025

**Version** 1.1.11

**Date** 2025-01-03

**Title** Bindings to C++ Libraries for Single-Cell Analysis

**Description** Implements R bindings to C++ code for analyzing single-cell (expression) data, mostly from various libscrans libraries. Each function performs an individual step in the single-cell analysis workflow, ranging from quality control to clustering and marker detection. It is mostly intended for other Bioconductor package developers to build more user-friendly end-to-end workflows.

**License** MIT + file LICENSE

**Imports** methods, Rcpp, beachmat (>= 2.21.6), DelayedArray, BiocNeighbors (>= 1.99.0), igraph, parallel

**Suggests** testthat, knitr, rmarkdown, BiocStyle, MatrixGenerics, sparseMatrixStats, Matrix, scRNAseq

**LinkingTo** Rcpp, assorthead, beachmat, BiocNeighbors

**biocViews** Normalization, RNASeq, Software, GeneExpression, Transcriptomics, SingleCell, BatchEffect, QualityControl, DifferentialExpression, FeatureExtraction, PrincipalComponent, Clustering

**VignetteBuilder** knitr

**Encoding** UTF-8

**RoxygenNote** 7.3.2

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

**git\_branch** devel

**git\_last\_commit** 644f28b

**git\_last\_commit\_date** 2025-01-03

**Repository** Bioconductor 3.21

**Date/Publication** 2025-01-03

**Author** Aaron Lun [cre, aut]

**Maintainer** Aaron Lun <[infinite.monkeys.with.keyboards@gmail.com](mailto:infinite.monkeys.with.keyboards@gmail.com)>

## Contents

adt_quality_control . . . . .	2
aggregateAcrossCells . . . . .	4
aggregateAcrossGenes . . . . .	6
buildSnnGraph . . . . .	7
centerSizeFactors . . . . .	8
chooseHighlyVariableGenes . . . . .	9
choosePseudoCount . . . . .	10
clusterGraph . . . . .	11
clusterKmeans . . . . .	13
combineFactors . . . . .	14
computeClrm1Factors . . . . .	15
correctMnn . . . . .	16
crispr_quality_control . . . . .	18
fitVarianceTrend . . . . .	20
modelGeneVariances . . . . .	21
normalizeCounts . . . . .	23
rna_quality_control . . . . .	24
runAllNeighborSteps . . . . .	26
runPca . . . . .	28
runTsne . . . . .	30
runUmap . . . . .	31
sanitizeSizeFactors . . . . .	33
scaleByNeighbors . . . . .	34
scoreGeneSet . . . . .	35
scoreMarkers . . . . .	37
subsampleByNeighbors . . . . .	39
summarizeEffects . . . . .	41
testEnrichment . . . . .	42
<b>Index</b>	<b>44</b>

---

adt\_quality\_control     *Quality control for ADT count data*

---

### Description

Compute per-cell QC metrics from an initialized matrix of ADT counts, and use the metrics to suggest filter thresholds to retain high-quality cells.

### Usage

```
computeAdtQcMetrics(x, subsets, num.threads = 1)

suggestAdtQcThresholds(
  metrics,
  block = NULL,
```

```

    min.detected.drop = 0.1,
    num.mads = 3
)

filterAdtQcMetrics(thresholds, metrics, block = NULL)

```

### Arguments

x	A matrix-like object where rows are ADTs and columns are cells. Values are expected to be counts.
subsets	List of vectors specifying tag subsets of interest, typically control tags like IgGs. Each vector may be logical (whether to keep each row), integer (row indices) or character (row names).
num.threads	Integer scalar specifying the number of threads to use.
metrics	List with the same structure as produced by <code>computeAdtQcMetrics</code> .
block	Factor specifying the block of origin (e.g., batch, sample) for each cell in <code>metrics</code> . Alternatively NULL if all cells are from the same block. For <code>filterAdtQcMetrics</code> , a blocking factor should be provided if <code>block</code> was used to construct thresholds.
min.detected.drop	Minimum drop in the number of detected features from the median, in order to consider a cell to be of low quality.
num.mads	Number of median from the median, to define the threshold for outliers in each metric.
thresholds	List with the same structure as produced by <code>suggestAdtQcThresholds</code> .

### Value

For `computeAdtQcMetrics`, a list is returned containing:

- `sum`, a numeric vector containing the total ADT count for each cell.
- `detected`, an integer vector containing the number of detected tags per cell.
- `subsets`, a list of numeric vectors containing the total count of each control subset.

Each vector is of length equal to the number of cells.

For `suggestAdtQcThresholds`, a named list is returned:

- If `block=NULL`, the list contains:
  - `detected`, a numeric scalar containing the lower bound on the number of detected tags.
  - `subsets`, a numeric vector containing the upper bound on the sum of counts in each control subset.
- Otherwise, if `block` is supplied, the list contains:
  - `detected`, a numeric vector containing the lower bound on the number of detected tags for each blocking level.
  - `subsets`, a list of numeric vectors containing the upper bound on the sum of counts in each control subset for each blocking level.

Each vector is of length equal to the number of levels in block and is named accordingly.

For `filterAdtQcMetrics`, a logical vector of length `ncol(x)` is returned indicating which cells are of high quality.

### Author(s)

Aaron Lun

### See Also

[https://libscran.github.io/scran\\_qc/](https://libscran.github.io/scran_qc/), for the rationale of QC filtering on ADT counts.

### Examples

```
# Mocking a matrix:
library(Matrix)
x <- round(abs(rsparsematrix(1000, 100, 0.1) * 100))

# Mocking up a control set.
sub <- list(IgG=rbinom(nrow(x), 1, 0.1) > 0)

qc <- computeAdtQcMetrics(x, sub)
str(qc)

filt <- suggestAdtQcThresholds(qc)
str(filt)

keep <- filterAdtQcMetrics(filt, qc)
summary(keep)
```

---

aggregateAcrossCells *Aggregate expression across cells*

---

### Description

Aggregate expression values across cells based on one or more grouping factors. This is primarily used to create pseudo-bulk profiles for each cluster/sample combination.

### Usage

```
aggregateAcrossCells(x, factors, num.threads = 1)
```

### Arguments

<code>x</code>	A matrix-like object where rows correspond to genes or genomic features and columns correspond to cells. Values are typically expected to be counts.
<code>factors</code>	A list or data frame containing one or more grouping factors, see <a href="#">combineFactors</a> .
<code>num.threads</code>	Integer specifying the number of threads to be used for aggregation.

**Value**

A list containing:

- `sums`, a numeric matrix where each row corresponds to a gene and each column corresponds to a unique combination of grouping levels. Each entry contains the summed expression across all cells with that combination.
- `detected`, an integer matrix where each row corresponds to a gene and each column corresponds to a unique combination of grouping levels. Each entry contains the number of cells with detected expression in that combination.
- `combinations`, a data frame describing the levels for each unique combination of factors. Rows of this data frame correspond to columns of `sums` and `detected`, while columns correspond to the factors in `factors`.
- `counts`, the number of cells associated with each combination. Each entry corresponds to a row of `combinations`.
- `index`, an integer vector of length equal to the number of cells in `x`. This specifies the combination in `combinations` to which each cell was assigned.

**Author(s)**

Aaron Lun

**See Also**

[aggregateAcrossGenes](#), to aggregate expression values across gene sets.

**Examples**

```
# Mocking a matrix:
library(Matrix)
x <- round(abs(rsparsematrix(1000, 100, 0.1) * 100))

# Simple aggregation:
clusters <- sample(LETTERS, 100, replace=TRUE)
agg <- aggregateAcrossCells(x, list(cluster=clusters))
str(agg)

# Multi-factor aggregation
samples <- sample(1:5, 100, replace=TRUE)
agg2 <- aggregateAcrossCells(x, list(cluster=clusters, sample=samples))
str(agg2)
```

---

aggregateAcrossGenes *Aggregate expression across genes*

---

### Description

Aggregate expression values across genes, potentially with weights. This is typically used to summarize expression values for gene sets into a single per-cell score.

### Usage

```
aggregateAcrossGenes(x, sets, average = FALSE, num.threads = 1)
```

### Arguments

x	A matrix-like object where rows correspond to genes or genomic features and columns correspond to cells. Values are typically expected to be counts.
sets	A list of integer vectors containing the row indices of genes in each set. Alternatively, each entry may be a list of length 2, containing an integer vector (row indices) and a numeric vector (weights).
average	Logical scalar indicating whether to compute the average rather than the sum.
num.threads	Integer specifying the number of threads to be used for aggregation.

### Value

A list of length equal to that of sets. Each entry is a numeric vector of length equal to the number of columns in x, containing the (weighted) sum/mean of expression values for the corresponding set across all cells.

### Author(s)

Aaron Lun

### See Also

[aggregateAcrossCells](#), to aggregate expression values across groups of cells.

### Examples

```
# Mocking a matrix:
library(Matrix)
x <- round(abs(rsparsematrix(1000, 100, 0.1) * 100))

# Unweighted aggregation:
sets <- list(
  foo = sample(nrow(x), 20),
  bar = sample(nrow(x), 10)
)
agg <- aggregateAcrossGenes(x, sets)
```

```

str(agg)

# Weighted aggregation:
sets <- list(
  foo = list(sample(nrow(x), 20), runif(20)),
  bar = list(sample(nrow(x), 10), runif(10))
)
agg2 <- aggregateAcrossGenes(x, sets, average = TRUE)
str(agg2)

```

---

buildSnnGraph      *Build a shared nearest neighbor graph*

---

### Description

Build a shared nearest neighbor (SNN) graph where each node is a cell. Edges are formed between cells that share one or more nearest neighbors, weighted by the number or importance of those shared neighbors.

### Usage

```

buildSnnGraph(
  x,
  num.neighbors = 10,
  weight.scheme = "ranked",
  num.threads = 1,
  BNPARAM = AnnoyParam()
)

```

### Arguments

x	For buildSnnGraph, a numeric matrix where rows are dimensions and columns are cells, typically containing a low-dimensional representation from, e.g., <a href="#">runPca</a> . Alternatively, a named list of nearest-neighbor search results. This should contain index, an integer matrix where rows are neighbors and columns are cells. Each column contains 1-based indices for the nearest neighbors of the corresponding cell, ordered by increasing distance. The number of neighbors for each cell should be equal to num.neighbors, otherwise a warning is raised. Alternatively, an index constructed by <a href="#">buildIndex</a> .
num.neighbors	Integer scalar specifying the number of neighbors to use to construct the graph.
weight.scheme	String specifying the weighting scheme to use for constructing the SNN graph. This can be "ranked" (default), "jaccard" or "number".
num.threads	Integer scalar specifying the number of threads to use. Only used if x is not a list of existing nearest-neighbor search results.
BNPARAM	A <a href="#">BiocNeighborParam</a> object specifying the algorithm to use. Only used if x is not a list of existing nearest-neighbor search results.

**Value**

If `as.pointer=FALSE`, a list is returned containing:

- `vertices`, an integer scalar specifying the number of vertices in the graph (i.e., cells in `x`).
- `edges`, an integer vector of 1-based indices for graph edges. Pairs of values represent the endpoints of an (undirected) edge, i.e., `edges[1:2]` form the first edge, `edges[3:4]` form the second edge and so on.
- `weights`, a numeric vector of weights for each edge. This has length equal to half the length of edges.

**Author(s)**

Aaron Lun

**See Also**

[https://libscran.github.io/scran\\_graph\\_cluster/](https://libscran.github.io/scran_graph_cluster/), for details on the underlying implementation.

**Examples**

```
data <- matrix(rnorm(10000), ncol=1000)
out <- buildSnnGraph(data)
str(out)

# We can use this to make an igraph::graph.
g <- igraph::make_undirected_graph(out$edges, n = out$vertices)
igraph::E(g)$weight <- out$weight
```

---

centerSizeFactors      *Center size factors*

---

**Description**

Scale the size factors so they are centered at unity, which ensures that the scale of the counts is preserved (on average) after normalization.

**Usage**

```
centerSizeFactors(size.factors, block = NULL, mode = c("lowest", "per-block"))
```



**Arguments**

size.factors	Numeric vector of size factors across cells.
block	Vector or factor of length equal to size.factors, specifying the block of origin for each cell. Alternatively NULL, in which case all cells are assumed to be in the same block.
mode	String specifying how to scale size factors across blocks. "lowest" will compute the average size factor in each block, identify the lowest average across all blocks, and then scale all size factors by that value. "per-block" will compute the average size factor in each block, and then scale each size factor by the average of block to which it belongs. Only used if block is provided.

**Value**

Numeric vector of length equal to size.factors, containing the centered size factors.

**Author(s)**

Aaron Lun

**See Also**

[https://libscrn.github.io/scrn\\_norm/](https://libscrn.github.io/scrn_norm/), for the rationale behind centering the size factors.

**Examples**

```
centerSizeFactors(runif(100))  
centerSizeFactors(runif(100), block=sample(3, 100, replace=TRUE))
```

---

chooseHighlyVariableGenes

*Choose highly variable genes*

---

**Description**

Choose highly variable genes (HVGs) based on a variance-related statistic.

**Usage**

```
chooseHighlyVariableGenes(  
  stats,  
  top = 4000,  
  larger = TRUE,  
  keep.ties = TRUE,  
  bound = NULL  
)
```

**Arguments**

stats	Numeric vector of variances (or a related statistic) across all genes. Typically the residuals from <code>modelGeneVariances</code> are used here.
top	Integer specifying the number of top genes to retain. Note that the actual number of retained genes may not be equal to top, depending on the other options.
larger	Logical scalar indicating whether larger values of stats correspond to more variable genes. If TRUE, HVGs are defined as those with the largest values of stats.
keep.ties	Logical scalar indicating whether to keep tied values of stats, even if top may be exceeded.
bound	Numeric scalar specifying the lower bound (if larger=TRUE) or upper bound (otherwise) to be applied to stats. Genes are not considered to be HVGs if they do not pass this bound, even if they are within the top genes. Ignored if NULL.

**Value**

Integer vector containing the indices of genes in stats that are considered to be highly variable.

**Author(s)**

Aaron Lun

**See Also**

[https://libscran.github.io/scran\\_variances/](https://libscran.github.io/scran_variances/), for the underlying implementation.

**Examples**

```
resids <- rexp(10000)
str(chooseHighlyVariableGenes(resids))
```

---

choosePseudoCount      *Choose a suitable pseudo-count*

---

**Description**

Choose a suitable pseudo-count to control the bias introduced by log-transformation of normalized counts.

**Usage**

```
choosePseudoCount(size.factors, quantile = 0.05, max.bias = 1, min.value = 1)
```

**Arguments**

size.factors	Numeric vector of size factors for all cells.
quantile	Numeric scalar specifying the quantile to use for defining extreme size factors.
max.bias	Numeric scalar specifying the maximum allowed bias.
min.value	Numeric scalar specifying the minimum value for the pseudo-count.

**Value**

A choice of pseudo-count for `normalizeCounts`.

**Author(s)**

Aaron Lun

**See Also**

[https://libscran.github.io/scran\\_norm/](https://libscran.github.io/scran_norm/), for the motivation behind calculating a larger pseudo-count.

**Examples**

```
sf <- runif(100)
choosePseudoCount(sf)
choosePseudoCount(sf, quantile=0.01)
choosePseudoCount(sf, max.bias=0.5)
```

---

clusterGraph

*Graph-based clustering of cells*

---

**Description**

Identify clusters of cells using a variety of community detection methods from a graph where similar cells are connected.

**Usage**

```
clusterGraph(
  x,
  method = c("multilevel", "leiden", "walktrap"),
  multilevel.resolution = 1,
  leiden.resolution = 1,
  leiden.objective = c("modularity", "cpm"),
  walktrap.steps = 4,
  seed = 42
)
```

**Arguments**

x	List containing graph information or an external pointer to a graph, as returned by <code>buildSnnGraph</code> . Alternatively, an <code>igraph</code> object with edge weights.
method	String specifying the algorithm to use.
multilevel.resolution	Numeric scalar specifying the resolution when method="multilevel".
leiden.resolution	Numeric scalar specifying the resolution when method="leiden".
leiden.objective	String specifying the objective function when method="leiden".
walktrap.steps	Integer scalar specifying the number of steps to use when method="walktrap".
seed	Integer scalar specifying the random seed to use for method="multilevel" or "leiden".

**Value**

A list containing membership, an integer vector containing the cluster assignment for each cell; and status, an integer scalar indicating whether the algorithm completed successfully (0) or not (non-zero). Additional fields may be present depending on the method:

- For method="multilevel", the levels list contains the clustering result at each level of the algorithm. A modularity numeric vector also contains the modularity at each level, the highest of which corresponds to the reported membership.
- For method="leiden", a quality numeric scalar containing the quality of the partitioning.
- For method="walktrap", a merges matrix specifies the pair of cells or clusters that were merged at each step of the algorithm. A modularity numeric scalar also contains the modularity of the final partitioning.

**Author(s)**

Aaron Lun

**See Also**

<https://igraph.org>, for the underlying implementation of each clustering method.

[https://libscran.github.io/scran\\_graph\\_cluster/](https://libscran.github.io/scran_graph_cluster/), for wrappers around the **igraph** code.

**Examples**

```
data <- matrix(rnorm(10000), ncol=1000)
gout <- buildSnnGraph(data)
str(gout)

str(clusterGraph(gout))
str(clusterGraph(gout, method="leiden"))
str(clusterGraph(gout, method="walktrap"))
```

---

clusterKmeans	<i>K-means clustering</i>
---------------	---------------------------

---

**Description**

Perform k-means clustering with a variety of different initialization and refinement algorithms.

**Usage**

```
clusterKmeans(
  x,
  k,
  init.method = c("var-part", "kmeans++", "random"),
  refine.method = c("hartigan-wong", "lloyd"),
  var.part.optimize.partition = TRUE,
  var.part.size.adjustment = 1,
  lloyd.iterations = 100,
  hartigan.wong.iterations = 10,
  hartigan.wong.quick.transfer.iterations = 50,
  hartigan.wong.quit.quick.transfer.failure = FALSE,
  seed = 5489L,
  num.threads = 1
)
```

**Arguments**

<code>x</code>	Numeric matrix where rows are dimensions and columns are cells.
<code>k</code>	Integer scalar specifying the number of clusters.
<code>init.method</code>	String specifying the initialization method: variance partitioning ("var-part"), kmeans++ ("kmeans++") or random initialization ("random").
<code>refine.method</code>	String specifying the refinement method: Lloyd's algorithm ("lloyd") or the Hartigan-Wong algorithm ("hartigan-wong").
<code>var.part.optimize.partition</code>	Logical scalar indicating whether each partition boundary should be optimized to reduce the sum of squares in the child partitions. Only used if <code>init.method = "var.part"</code> .
<code>var.part.size.adjustment</code>	Numeric scalar between 0 and 1, specifying the adjustment to the cluster size when prioritizing the next cluster to partition. Setting this to 0 will ignore the cluster size while setting this to 1 will generally favor larger clusters. Only used if <code>init.method = "var.part"</code> .
<code>lloyd.iterations</code>	Integer scalar specifying the maximum number of iterations for the Lloyd algorithm.

hartigan.wong.iterations	Integer scalar specifying the maximum number of iterations for the Hartigan-Wong algorithm.
hartigan.wong.quick.transfer.iterations	Integer scalar specifying the maximum number of quick transfer iterations for the Hartigan-Wong algorithm.
hartigan.wong.quit.quick.transfer.failure	Logical scalar indicating whether to quit the Hartigan-Wong algorithm upon convergence failure during quick transfer iterations.
seed	Integer scalar specifying the seed to use for random or kmeans++ initialization.
num.threads	Integer scalar specifying the number of threads to use.

**Value**

By default, a list is returned containing:

- `clusters`, a factor containing the cluster assignments for each cell.
- `centers`, a numeric matrix with the coordinates of the cluster centroids (dimensions in rows, centers in columns).
- `iterations`, an integer scalar specifying the number of refinement iterations that were performed.
- `status`, an integer scalar specifying the convergence status. Any non-zero value indicates a convergence failure though the exact meaning depends on the choice of `refine.method`.

**Author(s)**

Aaron Lun

**Examples**

```
x <- t(as.matrix(iris[,1:4]))
clustering <- clusterKmeans(x, k=3)
table(clustering$clusters, iris[, "Species"])
```

---

combineFactors

*Combine multiple factors*

---

**Description**

Combine multiple categorical factors based on the unique combinations of levels from each factor.

**Usage**

```
combineFactors(factors, keep.unused = FALSE)
```

**Arguments**

- `factors` List of vectors or factors of the same length. Corresponding elements across all vectors/factors represent the combination of levels for a single observation. For factors, any existing levels are respected. For other vectors, the sorted and unique values are used as levels.  
Alternatively, a data frame where each column is a vector or factor and each row corresponds to an observation.
- `keep.unused` Logical scalar indicating whether to report unused combinations of levels.

**Value**

List containing `levels`, a data frame containing the sorted and unique combinations of levels from `factors`; and `index`, an integer vector specifying the index into `levels` for each observation.

**Author(s)**

Aaron Lun

**Examples**

```
combineFactors(list(
  sample(LETTERS[1:5], 100, replace=TRUE),
  sample(3, 100, replace=TRUE)
))

combineFactors(list(
  factor(sample(LETTERS[1:5], 10, replace=TRUE), LETTERS[1:5]),
  factor(sample(5, 10, replace=TRUE), 1:5)
), keep.unused=TRUE)
```

---

`computeClrm1Factors`    *Compute size factors for ADT counts*

---

**Description**

Compute size factors from an ADT count matrix using the CLRm1 method.

**Usage**

```
computeClrm1Factors(x, num.threads = 1)
```

**Arguments**

- `x` A matrix-like object containing ADT count data. Rows correspond to tags and columns correspond to cells.
- `num.threads` Number of threads to use.

**Value**

Numeric vector containing the CLRM1 size factor for each cell. Note that these size factors are not centered and should be passed through, e.g., `centerSizeFactors` before normalization.

**Author(s)**

Aaron Lun

**See Also**

<https://github.com/libscrans/clrm1>, for a description of the CLRM1 method.

**Examples**

```
library(Matrix)
x <- abs(rsparsematrix(1000, 100, 0.1) * 10)
head(computeClrm1Factors(x))
```

---

correctMnn

*Batch correction with mutual nearest neighbors*

---

**Description**

Apply mutual nearest neighbor (MNN) correction to remove batch effects from a low-dimensional matrix.

**Usage**

```
correctMnn(
  x,
  block,
  num.neighbors = 15,
  num.mads = 3,
  robust.iterations = 2,
  robust.trim = 0.25,
  mass.cap = NULL,
  order = NULL,
  reference.policy = c("max-rss", "max-size", "max-variance", "input"),
  BNPARAM = AnnoyParam(),
  num.threads = 1
)
```



**Arguments**

<code>x</code>	Numeric matrix where rows are dimensions and columns are cells, typically containing low-dimensional coordinates (e.g., from <a href="#">runPca</a> ).
<code>block</code>	Factor specifying the block of origin (e.g., batch, sample) for each cell in <code>x</code> .
<code>num.neighbors</code>	Integer scalar specifying the number of neighbors to use when identifying MNN pairs.
<code>num.mads</code>	Numeric scalar specifying the number of median absolute deviations to use for removing outliers in the center-of-mass calculations.
<code>robust.iterations</code>	Integer scalar specifying the number of iterations for robust calculation of the center of mass.
<code>robust.trim</code>	Numeric scalar in $[0, 1)$ specifying the trimming proportion for robust calculation of the center of mass.
<code>mass.cap</code>	Integer scalar specifying the cap on the number of observations to use for center-of-mass calculations on the reference dataset. A value of 100,000 may be appropriate for speeding up correction of very large datasets. If NULL, no cap is used.
<code>order</code>	Vector containing levels of batch in the desired merge order. If NULL, a suitable merge order is automatically determined.
<code>reference.policy</code>	String specifying the policy to use to choose the first reference batch. This can be based on the largest batch (" <code>max-size</code> "), the most variable batch (" <code>max-variance</code> "), the batch with the largest residual sum of squares (" <code>max-rss</code> "), or the first specified input (" <code>input</code> "). Only used for automatic merges, i.e., when <code>order=NULL</code> .
<code>BNPARAM</code>	A <a href="#">BiocNeighborParam</a> object specifying the nearest-neighbor algorithm to use.
<code>num.threads</code>	Integer scalar specifying the number of threads to use.

**Value**

List containing:

- `corrected`, a numeric matrix of the same dimensions as `x`, containing the corrected values.
- `merge.order`, character vector containing the unique levels of batch in the automatically determined merge order. The first level in this vector is used as the reference batch; all other batches are iteratively merged to it.
- `num.pairs`, integer vector of length equal to the number of batches minus 1. This contains the number of MNN pairs at each merge.

**Author(s)**

Aaron Lun

**Examples**

```
# Mocking up a dataset with multiple batches.
x <- matrix(rnorm(10000), nrow=10)
b <- sample(3, ncol(x), replace=TRUE)
x[,b==2] <- x[,b==2] + 3
x[,b==3] <- x[,b==3] + 5
lapply(split(colMeans(x), b), mean) # indeed the means differ...

corrected <- correctMnn(x, b)
str(corrected)
lapply(split(colMeans(corrected$corrected), b), mean) # now merged.
```

---

```
crispr_quality_control
```

*Quality control for CRISPR count data*

---

**Description**

Compute per-cell QC metrics from an initialized matrix of CRISPR counts, and use the metrics to suggest filter thresholds to retain high-quality cells.

**Usage**

```
computeCrisprQcMetrics(x, num.threads = 1)

suggestCrisprQcThresholds(metrics, block = NULL, num.mads = 3)

filterCrisprQcMetrics(thresholds, metrics, block = NULL)
```

**Arguments**

<code>x</code>	A matrix-like object where rows are CRISPRs and columns are cells. Values are expected to be counts.
<code>num.threads</code>	Integer scalar specifying the number of threads to use.
<code>metrics</code>	List with the same structure as produced by <code>computeCrisprQcMetrics</code> .
<code>block</code>	Factor specifying the block of origin (e.g., batch, sample) for each cell in <code>metrics</code> . Alternatively NULL if all cells are from the same block. For <code>filterCrisprQcMetrics</code> , a blocking factor should be provided if <code>block</code> was used to construct thresholds.
<code>num.mads</code>	Number of median from the median, to define the threshold for outliers in each metric.
<code>thresholds</code>	List with the same structure as produced by <code>suggestCrisprQcThresholds</code> .

**Value**

For `computeCrisprQcMetrics`, a list is returned containing:

- `sum`, a numeric vector containing the total CRISPR count for each cell.
- `detected`, an integer vector containing the number of detected guides per cell.
- `max.value`, a numeric vector containing the count for the most abundant guide in cell.
- `max.index`, an integer vector containing the row index of the most abundant guide in cell.

Each vector is of length equal to the number of cells.

For `suggestCrisprQcThresholds`, a named list is returned.

- If `block=NULL`, the list contains:
  - `max.value`, a numeric scalar containing the lower bound on the maximum counts for each blocking level.
- Otherwise, if `block` is supplied, the list contains:
  - `max.value`, a numeric vector containing the lower bound on the maximum counts for each blocking level.

Each vector is of length equal to the number of levels in `block` and is named accordingly.

For `filterCrisprQcMetrics`, a logical vector of length `ncol(x)` is returned indicating which cells are of high quality.

**Author(s)**

Aaron Lun

**See Also**

[https://libscran.github.io/scran\\_qc/](https://libscran.github.io/scran_qc/), for the rationale of QC filtering on CRISPR counts.

**Examples**

```
# Mocking a matrix:
library(Matrix)
x <- round(abs(rsparsematrix(100, 100, 0.1) * 100))

qc <- computeCrisprQcMetrics(x)
str(qc)

filt <- suggestCrisprQcThresholds(qc)
str(filt)

keep <- filterCrisprQcMetrics(filt, qc)
summary(keep)
```

---

fitVarianceTrend      *Fit a mean-variance trend*

---

## Description

Fit a trend to the per-cell variances with respect to the mean.

## Usage

```
fitVarianceTrend(
  means,
  variances,
  mean.filter = TRUE,
  min.mean = 0.1,
  transform = TRUE,
  span = 0.3,
  use.min.width = FALSE,
  min.width = 1,
  min.window.count = 200,
  num.threads = 1
)
```

## Arguments

means	Numeric vector containing the mean (log-)expression for each gene.
variances	Numeric vector containing the variance in the (log-)expression for each gene.
mean.filter	Logical scalar indicating whether to filter on the means before trend fitting.
min.mean	Numeric scalar specifying the minimum mean of genes to use in trend fitting. Only used if mean.filter=TRUE.
transform	Logical scalar indicating whether a quarter-root transformation should be applied before trend fitting.
span	Numeric scalar specifying the span of the LOWESS smoother. Ignored if use.min.width=TRUE.
use.min.width	Logical scalar indicating whether a minimum width constraint should be applied to the LOWESS smoother. Useful to avoid overfitting in high-density intervals.
min.width	Minimum width of the window to use when use.min.width=TRUE.
min.window.count	Minimum number of observations in each window. Only used if use.min.width=TRUE.
num.threads	Number of threads to use.

## Value

List containing fitted, the fitted values of the trend for each gene; and residuals, the residuals from the trend.

**Author(s)**

Aaron Lun

**See Also**[https://libscrans.github.io/scrans\\_variances/](https://libscrans.github.io/scrans_variances/), for the underlying implementation.**Examples**

```
x <- runif(1000)
y <- 2^rnorm(1000)
out <- fitVarianceTrend(x, y)

plot(x, y)
o <- order(x)
lines(x[o], out$fitted[o], col="red")
```

---

modelGeneVariances      *Model per-gene variances in expression*

---

**Description**

Compute the variance in (log-)expression values for each gene, and model the trend in the variances with respect to the mean.

**Usage**

```
modelGeneVariances(
  x,
  block = NULL,
  block.weight.policy = c("variable", "equal", "none"),
  variable.block.weight = c(0, 1000),
  mean.filter = TRUE,
  min.mean = 0.1,
  transform = TRUE,
  span = 0.3,
  use.min.width = FALSE,
  min.width = 1,
  min.window.count = 200,
  num.threads = 1
)
```

**Arguments**

<code>x</code>	A matrix-like object where rows correspond to genes or genomic features and columns correspond to cells. It is typically expected to contain log-expression values, e.g., from <code>normalizeCounts</code> .
<code>block</code>	Factor specifying the block of origin (e.g., batch, sample) for each cell in <code>x</code> . Alternatively NULL if all cells are from the same block.
<code>block.weight.policy</code>	String specifying the policy to use for weighting different blocks when computing the average for each statistic Only used if <code>block</code> is not NULL.
<code>variable.block.weight</code>	Numeric vector of length 2, specifying the parameters for variable block weighting. The first and second values are used as the lower and upper bounds, respectively, for the variable weight calculation. Only used if <code>block</code> is not NULL and <code>block.weight.policy = "variable"</code> .
<code>mean.filter</code>	Logical scalar indicating whether to filter on the means before trend fitting.
<code>min.mean</code>	Numeric scalar specifying the minimum mean of genes to use in trend fitting. Only used if <code>mean.filter=TRUE</code> .
<code>transform</code>	Logical scalar indicating whether a quarter-root transformation should be applied before trend fitting.
<code>span</code>	Numeric scalar specifying the span of the LOWESS smoother. Ignored if <code>use.min.width=TRUE</code> .
<code>use.min.width</code>	Logical scalar indicating whether a minimum width constraint should be applied to the LOWESS smoother. Useful to avoid overfitting in high-density intervals.
<code>min.width</code>	Minimum width of the window to use when <code>use.min.width=TRUE</code> .
<code>min.window.count</code>	Minimum number of observations in each window. Only used if <code>use.min.width=TRUE</code> .
<code>num.threads</code>	Integer scalar specifying the number of threads to use.

**Value**

A list containing statistics. This is a data frame with the columns `means`, `variances`, `fitted` and `residuals`, each of which is a numeric vector containing the statistic of the same name across all genes.

If `block` is supplied, each of the column vectors described above contains the average across all blocks. The list will also contain `per.block`, a list of data frames containing the equivalent statistics for each block.

**Author(s)**

Aaron Lun

**See Also**

[https://libscran.github.io/scran\\_variances/](https://libscran.github.io/scran_variances/), for the variance modelling.

[https://libscran.github.io/scran\\_blocks/](https://libscran.github.io/scran_blocks/), for details on the blocking.

`fitVarianceTrend`, which fits the mean-variance trend.

**Examples**

```

library(Matrix)
x <- abs(rsparsematrix(1000, 100, 0.1) * 10)
out <- modelGeneVariances(x)
str(out)

# Throwing in some blocking.
block <- sample(letters[1:4], ncol(x), replace=TRUE)
out <- modelGeneVariances(x, block=block)
str(out)

```

---

normalizeCounts	<i>Normalize the count matrix</i>
-----------------	-----------------------------------

---

**Description**

Apply scaling normalization to a count matrix to obtain log-transformed normalized expression values.

**Usage**

```

normalizeCounts(
  x,
  size.factors,
  log = TRUE,
  pseudo.count = 1,
  log.base = 2,
  preserve.sparsity = FALSE
)

```

**Arguments**

x	A matrix-like object where rows correspond to genes or genomic features and columns correspond to cells. Values are typically expected to be counts. Alternatively, an external pointer created by <a href="#">initializeCpp</a> .
size.factors	A numeric vector of length equal to the number of cells in x, containing positive size factors for all cells.
log	Logical scalar indicating whether log-transformation should be performed.
pseudo.count	Numeric scalar specifying the positive pseudo-count to add before any log-transformation. Ignored if log=FALSE.
log.base	Numeric scalar specifying the base of the log-transformation. Ignored if log=FALSE.
preserve.sparsity	Logical scalar indicating whether to preserve sparsity for pseudo.count != 1. If TRUE, users should manually add $\log(\text{pseudo.count}, \text{log.base})$ to the returned matrix to obtain the desired log-transformed expression values. Ignored if log = FALSE or pseudo.count = 1.

**Value**

If `x` is a matrix-like object, a `DelayedArray` is returned containing the (log-transformed) normalized expression matrix.

If `x` is an external pointer produced by `initializeCpp`, a new external pointer is returned containing the normalized expression matrix.

**Author(s)**

Aaron Lun

**See Also**

[https://libscran.github.io/scran\\_norm/](https://libscran.github.io/scran_norm/), for the rationale behind normalization.

**Examples**

```
# Mocking a matrix:
library(Matrix)
x <- round(abs(rsparsematrix(1000, 100, 0.1) * 100))
sf <- centerSizeFactors(colSums(x))
normed <- normalizeCounts(x, size.factors=sf)
normed

# Passing a pointer.
ptr <- beachmat::initializeCpp(x)
optr <- normalizeCounts(ptr, sf)
optr
```

---

rna\_quality\_control    *Quality control for RNA count data*

---

**Description**

Compute per-cell QC metrics from an initialized matrix of RNA counts, and use the metrics to suggest filter thresholds to retain high-quality cells.

**Usage**

```
computeRnaQcMetrics(x, subsets, num.threads = 1)

suggestRnaQcThresholds(metrics, block = NULL, num.mads = 3)

filterRnaQcMetrics(thresholds, metrics, block = NULL)
```



**Arguments**

x	A matrix-like object where rows are genes and columns are cells. Values are expected to be counts.
subsets	List of vectors specifying gene subsets of interest, typically for control-like features like mitochondrial genes or spike-in transcripts. Each vector may be logical (whether to keep each row), integer (row indices) or character (row names).
num.threads	Integer scalar specifying the number of threads to use.
metrics	List with the same structure as produced by <code>computeRnaQcMetrics</code> .
block	Factor specifying the block of origin (e.g., batch, sample) for each cell in <code>metrics</code> . Alternatively NULL if all cells are from the same block. For <code>filterRnaQcMetrics</code> , a blocking factor should be provided if <code>block</code> was used to construct thresholds.
num.mads	Number of median from the median, to define the threshold for outliers in each metric.
thresholds	List with the same structure as produced by <code>suggestRnaQcThresholds</code> .

**Value**

For `computeRnaQcMetrics`, a list is returned containing:

- `sum`, a numeric vector containing the total RNA count for each cell.
- `detected`, an integer vector containing the number of detected genes per cell.
- `subsets`, a list of numeric vectors containing the proportion of counts in each feature subset.

Each vector is of length equal to the number of cells.

For `suggestRnaQcThresholds`, a named list is returned.

- If `block=NULL`, the list contains:
  - `sum`, a numeric scalar containing the lower bound on the sum.
  - `detected`, a numeric scalar containing the lower bound on the number of detected genes.
  - `subsets`, a numeric vector containing the upper bound on the sum of counts in each feature subset.
- Otherwise, if `block` is supplied, the list contains:
  - `sum`, a numeric vector containing the lower bound on the sum for each blocking level.
  - `detected`, a numeric vector containing the lower bound on the number of detected genes for each blocking level.
  - `subsets`, a list of numeric vectors containing the upper bound on the sum of counts in each feature subset for each blocking level.

Each vector is of length equal to the number of levels in `block` and is named accordingly.

For `filterRnaQcMetrics`, a logical vector of length `ncol(x)` is returned indicating which cells are of high quality.

**Author(s)**

Aaron Lun

**See Also**

[https://libscran.github.io/scran\\_qc/](https://libscran.github.io/scran_qc/), for the rationale of QC filtering on RNA counts.

**Examples**

```
# Mocking a matrix:
library(Matrix)
x <- round(abs(rsparsematrix(1000, 100, 0.1) * 100))

# Mocking up a control set.
sub <- list(mito=rbinom(nrow(x), 1, 0.1) > 0)

qc <- computeRnaQcMetrics(x, sub)
str(qc)

filt <- suggestRnaQcThresholds(qc)
str(filt)

keep <- filterRnaQcMetrics(filt, qc)
summary(keep)
```

---

runAllNeighborSteps    *Run all neighbor-related steps*

---

**Description**

Run all steps that require a nearest-neighbor search. This includes [runUmap](#), [runTsne](#) and [buildSnnGraph](#) with [clusterGraph](#). The idea is to build the index once, perform the neighbor search, and run each task in parallel to save time.

**Usage**

```
runAllNeighborSteps(
  x,
  runUmap.args = list(),
  runTsne.args = list(),
  buildSnnGraph.args = list(),
  clusterGraph.args = list(),
  BNPARAM = AnnoyParam(),
  return.graph = FALSE,
  collapse.search = FALSE,
  num.threads = 3
)
```

**Arguments**

<code>x</code>	Numeric matrix where rows are dimensions and columns are cells, typically containing a low-dimensional representation from, e.g., <code>runPca</code> . Alternatively, an index constructed by <code>buildIndex</code> .
<code>runUmap.args</code>	Named list of further arguments to pass to <code>runUmap</code> . This can be set to <code>NULL</code> to omit the UMAP.
<code>runTsne.args</code>	Named list of further arguments to pass to <code>runTsne</code> . This can be set to <code>NULL</code> to omit the t-SNE.
<code>buildSnnGraph.args</code>	Named list of further arguments to pass to <code>buildSnnGraph</code> .
<code>clusterGraph.args</code>	Named list of further arguments to pass to <code>clusterGraph</code> . This can be set to <code>NULL</code> to omit the clustering.
<code>BNPARAM</code>	A <code>BiocNeighborParam</code> instance specifying the nearest-neighbor search algorithm to use.
<code>return.graph</code>	Logical scalar indicating whether to return the output of <code>buildSnnGraph</code> . By default, only the output of <code>clusterGraph</code> is returned.
<code>collapse.search</code>	Logical scalar indicating whether to collapse the nearest-neighbor search for each step into a single search. Steps that need fewer neighbors will take a subset of the neighbors from the collapsed search. This is faster but may not give the same results as separate searches for some algorithms (e.g., approximate searches).
<code>num.threads</code>	Integer scalar specifying the number of threads to use. At least one thread should be available for each step.

**Value**

A named list containing the results of each step. See each individual function for the format of the results.

**Author(s)**

Aaron Lun

**Examples**

```
x <- t(as.matrix(iris[,1:4]))
# (Turning down the number of threads so that R CMD check is happy.)
res <- runAllNeighborSteps(x, num.threads=2)
str(res)
```

---

runPca *Principal components analysis*

---

### Description

Run a PCA on the gene-by-cell log-expression matrix to obtain a low-dimensional representation for downstream analyses.

### Usage

```
runPca(
  x,
  number = 25,
  scale = FALSE,
  block = NULL,
  block.weight.policy = c("variable", "equal", "none"),
  variable.block.weight = c(0, 1000),
  components.from.residuals = FALSE,
  extra.work = 7,
  iterations = 1000,
  seed = 5489,
  realized = TRUE,
  num.threads = 1
)
```

### Arguments

x	A matrix-like object where rows correspond to genes or genomic features and columns correspond to cells. Typically, the matrix is expected to contain log-expression values, and the rows should be filtered to relevant (e.g., highly variable) genes.
number	Integer scalar specifying the number of PCs to retain.
scale	Logical scalar indicating whether to scale all genes to have the same variance.
block	Factor specifying the block of origin (e.g., batch, sample) for each cell in x. Alternatively NULL if all cells are from the same block.
block.weight.policy	String specifying the policy to use for weighting different blocks when computing the average for each statistic. Only used if block is not NULL.
variable.block.weight	Numeric vector of length 2, specifying the parameters for variable block weighting. The first and second values are used as the lower and upper bounds, respectively, for the variable weight calculation. Only used if block is not NULL and block.weight.policy = "variable".
components.from.residuals	Logical scalar indicating whether to compute the PC scores from the residuals in the presence of a blocking factor. By default, the residuals are only used to

	compute the rotation matrix, and the original expression values of the cells are projected onto this new space. Only used if <code>block</code> is not <code>NULL</code> .
<code>extra.work</code>	Integer scalar specifying the extra dimensions for the IRLBA workspace.
<code>iterations</code>	Integer scalar specifying the maximum number of restart iterations for IRLBA.
<code>seed</code>	Integer scalar specifying the seed for the initial random vector in IRLBA.
<code>realized</code>	Logical scalar indicating whether to realize <code>x</code> into an optimal memory layout for IRLBA. This speeds up computation at the cost of increased memory usage.
<code>num.threads</code>	Number of threads to use.

### Value

List containing:

- `components`, a matrix of PC scores. Rows are dimensions (i.e., PCs) and columns are cells.
- `rotation`, the rotation matrix. Rows are genes and columns are dimensions.
- `variance.explained`, the vector of variances explained by each PC.
- `total.variance`, the total variance in the dataset.
- `center`, a numeric vector containing the mean for each gene. If `block` is provided, this is instead a matrix containing the mean for each gene (column) in each block (row).
- `scale`, a numeric vector containing the scaling for each gene. Only reported if `scale=TRUE`.

### Author(s)

Aaron Lun

### See Also

[https://libscran.github.io/scran\\_pca/](https://libscran.github.io/scran_pca/), for more details on the PCA.

[https://libscran.github.io/scran\\_blocks/](https://libscran.github.io/scran_blocks/), for more details on the block weighting.

### Examples

```
library(Matrix)
x <- abs(rsparsematrix(1000, 100, 0.1) * 10)
y <- normalizeCounts(x, size.factors=centerSizeFactors(colSums(x)))

# A simple PCA:
out <- runPca(y)
str(out)

# Blocking on uninteresting factors:
block <- sample(LETTERS[1:3], ncol(y), replace=TRUE)
bout <- runPca(y, block=block)
str(bout)
```

---

runTsne	<i>t-stochastic neighbor embedding</i>
---------	--

---

### Description

Compute t-SNE coordinates to visualize similarities between cells.

### Usage

```
runTsne(
  x,
  perplexity = 30,
  num.neighbors = tsnePerplexityToNeighbors(perplexity),
  max.depth = 20,
  leaf.approximation = FALSE,
  max.iterations = 500,
  seed = 42,
  num.threads = 1,
  BNPARAM = AnnoyParam()
)

tsnePerplexityToNeighbors(perplexity)
```

### Arguments

<code>x</code>	Numeric matrix where rows are dimensions and columns are cells, typically containing a low-dimensional representation from, e.g., <a href="#">runPca</a> . Alternatively, a named list of nearest-neighbor search results like that returned by <a href="#">findKNN</a> . This should contain <code>index</code> , an integer matrix where rows are neighbors and columns are cells; and <code>distance</code> , a numeric matrix of the same dimensions containing the distances to each neighbor. Each column contains 1-based indices for the nearest neighbors of the corresponding cell, ordered by increasing distance. The number of neighbors should be the same as <code>num.neighbors</code> , otherwise a warning is raised. Alternatively, an index constructed by <a href="#">buildIndex</a> .
<code>perplexity</code>	Numeric scalar specifying the perplexity to use in the t-SNE algorithm.
<code>num.neighbors</code>	Integer scalar specifying the number of neighbors, typically derived from <code>perplexity</code> .
<code>max.depth</code>	Integer scalar specifying the maximum depth of the Barnes-Hut quadtree. Smaller values (7-10) improve speed at the cost of accuracy.
<code>leaf.approximation</code>	Logical scalar indicating whether to use the “leaf approximation” approach, which sacrifices some accuracy for greater speed. Only effective when <code>max.depth</code> is small enough for multiple cells to be assigned to the same leaf node of the quadtree.
<code>max.iterations</code>	Integer scalar specifying the maximum number of iterations to perform.

seed	Integer scalar specifying the seed to use for generating the initial coordinates.
num.threads	Integer scalar specifying the number of threads to use.
BNPARAM	A <a href="#">BiocNeighborParam</a> object specifying the algorithm to use. Only used if x is not a prebuilt index or a list of existing nearest-neighbor search results.

**Value**

For runTsnE, a numeric matrix where rows are cells and columns are the two dimensions of the embedding.

For tsnePerplexityToNeighbors, an integer scalar specifying the number of neighbors to use for a given perplexity.

**Author(s)**

Aaron Lun

**See Also**

<https://libscran.github.io/qdtsne/>, for an explanation of the approximations.

**Examples**

```
x <- t(as.matrix(iris[,1:4]))
embedding <- runTsnE(x)
plot(embedding[,1], embedding[,2], col=iris[,5])
```

---

runUmap

*Uniform manifold approximation and projection*

---

**Description**

Compute UMAP coordinates to visualize similarities between cells.

**Usage**

```
runUmap(
  x,
  num.dim = 2,
  num.neighbors = 15,
  num.epochs = NULL,
  min.dist = 0.1,
  seed = 1234567890,
  num.threads = 1,
  parallel.optimization = FALSE,
  BNPARAM = AnnoyParam()
)
```

**Arguments**

<code>x</code>	Numeric matrix where rows are dimensions and columns are cells, typically containing a low-dimensional representation from, e.g., <code>runPca</code> . Alternatively, a named list of nearest-neighbor search results like that returned by <code>findKNN</code> . This should contain <code>index</code> , an integer matrix where rows are neighbors and columns are cells; and <code>distance</code> , a numeric matrix of the same dimensions containing the distances to each neighbor. Each column contains 1-based indices for the nearest neighbors of the corresponding cell, ordered by increasing distance. The number of neighbors should be the same as <code>num.neighbors</code> , otherwise a warning is raised. Alternatively, an index constructed by <code>buildIndex</code> .
<code>num.dim</code>	Integer scalar specifying the number of dimensions of the output embedding.
<code>num.neighbors</code>	Integer scalar specifying the number of neighbors to use in the UMAP algorithm.
<code>num.epochs</code>	Integer scalar specifying the number of epochs to perform. If set to <code>NULL</code> , an appropriate number of epochs is chosen based on <code>ncol(x)</code> .
<code>min.dist</code>	Numeric scalar specifying the minimum distance between points.
<code>seed</code>	Integer scalar specifying the seed to use.
<code>num.threads</code>	Integer scalar specifying the number of threads to use.
<code>parallel.optimization</code>	Logical scalar specifying whether to parallelize the optimization step.
<code>BNPARAM</code>	A <code>BiocNeighborParam</code> object specifying the algorithm to use. Only used if <code>x</code> is not a prebuilt index or a list of existing nearest-neighbor search results.

**Value**

A numeric matrix where rows are cells and columns are the two dimensions of the embedding.

**Author(s)**

Aaron Lun

**See Also**

<https://libscran.github.io/umapp/>, for details on the underlying implementation.

**Examples**

```
x <- t(as.matrix(iris[,1:4]))
embedding <- runUmap(x)
plot(embedding[,1], embedding[,2], col=iris[,5])
```



---

sanitizeSizeFactors     *Sanitize size factors*

---

### Description

Replace invalid size factors, i.e., zero, negative, infinite or NaN values.

### Usage

```
sanitizeSizeFactors(  
  size.factors,  
  replace.zero = TRUE,  
  replace.negative = TRUE,  
  replace.infinite = TRUE,  
  replace.nan = TRUE  
)
```

### Arguments

size.factors	Numeric vector of size factors across cells.
replace.zero	Logical scalar indicating whether to replace size factors of zero with the lowest positive factor. If FALSE, zeros are retained.
replace.negative	Logical scalar indicating whether to replace negative size factors with the lowest positive factor. If FALSE, negative values are retained.
replace.infinite	Logical scalar indicating whether to replace infinite size factors with the largest positive factor. If FALSE, infinite values are retained.
replace.nan	Logical scalar indicating whether to replace NaN size factors with unity. If FALSE, NaN values are retained.

### Value

Numeric vector of length equal to `size.factors`, containing the sanitized size factors.

### Author(s)

Aaron Lun

### See Also

[https://libscran.github.io/scran\\_norm/](https://libscran.github.io/scran_norm/), for more details on the sanitization.

**Examples**

```
sf <- 2^rnorm(100)
sf[1] <- 0
sf[2] <- -1
sf[3] <- Inf
sf[4] <- NaN
sanitizeSizeFactors(sf)
```

---

scaleByNeighbors      *Scale and combine multiple embeddings*

---

**Description**

Scale multiple embeddings (usually derived from different modalities across the same set of cells) so that their within-population variances are comparable, and then combine them into a single embedding matrix for combined downstream analysis.

**Usage**

```
scaleByNeighbors(
  x,
  num.neighbors = 20,
  num.threads = 1,
  weights = NULL,
  BNPARAM = AnnoyParam()
)
```

**Arguments**

<code>x</code>	List of numeric matrices of principal components or other embeddings, one for each modality. For each entry, rows are dimensions and columns are cells. All entries should have the same number of columns but may have different numbers of rows.
<code>num.neighbors</code>	Integer scalar specifying the number of neighbors to use to define the scaling factor.
<code>num.threads</code>	Integer scalar specifying the number of threads to use.
<code>weights</code>	Numeric vector of length equal to that of <code>x</code> , specifying the weights to apply to each modality. Each value represents a multiplier of the within-population variance of its modality, i.e., larger values increase the contribution of that modality in the combined output matrix. <code>NULL</code> is equivalent to an all-1 vector, i.e., all modalities are scaled to have the same within-population variance.
<code>BNPARAM</code>	A <a href="#">BiocNeighborParam</a> object specifying how to perform the neighbor search.

**Value**

List containing `scaling`, a vector of scaling factors to be applied to each embedding; and `combined`, a numeric matrix formed by scaling each entry of `x` and then `rbind`ing them together.

**Author(s)**

Aaron Lun

**Examples**

```
pcs <- list(
  gene = matrix(rnorm(10000), ncol=200),
  protein = matrix(rnorm(1000, sd=3), ncol=200),
  guide = matrix(rnorm(2000, sd=5), ncol=200)
)

out <- scaleByNeighbors(pcs)
out$scaling
dim(out$combined)
```

---

scoreGeneSet

*Score gene set activity for each cell*

---

**Description**

Compute per-cell scores for a gene set, defined as the column sums of a rank-1 approximation to the submatrix for the gene set. This uses the same approach implemented in the **GSDecon** package by Jason Hackney.

**Usage**

```
scoreGeneSet(
  x,
  set,
  rank = 1,
  scale = FALSE,
  block = NULL,
  block.weight.policy = c("variable", "equal", "none"),
  variable.block.weight = c(0, 1000),
  extra.work = 7,
  iterations = 1000,
  seed = 5489,
  realized = TRUE,
  num.threads = 1
)
```

**Arguments**

<code>x</code>	A matrix-like object where rows correspond to genes or genomic features and columns correspond to cells. Typically, the matrix is expected to contain log-expression values.
<code>set</code>	Integer, logical or character vector specifying the rows that belong to the gene set.
<code>rank</code>	Integer scalar specifying the rank of the approximation.
<code>scale</code>	Logical scalar indicating whether to scale all genes to have the same variance.
<code>block</code>	Factor specifying the block of origin (e.g., batch, sample) for each cell in <code>x</code> . Alternatively NULL if all cells are from the same block.
<code>block.weight.policy</code>	String specifying the policy to use for weighting different blocks when computing the average for each statistic. Only used if <code>block</code> is not NULL.
<code>variable.block.weight</code>	Numeric vector of length 2, specifying the parameters for variable block weighting. The first and second values are used as the lower and upper bounds, respectively, for the variable weight calculation. Only used if <code>block</code> is not NULL and <code>block.weight.policy = "variable"</code> .
<code>extra.work</code>	Integer scalar specifying the extra dimensions for the IRLBA workspace.
<code>iterations</code>	Integer scalar specifying the maximum number of restart iterations for IRLBA.
<code>seed</code>	Integer scalar specifying the seed for the initial random vector in IRLBA.
<code>realized</code>	Logical scalar indicating whether to realize <code>x</code> into an optimal memory layout for IRLBA. This speeds up computation at the cost of increased memory usage.
<code>num.threads</code>	Number of threads to use.

**Value**

List containing scores, a numeric vector of per-cell scores for each column in `x`; and weights, a numeric vector of per-gene weights for each gene in `set`.

**Author(s)**

Aaron Lun

**See Also**

<https://libscran.github.io/gsdecon/>, for more details on the underlying algorithm.

**Examples**

```
library(Matrix)
x <- round(abs(rsparsematrix(1000, 100, 0.1) * 100))
normed <- normalizeCounts(x, size.factors=centerSizeFactors(colSums(x)))
scoreGeneSet(normed, set=c(1,3,5,10,20,100))
```

---

scoreMarkers	<i>Score marker genes</i>
--------------	---------------------------

---

## Description

Score marker genes for each group using a variety of effect sizes from pairwise comparisons between groups. This includes Cohen's d, the area under the curve (AUC), the difference in the means (delta-mean) and the difference in the proportion of detected cells (delta-detected).

## Usage

```
scoreMarkers(
  x,
  groups,
  block = NULL,
  block.weight.policy = c("variable", "equal", "none"),
  variable.block.weight = c(0, 1000),
  compute.delta.mean = TRUE,
  compute.delta.detected = TRUE,
  compute.cohens.d = TRUE,
  compute.auc = TRUE,
  threshold = 0,
  all.pairwise = FALSE,
  num.threads = 1
)
```

## Arguments

<code>x</code>	A matrix-like object where rows correspond to genes or genomic features and columns correspond to cells. It is typically expected to contain log-expression values, e.g., from <a href="#">normalizeCounts</a> .
<code>groups</code>	A vector specifying the group assignment for each cell in <code>x</code> .
<code>block</code>	Factor specifying the block of origin (e.g., batch, sample) for each cell in <code>x</code> . Alternatively NULL if all cells are from the same block.
<code>block.weight.policy</code>	String specifying the policy to use for weighting different blocks when computing the average for each statistic Only used if <code>block</code> is not NULL.
<code>variable.block.weight</code>	Numeric vector of length 2, specifying the parameters for variable block weighting. The first and second values are used as the lower and upper bounds, respectively, for the variable weight calculation. Only used if <code>block</code> is not NULL and <code>block.weight.policy = "variable"</code> .
<code>compute.delta.mean</code>	Logical scalar indicating whether to compute the delta-means, i.e., the log-fold change when <code>x</code> contains log-expression values.

<code>compute.delta.detected</code>	Logical scalar indicating whether to compute the delta-detected, i.e., differences in the proportion of cells with detected expression.
<code>compute.cohens.d</code>	Logical scalar indicating whether to compute Cohen's d.
<code>compute.auc</code>	Logical scalar indicating whether to compute the AUC. Setting this to FALSE can improve speed and memory efficiency.
<code>threshold</code>	Non-negative numeric scalar specifying the minimum threshold on the differences in means (i.e., the log-fold change, if <code>x</code> contains log-expression values). This is incorporated into the effect sizes for Cohen's d and the AUC.
<code>all.pairwise</code>	Logical scalar indicating whether to report the full effects for every pairwise comparison between groups.
<code>num.threads</code>	Integer scalar specifying the number of threads to use.

### Value

If `all.pairwise=FALSE`, a named list is returned containing:

- `mean`, a numeric matrix containing the mean expression for each group. Each row is a gene and each column is a group.
- `detected`, a numeric matrix containing the proportion of detected cells in each group. Each row is a gene and each column is a group.
- `cohens.d`, a list of data frames where each data frame corresponds to a group. Each row of each data frame represents a gene, while each column contains a summary of Cohen's d from pairwise comparisons to all other groups. This includes the `min`, `mean`, `median`, `max` and `min.rank` - check out [?summarizeEffects](#) for details. Omitted if `compute.cohens.d=FALSE`.
- `auc`, a list like `cohens.d` but containing the summaries of the AUCs from each pairwise comparison. Omitted if `compute.auc=FALSE`.
- `delta.mean`, a list like `cohens.d` but containing the summaries of the delta-mean from each pairwise comparison. Omitted if `compute.delta.mean=FALSE`.
- `delta.detected`, a list like `cohens.d` but containing the summaries of the delta-detected from each pairwise comparison. Omitted if `compute.delta.detected=FALSE`.

If `all.pairwise=TRUE`, a list is returned containing:

- `mean`, a numeric matrix containing the mean expression for each group. Each row is a gene and each column is a group.
- `detected`, a numeric matrix containing the proportion of detected cells in each group. Each row is a gene and each column is a group.
- `cohens.d`, a 3-dimensional numeric array containing the Cohen's d from each pairwise comparison between groups. The extents of the first two dimensions are equal to the number of groups, while the extent of the final dimension is equal to the number of genes. The entry `[i, j, k]` represents Cohen's d from the comparison of group `j` over group `i` for gene `k`. Omitted if `compute.cohens.d=FALSE`.
- `auc`, an array like `cohens.d` but containing the AUCs from each pairwise comparison. Omitted if `compute.auc=FALSE`.

- `delta.mean`, an array like `cohens.d` but containing the delta-mean from each pairwise comparison. Omitted if `compute.delta.mean=FALSE`.
- `delta.detected`, an array like `cohens.d` but containing the delta-detected from each pairwise comparison. Omitted if `compute.delta.detected=FALSE`.

### See Also

[https://libscran.github.io/scran\\_markers/](https://libscran.github.io/scran_markers/), in particular the `score_markers_summary` function (for all `.pairwise=FALSE`), the `score_markers_pairwise` function (for all `.pairwise=TRUE`), and their blocked equivalents `score_markers_summary_blocked` and `score_markers_pairwise_blocked` (when `block` is not `NULL`).

[summarizeEffects](#), to summarize the pairwise effects returned when all `.pairwise=TRUE`.

### Examples

```
# Mocking a matrix:
library(Matrix)
x <- round(abs(rsparsematrix(1000, 100, 0.1) * 100))
normed <- normalizeCounts(x, size.factors=centerSizeFactors(colSums(x)))

g <- sample(letters[1:4], ncol(x), replace=TRUE)
scores <- scoreMarkers(normed, g)
names(scores)
head(scores$mean)
head(scores$cohens.d[["a"]])
```

---

subsampleByNeighbors    *Subsample cells based on their neighbors*

---

### Description

Subsample a dataset by selecting cells to represent all of their nearest neighbors.

### Usage

```
subsampleByNeighbors(
  x,
  num.neighbors = 20,
  min.remaining = 10,
  num.threads = 1,
  BNPARAM = AnnoyParam()
)
```

**Arguments**

<code>x</code>	<p>A numeric matrix where rows are dimensions and columns are cells, typically containing a low-dimensional representation from, e.g., <code>runPca</code>. Alternatively, an index constructed by <code>buildIndex</code>. Alternatively, a list containing existing nearest-neighbor search results. This should contain:</p> <ul style="list-style-type: none"> <li>• <code>index</code>, an integer matrix where rows are neighbors and columns are cells. Each column contains 1-based indices for the nearest neighbors of the corresponding cell, ordered by increasing distance.</li> <li>• <code>distance</code>, a numeric matrix of the same dimensions as <code>index</code>, containing the distances to each of the nearest neighbors.</li> </ul> <p>The number of neighbors should be equal to <code>num.neighbors</code>, otherwise a warning is raised.</p>
<code>num.neighbors</code>	Integer scalar specifying the number of neighbors to use. Larger values result in greater downsampling. Only used if <code>x</code> does not contain existing nearest-neighbor results.
<code>min.remaining</code>	Integer scalar specifying the minimum number of remaining (i.e., unselected) neighbors that a cell must have in order to be considered for selection. This should be less than or equal to <code>num.neighbors</code> .
<code>num.threads</code>	Integer scalar specifying the number of threads to use for the nearest-neighbor search. Only used if <code>x</code> does not contain existing nearest-neighbor results.
<code>BNPARAM</code>	A <code>BiocNeighborParam</code> object specifying the algorithm to use. Only used if <code>x</code> does not contain existing nearest-neighbor results.

**Value**

Integer vector with the indices of the selected cells in the subsample.

**Author(s)**

Aaron Lun

**See Also**

<https://libscan.github.io/nesub/>, for more details on the underlying algorithm.

**Examples**

```
x <- matrix(rnorm(10000), nrow=2)
keep <- subsampleByNeighbors(x, 10)
plot(x[1,], x[2,])
points(x[1,keep], x[2,keep], col="red")
legend('topright', col=c('black', 'red'), legend=c('all', 'subsample'), pch=1)
```



---

summarizeEffects	<i>Summarize pairwise effect sizes for each group</i>
------------------	---

---

### Description

For each group, summarize the effect sizes for all pairwise comparisons to other groups. This yields a set of summary statistics that can be used to rank marker genes for each group.

### Usage

```
summarizeEffects(effects, num.threads = 1)
```

### Arguments

effects	A 3-dimensional numeric containing the effect sizes from each pairwise comparison between groups. The extents of the first two dimensions are equal to the number of groups, while the extent of the final dimension is equal to the number of genes. The entry [i, j, k] represents the effect size from the comparison of group j against group i for gene k. See also the output of <a href="#">scoreMarkers</a> with <code>all.pairwise=TRUE</code> .
num.threads	Integer scalar specifying the number of threads to use.

### Details

Each summary statistic can be used to prioritize different sets of marker genes for the group of interest, by ranking them in decreasing order according to said statistic:

- `min` contains the minimum effect size across all comparisons involving the group of interest. Using this to define markers will focus on genes that are upregulated in all comparisons.
- `mean` contains the mean effect size across all comparisons involving the group of interest. Using this to define markers will focus on genes that are generally upregulated.
- `median` contains the median effect size across all comparisons involving the group of interest. Using this to define markers will focus on genes that are upregulated in most comparisons.
- `max` contains the maximum effect size across all comparisons involving the group of interest. Using this to define markers will focus on genes that are upregulated in at least one comparison.

The `min.rank` is a more exotic summary statistic, containing the minimum rank for each gene across all comparisons involving the group of interest. This is defined by ranking the effect sizes across genes within each comparison, and then taking the minimum of these ranks across comparisons. Taking all genes with `min.rank <= T` will yield a set containing the top T genes from each comparison; the idea is to ensure that there are at least T genes that can distinguish the group of interest from the others.

### Value

List of data frames containing summary statistics for the effect sizes. Each data frame corresponds to a group, each row corresponds to a gene, and each column contains a single summary.

**Author(s)**

Aaron Lun

**See Also**

[https://libscran.github.io/scran\\_markers/](https://libscran.github.io/scran_markers/), for more details on the statistics.  
[scoreMarkers](#), to compute the pairwise effects in the first place.

**Examples**

```
# Mocking a matrix:
library(Matrix)
x <- round(abs(rsparsematrix(1000, 100, 0.1) * 100))
normed <- normalizeCounts(x, size.factors=centerSizeFactors(colSums(x)))

g <- sample(letters[1:4], ncol(x), replace=TRUE)
effects <- scoreMarkers(normed, g, all.pairwise=TRUE)

summarized <- summarizeEffects(effects$cohens.d)
str(summarized)
```

---

testEnrichment

*Test for gene set enrichment*


---

**Description**

Perform a hypergeometric test for enrichment of interesting genes (e.g., markers) in one or more pre-defined gene sets.

**Usage**

```
testEnrichment(x, sets, universe, log = FALSE, num.threads = 1)
```

**Arguments**

x	Vector of identifiers for the interesting genes.
sets	List of vectors of identifiers for the pre-defined gene sets.
universe	Vector of identifiers for the universe of genes in the dataset. It is expected that x is a subset of universe. Alternatively, an integer scalar specifying the size of the universe.
log	Logical scalar indicating whether to report the log-transformed p-values.
num.threads	Integer scalar specifying the number of threads to use.

**Value**

Numeric vector of (log-transformed) p-values to test for significant enrichment of x in each entry of sets.

**Author(s)**

Aaron Lun

**See Also**

[phyper](#) and <https://github.com/libscran/phyper/>, which is the basis for the underlying calculation.

**Examples**

```
testEnrichment(  
  x=LETTERS[1:5],  
  sets=list(  
    first=LETTERS[1:10],  
    second=LETTERS[1:5 * 2],  
    third=LETTERS[10:20]  
  ),  
  universe=LETTERS  
)
```

# Index

adt\_quality\_control, 2  
aggregateAcrossCells, 4, 6  
aggregateAcrossGenes, 5, 6  
  
BiocNeighborParam, 7, 17, 27, 31, 32, 34, 40  
buildIndex, 7, 27, 30, 32, 40  
buildSnnGraph, 7, 12, 26, 27  
  
centerSizeFactors, 8, 16  
chooseHighlyVariableGenes, 9  
choosePseudoCount, 10  
clusterGraph, 11, 26, 27  
clusterKmeans, 13  
combineFactors, 4, 14  
computeAdtQcMetrics  
    (ad\_t\_quality\_control), 2  
computeClrm1Factors, 15  
computeCrisprQcMetrics  
    (crispr\_quality\_control), 18  
computeRnaQcMetrics  
    (rna\_quality\_control), 24  
correctMnn, 16  
crispr\_quality\_control, 18  
  
DelayedArray, 24  
  
filterAdtQcMetrics  
    (ad\_t\_quality\_control), 2  
filterCrisprQcMetrics  
    (crispr\_quality\_control), 18  
filterRnaQcMetrics  
    (rna\_quality\_control), 24  
findKNN, 30, 32  
fitVarianceTrend, 20, 22  
  
igraph, 12  
initializeCpp, 23, 24  
  
modelGeneVariances, 10, 21  
  
normalizeCounts, 11, 22, 23, 37  
  
phyper, 43  
  
rna\_quality\_control, 24  
runAllNeighborSteps, 26  
runPca, 7, 17, 27, 28, 30, 32, 40  
runTsne, 26, 27, 30  
runUmap, 26, 27, 31  
  
sanitizeSizeFactors, 33  
scaleByNeighbors, 34  
scoreGeneSet, 35  
scoreMarkers, 37, 41, 42  
subsampleByNeighbors, 39  
suggestAdtQcThresholds  
    (ad\_t\_quality\_control), 2  
suggestCrisprQcThresholds  
    (crispr\_quality\_control), 18  
suggestRnaQcThresholds  
    (rna\_quality\_control), 24  
summarizeEffects, 38, 39, 41  
  
testEnrichment, 42  
tsnePerplexityToNeighbors (runTsne), 30