

Package ‘scregclust’

December 6, 2024

Title Reconstructing the Regulatory Programs of Target Genes in
scRNA-Seq Data

Version 0.2.0

Description Implementation of the scregclust algorithm
described in Larsson, Held, et al. (2024) <[doi:10.1038/s41467-024-53954-3](https://doi.org/10.1038/s41467-024-53954-3)>
which reconstructs regulatory programs of target genes in scRNA-seq data.
Target genes are clustered into modules and each module is associated with a linear
model describing the regulatory program.

Encoding UTF-8

Depends R (>= 3.5.0)

Imports Matrix, stats, methods, utils, reshape, igraph, graphics,
grid, cli, prettyunits, ggplot2, rlang, Rcpp (>= 1.0.8)

Suggests knitr, rmarkdown, testthat (>= 3.0.0), Seurat (>= 5.0.0),
glmGamPoi, hdf5r

LinkingTo Rcpp, RcppEigen

VignetteBuilder knitr

RoxygenNote 7.3.2

License GPL (>= 3)

Config/testthat/edition 3

URL <https://scmethods.github.io/scregclust/>,
<https://github.com/scmethods/scregclust/>

BugReports <https://github.com/scmethods/scregclust/issues>

NeedsCompilation yes

Author Felix Held [aut, cre] (<<https://orcid.org/0000-0002-7679-7752>>),
Ida Larsson [aut] (<<https://orcid.org/0000-0001-5422-4243>>),
Sven Nelander [aut] (<<https://orcid.org/0000-0003-1758-1262>>),
André Armatowski [ctb]

Maintainer Felix Held <felix.held@gmail.com>

Repository CRAN

Date/Publication 2024-12-06 14:30:01 UTC

Contents

| | |
|------------------------------------|-----------|
| available_results | 2 |
| cluster_overlap | 3 |
| fast_cor | 3 |
| find_module_sizes | 4 |
| get_avg_num_regulators | 4 |
| get_num_final_configs | 5 |
| get_rand_indices | 5 |
| get_regulator_list | 6 |
| get_target_gene_modules | 6 |
| kmeanspp | 7 |
| plot_module_count_helper | 8 |
| plot_regulator_network | 8 |
| plot_silhouettes | 9 |
| scregclust | 10 |
| scregclust_format | 14 |
| Index | 15 |

| | |
|-------------------|---|
| available_results | <i>Extract final configurations into a data frame</i> |
|-------------------|---|

Description

Extract final configurations into a data frame

Usage

```
available_results(obj)
```

Arguments

obj An object of class scregclust

Value

A `data.frame` containing penalization parameters and final configurations for those penalizations.

| | |
|-----------------|---|
| cluster_overlap | <i>Create a table of module overlap for two clusterings</i> |
|-----------------|---|

Description

Compares two clusterings and creates a table of overlap between them. Module labels do not have to match.

Usage

```
cluster_overlap(k1, k2)
```

Arguments

| | |
|----|-------------------|
| k1 | First clustering |
| k2 | Second clustering |

Value

A matrix showing the module overlap with the labels of k1 in the columns and the labels of k2 in the rows.

| | |
|----------|--|
| fast_cor | <i>Fast computation of correlation</i> |
|----------|--|

Description

This uses a more memory-intensive but much faster algorithm than the built-in cor function.

Usage

```
fast_cor(x, y)
```

Arguments

| | |
|---|---------------------|
| x | first input matrix |
| y | second input matrix |

Details

Computes the correlation between the columns of x and y.

Value

Correlations matrix between the columns of x and y

find_module_sizes *Determine module sizes*

Description

Determine module sizes

Usage

```
find_module_sizes(module, n_modules)
```

Arguments

| | |
|-----------|--------------------------|
| module | Vector of module indices |
| n_modules | Total number of modules |

Value

A named vector containing the name of the module (its index or "Noise") and the number of elements in that module

get_avg_num_regulators
Get the average number of active regulators per module

Description

Get the average number of active regulators per module

Usage

```
get_avg_num_regulators(fit)
```

Arguments

| | |
|-----|-------------------------------|
| fit | An object of class scRegClust |
|-----|-------------------------------|

Value

A [data.frame](#) containing the average number of active regulators per module for each penalization parameter.

get_num_final_configs *Return the number of final configurations*

Description

Returns the number of final configurations per penalization parameter in an scRegClust object.

Usage

```
get_num_final_configs(fit)
```

Arguments

fit An object of class scRegClust

Value

An integer vector containing the number of final configurations for each penalization parameter.

get_rand_indices *Compute Rand indices*

Description

Compute Rand indices for fitted scregclust object

Usage

```
get_rand_indices(fit, groundtruth, adjusted = TRUE)
```

Arguments

fit An object of class scregclust
groundtruth A known clustering of the target genes (integer vector)
adjusted If TRUE, the Adjusted Rand index is computed. Otherwise the ordinary Rand index is computed.

Value

A `data.frame` containing the Rand indices. Since there can be more than one final configuration for some penalization parameters, Rand indices are averaged for each fixed penalization parameter. Returned are the mean, standard deviation and number of final configurations that were averaged.

References

- W. M. Rand (1971). "Objective criteria for the evaluation of clustering methods". *Journal of the American Statistical Association* 66 (336): 846–850. DOI:10.2307/2284239
- Lawrence Hubert and Phipps Arabie (1985). "Comparing partitions". *Journal of Classification*. 2 (1): 193–218. DOI:10.1007/BF01908075

get_regulator_list *Return list of regulator genes*

Description

Return list of regulator genes

Usage

```
get_regulator_list(mode = c("TF", "kinase"))
```

Arguments

mode Determines which genes are considered to be regulators. Currently supports TF=transcription factors and kinases.

Value

a list of gene symbols

See Also

[scregclust_format\(\)](#)

get_target_gene_modules
Extract target gene modules for given penalization parameters

Description

Extract target gene modules for given penalization parameters

Usage

```
get_target_gene_modules(fit, penalization = NULL)
```

Arguments

| | |
|--------------|--|
| fit | An object of class <code>scregclust</code> |
| penalization | A numeric vector of penalization parameters. The penalization parameters specified here must have been used during fitting of the <code>fit</code> object. |

Value

A list of lists of final target modules. One list for each parameter in `penalization`. The lists contain the modules of target genes for each final configuration.

| | |
|----------|--|
| kmeanspp | <i>Perform the k-means++ algorithm</i> |
|----------|--|

Description

Performs the k-means++ algorithm to cluster the rows of the input matrix.

Usage

```
kmeanspp(x, n_cluster, n_init_clusterings = 10L, n_max_iter = 10L)
```

Arguments

| | |
|--------------------|--|
| x | Input matrix (n x p) |
| n_cluster | Number of clusters |
| n_init_clusterings | Number of repeated random initializations to perform |
| n_max_iter | Number of maximum iterations to perform in the k-means algorithm |

Details

Estimation is repeated

Value

An object of class `stats::kmeans`.

References

David Arthur and Sergei Vassilvitskii. K-Means++: The advantages of careful seeding. In Proceedings of the Eighteenth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA '07, pages 1027—1035. Society for Industrial and Applied Mathematics, 2007.

plot_module_count_helper

Plot average silhouette scores and average predictive R^2

Description

Plot average silhouette scores and average predictive R^2

Usage

```
plot_module_count_helper(list_of_fits, penalization)
```

Arguments

| | |
|--------------|--|
| list_of_fits | A list of scregclust objects each fit to the same dataset across a variety of module counts (varying n_modules while running scregclust). |
| penalization | Either a single numeric value requesting the results for the same penalty parameter across all fits in list_of_fits, or one for each individual fit. |

Value

A ggplot2 plot showing the average silhouette score and the average predictive R^2

plot_regulator_network

Plotting the regulatory table from scregclust as a directed graph

Description

Plotting the regulatory table from scregclust as a directed graph

Usage

```
plot_regulator_network(
  output,
  arrow_size = 0.3,
  edge_scaling = 30,
  no_links = 6,
  col = c("gray80", "#FC7165", "#BD828C", "#9D8A9F", "#7D92B2", "#BDA88C", "#FCBD65",
    "#F2BB90", "#E7B9BA", "#BDB69C", "#92B27D", "#9B8BA5", "#9D7DB2", "#94A5BF")
)
```


Arguments

| | |
|--------------|--|
| output | Object of type scregclust_output from a fit of the scregclust algorithm. |
| arrow_size | Size of arrow head |
| edge_scaling | Scaling factor for edge width |
| no_links | Threshold value (0-10) for number of edges to show, higher value = more stringent threshold = less edges |
| col | color |

Value

Graph with gene modules and regulators as nodes

| | |
|------------------|--|
| plot_silhouettes | <i>Plot individual silhouette scores</i> |
|------------------|--|

Description

Plot individual silhouette scores

Usage

```
plot_silhouettes(list_of_fits, penalization, final_config = 1L)
```

Arguments

| | |
|--------------|---|
| list_of_fits | A list of scregclust objects each fit to the same dataset across a variety of module counts (varying n_modules when running scregclust). |
| penalization | Either a single numeric value requesting the results for the same penalty parameter across all fits in list_of_fits, or one for each individual fit. |
| final_config | The final configuration that should be visualized. Either a single number to be used for all fits in list_of_fits, or one for each individual fit. |

Value

A ggplot2 plot showing the the silhouette scores for each supplied fit.

| | |
|------------|---|
| scregclust | <i>Uncover gene modules and their regulatory programs from single-cell data</i> |
|------------|---|

Description

Use the scRegClust algorithm to determine gene modules and their regulatory programs from single-cell data.

Usage

```
scregclust(  
  expression,  
  genesymbols,  
  is_regulator,  
  penalization,  
  n_modules,  
  initial_target_modules = NULL,  
  sample_assignment = NULL,  
  center = TRUE,  
  split1_proportion = 0.5,  
  total_proportion = 1,  
  split_indices = NULL,  
  prior_indicator = NULL,  
  prior_genesymbols = NULL,  
  prior_baseline = 1e-06,  
  prior_weight = 0.5,  
  min_module_size = 0L,  
  allocate_per_obs = TRUE,  
  noise_threshold = 0.025,  
  n_cycles = 50L,  
  use_kmeanspp_init = TRUE,  
  n_initializations = 50L,  
  max_optim_iter = 10000L,  
  tol_coop_rel = 1e-08,  
  tol_coop_abs = 1e-12,  
  tol_nnl = 1e-04,  
  compute_predictive_r2 = TRUE,  
  compute_silhouette = FALSE,  
  nowarnings = FALSE,  
  verbose = TRUE,  
  quick_mode = FALSE,  
  quick_mode_percent = 0.1  
)
```

Arguments

| | |
|------------------------|---|
| expression | $p \times n$ matrix of pre-processed single cell expression data with p rows of genes and n columns of cells. |
| genesymbols | A vector of gene names corresponding to rows of expression. Has to be of length p . |
| is_regulator | An indicator vector where 1 indicates that the corresponding row in expression is a candidate regulator. All other rows represent target genes. Has to be of length p . |
| penalization | Sparsity penalty related to the amount of regulators associated with each module. Either a single positive number or a vector of positive numbers. |
| n_modules | Requested number of modules (integer). If this is provided without specifying <code>initial_target_modules</code> , then an initial module allocation is performed on the cross-correlation matrix of targets and genes on the first dataset after data splitting. |
| initial_target_modules | The initial assignment of target genes to modules of length <code>sum(is_regulator == 0L)</code> . If this is not specified, then see <code>n_modules</code> regarding module initialization. If provided, <code>use_kmeanspp_init</code> and <code>n_initializations</code> are ignored. |
| sample_assignment | A vector of sample assignment for each cell, can be used to perform the data splitting with stratification. Has to be of length n . No stratification if NULL is supplied. |
| center | Whether or not genes should be centered within each subgroup defined in <code>sample_assignment</code> . |
| split1_proportion | The proportion to use for the first dataset during data splitting. The proportion for the second dataset is $1 - \text{split1_proportion}$. If stratification with <code>sample_assignment</code> is used, then the proportion of each strata is controlled. |
| total_proportion | Can be used to only use a proportion of the supplied observations. The proportion of the first dataset during data splitting in relation to the full dataset will be $\text{total_proportion} * \text{split1_proportion}$. |
| split_indices | Can be used to provide an explicit data split. If this is supplied then <code>split1_proportion</code> , and <code>total_proportion</code> are ignored. Note that if <code>sample_assignment</code> is provided and <code>center == TRUE</code> , then subgroup centering will be performed as in the case of random splitting. A vector of length n containing entries 1 for cells in the first data split, 2 for cells in the second data split and NA for cells that should be excluded from the computations. |
| prior_indicator | An indicator matrix (sparse or dense) of size $q \times q$ that indicates whether there is a known functional relationship between two genes. Ideally, this is supplied as a sparse matrix (<code>sparseMatrix</code> in the <code>Matrix</code> package). If not, then the matrix is converted to one. |
| prior_genesymbols | A vector of gene names of length q corresponding to the rows/columns in <code>prior_indicator</code> . Does not have to be the same as <code>genesymbols</code> , but only useful if there is overlap. |

| | |
|------------------------------------|---|
| <code>prior_baseline</code> | A positive baseline for the network prior. The larger this parameter is, the less impact the network prior will have. |
| <code>prior_weight</code> | A number between 0 and 1 indicating the strength of the prior in relation to the data. 0 ignores the prior and makes the algorithm completely data-driven. 1 uses only the prior during module allocation. |
| <code>min_module_size</code> | Minimum required size of target genes in a module. Smaller modules are emptied. |
| <code>allocate_per_obs</code> | Whether module allocation should be performed for each observation in the second data split separately. If FALSE, target genes are allocated into modules on the aggregate sum of squares across all observations in the second data split. |
| <code>noise_threshold</code> | Threshold for the best R^2 of a target gene before it gets identified as noise. |
| <code>n_cycles</code> | Number of maximum algorithmic cycles. |
| <code>use_kmeanspp_init</code> | Use kmeans++ for module initialization if <code>initial_target_modules</code> is a single integer; otherwise use kmeans with random initial cluster centers |
| <code>n_initializations</code> | Number of kmeans(++) initialization runs. |
| <code>max_optim_iter</code> | Maximum number of iterations during optimization in the coop-Lasso and NNLS steps. |
| <code>tol_coop_rel</code> | Relative convergence tolerance during optimization in the coop-Lasso step. |
| <code>tol_coop_abs</code> | Absolute convergence tolerance during optimization in the coop-Lasso step. |
| <code>tol_nnls</code> | Convergence tolerance during optimization in the NNLS step. |
| <code>compute_predictive_r2</code> | Whether to compute predictive R^2 per module as well as regulator importance. |
| <code>compute_silhouette</code> | Whether to compute silhouette scores for each target gene. |
| <code>nowarnings</code> | When turned on then no warning messages are shown. |
| <code>verbose</code> | Whether to print progress. |
| <code>quick_mode</code> | Whether to use a reduced number of noise targets to speed up computations. |
| <code>quick_mode_percent</code> | A number in $[0, 1)$ indicating the amount of noise targets to use in the re-allocation process if <code>quick_mode = TRUE</code> . |

Value

A list with S3 class `scregclust` containing

| | |
|-------------------------------------|---|
| <code>penalization</code> | The supplied penalization parameters |
| <code>results</code> | A list of result lists (each with S3 class <code>scregclust_result</code>), one for each supplied penalization parameter. See below. |
| <code>initial_target_modules</code> | Initial allocation of target genes into modules. |

`split_indices` either verbatim the vector given as input or a vector encoding the splits as NA = not included, 1 = split 1 or 2 = split 2. Allows reproducibility of data splits.

For each supplied penalization parameter, `results` contains a list with

- the current penalization parameter,
- the supplied `genesymbols` after filtering (as used during fitting),
- the supplied `is_regulator` vector after filtering (as used during fitting),
- the number of fitted modules `n_modules`,
- whether the current run converged to a single configuration (as a boolean),
- as well as an output object containing the numeric results for each final configuration.

It is possible that the algorithm ends in a finite cycle of configurations instead of a unique final configuration. Therefore, output is a list with each element itself being a list with the following contents:

`reg_table` a regulator table, a matrix of weights for each regulator and module

`module` vector of same length as `genesymbols` containing the module assignments for all genes with regulators marked as NA. Genes considered noise are marked as -1.

`module_all` same as `module`, however, genes that were marked as noise (-1 in `module`) are assigned to the module in which it has the largest R^2 , even if it is below `noise_threshold`.

`r2` matrix of predictive R^2 value for each target gene and module

`best_r2` vector of best predictive R^2 for each gene (regulators marked with NA)

`best_r2_idx` module index corresponding to best predictive R^2 for each gene (regulators marked with NA)

`r2_module` a vector of predictive R^2 values for each module (included if `compute_predictive_r2 == TRUE`)

`importance` a matrix of importance values for each regulator (rows) and module (columns) (included if `compute_predictive_r2 == TRUE`)

`r2_cross_module_per_target` a matrix of cross module R^2 values for each target gene (rows) and each module (columns) (included if `compute_silhouette == TRUE`)

`silhouette` a vector of silhouette scores for each target gene (included if `compute_silhouette == TRUE`)

`models` regulator selection for each module as a matrix with regulators in rows and modules in columns

`signs` regulator signs for each module as a matrix with regulators in rows and modules in columns

`weights` average regulator coefficient for each module

`coeffs` list of regulator coefficient matrices for each module for all target genes as re-estimated in the NNLS step

`sigmas` matrix of residual variances, one per target gene in each module; derived from the residuals in NNLS step

scregclust_format *Package data before clustering*

Description

Package data before clustering

Usage

```
scregclust_format(expression_matrix, mode = c("TF", "kinase"))
```

Arguments

`expression_matrix` The $p \times n$ gene expression matrix with gene symbols as rownames.
`mode` Determines which genes are considered to be regulators.

Value

A list with

`genesymbols` The gene symbols extracted from the expression matrix
`sample_assignment` A vector filled with 1's of the same length as there are columns in the gene expression matrix.
`is_regulator` Whether a gene is considered to be a regulator or not, determined dependent on mode.

See Also

[get_regulator_list\(\)](#)

Index

* helpers

- available_results, 2
- cluster_overlap, 3
- fast_cor, 3
- find_module_sizes, 4
- kmeanspp, 7

* main

- scregclust, 10
- scregclust_format, 14

* plotting

- plot_module_count_helper, 8
- plot_regulator_network, 8
- plot_silhouettes, 9

* utilities

- get_avg_num_regulators, 4
- get_num_final_configs, 5
- get_rand_indices, 5
- get_regulator_list, 6
- get_target_gene_modules, 6

available_results, 2

cluster_overlap, 3

data.frame, 2, 4, 5

fast_cor, 3

find_module_sizes, 4

get_avg_num_regulators, 4

get_num_final_configs, 5

get_rand_indices, 5

get_regulator_list, 6

get_regulator_list(), 14

get_target_gene_modules, 6

kmeanspp, 7

plot_module_count_helper, 8

plot_regulator_network, 8

plot_silhouettes, 9

scregclust, 8, 9, 10

scregclust_format, 14

scregclust_format(), 6

stats::kmeans, 7