

Package ‘celaref’

December 11, 2018

Title Single-cell RNAseq cell cluster labelling by reference

Version 1.0.0

Description After the clustering step of a single-cell RNAseq experiment, this package aims to suggest labels/cell types for the clusters, on the basis of similarity to a reference dataset. It requires a table of read counts per cell per gene, and a list of the cells belonging to each of the clusters, (for both test and reference data).

Depends R (>= 3.5.0), SummarizedExperiment

Imports MAST, ggplot2, Matrix, dplyr, magrittr, stats, utils, rlang, BiocGenerics, S4Vectors, readr, tibble

Suggests limma, parallel, knitr, rmarkdown, ExperimentHub

biocViews SingleCell

VignetteBuilder knitr

License GPL-3

Encoding UTF-8

LazyData true

RoxygenNote 6.1.0

git_url <https://git.bioconductor.org/packages/celaref>

git_branch RELEASE_3_8

git_last_commit a2ffb05

git_last_commit_date 2018-10-30

Date/Publication 2018-12-11

Author Sarah Williams [aut, cre]

Maintainer Sarah Williams <sarah.williams1@monash.edu>

R topics documented:

contrast_each_group_to_the_rest	2
contrast_each_group_to_the_rest_for_norm_ma_with_limma	3
contrast_the_group_to_the_rest	5
contrast_the_group_to_the_rest_with_limma_for_microarray	6
convert_se_gene_ids	7
demo_cell_info_table	8
demo_counts_matrix	8

demo_gene_info_table	9
demo_microarray_expr	9
demo_microarray_sample_sheet	10
demo_query_se	10
demo_ref_se	11
de_table.demo_query	11
de_table.demo_ref	12
find_within_match_differences	12
get_inner_or_outer_ci	13
get_limma_top_table_with_ci	13
get_matched_stepped_mwtest_res_table	14
get_ranking_and_test_results	15
get_rankstat_table	16
get_reciprocal_matches	17
get_stepped_pvals_str	17
get_the_up_genes_for_all_possible_groups	18
get_the_up_genes_for_group	19
get_vs_random_pval	20
load_dataset_10Xdata	20
load_se_from_tables	22
make_ranking_violin_plot	24
make_ref_similarity_names	25
make_ref_similarity_names_for_group	28
run_pair_test_stats	29
trim_small_groups_and_low_expression_genes	30

Index 32

contrast_each_group_to_the_rest
contrast_each_group_to_the_rest

Description

Produces a table of within-experiment differential expression results (for either query or reference experiment), where each group (cluster) is compared to the rest of the cells.

Usage

```
contrast_each_group_to_the_rest(dataset_se, dataset_name,
  groups2test = NA, num_cores = 2)
```

Arguments

dataset_se	Summarised experiment object containing count data. Also requires 'ID' and 'group' to be set within the cell information (see colData())
dataset_name	Short, meaningful name for this dataset/experiment.
groups2test	An optional character vector specifying specific groups to check. By default (set to NA), all groups will be tested.
num_cores	Number of cores to use to run MAST jobs in parallel. Ignored if parallel package not available. Set to 1 to avoid parallelisation. Default = 2

Details

Note that this function is *slow*, because it runs the differential expression. It only needs to be run once per dataset though (unless group labels change). Having package **parallel** installed is highly recommended.

Both reference and query datasets should be processed with this function.

The tables produced by this function (usually named something like *de_table.datasetname*) contain summarised results of MAST results. Each group is compared versus cells in the group, versus not in the group, (Ie. always a 2-group contrast, other groups information is ignored). As per MAST recommendations, the proportion of genes seen in each cell is included in the model.

Value

A tibble the within-experiment *de_table* (differential expression table). This is a core summary of the individual experiment/dataset, which is used for the cross-dataset comparisons.

The table fields won't necessarily match across datasets, as they include cell annotations information. Important columns (used in downstream analysis) are:

ID Gene identifier

ci_inner Inner (conservative) 95% confidence interval of log2 fold-change.

fd Multiple hypothesis corrected p-value (using BH/FDR method)

group Cells from this group were compared to everything else

sig_up Significantly differentially expressed ($\text{fd} < 0.01$), with a positive fold change?

rank Rank position (within group), ranked by CI inner, highest to lowest.

rescaled_rank Rank scaled 0 (top most overrepresented genes in group) - 1 (top most not-present genes)

dataset Name of dataset/experiment

Examples

```
de_table.demo_query <- contrast_each_group_to_the_rest(
  demo_query_se, "a_demo_query")
```

```
## Not run:
```

```
de_table.demo_ref <- contrast_each_group_to_the_rest(
  demo_ref_se, "a_demo_ref", num_cores=4)
```

```
## End(Not run)
```

```
contrast_each_group_to_the_rest_for_norm_ma_with_limma
```

```
  contrast_each_group_to_the_rest_for_norm_ma_with_limma
```

Description

This function loads and processes microarray data (from purified cell populations) that can be used as a reference.

Usage

```
contrast_each_group_to_the_rest_for_norm_ma_with_limma(norm_expression_table,
  sample_sheet_table, dataset_name, sample_name, group_name = "group",
  groups2test = NA, extra_factor_name = NA, pval_threshold = 0.01)
```

Arguments

<code>norm_expression_table</code>	A logged, normalised expression table. Any filtering (removal of low-expression probes/genes)
<code>sample_sheet_table</code>	Tab-separated text file of sample information. Columns must have names. Sample/microarray ids should be listed under sample_name column. The cell-type (or 'group') of each sample should be listed under a column group_name .
<code>dataset_name</code>	Short, meaningful name for this dataset/experiment.
<code>sample_name</code>	Name of sample_sheet_table with sample ID
<code>group_name</code>	Name of sample_sheet_table with group/cell-type. Default = "group"
<code>groups2test</code>	An optional character vector specifying specific groups to check. By default (set to NA), all groups will be tested.
<code>extra_factor_name</code>	Optionally, an extra cross-group factor (as column name in sample_sheet_table) to include in the model used by limma. E.g. An individual/mouse id. Refer limma docs. Default = NA
<code>pval_threshold</code>	For reporting only, a p-value threshold. Default = 0.01

Details

Sometimes there are microarray studies measuring purified cell populations that would be measured together in a single-cell sequencing experiment. E.g. comparing PBMC scRNA to FACs-sorted blood cell populations. This function will process microarray data with limma and format it for comparisons.

The microarray data used should consist of purified cell types from *one* single study/experiment (due to batch effects). Ideally just those cell-types expected in the scRNAseq, but the method appears relatively robust to a few extra cell types.

Note that unlike the single-cell workflow there are no summarisedExperiment objects (they're not really comparable) - this function reads data and generates a table of within-dataset differential expression contrasts in one step. I.e. equivalent to the output of [contrast_each_group_to_the_rest](#).

Also, note that while downstream functions can accept the microarray-derived data as query datasets, it's not really intended and assumptions might not hold (Generally, it's known what got loaded onto a microarray!)

The (otherwise optional) 'limma' package must be installed to use this function.

Value

A tibble, the within-experiment `de_table` (differential expression table)

See Also

[contrast_each_group_to_the_rest](#) is the function that makes comparable output on the scRNAseq data (dataset_se objects).

Limma Limma package for differential expression.

Other Data loading functions: [load_dataset_10Xdata](#), [load_se_from_tables](#)

Examples

```
contrast_each_group_to_the_rest_for_norm_ma_with_limma(
  norm_expression_table=demo_microarray_expr,
  sample_sheet_table=demo_microarray_sample_sheet,
  dataset_name="DemoSimMicroarrayRef",
  sample_name="cell_sample", group_name="group")

## Not run:
contrast_each_group_to_the_rest_for_norm_ma_with_limma(
  norm_expression_table, sample_sheet_table=samples_table,
  dataset_name="Watkins2009PBMCs", extra_factor_name='description')

## End(Not run)
```

```
contrast_the_group_to_the_rest
      contrast_the_group_to_the_rest
```

Description

Internal function to calculate differential expression within an experiment between a specified group and cells not in that group.

Usage

```
contrast_the_group_to_the_rest(dataset_se, the_group,
  pvalue_threshold = 0.01)
```

Arguments

dataset_se	Datast summarisedExperiment object.
the_group	group to test
pvalue_threshold	Default = 0.01

Details

This function should only be called by `contrast_each_group_to_the_rest` (which can be passed a single group name if desired). Else 'pofgenes' will not be defined.

MAST is supplied with $\log_2(\text{counts} + 1.1)$, and `zlm` called with model `'~ TvsR + pofgenes'`. The p-values reported are from the hurdle model. FDR is with default `fdr/BH` method.

Value

A tibble, the within-experiment `de_table` (differential expression table), for the group specified.

See Also

[contrast_each_group_to_the_rest](#)

```
contrast_the_group_to_the_rest_with_limma_for_microarray
      contrast_the_group_to_the_rest_with_limma_for_microarray
```

Description

Private function used by `contrast_each_group_to_the_rest_for_norm_ma_with_limma`

Usage

```
contrast_the_group_to_the_rest_with_limma_for_microarray(norm_expression_table,
  sample_sheet_table, the_group, sample_name, extra_factor_name = NA,
  pval_threshold = 0.01)
```

Arguments

<code>norm_expression_table</code>	A logged, normalised expression table. Any filtering (removal of low-expression probes/genes)
<code>sample_sheet_table</code>	Tab-separated text file of sample information. Columns must have names. Sample/microarray ids should be listed under sample_name column. The cell-type (or 'group') of each sample should be listed under a column group_name .
<code>the_group</code>	Which query group is being tested.
<code>sample_name</code>	Name of sample_sheet_table with sample ID
<code>extra_factor_name</code>	Optionally, an extra cross-group factor (as column name in sample_sheet_table) to include in the model used by limma. E.g. An individual/mouse id. Refer limma docs. Default = NA
<code>pval_threshold</code>	For reporting only, a p-value threshold. Default = 0.01

Value

A tibble, the within-experiment `de_table` (differential expression table), for the group specified.

See Also

[contrast_each_group_to_the_rest_for_norm_ma_with_limma](#) public calling function

Limma Limma package for differential expression.

```
convert_se_gene_ids  convert_se_gene_ids
```

Description

Change the gene IDs in in the supplied dataset_se object to some other id already present in the gene info (as seen with rowData())

Usage

```
convert_se_gene_ids(dataset_se, new_id, eval_col, find_max = TRUE)
```

Arguments

dataset_se	Summarised experiment object containing count data. Also requires 'ID' and 'group' to be set within the cell information (see colData())
new_id	A column within the feature information (view colData(dataset_se)) of the dataset_se , which will become the new ID column. Non-uniqueness of this column is handled gracefully! Any <i>NAs</i> will be dropped.
eval_col	Which column to use to break ties of duplicate new_id . Must be a column within the feature information (view colData(dataset_se)) of the dataset_se . Total reads per gene feature is a good choice.
find_max	If false, this will choose the minimal eval_col instead of max. Default = TRUE

Value

A modified dataset_se - ID will now be **new_id**, and unique. It will have fewer genes if old ID to new ID was not a 1:1 mapping. The selected genes will be according to the eval col max (or min). *should* pick the alphabetical first on ties, but could change.

See Also

[SummarizedExperiment](#) For general doco on the SummarizedExperiment objects.

[load_se_from_files](#) For reading data from flat files (not 10X cellRanger output)

Examples

```
# The demo dataset doesn't have other names, so make some up
# (don't do this)
dataset_se <- demo_ref_se
rowData(dataset_se)$dummyname <- toupper(rowData(dataset_se)$ID)

# If not already present, define a column to evaluate,
# typically total reads/gene.
rowData(dataset_se)$total_count <- rowSums(assay(dataset_se))

dataset_se <- convert_se_gene_ids(dataset_se, new_id='dummyname', eval_col='total_count')
```

demo_cell_info_table *Demo cell info table*

Description

Sample sheet table listing each cell, its assigned cluster/group, and any other information that might be interesting (replicate, individual e.t.c)

Usage

demo_cell_info_table

Format

An object of class `data.frame` with 515 rows and 4 columns.

Value

An example cell info table

demo_counts_matrix *Demo count matrix*

Description

Counts matrix for a small, demo example datasets. Raw counts of reads per gene (row) per cell (column).

Usage

demo_counts_matrix

Format

An object of class `matrix` with 200 rows and 514 columns.

Value

An example counts matrix.

demo_gene_info_table *Demo gene info table*

Description

Extra table of gene-level information for the demo example dataset. Can contain anything as long as there's a unique gene id.

Usage

demo_gene_info_table

Format

An object of class `data.frame` with 200 rows and 2 columns.

Value

An example table of genes.

demo_microarray_expr *Demo microarray expression table*

Description

Microarray-style expression table for the demo example dataset. Rows are genes, columns are samples, as per counts matrix.

Usage

demo_microarray_expr

Format

An object of class `matrix` with 200 rows and 20 columns.

Value

An example table of (fake) microarray data.

demo_microarray_sample_sheet

Demo microarray sample sheet table

Description

Microarray sample sheet table for the demo example dataset. Contains array identifiers, their group and any other information that could be useful.

Usage

demo_microarray_sample_sheet

Format

An object of class grouped_df (inherits from tbl_df, tbl, data.frame) with 20 rows and 2 columns.

Value

An example microarray sample sheet

demo_query_se

Demo query se (summarizedExperiment)

Description

A summarisedExperiment object loaded from demo info tables, for a query set.

Usage

demo_query_se

Format

An object of class SummarizedExperiment with 200 rows and 485 columns.

Value

An example summarised experiment (for demo query dataset)

demo_ref_se	<i>Demo reference se (summarizedExperiment)</i>
-------------	---

Description

A summarisedExperiment object loaded from demo info tables, for a reference set.

Usage

```
demo_ref_se
```

Format

An object of class SummarizedExperiment with 200 rows and 515 columns.

Value

An example summarised experiment (for demo reference dataset)

de_table.demo_query	<i>Demo query de table</i>
---------------------	----------------------------

Description

Small example dataset that is the output of [contrast_each_group_to_the_rest](#). It contains the results of each group compared to the rest of the sample (ie within sample differential expression)

Usage

```
de_table.demo_query
```

Format

An object of class data.frame with 800 rows and 13 columns.

Value

An example de_table from [contrast_each_group_to_the_rest](#) (for demo query dataset)

```
de_table.demo_ref      Demo ref de table
```

Description

Small example dataset that is the output of [contrast_each_group_to_the_rest](#). It contains the results of each group compared to the rest of the sample (ie within sample differential expression)

Usage

```
de_table.demo_ref
```

Format

An object of class `data.frame` with 800 rows and 13 columns.

Value

An example `de_table` from [contrast_each_group_to_the_rest](#) (for demo ref dataset)

```
find_within_match_differences
      find_within_match_differences
```

Description

Internal function to find if there are significant difference between the distributions, when there are multiple match groups.

Usage

```
find_within_match_differences(de_table.ref.marked, matches, the_test_group,
                             the_test_dataset, the_ref_dataset, the_pval)
```

Arguments

```
de_table.ref.marked      see make_ref_similarity_names_for_group
matches                  see make_ref_similarity_names_for_group
the_test_group          see make_ref_similarity_names_for_group
the_test_dataset        see make_ref_similarity_names_for_group
the_ref_dataset          see make_ref_similarity_names_for_group
the_pval                see make_ref_similarity_names_for_group
```

Details

For use by `make_ref_similarity_names_for_group`

Value

String of within match differences

See Also

[make_ref_similarity_names_for_group](#)

get_inner_or_outer_ci *get_inner_or_outer_ci*

Description

Given a fold-change, and high and low confidence interval (where lower < higher), pick the inner-most/most conservative one.

Usage

```
get_inner_or_outer_ci(fc, ci.hi, ci.lo, get_inner = TRUE)
```

Arguments

fc	Fold-change
ci.hi	Higher fold-change CI (numerically)
ci.lo	smaller fold-change CI (numerically)
get_inner	If TRUE, get the more conservative inner CI, else the bigger outside one.

Value

inner or outer CI from **ci.hi** or **ci.low**

get_limma_top_table_with_ci
get_limma_top_table_with_ci

Description

Internal function that wraps limma topTable output but also adds upper and lower confidence intervals to the logFC. Calculated according to <https://support.bioconductor.org/p/36108/>

Usage

```
get_limma_top_table_with_ci(fit2, the_coef, ci = 0.95)
```

Arguments

fit2	The fit2 object after calling eBayes as per standard limma workflow. Ie object that topTable gets called on.
the_coef	Coefficient. As passed to topTable.
ci	Confidence interval. Number between 0 and 1, default 0.95 (95%)

Value

Output of topTable, but with the (95 for the logFC.

See Also

[contrast_the_group_to_the_rest_with_limma_for_microarray](#) Calling function.

`get_matched_stepped_mwtest_res_table`
get_matched_stepped_mwtest_res_table

Description

Internal function to grab a table of the matched group(s).

Usage

```
get_matched_stepped_mwtest_res_table(mwtest_res_table.this, the_pval)
```

Arguments

<code>mwtest_res_table.this</code>	Combined output of get_ranking_and_test_results
<code>the_pval</code>	Pvalue threshold

Details

For use by [make_ref_similarity_names_for_group](#)

Value

Stepped pvalues string

See Also

[make_ref_similarity_names_for_group](#)

get_ranking_and_test_results
get_ranking_and_test_results

Description

Internal function to get reference group similarity contrasts for an individual query group.

Usage

```
get_ranking_and_test_results(de_table.ref.marked, the_test_group,  
                             the_test_dataset, the_ref_dataset, num_steps, pval = 0.01)
```

Arguments

de_table.ref.marked
 see [make_ref_similarity_names_using_marked](#)

the_test_group The group to calculate the stats on.

the_test_dataset
 see [make_ref_similarity_names_using_marked](#)

the_ref_dataset
 see [make_ref_similarity_names_using_marked](#)

num_steps see [make_ref_similarity_names_using_marked](#)

pval see [make_ref_similarity_names_using_marked](#)

Details

For use by [make_ref_similarity_names_using_marked](#), see that function for parameter details.
This function just runs this for a single query group **the_test_group**

Value

Table of similarity contrast results/assigned names e.t.c for a single group. Used internally for populating mwtest_res_table tables.

See Also

[make_ref_similarity_names_using_marked](#) which calls this.

get_rankstat_table	<i>get_rankstat_table</i>
--------------------	---------------------------

Description

Summarise the comparison of the specified query group against in the comparison in **de_table.ref.marked** - number of 'top' genes and their median rank in each of the reference groups, with reference group rankings.

Usage

```
get_rankstat_table(de_table.ref.marked, the_test_group)
```

Arguments

`de_table.ref.marked`
The output of [get_the_up_genes_for_all_possible_groups](#) for the contrast of interest.

`the_test_group` Name of query group to test

Value

A tibble of query group name (`test_group`), number of 'top' genes (`n`), reference dataset group (`group`) with its ranking (`grouprank`) and the median (rescaled 0..1) ranking of 'top' genes (`median_rank`).

See Also

[get_the_up_genes_for_all_possible_groups](#) To prepare the **de_table.ref.marked** input.

Examples

```
# Make input
# de_table.demo_query <- contrast_each_group_to_the_rest(demo_query_se, "demo_query")
# de_table.demo_ref <- contrast_each_group_to_the_rest(demo_ref_se, "demo_ref")

de_table.marked.query_vs_ref <- get_the_up_genes_for_all_possible_groups(
  de_table.demo_query,
  de_table.demo_ref)

get_rankstat_table(de_table.marked.query_vs_ref, "Group3")
```

```
get_reciprocal_matches
    get_reciprocal_matches
```

Description

Internal function to run a binomial test of median test rank > 0.5 (random).

Usage

```
get_reciprocal_matches(mwtest_res_table.recip, de_table.recip.marked,
    the_pval)
```

Arguments

mwtest_res_table.recip	Combined output of get_ranking_and_test_results for reciprocal test - ref vs query.
de_table.recip.marked	Reciprocal ref vs query de_table.ref.marked
the_pval	See make_ref_similarity_names_using_marked

Details

For use by [make_ref_similarity_names_using_marked](#)

Value

List of table of reciprocal matches tested from reference to query.

See Also

[make_ref_similarity_names_using_marked](#)

```
get_stepped_pvals_str  get_stepped_pvals_str
```

Description

Internal function to construct the string of stepped pvalues reported by [make_ref_similarity_names_using_marked](#)

Usage

```
get_stepped_pvals_str(mwtest_res_table.this)
```

Arguments

mwtest_res_table.this	Combined output of get_ranking_and_test_results
-----------------------	---

Details

For use by `make_ref_similarity_names_for_group`

Value

Stepped pvalues string

See Also

[make_ref_similarity_names_for_group](#)

`get_the_up_genes_for_all_possible_groups`
get_the_up_genes_for_all_possible_groups

Description

For the most overrepresented genes of each group in the test dataset, get their rankings in all the groups of the reference dataset.

Usage

```
get_the_up_genes_for_all_possible_groups(de_table.test, de_table.ref,
  rankmetric = "TOP100_LOWER_CI_GTE1")
```

Arguments

<code>de_table.test</code>	A differential expression table of the query experiment, as generated from contrast_each_group_to
<code>de_table.ref</code>	A differential expression table of the reference dataset, as generated from contrast_each_group_to
<code>rankmetric</code>	Placeholder for support of different ranking methods, but only the default supported. Omit.

Details

This is effectively a subset of the reference data, 'marked' with the 'top' genes that represent the groups in the query data. The distribution of the *rescaled ranks* of these marked genes in each reference data group indicate how similar they are to the query group.

This function is simply a convenient wrapper for [get_the_up_genes_for_group](#) that merges output for each group in the query into one table.

Value

`de_table.marked` This will almost be a subset of `de_table.ref`, with an added column `test_group` set to the query groups, and `test_dataset` set to `test_dataset_name`.

If nothing passes the rankmetric criteria, a warning is thrown and NA is returned. (This can be a genuine inability to pick out the representative 'up' genes, or due to some problem in the analysis)

See Also

[get_the_up_genes_for_group](#) Function for testing a single group.

Examples

```
de_table.marked.query_vs_ref <- get_the_up_genes_for_all_possible_groups(
  de_table.test=de_table.demo_query ,
  de_table.ref=de_table.demo_ref )
```

```
get_the_up_genes_for_group
  get_the_up_genes_for_group
```

Description

For the most overrepresented genes of the specified group in the test dataset, get their rankings in all the groups of the reference dataset.

Usage

```
get_the_up_genes_for_group(the_group, de_table.test, de_table.ref,
  rankmetric = "TOP100_LOWER_CI_GTE1")
```

Arguments

<code>the_group</code>	The group (from the test/query experiment) to examine.
<code>de_table.test</code>	A differential expression table of the query experiment, as generated from contrast_each_group_to_the_rest .
<code>de_table.ref</code>	A differential expression table of the reference dataset, as generated from contrast_each_group_to_the_rest .
<code>rankmetric</code>	Placeholder for support of different ranking methods, but only the default supported. Omit.

Details

This is effectively a subset of the reference data, 'marked' with the 'top' genes that represent the group of interest in the query data. The distribution of the *rescaled ranks* of these marked genes in each reference data group indicate how similar they are to the query group.

Value

de_table.marked This will be a subset of **de_table.ref**, with an added column *test_group* set to **the_group**. If nothing passes the rankmetric criteria, NA.

See Also

[contrast_each_group_to_the_rest](#) For preparing the `de_table.*` tables. [get_the_up_genes_for_all_possible_groups](#) For running all query groups at once.

Examples

```
de_table.marked.Group3vsRef <- get_the_up_genes_for_group(
  the_group="Group3",
  de_table.test=de_table.demo_query,
  de_table.ref=de_table.demo_ref)
```

get_vs_random_pval *get_vs_random_pval*

Description

Internal function to run a binomial test of median test rank > 0.5 (random).

Usage

```
get_vs_random_pval(de_table.ref.marked, the_group, the_test_group)
```

Arguments

de_table.ref.marked
 see `make_ref_similarity_names_for_group`

the_group Reference group name

the_test_group Test group name #'

Details

For use by `make_ref_similarity_names_for_group`

Value

Pvalue result of a binomial test of each 'top gene' being greater than the theoretical random median rank of 0.5 (halfway).

See Also

[make_ref_similarity_names_for_group](#)

load_dataset_10Xdata *load_dataset_10Xdata*

Description

Convenience function to create a SummarizedExperiment object (`dataset_se`) from a the output of 10X cell ranger pipeline run.

Usage

```
load_dataset_10Xdata(dataset_path, dataset_genome, clustering_set,
  gene_id_cols_10X = c("ensembl_ID", "GeneSymbol"),
  id_to_use = gene_id_cols_10X[1])
```

Arguments

dataset_path	Path to the directory of 10X data, as generated by the cellRanger pipeline (versions 2.1.0 and 2.0.1). The directory should have subdirectories <i>analysis</i> , <i>filtered_gene_bc_matrices</i> and <i>raw_gene_bc_matrices</i> (only the first 2 are read).
dataset_genome	The genome that the reads were aligned against, e.g. GRCh38. Check for this as a directory name under the <i>filtered_gene_bc_matrices</i> subdirectory if unsure.
clustering_set	The 10X cellRanger pipeline produces several different cluster definitions per dataset. Specify which one to use e.g. kmeans_10_clusters Find them as directory names under <i>analysis/clustering/</i>
gene_id_cols_10X	Vector of the names of the columns in the gene description file (<i>filtered_gene_bc_matrices/GRCh38/gene</i>). The first element of this will become the ID. Default = c("ensembl_ID", "GeneSymbol")
id_to_use	Column from gene_id_cols_10X that defines the gene identifier to use as 'ID' in the returned SummarisedExperiment object. Many-to-one relationships between the assumed unique first element of gene_id_cols_10X and id_to_use will be handled gracefully by convert_se_gene_ids . Defaults to first element of gene_id_cols_10X

Details

This function makes a SummarizedExperiment object in a form that should work for celaref functions. Specifically, that means it will have an 'ID' feild for genes (view with `rowData(dataset_se)`), and both 'cell_sample' and 'group' feild for cells (view with `colData(dataset_se)`). See parameters for detail. Additionally, the counts will be an integer matrix (not a sparse matrix), and the *group* feild (but not *cell_sample* or *ID*) will be a factor.

The clustering information can be read from whichever cluster is specified, usually there will be several choices.

This funciton is designed to work with output of version 2.0.1 of the cellRanger pipeline, may not work with others (will not work for 1.x).

Value

A SummarisedExperiment object containing the count data, cell info and gene info.

See Also

[SummarizedExperiment](#) For general doco on the SummarizedExperiment objects.

[convert_se_gene_ids](#) describes method for converting IDs.

Other Data loading functions: [contrast_each_group_to_the_rest_for_norm_ma_with_limma](#), [load_se_from_tables](#)

Examples

```
example_10X_dir <- system.file("extdata", "sim_cr_dataset", package = "celaref")
dataset_se <- load_dataset_10Xdata(example_10X_dir, dataset_genome="GRCh38",
  clustering_set="kmeans_4_clusters", gene_id_cols_10X=c("gene"))

## Not run:
dataset_se <- load_dataset_10Xdata('~/.path/to/data/10X_pbmc4k',
  dataset_genome="GRCh38",
  clustering_set="kmeans_7_clusters")
```

```
## End(Not run)
```

```
load_se_from_tables  load_se_from_tables
```

Description

Create a SummarizedExperiment object (dataset_se) from a count matrix, cell information and optionally gene information.

load_se_from_files is a wrapper for load_se_from_tables that will read in tables from specified files.

Usage

```
load_se_from_tables(counts_matrix, cell_info_table, gene_info_table = NA,
  group_col_name = "group", cell_col_name = NA)
```

```
load_se_from_files(counts_file, cell_info_file, gene_info_file = NA,
  group_col_name = "group", cell_col_name = NA)
```

Arguments

counts_matrix A tab-separated matrix of read counts for each gene (row) and each cell (column). Columns and rows should be named.

cell_info_table Table of cell information. If there is a column labelled *cell_sample*, that will be used as the unique cell identifiers. If not, the first column is assumed to be cell identifiers, and will be copied to a new feild labelled *cell_sample*. Similarly - the clusters of these cells should be listed in one column - which can be called 'group' (case-sensitive) or specified with **group_col_name**. *Minimal data format: <cell_sample> <group>*

gene_info_table Optional table of gene information. If there is a column labelled *ID*, that will be used as the gene identifiers (they must be unique!). If not, the first column is assumed to be a gene identifier, and will be copied to a new feild labelled *ID*. Must match all rownames in **counts_matrix**. If omitted, ID will be generated from the rownames of counts_matrix. Default=NA

group_col_name Name of the column in **cell_info_table** containing the cluster/group that each cell belongs to. Case-sensitive. Default='group'

cell_col_name Name of the column in **cell_info_table** containing a cell id. Ignored if *cell_sample* column is already present. If omitted, (and no *cell_sample* column) will use first column. Case-sensitive. Default=NA

counts_file A tab-separated file of a matrix of read counts. As per **counts_matrix**. First column should be gene ID, and top row cell ids.

cell_info_file Tab-separated text file of cell information, as per **cell_info_table**. Columns must have names.

gene_info_file Optional tab-separated text file of gene information, as per **gene_info_file**. Columns must have names. Default=NA

```
make_ranking_violin_plot
      make_ranking_violin_plot
```

Description

Plot a panel of violin plots showing the distribution of the 'top' genes of each of query group, across the reference dataset.

Usage

```
make_ranking_violin_plot(de_table.marked = NA, de_table.test = NA,
  de_table.ref = NA, log10trans = FALSE)
```

Arguments

<code>de_table.marked</code>	The output of get_the_up_genes_for_all_possible_groups for the contrast of interest.
<code>de_table.test</code>	A differential expression table of the query experiment, as generated from contrast_each_group_to_the_rest .
<code>de_table.ref</code>	A differential expression table of the reference dataset, as generated from contrast_each_group_to_the_rest .
<code>log10trans</code>	Plot on a log scale? Useful for distinquishing multiple similar, yet distinct cell type that bunch at top of plot. Default=FALSE.

Details

In the plot output, each panel correponds to a different group/cluster in the query experiment. The x-axis has the groups in the reference dataset. The y-axis is the rescaled rank of each 'top' gene from the query group, within each reference group.

Only the 'top' genes for each query group are plotted, forming the violin plots - each individual gene is shown as a tickmark. Some groups have few top genes, and so their uncertainty can be seen on this plot.

The thick black lines represent the median gene rescaled ranking for each query group / reference group combination. Having this fall above the dotted median threshold marker is a quick indication of potential similarity. A complete lack of similarity would have a median rank around 0.5. Median rankings much less than 0.5 are common though (an 'anti-cell-groupA' signature), because genes overrepresented in one group in an experiment, are likely to be relatively 'underrepresented' in the other groups. Taken to an extreme, if there are only two reference groups, they'll be complete opposites.

Input can be either the precomputed `de_table.marked` object for the comparison, OR both `de_table.test` and `de_table.ref` differential expression results to compare from [contrast_each_group_to_the_rest](#)

Value

A ggplot object.

See Also

[get_the_up_genes_for_all_possible_groups](#) To make the input data.

Examples

```
# Make input
# de_table.demo_query <- contrast_each_group_to_the_rest(demo_query_se, "demo_query")
# de_table.demo_ref <- contrast_each_group_to_the_rest(demo_ref_se, "demo_ref")

# This:
make_ranking_violin_plot(de_table.test=de_table.demo_query,
                        de_table.ref=de_table.demo_ref )

# Is equivalent to this:
de_table.marked.query_vs_ref <-
  get_the_up_genes_for_all_possible_groups( de_table.test=de_table.demo_query,
                                           de_table.ref=de_table.demo_ref)
make_ranking_violin_plot(de_table.marked.query_vs_ref)
```

```
make_ref_similarity_names
      make_ref_similarity_names
```

Description

Construct some sensible labels or the groups/clusters in the query dataset, based on similarity the reference dataset.

This is a more low level/customisable version of [make_ref_similarity_names](#), (would usually use that instead). Suitable for rare cases to reuse an existing **de_table.ref.marked** object. Or use a **de_table.ref.marked** table with more than one dataset present (discouraged). Or to skip the reciprocal comparison step.

Usage

```
make_ref_similarity_names(de_table.test, de_table.ref, pval = 0.01,
                        num_steps = 5)
```

```
make_ref_similarity_names_using_marked(de_table.ref.marked,
                                       de_table.recip.marked = NA, the_test_dataset = NA,
                                       the_ref_dataset = NA, pval = 0.01, num_steps = 5)
```

Arguments

de_table.test	A differential expression table of the query experiment, as generated from contrast_each_group_to
de_table.ref	A differential expression table of the reference dataset, as generated from contrast_each_group_to
pval	Differences between the rescaled ranking distribution of 'top' genes on different reference groups are tested with a Mann-Whitney U test. If they are <i>significantly different</i> , only the top group(s) are reported. It isn't a simple cutoff threshold as it can change the number of similar groups reported. ie. A more stringent pval is more likely to decide that groups are similar - which would result in multiple group reporting, or no similarity at all. Unlikely that this parameter will ever need to change. Default = 0.01.

num_steps	After ranking reference groups according to median 'top' gene ranking, how many adjacent pairs to test for differences. Set to 1 to only compare each group to the next, or NA to perform an all-vs-all comparison. Setting too low may mean it is possible to miss groups with some similarity to the reported matches (<i>similar_non_match</i> column)). Too high (or NA) with a large number of reference groups could be slow. Default = 5.
de_table.ref.marked	The output of <code>get_the_up_genes_for_all_possible_groups</code> for the contrast of interest.
de_table.recip.marked	Optional. The (reciprocal) output of <code>get_the_up_genes_for_all_possible_groups</code> with the test and reference datasets swapped. If omitted a reciprocal test will not be done. Default = NA.
the_test_dataset	Optional. A short meaningful name for the experiment. (Should match <i>test_dataset</i> column in de_table.marked). Only needed in a table of more than one dataset. Default = NA.
the_ref_dataset	Optional. A short meaningful name for the experiment. (Should match <i>dataset</i> column in de_table.marked). Only needed in a table of more than one dataset. Default = NA.

Details

This function aims to report a) the top most similar reference group, if there's a clear frontrunner, b) A list of multiple similar groups if they have similar similarity, or c) 'No similarity', if there is none.

Each group is named according to the following rules. Testing for significant (smaller) differences with a one-directional Mann-Whitney U test on their rescaled ranks:

1. The first (as ranked by median rescaled rank) reference group is significantly more similar than the next: Report *first only*.
2. When comparing differences between groups stepwise ranked by median rescaled rank - no group is significantly different to its neighbour: Report *no similarity*
3. There's no significant differences in the stepwise comparisons of the first N reference groups - but there is a significant difference later on : Report *multiple group similarity*

There are some further heuristic caveats:

1. The distribution of top genes in the last (or only) match group is tested versus a theoretical random distribution around 0.5 (as reported in *pval_vs_random* column). If the distribution is not significantly above random (It is possible in edge cases where there is a skewed dataset and no/few matches), *no similarity* is reported. The significant *pval* column is left intact.
2. The comparison is repeated reciprocally - reference groups vs the query groups. This helps sensitivity of heterogeneous query groups - and investigating the reciprocal matches can be informative in these cases. If a query group doesn't 'match' a reference group, but the reference group does match that query group - it is reported in the group label in brackets. e.g. *c1:th_lymphocytes(tc_lymphocytes)*. It's even possible if there was no match (and *pval* = NA) e.g. *emphc2:(tc_lymphocytes)*

The similarity is formatted into a group label. Where there are multiple similar groups, they're listed from most to least similar by their median ranks.

For instance, a query dataset of clusters c1, c2, c3 and c4 againsts a cell-type labelled reference dataset might get names like: E.g.

- c1:macrophage
- c2:endothelialmesodermal
- c3:no_similarity
- c4:mesodermal(endothelial)

Function `make_ref_similarity_names` is a convenience wrapper function for `make_ref_similarity_names_from_ma`. It accepts two 'de_table' outputs of function `contrast_each_group_to_the_rest` to compare and handles running `get_the_up_genes_for_all_possible_groups`. Sister function `make_ref_similarity_names_from` may (rarely) be of use if the `de_table.marked` object has already been created, or if reciprocal tests are not wanted.

Value

A table of automagically-generated labels for each query group, given their similarity to reference groups.

The columns in this table:

- **test_group** : Query group e.g. "c1"
- **shortlab** : The cluster label described above e.g. "c1:macrophage"
- **pval** : If there is a similarity flagged, this is the P-value from a Mann-Whitney U test from the last 'matched' group to the adjacent 'non-matched' group. Ie. If only one label in shortlab, this will be the first of the stepped_pvals, if there are 2, it will be the second. If there is 'no_similarity' this will be NA (Because there is no confidence in what is the most appropriate of the all non-significant stepped pvalues.).
- **stepped_pvals** : P-values from Mann-Whitney U tests across adjacent pairs of reference groups ordered from most to least similar (ascending median rank). ie. 1st-2nd most similar first, 2nd-3rd, 3rd-4th e.t.c. The last value will always be NA (no more reference group). e.g. refA:8.44e-10,refB:2.37e-06,refC:0.000818,refD:0.435,refE:0.245,refF:NA
- **pval_to_random** : P-value of test of median rank (of last matched reference group) < random, from binomial test on top gene ranks (being < 0.5).
- **matches** : List of all reference groups that 'match', as described, except it also includes (rare) examples where pval_to_random is not significant. "|" delimited.
- **reciprocal_matches** : List of all reference groups that flagged test group as a match when direction of comparison is reversed. (significant pval and pval_to_random). "|" delimited.
- **similar_non_match**: This column lists any reference groups outside of shortlab that are not significantly different to a reported match group. Limited by `num_steps`, and will never find anything if `num_steps==1`. "|" delimited. Usually NA.
- **similar_non_match_detail** : P-values for any details about similar_non_match results. These p-values will always be non-significant. E.g. "A > C (p=0.0214,n.s)". "|" delimited. Usually NA.
- **differences_within** : This feild lists any pairs of reference groups in shortlab that are significantly different. "|" delimited. Usually NA.

Functions

- `make_ref_similarity_names_using_marked`: Construct some sensible cluster labels, but using a premade marked table.

See Also

[contrast_each_group_to_the_rest](#) For preparing de_table input

[get_the_up_genes_for_all_possible_groups](#) To prepare the **de_table.ref.marked** input.

Examples

```
# Make input
# de_table.demo_query <- contrast_each_group_to_the_rest(demo_query_se, "demo_query")
# de_table.demo_ref <- contrast_each_group_to_the_rest(demo_ref_se, "demo_ref")

make_ref_similarity_names(de_table.demo_query, de_table.demo_ref)
make_ref_similarity_names(de_table.demo_query, de_table.demo_ref, num_steps=3)
make_ref_similarity_names(de_table.demo_query, de_table.demo_ref, num_steps=NA)

# Make input
# de_table.demo_query <- contrast_each_group_to_the_rest(demo_query_se, "demo_query")
# de_table.demo_ref <- contrast_each_group_to_the_rest(demo_ref_se, "demo_ref")

de_table.marked.query_vs_ref <- get_the_up_genes_for_all_possible_groups(
  de_table.demo_query, de_table.demo_ref)
de_table.marked.reiprocal <- get_the_up_genes_for_all_possible_groups(
  de_table.demo_ref, de_table.demo_query)

make_ref_similarity_names_using_marked(de_table.marked.query_vs_ref,
  de_table.marked.reiprocal)

make_ref_similarity_names_using_marked(de_table.marked.query_vs_ref)
```

```
make_ref_similarity_names_for_group
  make_ref_similarity_names_for_group
```

Description

Internal function, called by `make_ref_similarity_names_using_marked` for each group.

Usage

```
make_ref_similarity_names_for_group(the_test_group, mwtest_res_table,
  de_table.ref.marked, reciprocal_matches = NA, the_test_dataset,
  the_ref_dataset, the_pval)
```

Arguments

`the_test_group` Query group to make name for

`mwtest_res_table`

Mann-whitney test results as constructed in [make_ref_similarity_names_using_marked](#)

de_table.ref.marked	The output of get_the_up_genes_for_all_possible_groups for the contrast of interest.
reciprocal_matches	Simplified table of reciprocal matches prepared within make_ref_similarity_names_using_marked . If omitted no reciprocal matching done. Default = NA.
the_test_dataset	A short meaningful name for the experiment. (Should match <i>test_dataset</i> column in de_table.marked)
the_ref_dataset	A short meaningful name for the experiment. (Should match <i>dataset</i> column in de_table.marked)
the_pval	pval as per make_ref_similarity_names_using_marked

Value

A tibble with just one group's labelling information, as per [make_ref_similarity_names_using_marked](#)

See Also

[make_ref_similarity_names_using_marked](#) Only place that uses this function, details there.

run_pair_test_stats *run_pair_test_stats*

Description

Internal function to compare the distribution of a query datasets 'top' genes between two different reference datasete groups with a Mann–Whitney U test. One directional test if groupA median < group B.

Usage

```
run_pair_test_stats(de_table.ref.marked, the_test_group, groupA, groupB,
  enforceAgtB = TRUE)
```

Arguments

de_table.ref.marked	The output of get_the_up_genes_for_all_possible_groups for the contrast of interest.
the_test_group	Name of the test group in query dataset.
groupA	One of the reference group names
groupB	Another of the reference group names
enforceAgtB	Do a one tailed test of A 'less' B (more similar)? Or two-tailed. Default = TRUE.

Details

For use by [make_ref_similarity_names_using_marked](#)

Value

A tibble of wilcox / man-whitneyU test results for this contrast.

See Also

[make_ref_similarity_names_using_marked](#)

```
trim_small_groups_and_low_expression_genes
  trim_small_groups_and_low_expression_genes
```

Description

Filter and return a SummarizedExperiment object (*dataset_se*) by several metrics:

- Cells with at least **min_lib_size** total reads.
- Genes expressed in at least **min_detected_by_min_samples** cells, at a threshold of **min_reads_in_sample** per cell.
- Remove entire groups (clusters) of cells where there are fewer than **min_group_membership** cells in that group.

Usage

```
trim_small_groups_and_low_expression_genes(dataset_se,
  min_lib_size = 1000, min_group_membership = 5,
  min_reads_in_sample = 1, min_detected_by_min_samples = 5)
```

Arguments

<i>dataset_se</i>	Summarised experiment object containing count data. Also requires 'ID' and 'group' to be set within the cell information (see <code>colData()</code>)
<i>min_lib_size</i>	Minimum library size. Cells with fewer than this many reads removed. Default = 1000
<i>min_group_membership</i>	Throw out groups/clusters with fewer than this many cells. May change with experiment size. Default = 5
<i>min_reads_in_sample</i>	Require this many reads to consider a gene detected in a sample. Default = 1
<i>min_detected_by_min_samples</i>	Keep genes detected in this many samples. May change with experiment size. Default = 5

Details

If it hasn't been done already, it is highly recommended to use this function to filter out genes with no/low total counts (especially in single cell data, there can be many) - without expression they are not useful and may reduce statistical power.

Likewise, very small groups (<5 cells) are unlikely to give useful results with this method. And cells with abnormally small library sizes may not be desirable.

Of course 'reasonable' thresholds for filtering cells/genes are subjective. Defaults are moderately sensible starting points.

Index

*Topic **datasets**

- de_table.demo_query, 11
- de_table.demo_ref, 12
- demo_cell_info_table, 8
- demo_counts_matrix, 8
- demo_gene_info_table, 9
- demo_microarray_expr, 9
- demo_microarray_sample_sheet, 10
- demo_query_se, 10
- demo_ref_se, 11

- load_dataset_10Xdata, 5, 20, 23
- load_se_from_files, 7
- load_se_from_files
(load_se_from_tables), 22
- load_se_from_tables, 5, 21, 22

- make_ranking_violin_plot, 24
- make_ref_similarity_names, 25, 25
- make_ref_similarity_names_for_group,
13, 14, 18, 20, 28
- make_ref_similarity_names_using_marked,
15, 17, 28–30
- make_ref_similarity_names_using_marked
(make_ref_similarity_names), 25

- run_pair_test_stats, 29

- trim_small_groups_and_low_expression_genes,
30

- contrast_each_group_to_the_rest, 2, 4–6,
11, 12, 18, 19, 24, 25, 28
- contrast_each_group_to_the_rest_for_norm_ma_with_limma,
3, 6, 21, 23
- contrast_the_group_to_the_rest, 5
- contrast_the_group_to_the_rest_with_limma_for_microarray,
6, 14
- convert_se_gene_ids, 7, 21

- de_table.demo_query, 11
- de_table.demo_ref, 12
- demo_cell_info_table, 8
- demo_counts_matrix, 8
- demo_gene_info_table, 9
- demo_microarray_expr, 9
- demo_microarray_sample_sheet, 10
- demo_query_se, 10
- demo_ref_se, 11

- find_within_match_differences, 12

- get_inner_or_outer_ci, 13
- get_limma_top_table_with_ci, 13
- get_matched_stepped_mwtest_res_table,
14
- get_ranking_and_test_results, 14, 15, 17
- get_rankstat_table, 16
- get_reciprocal_matches, 17
- get_stepped_pvals_str, 17
- get_the_up_genes_for_all_possible_groups,
16, 18, 19, 24, 26–29
- get_the_up_genes_for_group, 18, 19
- get_vs_random_pval, 20