

# Package ‘CIARA’

February 22, 2022

**Type** Package

**Title** Cluster Independent Algorithm for Rare Cell Types Identification

**Version** 0.1.0

**Author** Gabriele Lubatti

**Maintainer** Gabriele Lubatti

<gabriele.lubatti@helmholtz-muenchen.de>

**Description** Identification of markers of rare cell types by looking at genes whose expression is confined in small regions of the expression space <<https://github.com/ScialdoneLab>>.

**License** Artistic-2.0

**Depends** R (>= 4.0)

**Imports** Biobase, ggplot2, ggraph, magrittr

**Suggests** circlize, clustree, ComplexHeatmap, plotly, Seurat (>= 4.0), testthat, knitr, rmarkdown

**biocViews** software

**Config/testthat/edition** 3

**Encoding** UTF-8

**RoxygenNote** 7.1.1

**VignetteBuilder** knitr

**NeedsCompilation** no

**Repository** CRAN

**Date/Publication** 2022-02-22 20:00:02 UTC

## R topics documented:

CIARA . . . . .	2
CIARA_gene . . . . .	3
cluster_analysis_integrate_rare . . . . .	4
cluster_analysis_sub . . . . .	5
find_resolution . . . . .	6
get_background_full . . . . .	7

markers_cluster_seurat . . . . .	7
merge_cluster . . . . .	8
plot_balloon_marker . . . . .	9
plot_gene . . . . .	10
plot_genes_sum . . . . .	10
plot_heatmap_marker . . . . .	11
plot_interactive . . . . .	12
plot_umap . . . . .	13
selection_localized_genes . . . . .	14
test_hvg . . . . .	15
white_black_markers . . . . .	16

<b>Index</b>	<b>17</b>
--------------	-----------

---

CIARA

*CIARA*

---

## Description

It selects highly localized genes as specified in *CIARA\_gene*, starting from genes in *background*

## Usage

```
CIARA(
  norm_matrix,
  knn_matrix,
  background,
  cores_number = 1,
  p_value = 0.001,
  odds_ratio = 2,
  local_region = 1,
  approximation = FALSE
)
```

## Arguments

<code>norm_matrix</code>	Norm count matrix (n_genes X n_cells).
<code>knn_matrix</code>	K-nearest neighbors matrix (n_cells X n_cells).
<code>background</code>	Vector of genes for which the function <i>CIARA_gene</i> is run.
<code>cores_number</code>	Integer.Number of cores to use.
<code>p_value</code>	p value returned by the function <i>fisher.test</i> with parameter alternative = "g"
<code>odds_ratio</code>	odds_ratio returned by the function <i>fisher.test</i> with parameter alternative = "g"
<code>local_region</code>	Integer. Minimum number of local regions (cell with its knn neighbours) where the binarized gene expression is enriched in 1.
<code>approximation</code>	Logical.For a given gene, the fisher test is run in the local regions of only the cells where the binarized gene expression is 1.

**Value**

Dataframe with `n_rows` equal to the length of `background` . Each row is the output from `CIARA_gene`.

**Author(s)**

Gabriele Lubatti <gabriele.lubatti@helmholtz-muenchen.de>

---

 CIARA\_gene

*CIARA\_gene*


---

**Description**

The gene expression is binarized (1/0) if the value in a given cell is above/below the median. Each of cell with its first K nearest neighbors defined a local region. If there are at least `local_region` enriched in 1 according to `fisher.test`, then the gene is defined as highly localized and a final p value is assigned to it. The final p value is the minimum of the p values from all the enriched local regions. If there are no enriched local regions, then the p value by default is set to 1

**Usage**

```
CIARA_gene(
  norm_matrix,
  knn_matrix,
  gene_expression,
  p_value = 0.001,
  odds_ratio = 2,
  local_region = 1,
  approximation = FALSE
)
```

**Arguments**

<code>norm_matrix</code>	Norm count matrix ( <code>n_genes X n_cells</code> ).
<code>knn_matrix</code>	K-nearest neighbors matrix ( <code>n_cells X n_cells</code> ).
<code>gene_expression</code>	numeric vector with the gene expression (length equal to <code>n_cells</code> ). The gene expression is binarized (equal to 0/1 in the cells where the value is below/above the median)
<code>p_value</code>	p value returned by the function <code>fisher.test</code> with parameter <code>alternative = "g"</code>
<code>odds_ratio</code>	<code>odds_ratio</code> returned by the function <code>fisher.test</code> with parameter <code>alternative = "g"</code>
<code>local_region</code>	Integer. Minimum number of local regions (cell with its knn neighbours) where the binarized gene expression is enriched in 1.
<code>approximation</code>	Logical. For a given gene, the fisher test is run in the local regions of only the cells where the binarized gene expression is 1.

**Value**

List with one element corresponding to the p value of the gene.

**Author(s)**

Gabriele Lubatti <gabriele.lubatti@helmholtz-muenchen.de>

**See Also**

<https://www.rdocumentation.org/packages/stats/versions/3.6.2/topics/fisher.test>

---

```
cluster_analysis_integrate_rare
      cluster_analysis_integrate_rare
```

---

**Description**

cluster\_analysis\_integrate\_rare

**Usage**

```
cluster_analysis_integrate_rare(
  raw_counts,
  project_name,
  resolution,
  neighbors,
  max_dimension,
  feature_genes = NULL
)
```

**Arguments**

raw_counts	Raw count matrix (n_genes X n_cells).
project_name	Character name of the Seurat project.
resolution	Numeric value specifying the parameter <i>resolution</i> used in the Seurat function <i>FindClusters</i> .
neighbors	Numeric value specifying the parameter <i>k.param</i> in the Seurat function <i>FindNeighbors</i>
max_dimension	Numeric value specifying the maximum number of the PCA dimensions used in the parameter <i>dims</i> for the Seurat function <i>FindNeighbors</i>
feature_genes	vector of features specifying the argument <i>features</i> in the Seurat function <i>RunPCA</i> .

**Value**

Seurat object including raw and normalized counts matrices, UMAP coordinates and cluster result.

**Author(s)**

Gabriele Lubatti <gabriele.lubatti@helmholtz-muenchen.de>

**See Also**

<https://www.rdocumentation.org/packages/Seurat/versions/4.0.1/topics/FindClusters>  
<https://www.rdocumentation.org/packages/Seurat/versions/4.0.1/topics/FindNeighbors>  
<https://www.rdocumentation.org/packages/Seurat/versions/4.0.1/topics/RunPCA>

---

cluster\_analysis\_sub    *cluster\_analysis\_sub*

---

**Description**

cluster\_analysis\_sub

**Usage**

```
cluster_analysis_sub(  
  raw_counts,  
  resolution,  
  neighbors,  
  max_dimension,  
  name_cluster  
)
```

**Arguments**

raw_counts	Raw count matrix (n_genes X n_cells).
resolution	Numeric value specifying the parameter <i>resolution</i> used in the Seurat function <i>FindClusters</i> .
neighbors	Numeric value specifying the parameter <i>k.param</i> in the Seurat function <i>FindNeighbors</i> .
max_dimension	Numeric value specifying the maximum number of the PCA dimensions used in the parameter <i>dims</i> for the Seurat function <i>FindNeighbors</i> .
name_cluster	Character.Name of the original cluster for which the sub clustering is done.

**Value**

Seurat object including raw and normalized counts matrices and cluster result.

**Author(s)**

Gabriele Lubatti <gabriele.lubatti@helmholtz-muenchen.de>

**See Also**

<https://www.rdocumentation.org/packages/Seurat/versions/4.0.1/topics/RunPCA> <https://www.rdocumentation.org/packages/Seurat/versions/4.0.1/topics/FindVariableFeatures>

---

find\_resolution      *find\_resolution*

---

**Description**

find\_resolution

**Usage**

```
find_resolution(seurat_object, resolution_vector)
```

**Arguments**

`seurat_object`    Seurat object as returned by *cluster\_analysis\_integrate\_rare*  
`resolution_vector`  
                    vector with all values of resolution for which the Seurat function *FindClusters*  
                    is run

**Value**

Clustree object showing the connection between clusters obtained at different level of resolution as specified in *resolution\_vector*.

**Author(s)**

Gabriele Lubatti <gabriele.lubatti@helmholtz-muenchen.de>

**See Also**

<https://CRAN.R-project.org/package=clustree>

---

get\_background\_full    *get\_background\_full*

---

### Description

get\_background\_full

### Usage

```
get_background_full(  
  norm_matrix,  
  threshold = 1,  
  n_cells_low = 3,  
  n_cells_high = 20  
)
```

### Arguments

norm_matrix	Norm count matrix (n_genes X n_cells).
threshold	threshold in expression for a given gene
n_cells_low	minimum number of cells where a gene is expressed at a level above threshold
n_cells_high	maximum number of cells where a gene is expressed at a level above threshold

### Value

Character vector with all genes expressed at a level higher than *threshold* in a number of cells between *n\_cells* and *n\_cells\_high*.

### Author(s)

Gabriele Lubatti <gabriele.lubatti@helmholtz-muenchen.de>

---

markers\_cluster\_seurat  
*markers\_cluster\_seurat*

---

### Description

The Seurat function *FindMarkers* is used to identify general marker for each cluster (specific cluster vs all other cluster). This list of markers is then filtered keeping only the genes that appear as markers in a unique cluster.

### Usage

```
markers_cluster_seurat(seurat_object, cluster, cell_names, number_top)
```

**Arguments**

seurat_object	Seurat object as returned by <i>cluster_analysis_sub</i> or by <i>cluster_analysis_integrate_rare</i> .
cluster	Vector of length equal to the number of cells, with cluster assignment.
cell_names	Vector of length equal to the number of cells, with cell names.
number_top	Integer. Number of top marker genes to keep for each cluster.

**Value**

List of three elements. The first is a vector with *number\_top* marker genes for each cluster. The second is a vector with *number\_top* marker genes and corresponding cluster. The third element is a vector with all marker genes for each cluster.

**Author(s)**

Gabriele Lubatti <gabriele.lubatti@helmholtz-muenchen.de>

**See Also**

<https://www.rdocumentation.org/packages/Seurat/versions/4.0.1/topics/FindMarkers>

---

merge_cluster	<i>merge_cluster</i>
---------------	----------------------

---

**Description**

merge\_cluster

**Usage**

```
merge_cluster(old_cluster, new_cluster, max_number = NULL)
```

**Arguments**

old_cluster	original cluster assignment that need to be updated
new_cluster	new cluster assignment that need to be integrated with <i>old_cluster</i> .
max_number	Threshold in size for clusters in <i>new_cluster</i> . Only cluster with number of cells smaller than <i>max_number</i> will be integrated in <i>old cluster</i> . If <i>max_number</i> is NULL, then all the clusters in <i>new_cluster</i> are integrated in <i>old cluster</i> .

**Value**

Numeric vector of length equal to *old\_cluster* showing the merged cluster assignment between *old cluster* and *new\_cluster*.

**Author(s)**

Gabriele Lubatti <gabriele.lubatti@helmholtz-muenchen.de>

---

plot\_balloon\_marker    *plot\_balloon\_marker*

---

## Description

plot\_balloon\_marker

## Usage

```
plot_balloon_marker(  
  norm_counts,  
  cluster,  
  marker_complete,  
  max_number,  
  max_size = 5,  
  text_size = 7  
)
```

## Arguments

norm_counts	Norm count matrix (genes X cells).
cluster	Vector of length equal to the number of cells, with cluster assignment.
marker_complete	Third element of the output list as returned by the function <i>markers_cluster_seurat</i>
max_number	Integer. Maximum number of markers for each cluster for which we want to plot the expression.
max_size	Integer. Size of the dots to be plotted.
text_size	Size of the text in the heatmap plot.

## Value

ggplot2 object showing balloon plot.

## Author(s)

Gabriele Lubatti <gabriele.lubatti@helmholtz-muenchen.de>

---

 plot\_gene

*plot\_gene*


---

**Description**

Cells are coloured according to the expression of *gene\_id* and plotted according to *coordinate\_umap*.

**Usage**

```
plot_gene(norm_counts, coordinate_umap, gene_id, title_name)
```

**Arguments**

norm_counts	Norm count matrix (genes X cells).
coordinate_umap	Data frame with dimensionality reduction coordinates. Number of rows must be equal to the number of cells
gene_id	Character name of the gene.
title_name	Character name.

**Value**

ggplot2 object.

**Author(s)**

Gabriele Lubatti <gabriele.lubatti@helmholtz-muenchen.de>

**See Also**

<https://CRAN.R-project.org/package=ggplot2>

---

 plot\_genes\_sum

*plot\_genes\_sum*


---

**Description**

The sum of each gene in *genes\_relevant* across all cells is first normalized to 1. Then for each cell, the sum from the (normalized) genes expression is computed and shown in the output plot.

**Usage**

```
plot_genes_sum(coordinate_umap, norm_counts, genes_relevant, name_title)
```

**Arguments**

coordinate_umap	Data frame with dimensionality reduction coordinates. Number of rows must be equal to the number of cells
norm_counts	Norm count matrix (genes X cells).
genes_relevant	Vector with gene names for which we want to visualize the sum in each cell.
name_title	Character value.

**Value**

ggplot2 object.

**Author(s)**

Gabriele Lubatti <gabriele.lubatti@helmholtz-muenchen.de>

**See Also**

<https://CRAN.R-project.org/package=ggplot2>

---

plot\_heatmap\_marker     *plot\_heatmap\_marker*

---

**Description**

plot\_heatmap\_marker

**Usage**

```
plot_heatmap_marker(
  marker_top,
  marker_all_cluster,
  cluster,
  condition,
  norm_counts,
  text_size
)
```

**Arguments**

marker_top	First element returned by <i>markers_cluster_seurat</i>
marker_all_cluster	Second element returned by <i>markers_cluster_seurat</i>
cluster	Vector of length equal to the number of cells, with cluster assignment.
condition	Vector or length equal to the number of cells, specifying the condition of the cells (i.e. batch, dataset of origin..)
norm_counts	Norm count matrix (genes X cells).
text_size	Size of the text in the heatmap plot.

**Value**

Heatmap class object.

**Author(s)**

Gabriele Lubatti <gabriele.lubatti@helmholtz-muenchen.de>

**See Also**

<https://www.rdocumentation.org/packages/ComplexHeatmap/versions/1.10.2/topics/Heatmap>

---

plot_interactive	<i>plot_interactive</i>
------------------	-------------------------

---

**Description**

It shows in an interactive plot which are the highly localized genes in each cell. It is based on plotly library

**Usage**

```
plot_interactive(
  coordinate_umap,
  color,
  text,
  min_x = NULL,
  max_x = NULL,
  min_y = NULL,
  max_y = NULL
)
```

**Arguments**

coordinate_umap	Data frame with dimensionality reduction coordinates. Number of rows must be equal to the number of cells
color	vector of length equal to n_rows in coordinate_umap. Each cell will be coloured following a gradient according to the corresponding value of this vector.
text	Character vector specifying the highly localized genes in each cell. It is the output from <i>selection_localized_genes</i> .
min_x	Set the min limit on the x axis.
max_x	Set the max limit on the x axis.
min_y	Set the min limit on the y axis.
max_y	Set the min limit on the y axis.

**Value**

plotly object given by *plot\_ly function* (from library *plotly*).

**Author(s)**

Gabriele Lubatti <gabriele.lubatti@helmholtz-muenchen.de>

**See Also**

<https://plotly.com/r/>

---

plot\_umap

*plot\_umap*

---

**Description**

plot\_umap

**Usage**

```
plot_umap(coordinate_umap, cluster)
```

**Arguments**

coordinate\_umap

Data frame with dimensionality reduction coordinates. Number of rows must be equal to the number of cells

cluster

Vector of length equal to the number of cells, with cluster assignment.

**Value**

ggplot2 object.

**Author(s)**

Gabriele Lubatti <gabriele.lubatti@helmholtz-muenchen.de>

**See Also**

<https://CRAN.R-project.org/package=ggplot2>

---

```
selection_localized_genes  
    selection_localized_genes
```

---

**Description**

selection\_localized\_genes

**Usage**

```
selection_localized_genes(  
  norm_counts,  
  localized_genes,  
  min_number_cells = 4,  
  max_number_genes = 10  
)
```

**Arguments**

**norm\_counts** Norm count matrix (genes X cells).

**localized\_genes** vector of highly localized genes as provided by the last element of the list given as output from *CIARA\_mixing\_final*.

**min\_number\_cells** Minimum number of cells where a genes must be expressed (> 0).

**max\_number\_genes** Maximum number of genes to show for each cell in the interactive plot from *plot\_interactive*.

**Value**

Character vector where each entry contains the name of the top *max\_number\_genes* for the corresponding cell.

**Author(s)**

Gabriele Lubatti <gabriele.lubatti@helmholtz-muenchen.de>

---

test_hvg	<i>test_hvg</i>
----------	-----------------

---

### Description

For each cluster in *cluster*, HVGs are defined with Seurat function *FindVariableFeatures*. A Fisher test is performed to see if there is a statistically significant enrichment between the top *number\_hvg* and the *localized\_genes*

### Usage

```
test_hvg(
  raw_counts,
  cluster,
  localized_genes,
  background,
  number_hvg,
  min_p_value
)
```

### Arguments

<code>raw_counts</code>	Raw count matrix (n_genes X n_cells).
<code>cluster</code>	Vector of length equal to the number of cells, with cluster assignment.
<code>localized_genes</code>	Character vector with localized genes detected by CIARA.
<code>background</code>	Character vector with all the genes names to use as background for the Fisher test.
<code>number_hvg</code>	Integer value. Number of top HVGs provided by the Seurat function <i>FindVariableFeatures</i> .
<code>min_p_value</code>	Threshold on p values provided by Fisher test.

### Value

A list with two elements.

<code>first element</code>	The first one is a list with length equal to the number of clusters. Each entry is list of three elements. The first two elements contain the p value and the odds ration given by the Fisher test The third is a vector with genes names that are present both in <i>localized_genes</i> and in top <i>number_hvg</i> HVGs .
<code>second element</code>	a character vector with the name of the cluster that have a p value smaller than <i>min_p_value</i> .

### Author(s)

Gabriele Lubatti <gabriele.lubatti@helmholtz-muenchen.de>

**See Also**

<https://www.rdocumentation.org/packages/stats/versions/3.6.2/topics/fisher.test>

---

white\_black\_markers    *white\_black\_markers*

---

**Description**

A white-marker is a gene whose median expression across cells belong to *single\_cluster* is greater than *threshold* and in all the other clusters is equal to zero.

**Usage**

```
white_black_markers(  
  cluster,  
  single_cluster,  
  norm_counts,  
  marker_list,  
  threshold = 0  
)
```

**Arguments**

cluster	Vector of length equal to the number of cells, with cluster assignment.
single_cluster	Character. Label of one specify cluster
norm_counts	Norm count matrix (genes X cells).
marker_list	Third element of the output list as returned by the function <i>markers_cluster_seurat</i>
threshold	Numeric. The median of the genes across cells belong to <i>single_cluster</i> has to be greater than <i>threshold</i> in order to be consider as a white-black marker for <i>single_cluster</i>

**Value**

Logical vector of length equal to *marker\_list*, with TRUE/FALSE if the gene is/is not a white-black marker for *single\_cluster*.

**Author(s)**

Gabriele Lubatti <gabriele.lubatti@helmholtz-muenchen.de>

# Index

CIARA, [2](#)  
CIARA\_gene, [3](#)  
cluster\_analysis\_integrate\_rare, [4](#)  
cluster\_analysis\_sub, [5](#)  
  
find\_resolution, [6](#)  
  
get\_background\_full, [7](#)  
  
markers\_cluster\_seurat, [7](#)  
merge\_cluster, [8](#)  
  
plot\_balloon\_marker, [9](#)  
plot\_gene, [10](#)  
plot\_genes\_sum, [10](#)  
plot\_heatmap\_marker, [11](#)  
plot\_interactive, [12](#)  
plot\_umap, [13](#)  
  
selection\_localized\_genes, [14](#)  
  
test\_hvg, [15](#)  
  
white\_black\_markers, [16](#)