# Package 'ClustIRR'

October 24, 2025

Type Package

Title Clustering of immune receptor repertoires

**Version** 1.7.15

Description ClustIRR analyzes repertoires of B- and T-cell receptors. It starts by identifying communities of immune receptors with similar specificities, based on the sequences of their complementarity-determining regions (CDRs). Next, it employs a Bayesian probabilistic models to quantify differential community occupancy (DCO) between repertoires, allowing the identification of expanding or contracting communities in response to e.g. infection or cancer treatment.

License GPL-3 + file LICENSE

**Depends** R (>= 4.3.0)

**Imports** blaster, future, future.apply, grDevices, igraph, methods, Rcpp (>= 0.12.0), RcppParallel (>= 5.0.1), reshape2, rstan (>= 2.18.1), rstantools (>= 2.4.0), stats, stringdist, utils, posterior, visNetwork, dplyr, tidyr, ggplot2, ggforce, scales

Suggests BiocStyle, knitr, testthat, ggplot2, ggrepel, patchwork

SystemRequirements GNU make

**Encoding** UTF-8

LazyData false

NeedsCompilation yes

**biocViews** Clustering, ImmunoOncology, SingleCell, Software, Classification, Bayesian, BiomedicalInformatics, ImmunoOncology, MathematicalBiology

RoxygenNote 7.2.3 VignetteBuilder knitr

URL https://github.com/snaketron/ClustIRR

BugReports https://github.com/snaketron/ClustIRR/issues

Biarch true

2 BLOSUM62

# **Contents**

	BLOSUM62	2
	clustirr	3
	clust_irr-class	6
	Datasets	7
	dco	8
	decode_all_communities	10
	decode_community	12
	detect_communities	14
	get_ag_gene_hits	16
	get_ag_species_hits	17
	get_beta_violin	19
	get_community_purity	21
	get_honeycombs	22
	mcpas	23
	plot_graph	24
	tcr3d	25
	vdjdb	26
Index		28

BLOSUM62 BLOSUM62 matrix

# Description

Predefined scoring matrix for amino acid or nucleoitide alignments.

clustirr 3

### Usage

```
data("BLOSUM62")
```

### **Format**

BLOSUM62 is a square symmetric matrix. Rows and columns are identical single letters, representing nucleotide or amino acid. Elements are integer coefficients (substitution scores).

### **Details**

BLOSUM62 was obtained from NCBI (the same matrix used by the stand- alone BLAST software).

#### **Source**

See https://ftp.ncbi.nih.gov/blast/matrices/BLOSUM62

### References

See https://ftp.ncbi.nih.gov/blast/matrices/BLOSUM62

# **Examples**

```
data(BLOSUM62, package = "ClustIRR")
BLOSUM62
```

clustirr

Clustering of immune receptor repertoires (IRRs)

# Description

clustirr computes similarities between immune receptors (IRs = T-cell and B-cell receptors) based on their CDR3 sequences.

# Usage

4 clustirr

#### **Arguments**

s

a data.frame consisting of one or more IRRs. Each row is a clone of a given IRR with the following columns (clone features):

- sample: names of the repertoires (e.g. 'A', 'B', etc.)
- clone\_size: cell count in the clone (=clonal expansion)
- CDR3?: amino acid CDR3 sequence. Replace '?' with the appropriate name of the immmune receptor chain (e.g. CDR3a for CDR3s from TCR $\alpha$  chain; or CDR3d for CDR3s from TCR $\delta$  chain. Meanwhile, if paired CDR3s from both chains are available, then you can provide both in separate columns e.g.:
  - *CDR3b* and *CDR3a* [for  $\alpha\beta$  TCRs]
  - *CDR3g* and *CDR3d* [for  $\gamma \delta$  TCRs]
  - CDR3h and CDR3l [for heavy/light chain BCRs]

meta

data.frame with meta-data for each clone, which may contain clone-specific data, such as, V/J genes, cell-type (e.g. CD8+, CD4+), nut also repertoire-specific data, such as, biological condition, HLA type, age, etc. This data will be used to annotate the graph nodes and help downstream analyses.

cores

number of computer cores to use (default = 1)

control

auxiliary parameters to control the algorithm's behavior. See the details below:

- gmi: the minimum sequence identity between a pair of CDR3 sequences for them to even be considered for alignment and scoring (default = 0.8; 80 percent identity).
- trim\_flank\_aa: how many amino acids should be trimmed from the flanks of all CDR3 sequences to isolate the **CDR3 cores**. trim\_flank\_aa = 3 (default).
- db\_custom: additional database (data.frame) which allows us to annotate CDR3 sequences from the input (s) with their cognate antigens. The structure of db\_custom must be identical to that in data(vdjdb, package = "ClustIRR"). ClustIRR will use the internal databases if db\_custom=NULL (default). Three databases (data only from human CDR3) are integrated in ClustIRR: VDJdb, TCR3d and McPAS-TCR.
- db\_dist: we compute edit distances between CDR3 sequences from s and from a database (e.g. VDJdb). If a particular distance is smaller than or equal to db\_dist (default = 0), then we annotate the CDR3 from s with the specificity of the database CDR3 sequence.

### **Details**

ClustIRR performs the following steps.

- 1. Compute similarities between clones within each repertoire
- 2. Construct a graph from each TCR repertoire
- 3. Construct a joint similarity graph (J)
- 4. Detect communities in J
- 5. Analyze Differential Community Occupancy (DCO)

clustirr 5

- Between individual TCR repertoires with model M
- Between groups of TCR repertoires from biological conditions with model  $M_h$

6. Inspect results

the function clustirr performs the steps 1. to 3.

#### Value

The output is a list with the following elements.

- graph: the resulting igraph object
- clust\_irrs: list of clust\_irr objects for each repertoire (sample)
   Each element is an S4 object of class clust\_irr. This object contains two sublists:
  - clust, list, contains clustering results for each receptor chain. The results are stored as data.frame in separate sub-list named appropriately (e.g. CDR3a, CDR3b, CDR3g, etc.).
     Each row in the data.frames contains a pair of CDR3s.
    - The remaining columns contain similarity scores for the complete CDR3 sequences (column weight) or their cores (column cweight). The columns max\_len and max\_clen store the length of the longer CDR3 sequence and core in the pair, and these used to normalize the scores weight and cweight: the normalized scores are shown in the columns nweight and ncweight
  - inputs, list, contains all user provided inputs (see Arguments)
- multigraph: logical variable multigraph, which is set to TRUE if the graph is a joint graph made up of two or more repertoires (samples) and FALSE if the graph contains only one repertoire

```
# load package input data
data("CDR3ab", package = "ClustIRR")
s <- data.frame(CDR3b = CDR3ab[1:100, "CDR3b"], sample = "A", clone_size = 1)
# run analysis
c <- clustirr(s = s)
# output class
class(c)
# output structure
str(c)</pre>
```

6 clust\_irr-class

clust\_irr-class

clust\_irr class

### Description

Objects of the class clust\_irr are generated by the function cluster\_irr. These objects are used to store the clustering results in a structured way, such that they may be used as inputs of other ClustIRR functions (e.g. get\_graph, plot\_graph, etc.).

The output is an S4 object of class clust\_irr. This object contains two sublists:

• clust, list, contains clustering results for each IR chain. The results are stored as data.frame in separate sub-list named appropriately (e.g. CDR3a, CDR3b, CDR3g, etc.). Each row in the data.frames contains a pair of CDR3s.

The remaining columns contain similarity scores for the complete CDR3 sequences (column weight) or their cores (column cweight). The columns max\_len and max\_clen store the length of the longer CDR3 and CDR3 core sequence in the pair, and these used to normalize the scores weight and cweight: the normalized scores are shown in the columns nweight and ncweight

• inputs, list, contains all user provided inputs (see Arguments)

### **Arguments**

clust list, contains clustering results for each TCR/BCR chain. The results are stored

in separate sub-list named appropriately (e.g. CDR3a, CDR3b, CDR3g, etc.)

inputs list, contains all user provided inputs

### Value

The output is an S4 object of class clust\_irr

#### Accessors

To access the slots of clust\_irr object we have two accessor functions. In the description below, x is a clust\_irr object.

```
get_clustirr_clust get_clustirr_clust(x): Extract the clustering results (slot clust)
get_clustirr_inputs get_clustirr_inputs(x): Extract the processed inputs (slot inputs)
```

```
# load package input data
data("CDR3ab", package = "ClustIRR")
s <- data.frame(CDR3b = CDR3ab[1:100, "CDR3b"], sample = "A", clone_size = 1)
# run analysis
c <- clustirr(s = s)</pre>
```

Datasets 7

Datasets

Datasets CDR3ab and D1 with  $TCR\alpha\beta$  mock repertoires

### Description

 $TCR\alpha\beta$  repertoire with 10,000 T-cells (rows). Each T-cell has the following features: amino acid sequences of their complementarity determining region 3 (CDR3); and variable (V) and joining (J) gene names for TCR chains  $\alpha$  and  $\beta$ .

Important remark: this is a mock dataset, all CDR3 sequences and the genes were sampled from a larger set of CDR3 $\beta$  sequences and genes of naive CD8+ T cells in humans.

We used this data to create dataset D1: three  $TCR\alpha\beta$  repertoires a, b, and c, each with 500 TCR clones. We simulated clonal expansion with increasing degree in TCR repertoires b and c. The TCR repertoires as stores as element of a list. For each TCR repertoires we have a metadata: ma, mb, and mc.

### Usage

```
# For the raw data with 10,000 TCR clones
data(CDR3ab)
# For dataset D1
data(D1)
```

#### **Format**

data. frame with rows as TCR clones and 6 columns

- CDR3a: CDR3lpha amino acid sequence
- TRAV: variable (V) gene of TCR $\alpha$
- TRAV: joining (J) gene of TCR $\alpha$

8 dco

- CDR3b: CDR3 $\beta$  amino acid sequence
- TRBV: variable (V) gene of  $TCR\beta$
- TRBV: joining (J) gene of  $TCR\beta$

#### Value

data(CDR3ab) loads the object CDR3ab, which is a data.frame with six columns (3 for TCR $\alpha$  and 3 for TCR $\beta$ ) and rows for each TCR clone (see details).

### Source

GLIPH version 2

# **Examples**

```
data("CDR3ab")
data("D1")
```

dco

Model-based differential community occupancy (DCO) analysis

### **Description**

This algorithm takes as input a community matrix, and quantifies the relative enrichment/depletion of individual communities in each sample using a Bayesian hierarchical model.

### Usage

```
dco(community_occupancy_matrix, mcmc_control, compute_delta=TRUE, groups = NA)
```

# **Arguments**

community\_occupancy\_matrix

matrix, rows are communities, columns are repertoires, matrix entries are numbers of cells in each community and repertoire.

mcmc\_control list, configurations for the Markov Chain Monte Carlo (MCMC) simulation.

- mcmc\_warmup = 750; number of MCMC warmups
- mcmc iter = 1500; number of MCMC iterations
- mcmc\_chains = 4; number of MCMC chains
- mcmc\_cores = 1; number of computer cores
- mcmc\_algorithm = "NUTS"; which MCMC algorithm to use
- adapt\_delta = 0.95; MCMC step size
- max\_treedepth = 12; the max value, in exponents of 2, of what the binary tree size in NUTS should have.

dco 9

compute\_delta should delta be computed by the Stan model? This may be take up extra mem-

or

groups vector with integers  $\geq 1$ , one for each repertoire (column in community\_occupancy\_matrix).

This specifies the biological group of each repertoire (e.g. for cancer repertoire we may specify the index 1, and for normal repertoires the index 2). If this vector is specified, ClustIRR will employ a hierarchical model, modeling the dependence between the repertoires within each group. Else (which is the default setting in ClustIRR), ClustIRR will treat the repertoires as independent

samples by employing a simpler model.

#### Value

The output is a list with the folling elements:

fit model fit (stan object)

posterior\_summary

nested list with data.frames, summary of model parameters, including their means, medians, 95% credible intervals, etc. Predicted observations (y\_hat), which are useful for posterior predictive checks are also provided.

community\_occupancy\_matrix

matrix, rows are communities, columns are repertoires, matrix entries are numbers of cells in each community and repertoire.

mcmc\_control mcmc configuration inputs provided as list.

compute\_delta the input compute\_delta.

groups the input groups.

```
# load package input data
data("CDR3ab", package = "ClustIRR")
a <- data.frame(CDR3a = CDR3ab[1:500, "CDR3a"],</pre>
                   CDR3b = CDR3ab[1:500, "CDR3b"],
                   clone_size = 1,
                   sample = "a")
b <- data.frame(CDR3a = CDR3ab[401:900, "CDR3a"],</pre>
                   CDR3b = CDR3ab[401:900, "CDR3b"],
                   clone_size = 1,
                   sample = "b")
b$clone_size[1] <- 20
# run ClustIRR analysis
c <- clustirr(s = rbind(a, b))</pre>
# detect communities
gcd <- detect_communities(graph = c$graph,</pre>
                           algorithm = "leiden",
                           resolution = 1,
                           weight = "ncweight",
```

### **Description**

Given a graph based on which we have detected communities (with the function detect\_community), this function applies decode\_communities to **all** communities in the graph.

Each community is partitioned according to user-defined filters on edges (edge\_filter) and nodes (node\_filter). This allows extraction of subgraphs, cliques, and connected components across the full set of communities in the graph.

# Usage

```
decode_all_communities(graph, edge_filter, node_filter)
```

### **Arguments**

graph

igraph object that has been analyzed by graph-based community detection methods as implemented in detect\_community. Must have a vertex attribute named community.

edge\_filter

data.frame with edge filters. The data.frame has three columns:

- name: edge attribute name
- value: edge attribute value (threshold)
- operation: logical operation that tells ClustIRR which edge attribute values should pass the filter. Possible operations: "<", ">=", "<=", "==" and "!=".

node\_filter

data.frame with node filters. Groups of nodes that have the same attribute values among **ALL** provided attributes will be treated as a subcomponent.

#### **Details**

Internally, this function iterates over all unique community IDs found in the vertex attribute community, and applies decode\_communities.

- The edge\_filter controls which edges are retained based on their attributes.
- The node\_filter groups nodes with shared attributes into subcomponents.

#### Value

A named list of results, one per community. Each element of the list is the output of decode\_communities, containing:

- community\_graph: "filtered" igraph object for the community
- component\_stats: data.frame with summary about each connected component
- node\_summary: data.frame with summary about each node

```
# load package input data
data("CDR3ab", package = "ClustIRR")
a <- data.frame(CDR3a = CDR3ab[1:500, "CDR3a"],
                  CDR3b = CDR3ab[1:500, "CDR3b"],
                   clone_size = 1,
                   sample = "a")
b <- data.frame(CDR3a = CDR3ab[401:900, "CDR3a"],</pre>
                   CDR3b = CDR3ab[401:900, "CDR3b"],
                   clone_size = 1,
                  sample = "b")
b$clone_size[1] <- 20
# run ClustIRR analysis
c <- clustirr(s = rbind(a, b))</pre>
# detect communities
gcd <- detect_communities(graph = c$graph,</pre>
                           algorithm = "leiden",
                           resolution = 1,
                           weight = "ncweight",
                           iterations = 100,
                           chains = c("CDR3a", "CDR3b"))
# Construct edge and node filters
edge_filter <- rbind(data.frame(name = "nweight", value = 8, operation = ">="),
                      data.frame(name = "ncweight", value = 8, operation = ">="))
node_filter <- data.frame(name = "Ag_gene")</pre>
# Decode all communities at once
all_decoded <- decode_all_communities(graph = gcd$graph,</pre>
                                        edge_filter = edge_filter,
                                        node_filter = node_filter)
```

12 decode\_community

```
# Inspect one decoded community
names(all_decoded)
str(all_decoded[[1]]$component_stats)
```

decode\_community

Decode graph communities

### **Description**

Given a graph based on which we have detected communities (with the function detect\_communities), and a community ID, the function will try to partition the community nodes according to user-defined filters: edge and node filters.

# Usage

```
decode_community(community_id, graph, edge_filter, node_filter)
```

# **Arguments**

graph igraph object that has been analyzed by graph-based community detection meth-

ods as implemented in detect\_communities

community\_id which community should be decoded?

edge\_filter data.frame with edge filters. The deta.frame has three columns:

• name: edge attribute name

• value: edge attribute value (threshold)

• operation: logical operation that tells ClustIRR which edge attribute values should pass the filter. Possible operations: "<", ">=", "<=", "==" and "!=".

node\_filter

a vector with node attributes. Groups of nodes that have the same attribute values among **ALL** provided attributes will be treated as a subcomponent.

#### **Details**

How to decode a community?

For instance, the user may only be interested in retaining edges with core edge weight > 4; or making sure that nodes that have same 'cell\_type' (node meta datafrom) are grouped together. Or the user might want to treat all nodes that have the same V, D and J gene names and HLA types as subgroups, in which case all edges between nodes that do not share the same sets of attributes are dicarded.

Based on these filters, ClustIRR will reformat the edges in the selected community and then find **connected components** in the resulting graph.

decode\_community 13

#### Value

The output is a list with the following elements

- community\_graph: "filtered" igraph object
- component\_stats: data.frame with summary about each connected component
- node\_summary: data.frame with summary about each node

```
# load package input data
data("CDR3ab", package = "ClustIRR")
a <- data.frame(CDR3a = CDR3ab[1:500, "CDR3a"],</pre>
                  CDR3b = CDR3ab[1:500, "CDR3b"],
                  clone_size = 1,
                  sample = "a")
b <- data.frame(CDR3a = CDR3ab[401:900, "CDR3a"],</pre>
                  CDR3b = CDR3ab[401:900, "CDR3b"],
                  clone\_size = 1,
                  sample = "b")
b$clone_size[1] <- 20
# run ClustIRR analysis
c <- clustirr(s = rbind(a, b))</pre>
# detect communities
gcd <- detect_communities(graph = c$graph,</pre>
                          algorithm = "leiden",
                          resolution = 1,
                          weight = "ncweight",
                          iterations = 100,
                          chains = c("CDR3a", "CDR3b"))
# We "decompose" the communities in the gcd object using decode_community
# based on the attributes of the edges (edge_filter) and nodes (node_filter).
# We can pick from these edge attributes and create filters:
library(igraph)
edge_attr_names(graph = gcd$graph)
# For instance, the following edge-filter will instruct ClustIRR to keep
# edges with: edge attributes: nweight>=8 \bold{AND} ncweight>=8
edge_filter <- rbind(data.frame(name = "nweight", value = 8, operation = ">="),
                     data.frame(name = "ncweight", value = 8, operation = ">="))
# In addition, we can construct filters based on the following node attributes:
vertex_attr_names(graph = gcd$graph)
# The following node-filter will instruct ClustIRR to retain edges
# between nodes that have shared node attributed with respect to ALL
# of the following node attributes:
node_filter <- data.frame(name = "Ag_gene")</pre>
```

14 detect\_communities

```
# Lets inspect community with ID = 1.
 c1 <- decode_community(community_id = 1,</pre>
                           graph = gcd$graph,
                           edge_filter = edge_filter,
                           node_filter = node_filter)
 # Plot resulting igraph
 par(mar = c(0, 0, 0, 0))
 plot(c1$community_graph, vertex.size = 10)
 # Now look at node attributes
 as_data_frame(x = c1\$community_graph, what = "vertices")[,c("name", 
                                                                "component_id",
                                                                "CDR3b",
                                                                "CDR3a",
                                                                "Ag_gene")]
 str(c1$component_stats)
 str(c1$node_summary)
                         Graph-based community detection (GCD)
detect\_communities
```

# Description

Graph-based community detection in graphs constructed by get\_graph or get\_joint\_graph.

# Usage

#### **Arguments**

graph	igraph object
algorithm	graph-based community detection (GCD) method: leiden (default), louvain or infomap.
metric	possible metrics: "average" (default) or "max".
resolution	clustering resolution (default = 1) for GCD.
iterations	clustering iterations (default = 100) for GCD.
weight	which edge weight attribute (default = nweight) should be used for GCD
chains	which chains should be used for clustering? For instance: chains = "CDR3a"; or chains = "CDR3b"; or chains = c("CDR3a", "CDR3b").

detect\_communities 15

#### **Details**

ClustIRR employs graph-based community detection (GCD) algorithms, such as Louvain, Leiden or InfoMap, to identify communities of nodes that have high density of edges among each other, and low density of edges with nodes outside the community.

#### Value

The output is a list with the folling elements:

community\_occupancy\_matrix

matrix, rows are communities, columns are repertoires, matrix entries are numbers of cells in each community and repertoire.

community\_summary

data.frame, rows are communities and their properties are provided as columns.

node\_summary data.frame, rows are nodes (clones) and their properties are provided as column-

scontains all user provided.

graph igraph object, processed graph object.

graph\_structure\_quality

graph modularity and quality (only for Leiden) measure of the strength of divi-

sion of the graph into communities.

input\_config list, inputs provided as list.

### **Examples**

```
# load package input data
data("CDR3ab", package = "ClustIRR")
a <- data.frame(CDR3a = CDR3ab[1:500, "CDR3a"],
                  CDR3b = CDR3ab[1:500, "CDR3b"],
                  clone_size = 1,
                  sample = "a")
b <- data.frame(CDR3a = CDR3ab[401:900, "CDR3a"],</pre>
                  CDR3b = CDR3ab[401:900, "CDR3b"],
                  clone_size = 1,
                  sample = "b")
b$clone_size[1] <- 20
# run ClustIRR analysis
c <- clustirr(s = rbind(a, b))</pre>
# detect communities
gcd <- detect_communities(graph = c$graph,</pre>
                           algorithm = "leiden",
                           metric = "average",
                           resolution = 1,
                           weight = "ncweight",
                           iterations = 100,
                           chains = c("CDR3a", "CDR3b"))
```

# look at outputs

get\_ag\_gene\_hits

```
names(gcd)
# look at the community occupancymatrix
head(gcd$community_occupancy_matrix)
# look at the community summary
head(gcd$community_summary$wide)
# look at the node summary
head(gcd$node_summary)
```

get\_ag\_gene\_hits

Annotate antigen gene hits in node summary

### **Description**

Annotates clones (nodes) in the node summary output from the function detect\_communities with specificity to given antigen gene name. Adds new binary columns indicating hits and computes antigen-specific cellular statistics per community and sample.

# Usage

#### **Arguments**

node\_summary Node summary data.frame from detect\_communities

db Annotation database (e.g., "vdjdb", "mcpas", "tcr3d")

ag\_gene Antigen gene name (e.g., "MLANA", "gp100", "Spike")

### **Details**

Searches for antigen gene matches (e.g., "MLANA" or "gp100") in annotation columns from specified databases (e.g., "vdjdb"). For each match, adds a new column to node\_summary (1 = hit, 0 = no hit) and computes:

- Repertoire-level: Total cells/clones per sample
- Community-level: Cells/clones per community
- Antigen-specific: Cells/clones per antigen-community

Requires node\_summary from detect\_communities with database annotation columns).

get\_ag\_species\_hits 17

# Value

List containing:

node\_summary Input with added antigen hit columns. See new\_columns for the names of the

added columns

new\_columns Names of added columns

table\_summary Aggregated data.frame with columns:

• sample, community, ag\_key

• rep\_cells, rep\_clones (repertoire totals)

• com\_cells, com\_clones (community totals)

• ag\_cells, ag\_clones (antigen-specific counts)

```
# load package input data
data("CDR3ab", package = "ClustIRR")
a <- data.frame(CDR3a = CDR3ab[1:500, "CDR3a"],</pre>
                   CDR3b = CDR3ab[1:500, "CDR3b"],
                   clone\_size = 1,
                   sample = "a")
b <- data.frame(CDR3a = CDR3ab[401:900, "CDR3a"],</pre>
                   CDR3b = CDR3ab[401:900, "CDR3b"],
                   clone_size = 1,
                   sample = "b")
b$clone_size[1] <- 20
# run ClustIRR analysis
c <- clustirr(s = rbind(a, b))</pre>
# detect communities
gcd <- detect_communities(graph = c$graph,</pre>
                            algorithm = "leiden",
                            resolution = 1,
                            weight = "ncweight",
                            chains = c("CDR3a", "CDR3b"))
ag <- get_ag_gene_hits(node_summary = gcd$node_summary,</pre>
                        db = "vdjdb",
                        ag_gene = "Spike")
head(ag)
```

18 get\_ag\_species\_hits

# **Description**

Annotates clones (nodes) in the node summary output from the function detect\_communities with specificity to given antigen species. Adds new binary columns indicating hits and computes antigen-specific cellular statistics per community and sample.

# Usage

# Arguments

node\_summary Node summary data.frame from detect\_communities

db Annotation database (e.g., "vdjdb", "mcpas", "tcr3d")

ag\_species Antigen species (e.g., "EBV", "CMV", "SARS-CoV-2")

### **Details**

Searches for antigen species matches (e.g., "EBV" or "CMV") in annotation columns from specified databases (e.g., "vdjdb"). For each match, adds a new column to node\_summary (1 = hit, 0 = no hit) and computes:

• Repertoire-level: Total cells/clones per sample

• Community-level: Cells/clones per community

• Antigen-specific: Cells/clones per antigen-community

Requires node\_summary from detect\_communities with database annotation columns).

### Value

List containing:

node\_summary Input with added antigen hit columns. See new\_columns for the names of the

added columns

new\_columns Names of added columns

table\_summary Aggregated data.frame with columns:

• sample, community, ag\_key

• rep\_cells, rep\_clones (repertoire totals)

• com\_cells, com\_clones (community totals)

• ag\_cells, ag\_clones (antigen-specific counts)

get\_beta\_violin 19

### **Examples**

```
# load package input data
data("CDR3ab", package = "ClustIRR")
a <- data.frame(CDR3a = CDR3ab[1:500, "CDR3a"],</pre>
                   CDR3b = CDR3ab[1:500, "CDR3b"],
                   clone_size = 1,
                   sample = "a")
b <- data.frame(CDR3a = CDR3ab[401:900, "CDR3a"],</pre>
                   CDR3b = CDR3ab[401:900, "CDR3b"],
                   clone_size = 1,
                   sample = "b")
b$clone_size[1] <- 20
# run ClustIRR analysis
c <- clustirr(s = rbind(a, b))</pre>
# detect communities
gcd <- detect_communities(graph = c$graph,</pre>
                           algorithm = "leiden",
                           resolution = 1,
                           weight = "ncweight",
                           chains = c("CDR3a", "CDR3b"))
ag <- get_ag_species_hits(node_summary = gcd$node_summary,</pre>
                           db = "vdjdb",
                           ag_species = "EBV")
head(ag)
```

get\_beta\_violin

*Visualize distribution of*  $\beta$  *means in each repertoire as violin plots* 

### **Description**

Visualize the  $\beta$  means as violin plots, representing relative community occupancies for individual repertoires. At the same time, annotate the communities (dots) based on their specificity.

# Usage

20 get\_beta\_violin

### Arguments

beta	beta data.frame
node_summary	node_summary data.frame
ag	antigen species/gene, character, e.g. "EBV", "CMV", or "MLANA"
ag_species	is the antigen a species (TRUE) or gene (FALSE)
db	annotation database, character, e.g. "vdjdb"

#### **Details**

The user has to provide an antigen species (e.g. ag = "EBV" and ag\_species=TRUE) or an antigen gene (e.g. ag = "MLANA" and ag\_species=FALSE). Furthermore, the user has to provide nodes (node\_summary data.frame created by the function detect\_communities) and beta data.frame which is part of posterior\_summary generated by the function dco.

The user can also select an annotation database db, such as "vdjdb", "mcpas" or "tcr3d". We will look for perfect matches between CDR3 sequences in the input and in the annotation database for annotation.

#### Value

The output is a violin ggplot.

```
# load package input data
data("CDR3ab", package = "ClustIRR")
a <- data.frame(CDR3a = CDR3ab[1:500, "CDR3a"],
                  CDR3b = CDR3ab[1:500, "CDR3b"],
                  clone_size = 1,
                  sample = "a")
b <- data.frame(CDR3a = CDR3ab[401:900, "CDR3a"],</pre>
                  CDR3b = CDR3ab[401:900, "CDR3b"],
                  clone\_size = 1,
                  sample = "b")
b$clone_size[1] <- 20
# run ClustIRR analysis
c <- clustirr(s = rbind(a, b))</pre>
# detect communities
gcd <- detect_communities(graph = c$graph,</pre>
                           algorithm = "leiden",
                           resolution = 1,
                           weight = "ncweight",
                           chains = c("CDR3a", "CDR3b"))
# differential community occupancy analysis
dco <- dco(community_occupancy_matrix = gcd$community_occupancy_matrix)</pre>
```

get\_community\_purity 21

get\_community\_purity Compute community purity with respect to a node feature

# **Description**

Given an igraph object with detected communities (via a node attribute called community) and a node feature (e.g., the output graph of the function detect\_communities), this function computes the purity of each community with respect to that feature.

If the node feature is numeric (e.g., expression values, clone size, etc.), purity is measured by the coefficient of variation (CV). If the node feature is categorical (e.g., tissue type, sample name, cell type), purity is measured by Gini Impurity (GI).

### Usage

```
get_community_purity(graph, node_feature)
```

#### **Arguments**

graph igraph object with community membership stored as a vertex attribute named

community.

node\_feature Character scalar specifying the vertex attribute to evaluate (must be either cate-

gorical or numeric).

### **Details**

The function checks for the existence of a community vertex attribute in the graph and for a valid node\_feature. Depending on whether the feature is numeric or categorical, it computes CV or GI, respectively.

- For categorical features: communities with low GI contain mostly homogeneous feature values
- For numeric features: communities with low CV are more homogeneous in feature values.

#### Value

A data. frame summarizing purity for each community:

- For categorical features: columns GI (Gini impurity), n (community size), community.
- For numeric features: columns mean, sd, cv (coefficient of variation), n (community size), community.

22 get\_honeycombs

### **Examples**

```
library(igraph)

# create a simple graph
g <- make_ring(10)
V(g)$community <- rep(1:2, each = 5)
V(g)$tissue <- rep(c("A", "B"), times = 5)
V(g)$expression <- rnorm(10, mean = 5, sd = 1)

# categorical feature
get_community_purity(graph = g, node_feature = "tissue")

# numeric feature
get_community_purity(graph = g, node_feature = "expression")</pre>
```

get\_honeycombs

Generate honycomb plot: visualize community occupancy of pairs of immune receptor repertoires

# Description

Use the community\_occupancy\_matrix generated by the function detect\_communities to generate honeycomb plots for each pair of repertoires. In each plot, we will show communities (rows in the matric community\_occupancy\_matrix) as dots and their intensities in a pair of repertoires (x-axis and y-axis). The density of dots is encoded by the color of the honeycomb-like hexagons.

#### **Usage**

```
get_honeycombs(com)
```

#### **Arguments**

com

community\_occupancy\_matrix, matrix generated by detect\_communities

#### **Details**

Use the community\_occupancy\_matrix generated by the function detect\_communities to generate honeycomb plots for each pair of repertoires. In each plot, we will show communities (rows in the matric community\_occupancy\_matrix) as dots and their intensities in a pair of repertoires (x-axis and y-axis). The density of dots is encoded by the color of the honeycomb-like hexagons.

### Value

The output is a list with ggplots. Given n repertoires (columns in input community\_occupancy\_matrix), it will generate n\*(n-1)/2 plots. You can arrange the ggplots (or a portion of them) in any shape e.g. with the R-package patchwork.

mcpas 23

### **Examples**

```
# load package input data
data("CDR3ab", package = "ClustIRR")
a <- data.frame(CDR3a = CDR3ab[1:300, "CDR3a"],</pre>
                 CDR3b = CDR3ab[1:300, "CDR3b"],
                 clone\_size = 1,
                 sample = "a")
b <- data.frame(CDR3a = CDR3ab[201:400, "CDR3a"],</pre>
                   CDR3b = CDR3ab[201:400, "CDR3b"],
                   clone_size = 1,
                   sample = "b")
b$clone_size[1] <- 20
# run ClustIRR analysis
c <- clustirr(s = rbind(a, b))</pre>
# detect communities
gcd <- detect_communities(graph = c$graph,</pre>
                           algorithm = "leiden",
                            resolution = 1,
                           weight = "ncweight",
                            chains = c("CDR3a", "CDR3b"))
# get honeycombs
g <- get_honeycombs(com = gcd$community_occupancy_matrix)</pre>
g
```

mcpas

CDR3 sequences and their matching epitopes obtained from McPAS-TCR

# **Description**

data.frame with CDR3a and/or CDR3b sequences and their matching antigenic epitopes obtained from McPAS-TCR. The remaining CDR3 columns are set to NA. For data processing details see the script inst/script/get\_mcpastcr.R

# Usage

```
data(mcpas)
```

### **Format**

data.frame with columns:

- 1. CDR3a: CDR3a amino acid sequence
- 2. CDR3b: CDR3b amino acid sequence

plot\_graph

```
3. CDR3g: CDR3g amino acid sequence -> NA
```

4. CDR3d: CDR3d amino acid sequence -> NA

5. CDR3h: CDR3h amino acid sequence -> NA

6. CDR31: CDR31 amino acid sequence -> NA

7. CDR3\_species: CDR3 species (e.g. human, mouse, ...)

8. Antigen\_species: antigen species

9. Antigen\_gene: antigen gene

10. Reference: Reference (Pubmed ID)

#### Value

data(mcpas) loads the object McPAS-TCR

#### **Source**

```
McPAS-TCR, June 2024
```

### **Examples**

```
data(mcpas)
```

plot\_graph

Plot ClustIRR graph

### **Description**

This function visualizes a graph. The main input is g object created by the function clustirr.

### Usage

### **Arguments**

g Object returned by the function clustirr

as\_visnet

logical, if as\_visnet=TRUE we plot an interactive graph with visNetwork. If as\_visnet=FALSE, we plot a static graph with igraph.

tcr3d 25

select\_by

character string, two values are possible: "Ag\_species" or "Ag\_gene". This only has an effect if as\_visnet = TRUE, i.e. if the graph is interactive. It will allow the user to highligh clones (nodes) in the graph that are associated with a specific antigenic specie or gene. The mapping between CDR3 and antigens is extracted from databases, such as, VDJdb, McPAS-TCR and TCR3d. If none of the clones in the graph are matched to a CDR3, then the user will have no options to select/highlight.

show\_singletons

 $logical, if \verb|show_singletons=TRUE| we plot all vertices. If \verb|show_singletons=FALSE|,$ 

we plot only vertices connected by edges.

node\_opacity probability, controls the opacity of node colors. Lower values corresponding to

more transparent colors.

#### Value

The output is an igraph or visNetwork plot.

The size of the vertices increases linearly as the logarithm of the degree of the clonal expansion (number of cells per clone) in the corresponding clones.

### **Examples**

```
# load package input data
data("CDR3ab", package = "ClustIRR")
s <- data.frame(CDR3b = CDR3ab[1:100, "CDR3b"], sample = "A", clone_size = 1)
# run ClustIRR analysis
c <- clustirr(s = s)
# plot graph with vertices as clones
plot_graph(c, as_visnet=FALSE, show_singletons=TRUE, node_opacity = 0.8)</pre>
```

tcr3d

CDR3 sequences and their matching epitopes obtained from TCR3d

# **Description**

data.frame with paired CDR3a and CDR3b CDR3 sequences and their matching epitopes obtained from TCR3d. The remaining CDR3 columns are set to NA. The antigenic epitopes come from cancer antigens and from viral antigens. For data processing details see the script inst/script/get\_tcr3d.R

#### Usage

```
data(tcr3d)
```

26 vdjdb

### **Format**

data.frame with columns:

- 1. CDR3a: CDR3a amino acid sequence
- 2. CDR3b: CDR3b amino acid sequence
- 3. CDR3g: CDR3g amino acid sequence -> NA
- 4. CDR3d: CDR3d amino acid sequence -> NA
- 5. CDR3h: CDR3h amino acid sequence -> NA
- 6. CDR31: CDR31 amino acid sequence -> NA
- 7. CDR3\_species: CDR3 species (e.g. human, mouse, ...)
- 8. Antigen\_species: antigen species
- 9. Antigen\_gene: antigen gene
- 10. Reference: Reference ID

# Value

data(tcr3d) loads the object tcr3d

### **Source**

TCR3d, June 2024

# **Examples**

data("tcr3d")

vdjdb

CDR3 sequences and their matching epitopes obtained from VDJdb

# Description

data.frame with unpaired CDR3a or CDR3b sequences and their matching epitopes obtained from VDJdb. The remaining CDR3 columns are set to NA. For data processing details see the script inst/script/get\_vdjdb.R

### Usage

```
data(vdjdb)
```

vdjdb 27

# **Format**

data.frame with columns:

- 1. CDR3a: CDR3a amino acid sequence
- 2. CDR3b: CDR3b amino acid sequence
- 3. CDR3g: CDR3g amino acid sequence -> NA
- 4. CDR3d: CDR3d amino acid sequence -> NA
- 5. CDR3h: CDR3h amino acid sequence -> NA
- 6. CDR31: CDR31 amino acid sequence -> NA
- 7. CDR3\_species: CDR3 species (e.g. human, mouse, ...)
- 8. Antigen\_species: antigen species
- 9. Antigen\_gene: antigen gene
- 10. Reference: Reference (Pubmed ID)

# Value

data(vdjdb) loads the object vdjdb

### Source

VDJdb, December 2024

# **Examples**

data("vdjdb")

# **Index**

```
* datasets
    BLOSUM62, 2
    Datasets, 7
    mcpas, 23
    tcr3d, 25
    vdjdb, 26
BLOSUM62, 2
CDR3ab (Datasets), 7
class:clust_irr (clust_irr-class), 6
clust_irr (clust_irr-class), 6
clust\_irr-class, 6
clustirr, 3
D1 (Datasets), 7
Datasets, 7
dco. 8
decode_all_communities, 10
decode_community, 12
detect\_communities, 14
get_ag_gene_hits, 16
get_ag_species_hits, 17
get_beta_violin, 19
get_clustirr_clust (clust_irr-class), 6
{\tt get\_clustirr\_clust,clust\_irr-method}
        (clust_irr-class), 6
get_clustirr_inputs (clust_irr-class), 6
get_clustirr_inputs,clust_irr-method
        (clust_irr-class), 6
get_community_purity, 21
get_honeycombs, 22
mcpas, 23
plot_graph, 24
tcr3d, 25
vdjdb, 26
```