

Package ‘metasnf’

April 25, 2025

Title Meta Clustering with Similarity Network Fusion

Version 2.1.1

Description Framework to facilitate patient subtyping with similarity network fusion and meta clustering. The similarity network fusion (SNF) algorithm was introduced by Wang et al. (2014) in <[doi:10.1038/nmeth.2810](https://doi.org/10.1038/nmeth.2810)>. SNF is a data integration approach that can transform high-dimensional and diverse data types into a single similarity network suitable for clustering with minimal loss of information from each initial data source. The meta clustering approach was introduced by Caruana et al. (2006) in <[doi:10.1109/ICDM.2006.103](https://doi.org/10.1109/ICDM.2006.103)>. Meta clustering involves generating a wide range of cluster solutions by adjusting clustering hyperparameters, then clustering the solutions themselves into a manageable number of qualitatively similar solutions, and finally characterizing representative solutions to find ones that are best for the user's specific context. This package provides a framework to easily transform multimodal data into a wide range of similarity network fusion-derived cluster solutions as well as to visualize, characterize, and validate those solutions. Core package functionality includes easy customization of distance metrics, clustering algorithms, and SNF hyperparameters to generate diverse clustering solutions; calculation and plotting of associations between features, between patients, and between cluster solutions; and standard cluster validation approaches including resampled measures of cluster stability, standard metrics of cluster quality, and label propagation to evaluate generalizability in unseen data. Associated vignettes guide the user through using the package to identify patient subtypes while adhering to best practices for unsupervised learning.

License GPL (>= 3)

Encoding UTF-8

RoxygenNote 7.3.2

Imports cli, cluster, data.table, digest, dplyr, ggplot2, grDevices, MASS, mclust, methods, progressr, purrr, RColorBrewer, rlang, SNFtool, stats, tibble, tidyr, utils

Suggests circlize, ComplexHeatmap, InteractiveComplexHeatmap, clv, future, future.apply, knitr, rmarkdown, testthat (>= 3.0.0), ggalluvial, lifecycle, dbscan

Config/testthat/edition 3

Depends R (>= 4.1.0)

LazyData true

VignetteBuilder knitr

URL <https://branchlab.github.io/metasn/>,
<https://github.com/BRANCHlab/metasn/>

BugReports <https://github.com/BRANCHlab/metasn/issues>

NeedsCompilation no

Author Prashanth S Velayudhan [aut, cre],
 Xiaoqiao Xu [aut],
 Prajka Kallurkar [aut],
 Ana Patricia Balbon [aut],
 Maria T Secara [aut],
 Adam Taback [aut],
 Denise Sabac [aut],
 Nicholas Chan [aut],
 Shihao Ma [aut],
 Bo Wang [aut],
 Daniel Felsky [aut],
 Stephanie H Ameis [aut],
 Brian Cox [aut],
 Colin Hawco [aut],
 Lauren Erdman [aut],
 Anne L Wheeler [aut, ths]

Maintainer Prashanth S Velayudhan <psvelayu@gmail.com>

Repository CRAN

Date/Publication 2025-04-25 14:20:02 UTC

Contents

abcd_anxiety	6
abcd_colour	7
abcd_cort_sa	8
abcd_cort_t	9
abcd_depress	10
abcd_h_income	11
abcd_income	11
abcd_pubertal	12
abcd_subc_v	13
add_settings_df_rows	14
age_df	17
alluvial_cluster_plot	17
anxiety	19
as.data.frame.data_list	20
as.data.frame.ext_solutions_df	20
as.data.frame.settings_df	21

as.data.frame.snf_config	22
as.data.frame.solutions_df	22
as.data.frame.t_ext_solutions_df	23
as.data.frame.t_solutions_df	24
as.data.frame.weights_matrix	24
as.list.clust_fns_list	25
as.list.data_list	25
as.list.dist_fns_list	26
as.list.sim_mats_list	26
as.list.snf_config	27
as.matrix.ari_matrix	27
as.matrix.weights_matrix	28
assemble_data	28
assoc_pval_heatmap	29
as_ari_matrix	30
as_data_list	31
as_settings_df	31
as_sim_mats_list	32
as_snf_config	32
as_weights_matrix	33
auto_plot	33
bar_plot	34
batch_snf	35
batch_snf_subsamples	36
cache_a_complete_example_ext_sol_df	37
cache_a_complete_example_lp_ext_sol_df	38
cache_a_complete_example_sol_df	38
calculate_coclustering	39
calc_aris	40
calc_assoc_pval_matrix	41
calc_nmis	42
cancer_diagnosis_df	43
cell_significance_fn	44
check_dataless_annotations	45
check_hm_dependencies	45
check_similarity_matrices	46
clust_fns	46
clust_fns_list	47
cocluster_density	48
cocluster_heatmap	50
colour_scale	52
cort_sa	53
cort_t	53
data_list	54
depress	56
diagnosis_df	57
dist_fns	58
dist_fns_list	59

dlapply	60
dplyr_row_slice.ext_solutions_df	61
dplyr_row_slice.solutions_df	62
esm_manhattan_plot	62
estimate_nclust_given_graph	64
expression_df	65
extend_solutions	66
fav_colour	67
gender_df	68
get_complete_uids	68
get_heatmap_order	69
get_matrix_order	70
get_pvals	71
get_representative_solutions	72
income	73
is_data_list	74
jitter_plot	75
label_meta_clusters	75
label_propagate	77
linear_adjust	79
mc_manhattan_plot	80
merge.clust_fns_list	82
merge.data_list	83
merge.dist_fns_list	83
merge.ext_solutions_df	84
merge.settings_df	84
merge.sim_mats_list	85
merge.snf_config	85
merge.solutions_df	86
merge.t_ext_solutions_df	86
merge.t_solutions_df	87
merge.weights_matrix	87
merge_df_list	88
methylation_df	88
mock_ari_matrix	89
mock_clust_fns_list	89
mock_data_list	90
mock_dist_fns_list	90
mock_ext_solutions_df	91
mock_mc_solutions_df	91
mock_rep_solutions_df	92
mock_settings_df	92
mock_snf_config	93
mock_solutions_df	93
mock_t_solutions_df	94
mock_weights_matrix	94
new_solutions_df	95
plot.ari_matrix	95

plot.data_list	97
plot.ext_solutions_df	98
plot.snf_config	99
plot.solutions_df	102
print.ari_matrix	103
print.clust_fns_list	104
print.data_list	104
print.dist_fns_list	105
print.ext_solutions_df	105
print.settings_df	106
print.sim_mats_list	106
print.snf_config	107
print.solutions_df	107
print.t_ext_solutions_df	108
print.t_solutions_df	108
print.weights_matrix	109
pubertal	109
pval_heatmap	110
quality_measures	112
rbind.ext_solutions_df	113
rbind.solutions_df	114
rbind.t_solutions_df	114
rbind.weights_matrix	115
rename_dl	115
resample	116
save_heatmap	117
settings_df	117
shiny_annotator	121
similarity_matrix_heatmap	122
sim_mats_list	125
siw_euclidean_distance	125
snf_config	126
split_parser	130
str.ari_matrix	131
str.clust_fns_list	132
str.data_list	132
str.dist_fns_list	133
str.ext_solutions_df	133
str.settings_df	134
str.sim_mats_list	134
str.snf_config	135
str.solutions_df	135
str.t_ext_solutions_df	136
str.t_solutions_df	136
str.weights_matrix	137
subc_v	137
subsample_dl	138
subsample_pairwise_aris	139

summary.ari_matrix	140
summary.clust_fns_list	141
summary.data_list	141
summary.dist_fns_list	142
summary.ext_solutions_df	143
summary.settings_df	143
summary.sim_mats_list	144
summary.snf_config	144
summary.solutions_df	145
summary.t_ext_solutions_df	145
summary.t_solutions_df	146
summary.weights_matrix	146
train_test_assign	147
uids	147
validate_solutions_df	148
var_manhattan_plot	148
weights_matrix	149

Index **151**

abcd_anxiety	<i>Mock ABCD anxiety data</i>
--------------	-------------------------------

Description

A randomly shuffled and anonymized copy of anxiety data from the NIMH Data archive. The original file used was pdem02.txt. The file was pre-processed by the abcdutils package (<https://github.com/BRANCHlab/abcdutils>) function `get_cbcl_anxiety`.

Usage

```
abcd_anxiety
```

Format

`abcd_anxiety`:

A data frame with 275 rows and 2 columns:

patient The unique identifier of the ABCD dataset

cbcl_anxiety_r Ordinal value of impairment on CBCL anxiety, either 0 (no impairment), 1 (borderline clinical), or 2 (clinically impaired)

Source

Though this data is no longer "real" ABCD data, the reference for using ABCD as a data source is below:

Data used in the preparation of this article were obtained from the Adolescent Brain Cognitive Development (ABCD) Study (<https://abcdstudy.org>), held in the NIMH Data Archive (NDA). This is a multisite, longitudinal study designed to recruit more than 10,000 children age 9-10 and follow them over 10 years into early adulthood. The ABCD Study® is supported by the National Institutes of Health and additional federal partners under award numbers U01DA041048, U01DA050989, U01DA051016, U01DA041022, U01DA051018, U01DA051037, U01DA050987, U01DA041174, U01DA041106, U01DA041117, U01DA041028, U01DA041134, U01DA050988, U01DA051039, U01DA041156, U01DA041025, U01DA041120, U01DA051038, U01DA041148, U01DA041093, U01DA041089, U24DA041123, U24DA041147. A full list of supporters is available at <https://abcdstudy.org/federal-partners.html>. A listing of participating sites and a complete listing of the study investigators can be found at https://abcdstudy.org/consortium_members/. ABCD consortium investigators designed and implemented the study and/or provided data but did not necessarily participate in the analysis or writing of this report. This manuscript reflects the views of the authors and may not reflect the opinions or views of the NIH or ABCD consortium investigators.

abcd_colour

Mock ABCD "colour" data

Description

A randomly shuffled and anonymized copy of depression data from the NIMH Data archive. The original file used was pdem02.txt. The file was pre-processed by the abcdutils package (<https://github.com/BRANCHlab/abcdutils>) function `get_cbc1_depress`. The data was transformed into categorical colour values to demonstrate the Chi-squared test capabilities of `extend_solutions`.

Usage

```
abcd_colour
```

Format

`abcd_colour`:

A data frame with 275 rows and 2 columns:

patient The unique identifier of the ABCD dataset

colour Categorical transformation of `cbc1_depress`.

Source

Though this data is no longer "real" ABCD data, the reference for using ABCD as a data source is below:

Data used in the preparation of this article were obtained from the Adolescent Brain Cognitive Development (ABCD) Study (<https://abcdstudy.org>), held in the NIMH Data Archive (NDA). This is a multisite, longitudinal study designed to recruit more than 10,000 children age 9-10 and follow

them over 10 years into early adulthood. The ABCD Study® is supported by the National Institutes of Health and additional federal partners under award numbers U01DA041048, U01DA050989, U01DA051016, U01DA041022, U01DA051018, U01DA051037, U01DA050987, U01DA041174, U01DA041106, U01DA041117, U01DA041028, U01DA041134, U01DA050988, U01DA051039, U01DA041156, U01DA041025, U01DA041120, U01DA051038, U01DA041148, U01DA041093, U01DA041089, U24DA041123, U24DA041147. A full list of supporters is available at <https://abcdstudy.org/federal-partners.html>. A listing of participating sites and a complete listing of the study investigators can be found at https://abcdstudy.org/consortium_members/. ABCD consortium investigators designed and implemented the study and/or provided data but did not necessarily participate in the analysis or writing of this report. This manuscript reflects the views of the authors and may not reflect the opinions or views of the NIH or ABCD consortium investigators.

abcd_cort_sa

Mock ABCD cortical surface area data

Description

A randomly shuffled and anonymized copy of cortical surface area data from the NIMH Data archive. The original file used was `mrisdp10201.txt`. The file was pre-processed by the `abcdutils` package (<https://github.com/BRANCHlab/abcdutils>) function `get_cort_t`.

Usage

`abcd_cort_sa`

Format

`abcd_cort_sa`:

A data frame with 188 rows and 152 columns:

patient The unique identifier of the ABCD dataset

... Cortical surface areas of various ROIs (mm², I think)

Source

Though this data is no longer "real" ABCD data, the reference for using ABCD as a data source is below:

Data used in the preparation of this article were obtained from the Adolescent Brain Cognitive Development (ABCD) Study (<https://abcdstudy.org>), held in the NIMH Data Archive (NDA). This is a multisite, longitudinal study designed to recruit more than 10,000 children age 9-10 and follow them over 10 years into early adulthood. The ABCD Study® is supported by the National Institutes of Health and additional federal partners under award numbers U01DA041048, U01DA050989, U01DA051016, U01DA041022, U01DA051018, U01DA051037, U01DA050987, U01DA041174, U01DA041106, U01DA041117, U01DA041028, U01DA041134, U01DA050988, U01DA051039, U01DA041156, U01DA041025, U01DA041120, U01DA051038, U01DA041148, U01DA041093, U01DA041089, U24DA041123, U24DA041147. A full list of supporters is available at <https://abcdstudy.org/federal-partners.html>. A listing of participating sites and a complete listing of the study investigators can

be found at https://abcdstudy.org/consortium_members/. ABCD consortium investigators designed and implemented the study and/or provided data but did not necessarily participate in the analysis or writing of this report. This manuscript reflects the views of the authors and may not reflect the opinions or views of the NIH or ABCD consortium investigators.

 abcd_cort_t

Mock ABCD cortical thickness data

Description

A randomly shuffled and anonymized copy of cortical thickness data from the NIMH Data archive. The original file used was `mrisd10201.txt`. The file was pre-processed by the `abcdutils` package (<https://github.com/BRANCHlab/abcdutils>) function `get_cort_t`.

Usage

```
abcd_cort_t
```

Format

```
abcd_cort_t:
```

A data frame with 188 rows and 152 columns:

patient The unique identifier of the ABCD dataset

... Cortical thicknesses of various ROIs (mm³, I think)

Source

Though this data is no longer "real" ABCD data, the reference for using ABCD as a data source is below:

Data used in the preparation of this article were obtained from the Adolescent Brain Cognitive Development (ABCD) Study (<https://abcdstudy.org>), held in the NIMH Data Archive (NDA). This is a multisite, longitudinal study designed to recruit more than 10,000 children age 9-10 and follow them over 10 years into early adulthood. The ABCD Study® is supported by the National Institutes of Health and additional federal partners under award numbers U01DA041048, U01DA050989, U01DA051016, U01DA041022, U01DA051018, U01DA051037, U01DA050987, U01DA041174, U01DA041106, U01DA041117, U01DA041028, U01DA041134, U01DA050988, U01DA051039, U01DA041156, U01DA041025, U01DA041120, U01DA051038, U01DA041148, U01DA041093, U01DA041089, U24DA041123, U24DA041147. A full list of supporters is available at <https://abcdstudy.org/federal-partners.html>. A listing of participating sites and a complete listing of the study investigators can be found at https://abcdstudy.org/consortium_members/. ABCD consortium investigators designed and implemented the study and/or provided data but did not necessarily participate in the analysis or writing of this report. This manuscript reflects the views of the authors and may not reflect the opinions or views of the NIH or ABCD consortium investigators.

`abcd_depress`*Mock ABCD depression data*

Description

A randomly shuffled and anonymized copy of depression data from the NIMH Data archive. The original file used was pdem02.txt. The file was pre-processed by the abcdutils package (<https://github.com/BRANCHlab/abcdutils>) function `get_cbcl_depress`.

Usage

```
abcd_depress
```

Format

`abcd_depress`:

A data frame with 275 rows and 2 columns:

patient The unique identifier of the ABCD dataset

cbcl_depress_r Ordinal value of impairment on CBCL anxiety, either 0 (no impairment), 1 (borderline clinical), or 2 (clinically impaired)

Source

Though this data is no longer "real" ABCD data, the reference for using ABCD as a data source is below:

Data used in the preparation of this article were obtained from the Adolescent Brain Cognitive Development (ABCD) Study (<https://abcdstudy.org>), held in the NIMH Data Archive (NDA). This is a multisite, longitudinal study designed to recruit more than 10,000 children age 9-10 and follow them over 10 years into early adulthood. The ABCD Study® is supported by the National Institutes of Health and additional federal partners under award numbers U01DA041048, U01DA050989, U01DA051016, U01DA041022, U01DA051018, U01DA051037, U01DA050987, U01DA041174, U01DA041106, U01DA041117, U01DA041028, U01DA041134, U01DA050988, U01DA051039, U01DA041156, U01DA041025, U01DA041120, U01DA051038, U01DA041148, U01DA041093, U01DA041089, U24DA041123, U24DA041147. A full list of supporters is available at <https://abcdstudy.org/federal-partners.html>. A listing of participating sites and a complete listing of the study investigators can be found at https://abcdstudy.org/consortium_members/. ABCD consortium investigators designed and implemented the study and/or provided data but did not necessarily participate in the analysis or writing of this report. This manuscript reflects the views of the authors and may not reflect the opinions or views of the NIH or ABCD consortium investigators.

abcd_h_income	<i>Mock ABCD income data</i>
---------------	------------------------------

Description

Like abcd_income, but with no NAs in patient column

Usage

abcd_h_income

Format

abcd_income:

A data frame with 300 rows and 2 columns:

patient The unique identifier of the ABCD dataset

household_income Household income in 3 category levels (low = 1, medium = 2, high = 3)

Source

Though this data is no longer "real" ABCD data, the reference for using ABCD as a data source is below:

Data used in the preparation of this article were obtained from the Adolescent Brain Cognitive Development (ABCD) Study (<https://abcdstudy.org>), held in the NIMH Data Archive (NDA). This is a multisite, longitudinal study designed to recruit more than 10,000 children age 9-10 and follow them over 10 years into early adulthood. The ABCD Study® is supported by the National Institutes of Health and additional federal partners under award numbers U01DA041048, U01DA050989, U01DA051016, U01DA041022, U01DA051018, U01DA051037, U01DA050987, U01DA041174, U01DA041106, U01DA041117, U01DA041028, U01DA041134, U01DA050988, U01DA051039, U01DA041156, U01DA041025, U01DA041120, U01DA051038, U01DA041148, U01DA041093, U01DA041089, U24DA041123, U24DA041147. A full list of supporters is available at <https://abcdstudy.org/federal-partners.html>. A listing of participating sites and a complete listing of the study investigators can be found at https://abcdstudy.org/consortium_members/. ABCD consortium investigators designed and implemented the study and/or provided data but did not necessarily participate in the analysis or writing of this report. This manuscript reflects the views of the authors and may not reflect the opinions or views of the NIH or ABCD consortium investigators.

abcd_income	<i>Mock ABCD income data</i>
-------------	------------------------------

Description

A randomly shuffled and anonymized copy of income data from the NIMH Data archive. The original file used was pdem02.txt The file was pre-processed by the abcdutils package (<https://github.com/BRANCHlab/abcdutils>) function get_income.

Usage

abcd_income

Format

abcd_income:

A data frame with 300 rows and 2 columns:

patient The unique identifier of the ABCD dataset

household_income Household income in 3 category levels (low = 1, medium = 2, high = 3)

Source

Though this data is no longer "real" ABCD data, the reference for using ABCD as a data source is below:

Data used in the preparation of this article were obtained from the Adolescent Brain Cognitive Development (ABCD) Study (<https://abcdstudy.org>), held in the NIMH Data Archive (NDA). This is a multisite, longitudinal study designed to recruit more than 10,000 children age 9-10 and follow them over 10 years into early adulthood. The ABCD Study® is supported by the National Institutes of Health and additional federal partners under award numbers U01DA041048, U01DA050989, U01DA051016, U01DA041022, U01DA051018, U01DA051037, U01DA050987, U01DA041174, U01DA041106, U01DA041117, U01DA041028, U01DA041134, U01DA050988, U01DA051039, U01DA041156, U01DA041025, U01DA041120, U01DA051038, U01DA041148, U01DA041093, U01DA041089, U24DA041123, U24DA041147. A full list of supporters is available at <https://abcdstudy.org/federal-partners.html>. A listing of participating sites and a complete listing of the study investigators can be found at https://abcdstudy.org/consortium_members/. ABCD consortium investigators designed and implemented the study and/or provided data but did not necessarily participate in the analysis or writing of this report. This manuscript reflects the views of the authors and may not reflect the opinions or views of the NIH or ABCD consortium investigators.

abcd_pubertal

Mock ABCD pubertal status data

Description

A randomly shuffled and anonymized copy of pubertal status data from the NIMH Data archive. The original files used were abcd_ssphp01.txt and abcd_ssphy01.txt. The file was pre-processed by the abcdutils package (<https://github.com/BRANCHlab/abcdutils>) function `get_pubertal_status`.

Usage

abcd_pubertal

Format

abcd_pubertal:

A data frame with 275 rows and 2 columns:

patient The unique identifier of the ABCD dataset

pubertal_status Average reported pubertal status between child and parent (1-5 categorical scale)

Source

Though this data is no longer "real" ABCD data, the reference for using ABCD as a data source is below:

Data used in the preparation of this article were obtained from the Adolescent Brain Cognitive Development (ABCD) Study (<https://abcdstudy.org>), held in the NIMH Data Archive (NDA). This is a multisite, longitudinal study designed to recruit more than 10,000 children age 9-10 and follow them over 10 years into early adulthood. The ABCD Study® is supported by the National Institutes of Health and additional federal partners under award numbers U01DA041048, U01DA050989, U01DA051016, U01DA041022, U01DA051018, U01DA051037, U01DA050987, U01DA041174, U01DA041106, U01DA041117, U01DA041028, U01DA041134, U01DA050988, U01DA051039, U01DA041156, U01DA041025, U01DA041120, U01DA051038, U01DA041148, U01DA041093, U01DA041089, U24DA041123, U24DA041147. A full list of supporters is available at <https://abcdstudy.org/federal-partners.html>. A listing of participating sites and a complete listing of the study investigators can be found at https://abcdstudy.org/consortium_members/. ABCD consortium investigators designed and implemented the study and/or provided data but did not necessarily participate in the analysis or writing of this report. This manuscript reflects the views of the authors and may not reflect the opinions or views of the NIH or ABCD consortium investigators.

abcd_subc_v

Mock ABCD subcortical volumes data

Description

A randomly shuffled and anonymized copy of subcortical volume data from the NIMH Data archive. The original file used was smrip10201.txt The file was pre-processed by the abcdutils package (<https://github.com/BRANCHlab/abcdutils>) function `get_subc_v`.

Usage

`abcd_subc_v`

Format

`abcd_subc_v`:

A data frame with 174 rows and 31 columns:

patient The unique identifier of the ABCD dataset

... Subcortical volumes of various ROIs (mm³, I think)

Source

Though this data is no longer "real" ABCD data, the reference for using ABCD as a data source is below:

Data used in the preparation of this article were obtained from the Adolescent Brain Cognitive Development (ABCD) Study (<https://abcdstudy.org>), held in the NIMH Data Archive (NDA). This is a multisite, longitudinal study designed to recruit more than 10,000 children age 9-10 and follow them over 10 years into early adulthood. The ABCD Study® is supported by the National Institutes

of Health and additional federal partners under award numbers U01DA041048, U01DA050989, U01DA051016, U01DA041022, U01DA051018, U01DA051037, U01DA050987, U01DA041174, U01DA041106, U01DA041117, U01DA041028, U01DA041134, U01DA050988, U01DA051039, U01DA041156, U01DA041025, U01DA041120, U01DA051038, U01DA041148, U01DA041093, U01DA041089, U24DA041123, U24DA041147. A full list of supporters is available at <https://abcdstudy.org/federal-partners.html>. A listing of participating sites and a complete listing of the study investigators can be found at https://abcdstudy.org/consortium_members/. ABCD consortium investigators designed and implemented the study and/or provided data but did not necessarily participate in the analysis or writing of this report. This manuscript reflects the views of the authors and may not reflect the opinions or views of the NIH or ABCD consortium investigators.

add_settings_df_rows *Add rows to a settings_df*

Description

Add rows to a settings_df

Usage

```
add_settings_df_rows(
  sdf,
  n_solutions = 0,
  min_removed_inputs = 0,
  max_removed_inputs = sum(startsWith(colnames(sdf), "inc_")) - 1,
  dropout_dist = "exponential",
  min_alpha = NULL,
  max_alpha = NULL,
  min_k = NULL,
  max_k = NULL,
  min_t = NULL,
  max_t = NULL,
  alpha_values = NULL,
  k_values = NULL,
  t_values = NULL,
  possible_snf_schemes = c(1, 2, 3),
  clustering_algorithms = NULL,
  continuous_distances = NULL,
  discrete_distances = NULL,
  ordinal_distances = NULL,
  categorical_distances = NULL,
  mixed_distances = NULL,
  dfl = NULL,
  snf_input_weights = NULL,
  snf_domain_weights = NULL,
  retry_limit = 10,
  allow_duplicates = FALSE
)
```

Arguments

sdf	The existing settings data frame
n_solutions	Number of rows to generate for the settings data frame.
min_removed_inputs	The smallest number of input data frames that may be randomly removed. By default, 0.
max_removed_inputs	The largest number of input data frames that may be randomly removed. By default, this is 1 less than all the provided input data frames in the data list.
dropout_dist	Parameter controlling how the random removal of input data frames should occur. Can be "none" (no input data frames are randomly removed), "uniform" (uniformly sample between min_removed_inputs and max_removed_inputs to determine number of input data frames to remove), or "exponential" (pick number of input data frames to remove by sampling from min_removed_inputs to max_removed_inputs with an exponential distribution; the default).
min_alpha	The minimum value that the alpha hyperparameter can have. Random assigned value of alpha for each row will be obtained by uniformly sampling numbers between min_alpha and max_alpha at intervals of 0.1. Cannot be used in conjunction with the alpha_values parameter.
max_alpha	The maximum value that the alpha hyperparameter can have. See min_alpha parameter. Cannot be used in conjunction with the alpha_values parameter.
min_k	The minimum value that the k hyperparameter can have. Random assigned value of k for each row will be obtained by uniformly sampling numbers between min_k and max_k at intervals of 1. Cannot be used in conjunction with the k_values parameter.
max_k	The maximum value that the k hyperparameter can have. See min_k parameter. Cannot be used in conjunction with the k_values parameter.
min_t	The minimum value that the t hyperparameter can have. Random assigned value of t for each row will be obtained by uniformly sampling numbers between min_t and max_t at intervals of 1. Cannot be used in conjunction with the t_values parameter.
max_t	The maximum value that the t hyperparameter can have. See min_t parameter. Cannot be used in conjunction with the t_values parameter.
alpha_values	A number or numeric vector of a set of possible values that alpha can take on. Value will be obtained by uniformly sampling the vector. Cannot be used in conjunction with the min_alpha or max_alpha parameters.
k_values	A number or numeric vector of a set of possible values that k can take on. Value will be obtained by uniformly sampling the vector. Cannot be used in conjunction with the min_k or max_k parameters.
t_values	A number or numeric vector of a set of possible values that t can take on. Value will be obtained by uniformly sampling the vector. Cannot be used in conjunction with the min_t or max_t parameters.
possible_snf_schemes	A vector containing the possible snf_schemes to uniformly randomly select from. By default, the vector contains all 3 possible schemes: c(1, 2, 3). 1

corresponds to the "individual" scheme, 2 corresponds to the "domain" scheme, and 3 corresponds to the "two-step" scheme.

clustering_algorithms

A list of clustering algorithms to uniformly randomly pick from when clustering. When not specified, randomly select between spectral clustering using the eigen-gap heuristic and spectral clustering using the rotation cost heuristic. See `?clust_fns_list` for more details on running custom clustering algorithms.

continuous_distances

A vector of continuous distance metrics to use when a custom `dist_fns_list` is provided.

discrete_distances

A vector of categorical distance metrics to use when a custom `dist_fns_list` is provided.

ordinal_distances

A vector of categorical distance metrics to use when a custom `dist_fns_list` is provided.

categorical_distances

A vector of categorical distance metrics to use when a custom `dist_fns_list` is provided.

mixed_distances

A vector of mixed distance metrics to use when a custom `dist_fns_list` is provided.

df1

List containing distance metrics to vary over. See `?generate_dist_fns_list`.

snf_input_weights

Nested list containing weights for when SNF is used to merge individual input measures (see `?generate_snf_weights`)

snf_domain_weights

Nested list containing weights for when SNF is used to merge domains (see `?generate_snf_weights`)

retry_limit

The maximum number of attempts to generate a novel row. This function does not return matrices with identical rows. As the range of requested possible settings tightens and the number of requested rows increases, the risk of randomly generating a row that already exists increases. If a new random row has matched an existing row `retry_limit` number of times, the function will terminate.

allow_duplicates

If TRUE, enables creation of a settings data frame with duplicate non-feature weighting related hyperparameters. This function should only be used when paired with a custom weights matrix that has non-duplicate rows.

Value

A settings data frame

age_df	<i>Mock age data</i>
--------	----------------------

Description

Mock age data

Usage

```
age_df
```

Format

age_df:

A data frame with 200 rows and 2 columns:

patient_id Random three-digit number uniquely identifying the patient

age Mock age feature

Source

This data came from the SNFtool package, with slight modifications.

alluvial_cluster_plot	<i>Alluvial plot of patients across cluster counts and important features</i>
-----------------------	---

Description

This function creates an alluvial plot that shows how observations in a similarity matrix could have been clustered over a set of clustering functions.

Usage

```
alluvial_cluster_plot(  
  cluster_sequence,  
  similarity_matrix,  
  dl = NULL,  
  data = NULL,  
  key_outcome,  
  key_label = key_outcome,  
  extra_outcomes = NULL,  
  title = NULL  
)
```

Arguments

<code>cluster_sequence</code>	A list of clustering algorithms.
<code>similarity_matrix</code>	A similarity matrix.
<code>d1</code>	A data list.
<code>data</code>	A data frame that contains any features to include in the plot.
<code>key_outcome</code>	The name of the feature that determines how each patient stream is coloured in the alluvial plot.
<code>key_label</code>	Name of key outcome to be used for the plot legend.
<code>extra_outcomes</code>	Names of additional features to add to the plot.
<code>title</code>	Title of the plot.

Value

An alluvial plot (class "gg" and "ggplot") showing distribution of a feature across varying number cluster solutions.

Examples

```
input_dl <- data_list(
  list(gender_df, "gender", "demographics", "categorical"),
  list(diagnosis_df, "diagnosis", "clinical", "categorical"),
  uid = "patient_id"
)

sc <- snf_config(input_dl, n_solutions = 1)

sol_df <- batch_snf(input_dl, sc, return_sim_mats = TRUE)

sim_mats <- sim_mats_list(sol_df)

clust_fn_sequence <- list(spectral_two, spectral_four)

alluvial_cluster_plot(
  cluster_sequence = clust_fn_sequence,
  similarity_matrix = sim_mats[[1]],
  d1 = input_dl,
  key_outcome = "gender",
  key_label = "Gender",
  extra_outcomes = "diagnosis",
  title = "Gender Across Cluster Counts"
)
```

anxiety	<i>Mock ABCD anxiety data</i>
---------	-------------------------------

Description

Like the mock data frame "abcd_colour", but with "unique_id" as the "uid".

Usage

```
anxiety
```

Format

`anxiety`:

A data frame with 275 rows and 2 columns:

unique_id The unique identifier of the ABCD dataset

cbcl_anxiety_r Ordinal value of impairment on CBCL anxiety, either 0 (no impairment), 1 (borderline clinical), or 2 (clinically impaired)

Source

Though this data is no longer "real" ABCD data, the reference for using ABCD as a data source is below:

Data used in the preparation of this article were obtained from the Adolescent Brain Cognitive Development (ABCD) Study (<https://abcdstudy.org>), held in the NIMH Data Archive (NDA). This is a multisite, longitudinal study designed to recruit more than 10,000 children age 9-10 and follow them over 10 years into early adulthood. The ABCD Study® is supported by the National Institutes of Health and additional federal partners under award numbers U01DA041048, U01DA050989, U01DA051016, U01DA041022, U01DA051018, U01DA051037, U01DA050987, U01DA041174, U01DA041106, U01DA041117, U01DA041028, U01DA041134, U01DA050988, U01DA051039, U01DA041156, U01DA041025, U01DA041120, U01DA051038, U01DA041148, U01DA041093, U01DA041089, U24DA041123, U24DA041147. A full list of supporters is available at <https://abcdstudy.org/federal-partners.html>. A listing of participating sites and a complete listing of the study investigators can be found at https://abcdstudy.org/consortium_members/. ABCD consortium investigators designed and implemented the study and/or provided data but did not necessarily participate in the analysis or writing of this report. This manuscript reflects the views of the authors and may not reflect the opinions or views of the NIH or ABCD consortium investigators.

```
as.data.frame.data_list
```

Coerce a data_list class object into a data.frame class object

Description

Horizontally joins data frames within a data list into a single data frame, using the uid attribute as the joining key.

Usage

```
## S3 method for class 'data_list'
as.data.frame(x, row.names = NULL, optional = FALSE, ...)
```

Arguments

x	A data_list class object.
row.names	Additional parameter passed to as.data.frame().
optional	Additional parameter passed to as.data.frame().
...	Additional parameter passed to as.data.frame().

Value

dl_df A data.frame class object with all the features and observations of dl.

```
as.data.frame.ext_solutions_df
```

Coerce a ext_solutions_df class object into a data.frame class object

Description

Coerce a ext_solutions_df class object into a data.frame class object

Usage

```
## S3 method for class 'ext_solutions_df'
as.data.frame(
  x,
  row.names = NULL,
  optional = FALSE,
  keep_attributes = FALSE,
  ...
)
```

Arguments

x	A ext_solutions_df class object.
row.names	Additional parameter passed to as.data.frame().
optional	Additional parameter passed to as.data.frame().
keep_attributes	If TRUE, resulting data frame includes settings data frame and weights matrix.
...	Additional parameter passed to as.data.frame().

Value

A data.frame class object with all the columns of x and its contained solutions data frame.

```
as.data.frame.settings_df
```

Coerce a settings_df class object into a data.frame class object

Description

Coerce a settings_df class object into a data.frame class object

Usage

```
## S3 method for class 'settings_df'
as.data.frame(x, ...)
```

Arguments

x	A settings_df class object.
...	Additional parameter passed to as.data.frame().

Value

A data.frame class object with all the columns of x and its contained solutions data frame.

```
as.data.frame.snf_config
```

Coerce a settings_df class object into a data.frame class object

Description

Coerce a settings_df class object into a data.frame class object

Usage

```
## S3 method for class 'snf_config'  
as.data.frame(x, ...)
```

Arguments

x A settings_df class object.
... Additional parameter passed to as.data.frame().

Value

A data.frame class object with all the columns of x and its contained solutions data frame.

```
as.data.frame.solutions_df
```

Coerce a solutions_df class object into a data.frame class object

Description

Coerce a solutions_df class object into a data.frame class object

Usage

```
## S3 method for class 'solutions_df'  
as.data.frame(  
  x,  
  row.names = NULL,  
  optional = FALSE,  
  keep_attributes = FALSE,  
  ...  
)
```

Arguments

x	A solutions_df class object.
row.names	Additional parameter passed to as.data.frame().
optional	Additional parameter passed to as.data.frame().
keep_attributes	If TRUE, resulting data frame includes settings data frame and weights matrix.
...	Additional parameter passed to as.data.frame().

Value

A data.frame class object with all the columns of x and its contained solutions data frame.

```
as.data.frame.t_ext_solutions_df
      Coerce a t_ext_solutions_df class object into a data.frame class
      object
```

Description

Coerce a t_ext_solutions_df class object into a data.frame class object

Usage

```
## S3 method for class 't_ext_solutions_df'
as.data.frame(x, ...)
```

Arguments

x	A t_ext_solutions_df class object.
...	Additional parameter passed to as.data.frame().

Value

A data.frame class object with all the columns of x and its contained solutions data frame.

as.data.frame.t_solutions_df

Coerce a t_solutions_df class object into a data.frame class object

Description

Coerce a t_solutions_df class object into a data.frame class object

Usage

```
## S3 method for class 't_solutions_df'  
as.data.frame(x, ...)
```

Arguments

x A t_solutions_df class object.
... Additional parameter passed to as.data.frame().

Value

A data.frame class object with all the columns of x and its contained solutions data frame.

as.data.frame.weights_matrix

Coerce a weights_matrix class object into a data.frame class object

Description

Coerce a weights_matrix class object into a data.frame class object

Usage

```
## S3 method for class 'weights_matrix'  
as.data.frame(x, ...)
```

Arguments

x A weights_matrix class object.
... Additional parameter passed to as.data.frame().

Value

A data.frame class object with all the columns of x and its contained solutions data frame.

`as.list.clust_fns_list`*Coerce a clust_fns_list class object into a list class object*

Description

Coerce a clust_fns_list class object into a list class object

Usage

```
## S3 method for class 'clust_fns_list'  
as.list(x, ...)
```

Arguments

x A clust_fns_list class object.
... Additional parameter passed to as.list().

Value

A list class object with all the functions of x.

`as.list.data_list`*Coerce a data_list class object into a list class object*

Description

Coerce a data_list class object into a list class object

Usage

```
## S3 method for class 'data_list'  
as.list(x, ...)
```

Arguments

x A data_list class object.
... Additional parameter passed to as.list().

Value

A list class object with all the objects of x.

as.list.dist_fns_list *Coerce a dist_fns_list class object into a list class object*

Description

Coerce a dist_fns_list class object into a list class object

Usage

```
## S3 method for class 'dist_fns_list'  
as.list(x, ...)
```

Arguments

x A dist_fns_list class object.
... Additional parameter passed to as.list().

Value

A list class object with all the functions of x.

as.list.sim_mats_list *Coerce a sim_mats_list class object into a list class object*

Description

Coerce a sim_mats_list class object into a list class object

Usage

```
## S3 method for class 'sim_mats_list'  
as.list(x, ...)
```

Arguments

x A sim_mats_list class object.
... Additional parameter passed to as.list().

Value

A list class object with all the functions of x.

as.list.snf_config *Coerce a snf_config class object into a list class object*

Description

Coerce a snf_config class object into a list class object

Usage

```
## S3 method for class 'snf_config'  
as.list(x, ...)
```

Arguments

x A snf_config class object.
... Additional parameter passed to as.list().

Value

A list class object with all the functions of x.

as.matrix.ari_matrix *Coerce a ari_matrix class object into a matrix class object*

Description

Coerce a ari_matrix class object into a matrix class object

Usage

```
## S3 method for class 'ari_matrix'  
as.matrix(x, ...)
```

Arguments

x A ari_matrix class object.
... Additional parameter passed to as.matrix().

Value

A matrix and array class object.

```
as.matrix.weights_matrix
```

Coerce a weights_matrix class object into a matrix class object

Description

Coerce a weights_matrix class object into a matrix class object

Usage

```
## S3 method for class 'weights_matrix'  
as.matrix(x, ...)
```

Arguments

x A weights_matrix class object.
... Additional parameter passed to as.matrix().

Value

A matrix and array class object.

```
assemble_data
```

Collapse a data frame and/or a data list into a single data frame

Description

Collapse a data frame and/or a data list into a single data frame

Usage

```
assemble_data(data, dl)
```

Arguments

data A data frame.
dl A nested list of input data from data_list().

Value

A class "data.frame" object containing all the features of the provided data frame and/or data list.

assoc_pval_heatmap *Heatmap of pairwise associations between features*

Description

Heatmap of pairwise associations between features

Usage

```
assoc_pval_heatmap(
  correlation_matrix,
  scale_diag = "max",
  cluster_rows = TRUE,
  cluster_columns = TRUE,
  show_row_names = TRUE,
  show_column_names = TRUE,
  show_heatmap_legend = FALSE,
  confounders = NULL,
  out_of_models = NULL,
  annotation_colours = NULL,
  labels_colour = NULL,
  split_by_domain = FALSE,
  dl = NULL,
  significance_stars = TRUE,
  slice_font_size = 8,
  ...
)
```

Arguments

correlation_matrix	Matrix containing all pairwise association p-values. The recommended way to obtain this matrix is through the <code>calc_assoc_pval</code> function.
scale_diag	Parameter that controls how the diagonals of the <code>correlation_matrix</code> are adjusted in the heatmap. For best viewing, this is set to "max", which will match the diagonals to whichever pairwise association has the highest p-value.
cluster_rows	Parameter for <code>ComplexHeatmap::Heatmap</code> . Will be ignored if <code>split_by_domain</code> is also provided.
cluster_columns	Parameter for <code>ComplexHeatmap::Heatmap</code> . Will be ignored if <code>split_by_domain</code> is also provided.
show_row_names	Parameter for <code>ComplexHeatmap::Heatmap</code> .
show_column_names	Parameter for <code>ComplexHeatmap::Heatmap</code> .
show_heatmap_legend	Parameter for <code>ComplexHeatmap::Heatmap</code> .

confounders	A named list where the elements are columns in the correlation_matrix and the names are the corresponding display names.
out_of_models	Like confounders, but a named list of out of model measures (who are also present as columns in the correlation_matrix).
annotation_colours	Named list of heatmap annotations and their colours.
labels_colour	Vector of colours to use for the columns and rows of the heatmap.
split_by_domain	Visually slice the heatmap based on feature domains.
dl	A nested list of input data from data_list().
significance_stars	If TRUE (default), plots significance stars on heatmap cells
slice_font_size	Font size for domain separating labels.
...	Additional parameters passed into ComplexHeatmap::Heatmap.

Value

Returns a heatmap (class "Heatmap" from package ComplexHeatmap) that displays the pairwise associations between features from the provided correlation_matrix.

Examples

```
#data_list <- data_list(
#   list(income, "household_income", "demographics", "ordinal"),
#   list(pubertal, "pubertal_status", "demographics", "continuous"),
#   list(fav_colour, "favourite_colour", "demographics", "categorical"),
#   list(anxiety, "anxiety", "behaviour", "ordinal"),
#   list(depress, "depressed", "behaviour", "ordinal"),
#   uid = "unique_id"
#)
#
#assoc_pval_matrix <- calc_assoc_pval_matrix(data_list)
#ap_heatmap <- assoc_pval_heatmap(assoc_pval_matrix)
```

as_ari_matrix

Convert an object to an ARI matrix

Description

This function coerces non-ari_matrix class objects into ari_matrix class objects.

Usage

```
as_ari_matrix(x)
```

Arguments

x The object to convert into a weights matrix.

Value

An `ari_matrix` class object.

as_data_list *Convert an object to a data list*

Description

This function coerces non-`data_list` class objects into `data_list` class objects.

Usage

```
as_data_list(x)
```

Arguments

x The object to convert into a data list.

Value

A `data_list` class object.

as_settings_df *Convert an object to a settings data frame*

Description

This function coerces non-`settings_df` class objects into `settings_df` class objects.

Usage

```
as_settings_df(x)
```

Arguments

x The object to convert into a data list.

Value

A `settings_df` class object.

as_sim_mats_list	<i>Convert an object to a similarity matrix list</i>
------------------	--

Description

This function converts non-sim_mats_list class objects into sim_mats_list class objects.

Usage

```
as_sim_mats_list(x)
```

Arguments

x The object to convert into a sim_mats_list. Must be a list of square matrices with identical column and row names.

Value

A sim_mats_list class object.

as_snf_config	<i>Convert an object to a snf config</i>
---------------	--

Description

This function coerces non-snf_config class objects into snf_config class objects.

Usage

```
as_snf_config(x)
```

Arguments

x The object to convert into a snf config.

Value

A snf_config class object.

as_weights_matrix	<i>Convert an object to a weights matrix</i>
-------------------	--

Description

This function converts non-weights_matrix objects into weights_matrix class objects.

Usage

```
as_weights_matrix(x)
```

Arguments

x The object to convert into a data list.

Value

A weights_matrix class object.

auto_plot	<i>Automatically plot features across clusters</i>
-----------	--

Description

Given a single row of a solutions data frame and data provided through a data list, this function will return a series of bar and/or jitter plots based on feature types.

Usage

```
auto_plot(  
  sol_df_row = NULL,  
  dl = NULL,  
  cluster_df = NULL,  
  return_plots = TRUE,  
  save = NULL,  
  jitter_width = 6,  
  jitter_height = 6,  
  bar_width = 6,  
  bar_height = 6,  
  verbose = FALSE  
)
```

Arguments

sol_df_row	A single row of a solutions data frame.
dl	A data list containing data to plot.
cluster_df	Directly provide a cluster_df rather than a solutions matrix. Useful if plotting data from label propagated results.
return_plots	If TRUE, the function will return a list of plots. If FALSE, the function will instead return the full data frame used for plotting.
save	If a string is provided, plots will be saved and this string will be used to prefix plot names.
jitter_width	Width of jitter plots if save is specified.
jitter_height	Height of jitter plots if save is specified.
bar_width	Width of bar plots if save is specified.
bar_height	Height of bar plots if save is specified.
verbose	If TRUE, output progress to console.

Value

By default, returns a list of plots (class "gg", "ggplot") with one plot for every feature in the provided data list and/or target list. If return_plots is FALSE, will instead return a single "data.frame" object containing every provided feature for every observation in long format.

bar_plot	<i>Bar plot separating a feature by cluster</i>
----------	---

Description

Bar plot separating a feature by cluster

Usage

```
bar_plot(df, feature)
```

Arguments

df	A data.frame containing cluster column and the feature to plot.
feature	The feature to plot.

Value

A bar plot (class "gg", "ggplot") showing the distribution of a feature across clusters.

batch_snf	<i>Run variations of SNF</i>
-----------	------------------------------

Description

This is the core function of the `metasnf` package. Using the information stored in a `settings_df` (see `?settings_df`) and a data list (see `?data_list`), run repeated complete SNF pipelines to generate a broad space of post-SNF cluster solutions.

Usage

```
batch_snf(dl, sc, processes = 1, return_sim_mats = FALSE, sim_mats_dir = NULL)
```

Arguments

<code>dl</code>	A nested list of input data from <code>data_list()</code> .
<code>sc</code>	An <code>snf_config</code> class object which stores all sets of hyperparameters used to transform data in <code>dl</code> into a cluster solutions. See <code>?settings_df</code> or https://branchlab.github.io/metasnf/arti for more details.
<code>processes</code>	Specify number of processes used to complete SNF iterations <ul style="list-style-type: none"> • 1 (default) Sequential processing: function will iterate through the <code>settings_df</code> one row at a time with a for loop. This option will not make use of multiple CPU cores, but will show a progress bar. • 2 or higher: Parallel processing will use the <code>future.apply::future_apply</code> to distribute the SNF iterations across the specified number of CPU cores. If higher than the number of available cores, a warning will be raised and the maximum number of cores will be used. • max: All available cores will be used.
<code>return_sim_mats</code>	If TRUE, function will return a list where the first element is the solutions data frame and the second element is a list of similarity matrices for each row in the <code>sol_df</code> . Default FALSE.
<code>sim_mats_dir</code>	If specified, this directory will be used to save all generated similarity matrices.

Value

By default, returns a solutions data frame (class "data.frame"), a a data frame containing one row for every row of the provided settings matrix, all the original columns of that settings data frame, and new columns containing the assigned cluster of each observation from the cluster solution derived by that row's settings. If `return_sim_mats` is TRUE, the function will instead return a list containing the solutions data frame as well as a list of the final similarity matrices (class "matrix") generated by SNF for each row of the settings data frame. If `suppress_clustering` is TRUE, the solutions data frame will not be returned in the output.

Examples

```

input_dl <- data_list(
  list(gender_df, "gender", "demographics", "categorical"),
  list(diagnosis_df, "diagnosis", "clinical", "categorical"),
  uid = "patient_id"
)

sc <- snf_config(input_dl, n_solutions = 3)

# A solutions data frame without similarity matrices:
sol_df <- batch_snf(input_dl, sc)

# A solutions data frame with similarity matrices:
sol_df <- batch_snf(input_dl, sc, return_sim_mats = TRUE)
sim_mats_list(sol_df)

```

batch_snf_subsamples *Run SNF clustering pipeline on a list of subsampled data lists*

Description

Run SNF clustering pipeline on a list of subsampled data lists

Usage

```

batch_snf_subsamples(
  dl_subsamples,
  sc,
  processes = 1,
  return_sim_mats = FALSE,
  sim_mats_dir = NULL
)

```

Arguments

- | | |
|---------------|--|
| dl_subsamples | A list of subsampled data lists. This object is generated by the function <code>batch_snf_subsamples()</code> . |
| sc | An <code>snf_config</code> class object which stores all sets of hyperparameters used to transform data in <code>dl</code> into a cluster solutions. See <code>?settings_df</code> or https://branchlab.github.io/metasnfn/artic for more details. |
| processes | Specify number of processes used to complete SNF iterations <ul style="list-style-type: none"> • 1 (default) Sequential processing: function will iterate through the <code>settings_df</code> one row at a time with a for loop. This option will not make use of multiple CPU cores, but will show a progress bar. • 2 or higher: Parallel processing will use the <code>future.apply::future_apply</code> to distribute the SNF iterations across the specified number of CPU cores. If higher than the number of available cores, a warning will be raised and the maximum number of cores will be used. |

- max: All available cores will be used.
- return_sim_mats If TRUE, function will return a list where the first element is the solutions data frame and the second element is a list of similarity matrices for each row in the sol_df. Default FALSE.
- sim_mats_dir If specified, this directory will be used to save all generated similarity matrices.

Value

By default, returns a one-element list: `cluster_solutions`, which is itself a list of cluster solution data frames corresponding to each of the provided data list subsamples. Setting the parameters `return_sim_mats` and `return_solutions` to TRUE will turn the result of the function to a three-element list containing the corresponding solutions data frames and final fused similarity matrices of those cluster solutions, should you require these objects for your own stability calculations.

Examples

```
my_dl <- data_list(
  list(subc_v, "subcortical_volume", "neuroimaging", "continuous"),
  list(income, "household_income", "demographics", "continuous"),
  list(pubertal, "pubertal_status", "demographics", "continuous"),
  uid = "unique_id"
)

sc <- snf_config(my_dl, n_solutions = 5, max_k = 40)

my_dl_subsamples <- subsample_dl(
  my_dl,
  n_subsamples = 20,
  subsample_fraction = 0.85
)

batch_subsample_results <- batch_snf_subsamples(
  my_dl_subsamples,
  sc
)
```

cache_a_complete_example_ext_sol_df

Cached example extended solutions data frame

Description

An extended solutions data frame used as a cached example in the "a_complete_example.Rmd" vignette.

Usage

```
cache_a_complete_example_ext_sol_df
```

Format

cache_a_complete_example_ext_sol_df:
Contains 20 cluster solutions, 87 observations, and p-values for 336 features.

Source

This data came from the metasnf package.

cache_a_complete_example_lp_ext_sol_df
Cached example extended solutions data frame

Description

An extended solutions data frame used as a cached example in the "a_complete_example.Rmd" vignette.

Usage

cache_a_complete_example_lp_ext_sol_df

Format

cache_a_complete_example_lp_ext_sol_df:
Contains 5 cluster solutions, 74 observations, and p-values for 2 features.

Source

This data comes from the metasnf package.

cache_a_complete_example_sol_df
Cached example solutions data frame

Description

An solutions data frame used as a cached example in the "a_complete_example.Rmd" vignette.

Usage

cache_a_complete_example_sol_df

Format

cache_a_complete_example_sol_df:
A solutions data frame with 20 cluster solutions and 87 observations.

Source

This data came from the metasnf package.

```
calculate_coclustering
```

Calculate co-clustering data

Description

Calculate co-clustering data

Usage

```
calculate_coclustering(subsample_solutions, sol_df, verbose = FALSE)
```

Arguments

subsample_solutions	A list of containing cluster solutions from distinct subsamples of the data. This object is generated by the function <code>batch_snf_subsamples()</code> . These solutions should correspond to the ones in the solutions data frame.
sol_df	A solutions data frame. This object is generated by the function <code>batch_snf()</code> . The solutions in the solutions data frame should correspond to those in the subsample solutions.
verbose	If TRUE, output time remaining estimates to console.

Value

A list containing the following components:

- `cocluster_dfs`: A list of data frames, one per cluster solution, that shows the number of times that every pair of observations in the original cluster solution occurred in the same subsample, the number of times that every pair clustered together in a subsample, and the corresponding fraction of times that every pair clustered together in a subsample.
- `cocluster_ss_mats`: The number of times every pair of observations occurred in the same subsample, formatted as a pairwise matrix.
- `cocluster_sc_mats`: The number of times every pair of observations occurred in the same cluster, formatted as a pairwise matrix.
- `cocluster_cf_mats`: The fraction of times every pair of observations occurred in the same cluster, formatted as a pairwise matrix.
- `cocluster_summary`: Specifically among pairs of observations that clustered together in the original full cluster solution, what fraction of those pairs remained clustered together throughout the subsample solutions. This information is formatted as a data frame with one row per cluster solution.

Examples

```

my_dl <- data_list(
  list(subc_v, "subcortical_volume", "neuroimaging", "continuous"),
  list(income, "household_income", "demographics", "continuous"),
  list(pubertal, "pubertal_status", "demographics", "continuous"),
  uid = "unique_id"
)

sc <- snf_config(my_dl, n_solutions = 5, max_k = 40)

sol_df <- batch_snf(my_dl, sc)

my_dl_subsamples <- subsample_dl(
  my_dl,
  n_subsamples = 20,
  subsample_fraction = 0.85
)

batch_subsample_results <- batch_snf_subsamples(
  my_dl_subsamples,
  sc
)

coclustering_results <- calculate_coclustering(
  batch_subsample_results,
  sol_df,
  verbose = TRUE
)

```

calc_aris

Construct an ARI matrix storing inter-solution similarities

Description

This function constructs an `ari_matrix` class object from a `solutions_df` class object. The ARI matrix stores pairwise adjusted Rand indices for all cluster solutions as well as a numeric order for the solutions data frame based on the hierarchical clustering of the ARI matrix.

Usage

```

calc_aris(
  sol_df,
  processes = 1,
  verbose = FALSE,
  dist_method = "euclidean",
  hclust_method = "complete"
)

```


Arguments

sol_df	Solutions data frame containing cluster solutions to calculate pairwise ARIs for.
processes	Specify number of processes used to complete calculations <ul style="list-style-type: none"> • 1 (default) Sequential processing • 2 or higher: Parallel processing will use the <code>future.apply::future_apply</code> to distribute the calculations across the specified number of CPU cores. If higher than the number of available cores, a warning will be raised and the maximum number of cores will be used. • max: All available cores will be used. Note that no progress indicator is available during multi-core processing.
verbose	If TRUE, output progress to console.
dist_method	Distance method to use when calculating sorting order to of the matrix. Argument is directly passed into <code>stats::dist</code> . Options include "euclidean", "maximum", "manhattan", "canberra", "binary", or "minkowski".
hclust_method	Agglomerative method to use when calculating sorting order by <code>stats::hclust</code> . Options include "ward.D", "ward.D2", "single", "complete", "average", "mcquitty", "median", or "centroid".

Value

om_aris ARIs between clustering solutions of an solutions data frame

Examples

```
dl <- data_list(
  list(subc_v, "subcortical_volume", "neuroimaging", "continuous"),
  list(pubertal, "pubertal_status", "demographics", "continuous"),
  uid = "unique_id"
)

sc <- snf_config(dl, n_solutions = 3)
sol_df <- batch_snf(dl, sc)
calc_aris(sol_df)
```

calc_assoc_pval_matrix

Calculate p-values for all pairwise associations of features in a data list

Description

Calculate p-values for all pairwise associations of features in a data list

Usage

```
calc_assoc_pval_matrix(dl, verbose = FALSE, cat_test = "chi_squared")
```

Arguments

dl	A nested list of input data from <code>data_list()</code> .
verbose	If TRUE, output progress to the console.
cat_test	String indicating which statistical test will be used to associate cluster with a categorical feature. Options are "chi_squared" for the Chi-squared test and "fisher_exact" for Fisher's exact test.

Value

A "matrix" class object containing pairwise association p-values between the features in the provided data list.

Examples

```
data_list <- data_list(  
  list(income, "household_income", "demographics", "ordinal"),  
  list(pubertal, "pubertal_status", "demographics", "continuous"),  
  list(anxiety, "anxiety", "behaviour", "ordinal"),  
  list(depress, "depressed", "behaviour", "ordinal"),  
  uid = "unique_id"  
)  
  
assoc_pval_matrix <- calc_assoc_pval_matrix(data_list)
```

`calc_nmis`*Calculate feature NMIs for a data list and a solutions data frame*

Description

Normalized mutual information scores can be used to indirectly measure how important a feature may have been in producing a cluster solution. This function will calculate the normalized mutual information between cluster solutions in a solutions data frame as well as cluster solutions created by including only a single feature from a provided data list, but otherwise using all the same hyper-parameters as specified in the original SNF config. Note that NMIs can be calculated between two cluster solutions regardless of what features were actually used to create those cluster solutions. For example, a feature that was not involved in producing a particular cluster solution may still have a high NMI with that cluster solution (typically because it was highly correlated with a different feature that was used).

Usage

```
calc_nmis(  
  dl,  
  sol_df,  
  transpose = TRUE,  
  ignore_inclusions = TRUE,  
  processes = 1  
)
```

Arguments

<code>dl</code>	A nested list of input data from <code>data_list()</code> .
<code>sol_df</code>	Result of <code>batch_snf</code> storing cluster solutions and the settings that were used to generate them. Use the same value as was used in the original call to <code>batch_snf()</code> .
<code>transpose</code>	If TRUE, will transpose the output data frame.
<code>ignore_inclusions</code>	If TRUE, will ignore the inclusion columns in the solutions data frame and calculate NMIs for all features. If FALSE, will give NAs for features that were dropped on a given <code>settings_df</code> row.
<code>processes</code>	Specify number of processes used to complete SNF iterations <ul style="list-style-type: none"> • 1 (default) Sequential processing: function will iterate through the <code>settings_df</code> one row at a time with a for loop. This option will not make use of multiple CPU cores, but will show a progress bar. • 2 or higher: Parallel processing will use the <code>future.apply::future_apply</code> to distribute the SNF iterations across the specified number of CPU cores. If higher than the number of available cores, a warning will be raised and the maximum number of cores will be used. • max: All available cores will be used.

Value

A "data.frame" class object containing one row for every feature in the provided data list and one column for every solution in the provided solutions data frame. Populated values show the calculated NMI score for each feature-solution combination.

Examples

```
input_dl <- data_list(
  list(gender_df, "gender", "demographics", "categorical"),
  list(diagnosis_df, "diagnosis", "clinical", "categorical"),
  uid = "patient_id"
)

sc <- snf_config(input_dl, n_solutions = 2)

sol_df <- batch_snf(input_dl, sc)

calc_nmis(input_dl, sol_df)
```

`cancer_diagnosis_df` *Mock diagnosis data*

Description

This is the same data as `diagnosis_df`, with renamed features and columns.

Usage

```
cancer_diagnosis_df
```

Format

```
cancer_diagnosis_df:
```

A data frame with 200 rows and 2 columns:

patient_id Random three-digit number uniquely identifying the patient

diagnosis Mock cancer diagnosis feature (1, 2, or 3)

Source

This data came from the SNFtool package, with slight modifications.

```
cell_significance_fn
```

Place significance stars on ComplexHeatmap cells

Description

This is an internal function meant to be used to by the `assoc_pval_heatmap` function.

Usage

```
cell_significance_fn(data)
```

Arguments

`data` The matrix containing the cells to base the significance stars on.

Value

`cell_fn` Another function that is well-formatted for usage as the `cell_fun` argument in `ComplexHeatmap::Heatmap`.

`check_dataless_annotations`

Helper function to stop annotation building when no data was provided

Description

Helper function to stop annotation building when no data was provided

Usage

```
check_dataless_annotations(annotation_requests, data)
```

Arguments

`annotation_requests`

A list of requested annotations

`data`

A data frame with data to build annotations

Value

Does not return any value. This function just raises an error when annotations are requested without any provided data for a heatmap.

`check_hm_dependencies` *Check for ComplexHeatmap and circlize dependencies*

Description

Check for ComplexHeatmap and circlize dependencies

Usage

```
check_hm_dependencies()
```

Value

Does not return any value. This function just checks that the ComplexHeatmap and circlize packages are installed.

`check_similarity_matrices`*Check validity of similarity matrices*

Description

Check to see if similarity matrices in a list have the following properties:

1. The maximum value in the entire matrix is 0.5
2. Every value in the diagonal is 0.5

Usage

```
check_similarity_matrices(similarity_matrices)
```

Arguments`similarity_matrices`

A list of similarity matrices

Value

`valid_matrices` Boolean indicating if properties are met by all similarity matrices

`clust_fns`*Built-in clustering algorithms*

Description

These functions can be used when building a `metasnf` clustering functions list. Each function converts a similarity matrix (matrix class object) to a cluster solution (numeric vector). Note that these functions (or custom clustering functions) cannot accept number of clusters as a parameter; this value must be built into the function itself if necessary.

Usage

```
spectral_eigen(similarity_matrix)
```

```
spectral_rot(similarity_matrix)
```

```
spectral_eigen_classic(similarity_matrix)
```

```
spectral_rot_classic(similarity_matrix)
```

```
spectral_two(similarity_matrix)
```

```
spectral_three(similarity_matrix)
spectral_four(similarity_matrix)
spectral_five(similarity_matrix)
spectral_six(similarity_matrix)
spectral_seven(similarity_matrix)
spectral_eight(similarity_matrix)
spectral_nine(similarity_matrix)
spectral_ten(similarity_matrix)
```

Arguments

`similarity_matrix`
A similarity matrix.

Details

- `spectral_eigen`: Spectral clustering where the number of clusters is based on the eigen-gap heuristic
- `spectral_rot`: Spectral clustering where the number of clusters is based on the rotation-cost heuristic
- `spectral_(C)`: Spectral clustering for a C-cluster solution.

Value

`solution_data` A vector of cluster assignments

`clust_fns_list` *Build a clustering algorithms list*

Description

This function can be used to specify custom clustering algorithms to apply to the final similarity matrices produced by each run of the `batch_snf` function.

Usage

```
clust_fns_list(clust_fns = NULL, use_default_clust_fns = FALSE)
```

Arguments

`clust_fns` A list of named clustering functions

`use_default_clust_fns`
 If TRUE, prepend the base clustering algorithms (`spectral_eigen` and `spectral_rot`, which apply spectral clustering and use the eigen-gap and rotation cost heuristics respectively for determining the number of clusters in the graph) to `clust_fns`.

Value

A list of clustering algorithm functions that can be passed into the `batch_snf` and `generate_settings_list` functions.

Examples

```
# Using just the base clustering algorithms -----
# This will just contain spectral_eigen and spectral_rot
cfl <- clust_fns_list(use_default_clust_fns = TRUE)

# Adding algorithms provided by the package -----
# This will contain the base clustering algorithms (spectral_eigen,
# spectral_rot) as well as two pre-defined spectral clustering functions
# that force the number of clusters to be two or five
cfl <- clust_fns_list(
  clust_fns = list(
    "two_cluster_spectral" = spectral_two,
    "five_cluster_spectral" = spectral_five
  )
)

# Adding your own algorithms -----
# This will contain the base and user-provided clustering algorithms
my_clustering_algorithm <- function(similarity_matrix) {
  # your code that converts similarity matrix to clusters here...
}

# Suppress the base algorithms-----
# This will contain only user-provided clustering algorithms
cfl <- clust_fns_list(
  clust_fns = list(
    "two_cluster_spectral" = spectral_two,
    "five_cluster_spectral" = spectral_five
  )
)
```


Description

This function creates a density plot that shows, for all pairs of observations that originally clustered together, the distribution of the the fractions that those pairs clustered together across subsampled data.

Usage

```
cocluster_density(cocluster_df)
```

Arguments

`cocluster_df` A data frame containing co-clustering data for a single cluster solution. This object is generated by the `calculate_coclustering` function.

Value

Density plot (class "gg", "ggplot") of the distribution of co-clustering across pairs and subsamples of the data.

Examples

```
my_dl <- data_list(  
  list(subc_v, "subcortical_volume", "neuroimaging", "continuous"),  
  list(income, "household_income", "demographics", "continuous"),  
  list(pubertal, "pubertal_status", "demographics", "continuous"),  
  uid = "unique_id"  
)  
  
sc <- snf_config(my_dl, n_solutions = 5, max_k = 40)  
  
sol_df <- batch_snf(my_dl, sc)  
  
my_dl_subsamples <- subsample_dl(  
  my_dl,  
  n_subsamples = 20,  
  subsample_fraction = 0.85  
)  
  
batch_subsample_results <- batch_snf_subsamples(  
  my_dl_subsamples,  
  sc  
)  
  
coclustering_results <- calculate_coclustering(  
  batch_subsample_results,  
  sol_df,  
  verbose = TRUE  
)  
  
cocluster_dfs <- coclustering_results$"cocluster_dfs"
```

```
cocluster_density(cocluster_dfs[[1]])
```

cocluster_heatmap	<i>Heatmap of observation co-clustering across resampled data</i>
-------------------	---

Description

Create a heatmap that shows the distribution of observation co-clustering across resampled data.

Usage

```
cocluster_heatmap(
  cocluster_df,
  cluster_rows = TRUE,
  cluster_columns = TRUE,
  show_row_names = FALSE,
  show_column_names = FALSE,
  dl = NULL,
  data = NULL,
  left_bar = NULL,
  right_bar = NULL,
  top_bar = NULL,
  bottom_bar = NULL,
  left_hm = NULL,
  right_hm = NULL,
  top_hm = NULL,
  bottom_hm = NULL,
  annotation_colours = NULL,
  min_colour = NULL,
  max_colour = NULL,
  ...
)
```

Arguments

cocluster_df	A data frame containing co-clustering data for a single cluster solution. This object is generated by the <code>calculate_coclustering</code> function.
cluster_rows	Argument passed to <code>ComplexHeatmap::Heatmap()</code> .
cluster_columns	Argument passed to <code>ComplexHeatmap::Heatmap()</code> .
show_row_names	Argument passed to <code>ComplexHeatmap::Heatmap()</code> .
show_column_names	Argument passed to <code>ComplexHeatmap::Heatmap()</code> .
dl	See <code>?similarity_matrix_heatmap</code> .
data	See <code>?similarity_matrix_heatmap</code> .

left_bar	See ?similarity_matrix_heatmap.
right_bar	See ?similarity_matrix_heatmap.
top_bar	See ?similarity_matrix_heatmap.
bottom_bar	See ?similarity_matrix_heatmap.
left_hm	See ?similarity_matrix_heatmap.
right_hm	See ?similarity_matrix_heatmap.
top_hm	See ?similarity_matrix_heatmap.
bottom_hm	See ?similarity_matrix_heatmap.
annotation_colours	See ?similarity_matrix_heatmap.
min_colour	See ?similarity_matrix_heatmap.
max_colour	See ?similarity_matrix_heatmap.
...	Arguments passed to ComplexHeatmap::Heatmap().

Value

Heatmap (class "Heatmap" from ComplexHeatmap) object showing the distribution of observation co-clustering across resampled data.

Examples

```
my_dl <- data_list(
  list(subc_v, "subcortical_volume", "neuroimaging", "continuous"),
  list(income, "household_income", "demographics", "continuous"),
  list(pubertal, "pubertal_status", "demographics", "continuous"),
  uid = "unique_id"
)

sc <- snf_config(my_dl, n_solutions = 5, max_k = 40)

sol_df <- batch_snf(my_dl, sc)

my_dl_subsamples <- subsample_dl(
  my_dl,
  n_subsamples = 20,
  subsample_fraction = 0.85
)

batch_subsample_results <- batch_snf_subsamples(
  my_dl_subsamples,
  sc
)

coclustering_results <- calculate_coclustering(
  batch_subsample_results,
  sol_df,
  verbose = TRUE
)
```

```

cocluster_dfs <- coclustering_results$"cocluster_dfs"

cocluster_heatmap(
  cocluster_dfs[[1]],
  dl = my_dl,
  top_hm = list(
    "Income" = "household_income",
    "Pubertal Status" = "pubertal_status"
  ),
  annotation_colours = list(
    "Pubertal Status" = colour_scale(
      c(1, 4),
      min_colour = "black",
      max_colour = "purple"
    ),
    "Income" = colour_scale(
      c(0, 4),
      min_colour = "black",
      max_colour = "red"
    )
  )
)

```

 colour_scale

Return a colour ramp for a given vector

Description

Given a numeric vector and min and max colour values, return a colour ramp that assigns a colour to each element in the vector. This function is a wrapper for `circlize::colorRamp2`.

Usage

```
colour_scale(data, min_colour, max_colour)
```

Arguments

data	Vector of numeric values.
min_colour	Minimum colour value.
max_colour	Maximum colour value.

Value

A "function" class object that can build a circlize-style colour ramp.

cort_sa	<i>Mock ABCD cortical surface area data</i>
---------	---

Description

Like the mock data frame "abcd_cort_sa", but with "unique_id" as the "uid".

Usage

```
cort_sa
```

Format

cort_sa:

A data frame with 188 rows and 152 columns:

unique_id The unique identifier of the ABCD dataset

... Cortical surface areas of various ROIs (mm², I think)

Source

Though this data is no longer "real" ABCD data, the reference for using ABCD as a data source is below:

Data used in the preparation of this article were obtained from the Adolescent Brain Cognitive Development (ABCD) Study (<https://abcdstudy.org>), held in the NIMH Data Archive (NDA). This is a multisite, longitudinal study designed to recruit more than 10,000 children age 9-10 and follow them over 10 years into early adulthood. The ABCD Study® is supported by the National Institutes of Health and additional federal partners under award numbers U01DA041048, U01DA050989, U01DA051016, U01DA041022, U01DA051018, U01DA051037, U01DA050987, U01DA041174, U01DA041106, U01DA041117, U01DA041028, U01DA041134, U01DA050988, U01DA051039, U01DA041156, U01DA041025, U01DA041120, U01DA051038, U01DA041148, U01DA041093, U01DA041089, U24DA041123, U24DA041147. A full list of supporters is available at <https://abcdstudy.org/federal-partners.html>. A listing of participating sites and a complete listing of the study investigators can be found at https://abcdstudy.org/consortium_members/. ABCD consortium investigators designed and implemented the study and/or provided data but did not necessarily participate in the analysis or writing of this report. This manuscript reflects the views of the authors and may not reflect the opinions or views of the NIH or ABCD consortium investigators.

cort_t	<i>Mock ABCD cortical thickness data</i>
--------	--

Description

Like the mock data frame "abcd_cort_t", but with "unique_id" as the "uid".

Usage

```
cort_t
```

Format

```
cort_t:
```

A data frame with 188 rows and 152 columns:

unique_id The unique identifier of the ABCD dataset

... Cortical thicknesses of various ROIs (mm³, I think)

Source

Though this data is no longer "real" ABCD data, the reference for using ABCD as a data source is below:

Data used in the preparation of this article were obtained from the Adolescent Brain Cognitive Development (ABCD) Study (<https://abcdstudy.org>), held in the NIMH Data Archive (NDA). This is a multisite, longitudinal study designed to recruit more than 10,000 children age 9-10 and follow them over 10 years into early adulthood. The ABCD Study® is supported by the National Institutes of Health and additional federal partners under award numbers U01DA041048, U01DA050989, U01DA051016, U01DA041022, U01DA051018, U01DA051037, U01DA050987, U01DA041174, U01DA041106, U01DA041117, U01DA041028, U01DA041134, U01DA050988, U01DA051039, U01DA041156, U01DA041025, U01DA041120, U01DA051038, U01DA041148, U01DA041093, U01DA041089, U24DA041123, U24DA041147. A full list of supporters is available at <https://abcdstudy.org/federal-partners.html>. A listing of participating sites and a complete listing of the study investigators can be found at https://abcdstudy.org/consortium_members/. ABCD consortium investigators designed and implemented the study and/or provided data but did not necessarily participate in the analysis or writing of this report. This manuscript reflects the views of the authors and may not reflect the opinions or views of the NIH or ABCD consortium investigators.

```
data_list
```

```
Build a data_list class object
```

Description

`data_list()` constructs a data list object which inherits from classes `data_list` and `list`. This object is the primary way in which features to be used along the `metasnf` clustering pipeline are stored. The data list is fundamentally a 2-level nested list object where each inner list contains a data frame and associated metadata for that data frame. The metadata includes the name of the data frame, the 'domain' of that data frame (the broader source of information that the input data frame is capturing, determined by user's domain knowledge), and the type of feature stored in the data frame (continuous, discrete, ordinal, categorical, or mixed).

Usage

```
data_list(..., uid)
```

Arguments

... Any number of lists formatted as (df, "df_name", "df_domain", "df_type") and/or any number of lists of lists formatted as (df, "df_name", "df_domain", "df_type").

uid (character) the name of the uid column currently used data. data frame.

Examples

```
heart_rate_df <- data.frame(
  patient_id = c("1", "2", "3"),
  var1 = c(0.04, 0.1, 0.3),
  var2 = c(30, 2, 0.3)
)

personality_test_df <- data.frame(
  patient_id = c("1", "2", "3"),
  var3 = c(900, 1990, 373),
  var4 = c(509, 2209, 83)
)

survey_response_df <- data.frame(
  patient_id = c("1", "2", "3"),
  var5 = c(1, 3, 3),
  var6 = c(2, 3, 3)
)

city_df <- data.frame(
  patient_id = c("1", "2", "3"),
  var7 = c("toronto", "montreal", "vancouver")
)

# Explicitly (Name each nested list element):
dl <- data_list(
  list(
    data = heart_rate_df,
    name = "heart_rate",
    domain = "clinical",
    type = "continuous"
  ),
  list(
    data = personality_test_df,
    name = "personality_test",
    domain = "surveys",
    type = "continuous"
  ),
  list(
    data = survey_response_df,
    name = "survey_response",
    domain = "surveys",
    type = "ordinal"
  ),
  list(
    data = city_df,
```

```

      name = "city",
      domain = "location",
      type = "categorical"
    ),
    uid = "patient_id"
  )

# Compact loading
dl <- data_list(
  list(heart_rate_df, "heart_rate", "clinical", "continuous"),
  list(personality_test_df, "personality_test", "surveys", "continuous"),
  list(survey_response_df, "survey_response", "surveys", "ordinal"),
  list(city_df, "city", "location", "categorical"),
  uid = "patient_id"
)

# Printing data list summaries
summary(dl)

# Alternative loading: providing a single list of lists
list_of_lists <- list(
  list(heart_rate_df, "data1", "domain1", "continuous"),
  list(personality_test_df, "data2", "domain2", "continuous")
)

dl <- data_list(
  list_of_lists,
  uid = "patient_id"
)

```

depress

Mock ABCD depression data

Description

Like the mock data frame "abcd_depress", but with "unique_id" as the "uid".

Usage

```
depress
```

Format

depress:

A data frame with 275 rows and 2 columns:

unique_id The unique identifier of the ABCD dataset

cbcl_depress_r Ordinal value of impairment on CBCL anxiety, either 0 (no impairment), 1 (borderline clinical), or 2 (clinically impaired)

Source

Though this data is no longer "real" ABCD data, the reference for using ABCD as a data source is below:

Data used in the preparation of this article were obtained from the Adolescent Brain Cognitive Development (ABCD) Study (<https://abcdstudy.org>), held in the NIMH Data Archive (NDA). This is a multisite, longitudinal study designed to recruit more than 10,000 children age 9-10 and follow them over 10 years into early adulthood. The ABCD Study® is supported by the National Institutes of Health and additional federal partners under award numbers U01DA041048, U01DA050989, U01DA051016, U01DA041022, U01DA051018, U01DA051037, U01DA050987, U01DA041174, U01DA041106, U01DA041117, U01DA041028, U01DA041134, U01DA050988, U01DA051039, U01DA041156, U01DA041025, U01DA041120, U01DA051038, U01DA041148, U01DA041093, U01DA041089, U24DA041123, U24DA041147. A full list of supporters is available at <https://abcdstudy.org/federal-partners.html>. A listing of participating sites and a complete listing of the study investigators can be found at https://abcdstudy.org/consortium_members/. ABCD consortium investigators designed and implemented the study and/or provided data but did not necessarily participate in the analysis or writing of this report. This manuscript reflects the views of the authors and may not reflect the opinions or views of the NIH or ABCD consortium investigators.

diagnosis_df

Mock diagnosis data

Description

This is the same data as cancer_diagnosis_df, with renamed features and columns.

Usage

```
diagnosis_df
```

Format

```
diagnosis_df:
```

A data frame with 200 rows and 2 columns:

patient_id Random three-digit number uniquely identifying the patient

diagnosis Mock diagnosis feature

Source

This data came from the SNFtool package, with slight modifications.

dist_fns	<i>Built-in distance functions</i>
----------	------------------------------------

Description

These functions can be used when building a `metasnf` distance functions list. Each function converts a data frame into to a distance matrix.

Usage

```
euclidean_distance(df, weights_row)
```

```
gower_distance(df, weights_row)
```

```
sn_euclidean_distance(df, weights_row)
```

```
sew_euclidean_distance(df, weights_row)
```

```
hamming_distance(df, weights_row)
```

Arguments

<code>df</code>	Data frame containing at least 1 data column
<code>weights_row</code>	Single-row data frame where the column names contain the column names in <code>df</code> and the row contains the corresponding <code>weights_row</code> .

Details

Functions that work for numeric data:

- `euclidean_distance`: typical Euclidean distance
- `sn_euclidean_distance`: Data frame is first standardized and normalized before typical Euclidean distance is applied
- `siw_euclidean_distance`: Squared (including weights) Euclidean distance, where the weights are also squared
- `sew_euclidean_distance`: Squared (excluding weights) Euclidean distance, where the weights are not also squared

Functions that work for binary data:

- `hamming_distance`: typical Hamming distance

Functions that work for any type of data:

- `gower_distance`: Gower distance (`cluster::daisy`)

Value

A matrix class object containing pairwise distances.

dist_fns_list	<i>Build a distance metrics list</i>
---------------	--------------------------------------

Description

The distance metrics list object (inherits classes `dist_fns_list` and `list`) is a list that stores R functions which can convert a data frame of features into a matrix of pairwise distances. The list is a nested one, where the first layer of the list can hold up to 5 items (one for each of the `metasnf` recognized feature types, continuous, discrete, ordinal, categorical, and mixed), and the second layer can hold an arbitrary number of distance functions for each of those types.

Usage

```
dist_fns_list(
  cnt_dist_fns = NULL,
  dsc_dist_fns = NULL,
  ord_dist_fns = NULL,
  cat_dist_fns = NULL,
  mix_dist_fns = NULL,
  automatic_standard_normalize = FALSE,
  use_default_dist_fns = FALSE
)
```

Arguments

<code>cnt_dist_fns</code>	A named list of continuous distance metric functions.
<code>dsc_dist_fns</code>	A named list of discrete distance metric functions.
<code>ord_dist_fns</code>	A named list of ordinal distance metric functions.
<code>cat_dist_fns</code>	A named list of categorical distance metric functions.
<code>mix_dist_fns</code>	A named list of mixed distance metric functions.
<code>automatic_standard_normalize</code>	If TRUE, will automatically use standard normalization prior to calculation of any numeric distances. This parameter overrides all other distance functions list-related parameters.
<code>use_default_dist_fns</code>	If TRUE, prepend the base distance metrics (euclidean distance for continuous, discrete, and ordinal data and gower distance for categorical and mixed data) to the resulting distance metrics list.

Details

Call `?distance_metrics` to see all distance metric functions provided in `metasnf`.

Value

A distance metrics list object.

Examples

```

# Using just the base distance metrics -----
dist_fns_list <- dist_fns_list()

# Adding your own metrics -----
# This will contain only the and user-provided distance function:
cubed_euclidean <- function(df, weights_row) {
  # (your code that converts a data frame to a distance metric here...)
  weights <- diag(weights_row, nrow = length(weights_row))
  weighted_df <- as.matrix(df) %*% weights
  distance_matrix <- weighted_df |>
    stats::dist(method = "euclidean") |>
    as.matrix()
  distance_matrix <- distance_matrix^3
  return(distance_matrix)
}

dist_fns_list <- dist_fns_list(
  cnt_dist_fns = list(
    "my_cubed_euclidean" = cubed_euclidean
  )
)

# Using default base metrics-----
# Call ?distance_metrics to see all distance metric functions provided in
# metasnf. The code below will contain a mix of user-provided and built-in
# distance metric functions.
dist_fns_list <- dist_fns_list(
  cnt_dist_fns = list(
    "my_distance_metric" = cubed_euclidean
  ),
  dsc_dist_fns = list(
    "my_distance_metric" = cubed_euclidean
  ),
  ord_dist_fns = list(
    "my_distance_metric" = cubed_euclidean
  ),
  cat_dist_fns = list(
    "my_distance_metric" = gower_distance
  ),
  mix_dist_fns = list(
    "my_distance_metric" = gower_distance
  ),
  use_default_dist_fns = TRUE
)

```

Description

This function enables manipulating a `data_list` class object with `lapply` syntax without removing that object's `data_list` class attribute. The function will only preserve this attribute if the result of the `apply` call has a valid data list structure.

Usage

```
dlapply(X, FUN, ...)
```

Arguments

<code>X</code>	A <code>data_list</code> class object.
<code>FUN</code>	The function to be applied to each data list component.
<code>...</code>	Optional arguments to <code>FUN</code> .

Value

If `FUN` applied to each component of `X` yields a valid data list, a data list. Otherwise, a list.

Examples

```
# Convert all UID values to lowercase
dl <- data_list(
  list(abcd_income, "income", "demographics", "discrete"),
  list(abcd_colour, "colour", "likes", "categorical"),
  uid = "patient"
)

dl_lower <- dlapply(
  dl,
  function(x) {
    x$"data$"uid <- tolower(x$"data$"uid)
    return(x)
  }
)
```

```
dplyr_row_slice.ext_solutions_df
```

Function to extend dplyr to extended solutions data frame objects

Description

Function to extend `dplyr` to extended solutions data frame objects

Usage

```
dplyr_row_slice.ext_solutions_df(data, i, ...)
```

Arguments

<code>data</code>	An extended solutions data frame.
<code>i</code>	A vector of row indices.
<code>...</code>	Additional arguments.

Value

Row sliced object with appropriately preserved attributes.

```
dplyr_row_slice.solutions_df
```

Function to extend dplyr to solutions data frame objects

Description

Function to extend dplyr to solutions data frame objects

Usage

```
dplyr_row_slice.solutions_df(data, i, ...)
```

Arguments

<code>data</code>	A solutions data frame.
<code>i</code>	A vector of row indices.
<code>...</code>	Additional arguments.

Value

Row sliced object with appropriately preserved attributes.

```
esm_manhattan_plot
```

Manhattan plot of feature-cluster association p-values

Description

Manhattan plot of feature-cluster association p-values

Usage

```
esm_manhattan_plot(
  esm,
  neg_log_pval_thresh = 5,
  threshold = NULL,
  point_size = 5,
  jitter_width = 0.1,
  jitter_height = 0.1,
  text_size = 15,
  plot_title = NULL,
  hide_x_labels = FALSE,
  bonferroni_line = FALSE
)
```

Arguments

esm	Extended solutions data frame storing associations between features and cluster assignments. See <code>?extend_solutions</code> .
neg_log_pval_thresh	Threshold for negative log p-values.
threshold	P-value threshold to plot dashed line at.
point_size	Size of points in the plot.
jitter_width	Width of jitter.
jitter_height	Height of jitter.
text_size	Size of text in the plot.
plot_title	Title of the plot.
hide_x_labels	If TRUE, hides x-axis labels.
bonferroni_line	If TRUE, plots a dashed black line at the Bonferroni-corrected equivalent of the p-value threshold.

Value

A Manhattan plot (class "gg", "ggplot") showing the association p-values of features against each solution in the provided solutions data frame.

Examples

```
# full_dl <- data_list(
#   list(subc_v, "subcortical_volume", "neuroimaging", "continuous"),
#   list(income, "household_income", "demographics", "continuous"),
#   list(pubertal, "pubertal_status", "demographics", "continuous"),
#   list(anxiety, "anxiety", "behaviour", "ordinal"),
#   list(depress, "depressed", "behaviour", "ordinal"),
#   uid = "unique_id"
# )
#
```

```

# dl <- full_dl[1:3]
# target_dl <- full_dl[4:5]
#
# set.seed(42)
# sc <- snf_config(
#   dl = dl,
#   n_solutions = 20,
#   min_k = 20,
#   max_k = 50
# )
#
# sol_df <- batch_snf(dl, sc)
#
# ext_sol_df <- extend_solutions(
#   sol_df,
#   dl = dl,
#   target = target_dl,
#   min_pval = 1e-10 # p-values below 1e-10 will be thresholded to 1e-10
# )
#
# esm_manhattan <- esm_manhattan_plot(
#   ext_sol_df[1:5, ],
#   neg_log_pval_thresh = 5,
#   threshold = 0.05,
#   point_size = 3,
#   jitter_width = 0.1,
#   jitter_height = 0.1,
#   plot_title = "Feature-Solution Associations",
#   text_size = 14,
#   bonferroni_line = TRUE
# )

```

estimate_nclust_given_graph

Estimate number of clusters for a similarity matrix

Description

Calculate eigengap and rotation-cost estimates of the number of clusters to use when clustering a similarity matrix. This function was adapted from `SNFtool::estimateClustersGivenGraph`, but scales up the Laplacian operator prior to eigenvalue calculations to minimize the risk of floating point-related errors.

Usage

```
estimate_nclust_given_graph(W, NUMC = 2:10)
```


Arguments

W	Similarity matrix to calculate number of clusters for.
NUMC	Range of cluster counts to consider among when picking best number of clusters.

Value

A list containing the top two eigengap and rotation-cost estimates for the number of clusters in a given similarity matrix.

Examples

```
input_dl <- data_list(
  list(gender_df, "gender", "demographics", "categorical"),
  list(diagnosis_df, "diagnosis", "clinical", "categorical"),
  uid = "patient_id"
)

sc <- snf_config(input_dl, n_solutions = 1)
sol_df <- batch_snf(input_dl, sc, return_sim_mats = TRUE)
sim_mat <- sim_mats_list(sol_df)[[1]]
estimate_nclust_given_graph(sim_mat)
```

 expression_df

Modification of SNFtool mock data frame "Data1"

Description

Modification of SNFtool mock data frame "Data1"

Usage

```
expression_df
```

Format

expression_df:

A data frame with 200 rows and 3 columns:

gene_1_expression Mock gene expression feature

gene_2_expression Mock gene expression feature

patient_id Random three-digit number uniquely identifying the patient

Source

This data came from the SNFtool package, with slight modifications.

extend_solutions	<i>Extend a solutions data frame to include outcome evaluations</i>
------------------	---

Description

Extend a solutions data frame to include outcome evaluations

Usage

```
extend_solutions(
  sol_df,
  target_dl = NULL,
  dl = NULL,
  cat_test = "chi_squared",
  min_pval = 1e-10,
  processes = 1,
  verbose = FALSE
)
```

Arguments

sol_df	Result of batch_snf storing cluster solutions and the settings that were used to generate them.
target_dl	A data list with features to calculate p-values for. Features in the target list will be included during p-value summary measure calculations.
dl	A data list with features to calculate p-values for, but that should not be incorporated into p-value summary measure columns (i.e., min/mean/max p-value columns).
cat_test	String indicating which statistical test will be used to associate cluster with a categorical feature. Options are "chi_squared" for the Chi-squared test and "fisher_exact" for Fisher's exact test.
min_pval	If assigned a value, any p-value less than this will be replaced with this value.
processes	The number of processes to use for parallelization. Progress is only reported for sequential processing (processes = 1).
verbose	If TRUE, output progress to console.

Value

An extended solutions data frame (ext_sol_df class object) that contains p-value columns for each outcome in the provided data lists

Examples

```
## Not run:
input_dl <- data_list(
  list(gender_df, "gender", "demographics", "categorical"),
  list(diagnosis_df, "diagnosis", "clinical", "categorical"),
  uid = "patient_id"
)

sc <- snf_config(input_dl, n_solutions = 2)

sol_df <- batch_snf(input_dl, sc)

ext_sol_df <- extend_solutions(sol_df, input_dl)

## End(Not run)
```

fav_colour	<i>Mock ABCD "colour" data</i>
------------	--------------------------------

Description

Like the mock data frame "abcd_colour", but with "unique_id" as the "uid".

Usage

```
fav_colour
```

Format

fav_colour:

A data frame with 275 rows and 2 columns:

unique_id The unique identifier of the ABCD dataset

colour Categorical transformation of cbcl_depress.

Source

Though this data is no longer "real" ABCD data, the reference for using ABCD as a data source is below:

Data used in the preparation of this article were obtained from the Adolescent Brain Cognitive Development (ABCD) Study (<https://abcdstudy.org>), held in the NIMH Data Archive (NDA). This is a multisite, longitudinal study designed to recruit more than 10,000 children age 9-10 and follow them over 10 years into early adulthood. The ABCD Study® is supported by the National Institutes of Health and additional federal partners under award numbers U01DA041048, U01DA050989, U01DA051016, U01DA041022, U01DA051018, U01DA051037, U01DA050987, U01DA041174, U01DA041106, U01DA041117, U01DA041028, U01DA041134, U01DA050988, U01DA051039, U01DA041156, U01DA041025, U01DA041120, U01DA051038, U01DA041148, U01DA041093,

U01DA041089, U24DA041123, U24DA041147. A full list of supporters is available at <https://abcdstudy.org/federal-partners.html>. A listing of participating sites and a complete listing of the study investigators can be found at https://abcdstudy.org/consortium_members/. ABCD consortium investigators designed and implemented the study and/or provided data but did not necessarily participate in the analysis or writing of this report. This manuscript reflects the views of the authors and may not reflect the opinions or views of the NIH or ABCD consortium investigators.

gender_df	<i>Mock gender data</i>
-----------	-------------------------

Description

Mock gender data

Usage

gender_df

Format

gender_df:

A data frame with 200 rows and 2 columns:

patient_id Random three-digit number uniquely identifying the patient

gender_df Mock gene methylation feature

Source

This data came from the SNFtool package, with slight modifications.

get_complete_uids	<i>Pull complete-data UIDs from a list of data frames</i>
-------------------	---

Description

This function identifies all observations within a list of data frames that have no missing data across all data frames. This function is useful when constructing data lists of distinct feature sets from the same sample of observations. As `data_list()` strips away observations with any missing data, distinct sets of observations may be generated by building a data list from the same group of observations over different sets of features. Reducing the pool of observations to only those with complete UIDs first will avoid downstream generation of data lists of differing sizes.

Usage

```
get_complete_uids(list_of_dfs, uid)
```

Arguments

`list_of_dfs` List of data frames.
`uid` Name of column across data frames containing UIDs

Value

A character vector of the UIDs of observations that have complete data across the provided list of data frames.

Examples

```
complete_uids <- get_complete_uids(
  list(income, pubertal, anxiety, depress),
  uid = "unique_id"
)

income <- income[income$"unique_id" %in% complete_uids, ]
pubertal <- pubertal[pubertal$"unique_id" %in% complete_uids, ]
anxiety <- anxiety[anxiety$"unique_id" %in% complete_uids, ]
depress <- depress[depress$"unique_id" %in% complete_uids, ]

input_dl <- data_list(
  list(income, "income", "demographics", "ordinal"),
  list(pubertal, "pubertal", "demographics", "continuous"),
  uid = "unique_id"
)

target_dl <- data_list(
  list(anxiety, "anxiety", "behaviour", "ordinal"),
  list(depress, "depressed", "behaviour", "ordinal"),
  uid = "unique_id"
)
```

`get_heatmap_order` *Return the row or column ordering present in a heatmap*

Description

Return the row or column ordering present in a heatmap

Usage

```
get_heatmap_order(heatmap, type = "rows")
```

Arguments

`heatmap` A heatmap object to collect ordering from.
`type` The type of ordering to return. Either "rows" or "columns".

Value

A numeric vector of the ordering used within the provided ComplexHeatmap "Heatmap" object.

get_matrix_order	<i>Return the hierarchical clustering order of a matrix</i>
------------------	---

Description

Return the hierarchical clustering order of a matrix

Usage

```
get_matrix_order(matrix, dist_method = "euclidean", hclust_method = "complete")
```

Arguments

matrix	Matrix to cluster.
dist_method	Distance method to use when calculating sorting order to of the matrix. Argument is directly passed into stats::dist. Options include "euclidean", "maximum", "manhattan", "canberra", "binary", or "minkowski".
hclust_method	Agglomerative method to use when calculating sorting order by stats::hclust. Options include "ward.D", "ward.D2", "single", "complete", "average", "mcquitty", "median", or "centroid".

Value

A numeric vector of the ordering derived by the specified hierarchical clustering method applied to the provided matrix.

Examples

```
# dl <- data_list(
#   list(subc_v, "subcortical_volume", "neuroimaging", "continuous"),
#   list(income, "household_income", "demographics", "continuous"),
#   list(pubertal, "pubertal_status", "demographics", "continuous"),
#   list(anxiety, "anxiety", "behaviour", "ordinal"),
#   list(depress, "depressed", "behaviour", "ordinal"),
#   uid = "unique_id"
# )
#
# sc <- snf_config(
#   dl = dl,
#   n_solutions = 20,
#   min_k = 20,
#   max_k = 50
# )
#
# sol_df <- batch_snf(dl, sc)
```

```
#
# ext_sol_df <- extend_solutions(
#   sol_df,
#   dl = dl,
#   min_pval = 1e-10 # p-values below 1e-10 will be thresholded to 1e-10
# )
#
# # Calculate pairwise similarities between cluster solutions
# sol_aris <- calc_aris(sol_df)
#
# # Extract hierarchical clustering order of the cluster solutions
# meta_cluster_order <- get_matrix_order(sol_aris)
```

get_pvals

Get p-values from an extended solutions data frame

Description

This function can be used to neatly format the p-values associated with an extended solutions data frame. It can also calculate the negative logs of those p-values to make it easier to interpret large-scale differences.

Usage

```
get_pvals(ext_sol_df, negative_log = FALSE, keep_summaries = TRUE)
```

Arguments

`ext_sol_df` The output of `extend_solutions`. A data frame that contains at least one p-value column ending in `"_pval"`.

`negative_log` If TRUE, will replace p-values with negative log p-values.

`keep_summaries` If FALSE, will remove the mean, min, and max p-value.

Value

A "data.frame" class object Of only the p-value related columns of the provided `ext_sol_df`.

```
get_representative_solutions
```

Extract representative solutions from a matrix of ARIs

Description

Following clustering with `batch_snf`, a matrix of pairwise ARIs that show how related each cluster solution is to each other can be generated by the `calc_aris` function. Partitioning of the ARI matrix can be done by visual inspection of `meta_cluster_heatmap()` results or by `shiny_annotator`. Given the indices of meta cluster boundaries, this function will return a single representative solution from each meta cluster based on maximum average ARI to all other solutions within that meta cluster.

Usage

```
get_representative_solutions(aris, sol_df, filter_fn = NULL)
```

Arguments

<code>aris</code>	Matrix of adjusted rand indices from <code>calc_aris()</code>
<code>sol_df</code>	Output of <code>batch_snf</code> containing cluster solutions.
<code>filter_fn</code>	Optional function to filter the meta-cluster by prior to maximum average ARI determination. This can be useful if you are explicitly trying to select a solution that meets a certain condition, such as only picking from the 4 cluster solutions within a meta cluster. An example valid function could be <code>fn <- function(x) x[x\$"nclust" == 4,]</code> .

Value

The provided solutions data frame reduced to just one row per meta cluster defined by the split vector.

Examples

```
# dl <- data_list(
#   list(subc_v, "subcortical_volume", "neuroimaging", "continuous"),
#   list(income, "household_income", "demographics", "continuous"),
#   list(pubertal, "pubertal_status", "demographics", "continuous"),
#   list(anxiety, "anxiety", "behaviour", "ordinal"),
#   list(depress, "depressed", "behaviour", "ordinal"),
#   uid = "unique_id"
# )
#
# sc <- snf_config(
#   dl = dl,
#   n_solutions = 20,
#   min_k = 20,
#   max_k = 50
```



```

# )
#
# sol_df <- batch_snf(dl, sc)
#
# ext_sol_df <- extend_solutions(
#   sol_df,
#   dl = dl,
#   min_pval = 1e-10 # p-values below 1e-10 will be thresholded to 1e-10
# )
#
# # Calculate pairwise similarities between cluster solutions
# sol_aris <- calc_aris(sol_df)
#
# # Extract hierarchical clustering order of the cluster solutions
# meta_cluster_order <- get_matrix_order(sol_aris)
#
# # Identify meta cluster boundaries with shiny app or trial and error
# # ari_hm <- meta_cluster_heatmap(sol_aris, order = meta_cluster_order)
# # shiny_annotator(ari_hm)
#
# # Result of meta cluster examination
# split_vec <- c(2, 5, 12, 17)
#
# ext_sol_df <- label_meta_clusters(ext_sol_df, split_vec, meta_cluster_order)
#
# # Extracting representative solutions from each defined meta cluster
# rep_solutions <- get_representative_solutions(sol_aris, ext_sol_df)

```

income

Mock ABCD income data

Description

Like the mock data frame "abcd_h_income", but with "unique_id" as the "uid".

Like the mock data frame "abcd_cort_sa", but with "unique_id" as the "uid".

Usage

```
income
```

```
income
```

Format

income:

A data frame with 300 rows and 2 columns:

unique_id The unique identifier of the ABCD dataset

household_income Household income in 3 category levels (low = 1, medium = 2, high = 3)

income:

A data frame with 300 rows and 2 columns:

unique_id The unique identifier of the ABCD dataset

household_income Household income in 3 category levels (low = 1, medium = 2, high = 3)

Source

Though this data is no longer "real" ABCD data, the reference for using ABCD as a data source is below:

Data used in the preparation of this article were obtained from the Adolescent Brain Cognitive Development (ABCD) Study (<https://abcdstudy.org>), held in the NIMH Data Archive (NDA). This is a multisite, longitudinal study designed to recruit more than 10,000 children age 9-10 and follow them over 10 years into early adulthood. The ABCD Study® is supported by the National Institutes of Health and additional federal partners under award numbers U01DA041048, U01DA050989, U01DA051016, U01DA041022, U01DA051018, U01DA051037, U01DA050987, U01DA041174, U01DA041106, U01DA041117, U01DA041028, U01DA041134, U01DA050988, U01DA051039, U01DA041156, U01DA041025, U01DA041120, U01DA051038, U01DA041148, U01DA041093, U01DA041089, U24DA041123, U24DA041147. A full list of supporters is available at <https://abcdstudy.org/federal-partners.html>. A listing of participating sites and a complete listing of the study investigators can be found at https://abcdstudy.org/consortium_members/. ABCD consortium investigators designed and implemented the study and/or provided data but did not necessarily participate in the analysis or writing of this report. This manuscript reflects the views of the authors and may not reflect the opinions or views of the NIH or ABCD consortium investigators.

Though this data is no longer "real" ABCD data, the reference for using ABCD as a data source is below:

Data used in the preparation of this article were obtained from the Adolescent Brain Cognitive Development (ABCD) Study (<https://abcdstudy.org>), held in the NIMH Data Archive (NDA). This is a multisite, longitudinal study designed to recruit more than 10,000 children age 9-10 and follow them over 10 years into early adulthood. The ABCD Study® is supported by the National Institutes of Health and additional federal partners under award numbers U01DA041048, U01DA050989, U01DA051016, U01DA041022, U01DA051018, U01DA051037, U01DA050987, U01DA041174, U01DA041106, U01DA041117, U01DA041028, U01DA041134, U01DA050988, U01DA051039, U01DA041156, U01DA041025, U01DA041120, U01DA051038, U01DA041148, U01DA041093, U01DA041089, U24DA041123, U24DA041147. A full list of supporters is available at <https://abcdstudy.org/federal-partners.html>. A listing of participating sites and a complete listing of the study investigators can be found at https://abcdstudy.org/consortium_members/. ABCD consortium investigators designed and implemented the study and/or provided data but did not necessarily participate in the analysis or writing of this report. This manuscript reflects the views of the authors and may not reflect the opinions or views of the NIH or ABCD consortium investigators.

is_data_list

Test if the object is a data list

Description

Given an object, returns TRUE if that object inherits from the data_list class.

Usage

```
is_data_list(x)
```

Arguments

x An object.

Value

TRUE if the object inherits from the `data_list` class.

jitter_plot	<i>Jitter plot separating a feature by cluster</i>
-------------	--

Description

Jitter plot separating a feature by cluster

Usage

```
jitter_plot(df, feature)
```

Arguments

df A data.frame containing cluster column and the feature to plot.
feature The feature to plot.

Value

A jitter+violin plot (class "gg", "ggplot") showing the distribution of a feature across clusters.

label_meta_clusters	<i>Assign meta cluster labels to rows of a solutions data frame or extended solutions data frame</i>
---------------------	--

Description

Given a solutions data frame or extended solutions data frame class object and a numeric vector indicating which rows correspond to which meta clusters, assigns meta clustering information to the "meta_clusters" attribute of the data frame.

Usage

```
label_meta_clusters(sol_df, split_vector, order = NULL)
```

Arguments

sol_df	A solutions data frame or extended solutions data frame to assign meta clusters to.
split_vector	A numeric vector indicating which rows of sol_df should be the split points for meta cluster labeling.
order	An optional numeric vector indicating how the solutions data frame should be reordered prior to meta cluster labeling. This vector can be obtained by running <code>get_matrix_order()</code> on an ARI matrix, which itself can be obtained by calling <code>calc_aris()</code> on a solutions data frame.

Value

A solutions data frame with a populated "meta_clusters" attribute.

Examples

```
#dl <- data_list(
#   list(cort_sa, "cortical_surface_area", "neuroimaging", "continuous"),
#   list(subc_v, "subcortical_volume", "neuroimaging", "continuous"),
#   list(income, "household_income", "demographics", "continuous"),
#   list(pubertal, "pubertal_status", "demographics", "continuous"),
#   uid = "unique_id"
#)
#
#set.seed(42)
#my_sc <- snf_config(
#   dl = dl,
#   n_solutions = 20,
#   min_k = 20,
#   max_k = 50
#)
#
#sol_df <- batch_snf(dl, my_sc)
#
#sol_df
#
#sol_aris <- calc_aris(sol_df)
#
#meta_cluster_order <- get_matrix_order(sol_aris)
#
## `split_vec` found by iteratively plotting ari_hm or by ?shiny_annotator()
#split_vec <- c(6, 10, 16)
#ari_hm <- meta_cluster_heatmap(
#   sol_aris,
#   order = meta_cluster_order,
#   split_vector = split_vec
#)
#
#mc_sol_df <- label_meta_clusters(
#   sol_df,
#   order = meta_cluster_order,
```

```
#   split_vector = split_vec
#)
#
#mc_sol_df
```

label_propagate	<i>Label propagate cluster solutions to non-clustered observations</i>
-----------------	--

Description

Given a solutions data frame containing clustered observations and a data list containing those clustered observations as well as additional to-be-clustered observations, this function will re-run SNF to generate a similarity matrix of all observations and use the label propagation algorithm to assigned predicted clusters to the non-clustered observations.

Usage

```
label_propagate(partial_sol_df, full_dl, verbose = FALSE)
```

Arguments

`partial_sol_df` A solutions data frame derived from the training set.
`full_dl` A data list containing observations from both the training and testing sets.
`verbose` If TRUE, output progress to console.

Value

A data frame with one row per observation containing a column for UIDs, a column for whether the observation was in the train (original) or test (held out) set, and one column per row of the solutions data frame indicating the original and propagated clusters.

Examples

```
## Function to identify observations with complete data
#uids_with_complete_obs <- get_complete_uids(
#   list(subc_v, income, pubertal, anxiety, depress),
#   uid = "unique_id"
#)
#
## Dataframe assigning 80% of observations to train and 20% to test
#train_test_split <- train_test_assign(
#   train_frac = 0.8,
#   uids = uids_with_complete_obs
#)
#
## Pulling the training and testing observations specifically
#train_obs <- train_test_split$"train"
#test_obs <- train_test_split$"test"
#
```

```

## Partition a training set
#train_subc_v <- subc_v[subc_v$"unique_id" %in% train_obs, ]
#train_income <- income[income$"unique_id" %in% train_obs, ]
#train_pubertal <- pubertal[pubertal$"unique_id" %in% train_obs, ]
#train_anxiety <- anxiety[anxiety$"unique_id" %in% train_obs, ]
#train_depress <- depress[depress$"unique_id" %in% train_obs, ]
#
## Partition a test set
#test_subc_v <- subc_v[subc_v$"unique_id" %in% test_obs, ]
#test_income <- income[income$"unique_id" %in% test_obs, ]
#test_pubertal <- pubertal[pubertal$"unique_id" %in% test_obs, ]
#test_anxiety <- anxiety[anxiety$"unique_id" %in% test_obs, ]
#test_depress <- depress[depress$"unique_id" %in% test_obs, ]
#
## Find cluster solutions in the training set
#train_dl <- data_list(
#  list(train_subc_v, "subc_v", "neuroimaging", "continuous"),
#  list(train_income, "household_income", "demographics", "continuous"),
#  list(train_pubertal, "pubertal_status", "demographics", "continuous"),
#  uid = "unique_id"
#)
#
## We'll pick a solution that has good separation over our target features
#train_target_dl <- data_list(
#  list(train_anxiety, "anxiety", "behaviour", "ordinal"),
#  list(train_depress, "depressed", "behaviour", "ordinal"),
#  uid = "unique_id"
#)
#
#sc <- snf_config(
#  train_dl,
#  n_solutions = 5,
#  min_k = 10,
#  max_k = 30
#)
#
#train_sol_df <- batch_snf(
#  train_dl,
#  sc,
#  return_sim_mats = TRUE
#)
#
#ext_sol_df <- extend_solutions(
#  train_sol_df,
#  train_target_dl
#)
#
## Determining solution with the lowest minimum p-value
#lowest_min_pval <- min(ext_sol_df$"min_pval")
#which(ext_sol_df$"min_pval" == lowest_min_pval)
#top_row <- ext_sol_df[1, ]
#
## Propagate that solution to the observations in the test set

```

```

## data list below has both training and testing observations
#full_dl <- data_list(
#  list(subc_v, "subc_v", "neuroimaging", "continuous"),
#  list(income, "household_income", "demographics", "continuous"),
#  list(pubertal, "pubertal_status", "demographics", "continuous"),
#  uid = "unique_id"
#)
#
## Use the solutions data frame from the training observations and the data list
## from the training and testing observations to propagate labels to the test observations
#propagated_labels <- label_propagate(top_row, full_dl)
#
#propagated_labels_all <- label_propagate(ext_sol_df, full_dl)
#
#head(propagated_labels_all)
#tail(propagated_labels_all)

```

linear_adjust

Linearly correct data list by features with unwanted signal

Description

Given a data list to correct and another data list of categorical features to linearly adjust for, corrects the first data list based on the residuals of the linear model relating the numeric features in the first data list to the unwanted signal features in the second data list.

Usage

```
linear_adjust(dl, unwanted_signal_list, sig_digs = NULL)
```

Arguments

dl A nested list of input data from `data_list()`.

unwanted_signal_list A data list of categorical features that should have their mean differences removed in the first data list.

sig_digs Number of significant digits to round the residuals to.

Value

A data list ("list") in which each data component has been converted to contain residuals off of the linear model built against the features in the `unwanted_signal_list`.

Examples

```

has_tutor <- sample(c(1, 0), size = 9, replace = TRUE)
math_score <- 70 + 30 * has_tutor + rnorm(9, mean = 0, sd = 5)

math_df <- data.frame(uid = paste0("id_", 1:9), math = math_score)
tutor_df <- data.frame(uid = paste0("id_", 1:9), tutor = has_tutor)

dl <- data_list(
  list(math_df, "math_score", "school", "continuous"),
  uid = "uid"
)

adjustment_dl <- data_list(
  list(tutor_df, "tutoring", "school", "categorical"),
  uid = "uid"
)

adjusted_dl <- linear_adjust(dl, adjustment_dl)

adjusted_dl[[1]]$"data$"math"

# Equivalent to:
as.numeric(resid(lm(math_score ~ has_tutor)))

```

mc_manhattan_plot *Manhattan plot of feature-meta cluster association p-values*

Description

Given a data frame of representative meta cluster solutions (see `get_representative_solutions()`), returns a Manhattan plot for showing feature separation across all features in provided data/target lists.

Usage

```

mc_manhattan_plot(
  ext_sol_df,
  dl = NULL,
  target_dl = NULL,
  variable_order = NULL,
  neg_log_pval_thresh = 5,
  threshold = NULL,
  point_size = 5,
  text_size = 20,
  plot_title = NULL,
  xints = NULL,
  hide_x_labels = FALSE,
  domain_colours = NULL
)

```


Arguments

ext_sol_df	A sol_df that contains "_pval" columns containing the values to be plotted. This object is the output of extend_solutions().
dl	List of data frames containing data information.
target_dl	List of data frames containing target information.
variable_order	Order of features to be displayed in the plot.
neg_log_pval_thresh	Threshold for negative log p-values.
threshold	p-value threshold to plot horizontal dashed line at.
point_size	Size of points in the plot.
text_size	Size of text in the plot.
plot_title	Title of the plot.
xints	Either "outcomes" or a vector of numeric values to plot vertical lines at.
hide_x_labels	If TRUE, hides x-axis labels.
domain_colours	Named vector of colours for domains.

Value

A Manhattan plot (class "gg", "ggplot") showing the association p-values of features against each solution in the provided solutions data frame, stratified by meta cluster label.

Examples

```
# dl <- data_list(
#   list(subc_v, "subcortical_volume", "neuroimaging", "continuous"),
#   list(income, "household_income", "demographics", "continuous"),
#   list(pubertal, "pubertal_status", "demographics", "continuous"),
#   list(anxiety, "anxiety", "behaviour", "ordinal"),
#   list(depress, "depressed", "behaviour", "ordinal"),
#   uid = "unique_id"
# )
#
# sc <- snf_config(
#   dl = dl,
#   n_solutions = 20,
#   min_k = 20,
#   max_k = 50
# )
#
# sol_df <- batch_snf(dl, sc)
#
# ext_sol_df <- extend_solutions(
#   sol_df,
#   dl = dl,
#   min_pval = 1e-10 # p-values below 1e-10 will be thresholded to 1e-10
# )
#
```

```

## Calculate pairwise similarities between cluster solutions
# sol_aris <- calc_aris(sol_df)
#
## Extract hierarchical clustering order of the cluster solutions
# meta_cluster_order <- get_matrix_order(sol_aris)
#
## Identify meta cluster boundaries with shiny app or trial and error
## ari_hm <- meta_cluster_heatmap(sol_aris, order = meta_cluster_order)
## shiny_annotator(ari_hm)
#
## Result of meta cluster examination
# split_vec <- c(2, 5, 12, 17)
#
# ext_sol_df <- label_meta_clusters(ext_sol_df, split_vec, meta_cluster_order)
#
## Extracting representative solutions from each defined meta cluster
# rep_solutions <- get_representative_solutions(sol_aris, ext_sol_df)
#
# mc_manhattan <- mc_manhattan_plot(
#   rep_solutions,
#   dl = dl,
#   point_size = 3,
#   text_size = 12,
#   plot_title = "Feature-Meta Cluster Associations",
#   threshold = 0.05,
#   neg_log_pval_thresh = 5
# )
#
# mc_manhattan

```

merge.clust_fns_list *Merge clust_fns_list objects*

Description

Merge clust_fns_list objects

Usage

```

## S3 method for class 'clust_fns_list'
merge(x, y, ...)

```

Arguments

x	The first clust_fns_list object to merge.
y	The second clust_fns_list object to merge.
...	Additional arguments (not used).

Value

A new `clust_fns_list` object containing the merged clustering functions.

<code>merge.data_list</code>	<i>Merge observations between two compatible data lists</i>
------------------------------	---

Description

Join two data lists with the same components (data frames) but separate observations. To instead merge two data lists that have the same observations but different components, simply use `c()`.

Usage

```
## S3 method for class 'data_list'
merge(x, y, ...)
```

Arguments

<code>x</code>	The first data list to merge.
<code>y</code>	The second data list to merge.
<code>...</code>	Additional arguments passed into merge function.

Value

A data list ("list"-class object) containing the observations of both provided data lists.

<code>merge.dist_fns_list</code>	<i>Merge dist_fns_list objects</i>
----------------------------------	------------------------------------

Description

Merge `dist_fns_list` objects

Usage

```
## S3 method for class 'dist_fns_list'
merge(x, y, ...)
```

Arguments

<code>x</code>	The first <code>clust_fns_list</code> object to merge.
<code>y</code>	The second <code>clust_fns_list</code> object to merge.
<code>...</code>	Additional arguments (not used).

Value

A new `clust_fns_list` object containing the merged clustering functions.

merge.ext_solutions_df *Merge ext_solutions_df objects*

Description

Merge ext_solutions_df objects

Usage

```
## S3 method for class 'ext_solutions_df'  
merge(x, y, ...)
```

Arguments

x	The first ext_solutions_df object to merge.
y	The second ext_solutions_df object to merge.
...	Additional arguments (not used).

Value

Error message indicating that the merge function is not applicable to ext_solutions_df objects.

merge.settings_df *Merge settings_df objects*

Description

Merge settings_df objects

Usage

```
## S3 method for class 'settings_df'  
merge(x, y, ...)
```

Arguments

x	The first settings_df object to merge.
y	The second settings_df object to merge.
...	Additional arguments (not used).

Value

Error message indicating that the merge function is not applicable to settings_df objects.

merge.sim_mats_list *Merge sim_mats_list objects*

Description

Merge sim_mats_list objects

Usage

```
## S3 method for class 'sim_mats_list'  
merge(x, y, ...)
```

Arguments

x	The first sim_mats_list object to merge.
y	The second sim_mats_list object to merge.
...	Additional arguments (not used).

Value

A merged sim_mats_list object containing the similarity matrices from both input objects.

merge.snf_config *Merge method for SNF config objects*

Description

Merge method for SNF config objects

Usage

```
## S3 method for class 'snf_config'  
merge(x, y, reset_indices = TRUE, ...)
```

Arguments

x	SNF config to merge.
y	SNF config to merge.
reset_indices	If TRUE (default), re-labels the "solutions" indices in the config from 1 to the number of defined settings.
...	Additional arguments passed into merge function.

Value

An SNF config combining the rows of both prior configurations.

merge.solutions_df *Merge solutions_df objects*

Description

Merge solutions_df objects

Usage

```
## S3 method for class 'solutions_df'  
merge(x, y, ...)
```

Arguments

x The first solutions_df object to merge.
y The second solutions_df object to merge.
... Additional arguments (not used).

Value

Error message indicating that the merge function is not applicable to solutions_df objects.

merge.t_ext_solutions_df
 Merge t_ext_solutions_df objects

Description

Merge t_ext_solutions_df objects

Usage

```
## S3 method for class 't_ext_solutions_df'  
merge(x, y, ...)
```

Arguments

x The first t_ext_solutions_df object to merge.
y The second t_ext_solutions_df object to merge.
... Additional arguments (not used).

Value

Error message indicating that the merge function is not applicable to t_ext_solutions_df objects.

merge.t_solutions_df *Merge t_solutions_df objects*

Description

Merge t_solutions_df objects

Usage

```
## S3 method for class 't_solutions_df'  
merge(x, y, ...)
```

Arguments

x	The first t_solutions_df object to merge.
y	The second t_solutions_df object to merge.
...	Additional arguments (not used).

Value

Error message indicating that the merge function is not applicable to t_solutions_df objects.

merge.weights_matrix *Merge weights_matrix objects*

Description

Merge weights_matrix objects

Usage

```
## S3 method for class 'weights_matrix'  
merge(x, y, ...)
```

Arguments

x	The first weights_matrix object to merge.
y	The second weights_matrix object to merge.
...	Additional arguments (not used).

Value

Error message indicating that the merge function is not applicable to weights_matrix objects.

merge_df_list	<i>Merge list of data frames into a single data frame</i>
---------------	---

Description

This helper function combines all data frames in a single-level list into a single data frame.

Usage

```
merge_df_list(df_list, join = "inner", uid = "uid", no_na = FALSE)
```

Arguments

df_list	list of data frames.
join	String indicating if join should be "inner" or "full".
uid	Column name to join on. Default is "uid".
no_na	Whether to remove NA values from the merged data frame.

Value

Inner join of all data frames in list.

Examples

```
merge_df_list(list(income, pubertal), uid = "unique_id")
```

methylation_df	<i>Modification of SNFtool mock data frame "Data2"</i>
----------------	--

Description

Modification of SNFtool mock data frame "Data2"

Usage

```
methylation_df
```

Format

methylation_df:

A data frame with 200 rows and 3 columns:

gene_1_expression Mock gene methylation feature

gene_2_expression Mock gene methylation feature

patient_id Random three-digit number uniquely identifying the patient

Source

This data came from the SNFtool package, with slight modifications.

mock_ari_matrix *Mock example of an ari_matrix metasnfn object*

Description

An ari_matrix class object containing adjusted Rand indices (ARIs) between 20 cluster solutions. Used as an example of an ari_matrix metasnfn object.

Usage

```
mock_ari_matrix
```

Format

```
mock_ari_matrix:  
A 20 by 20 ARI matrix.
```

Source

This data comes from the metasnfn package.

mock_clust_fns_list *Mock example of a clust_fns_list metasnfn object*

Description

Mock example of a clust_fns_list metasnfn object

Usage

```
mock_clust_fns_list
```

Format

```
mock_clust_fns_list:  
A clust_fns_list object containing two clustering functions covering 2 and 5 five cluster solution versions of spectral clustering. Extracted from mock_snfn_config.
```

Source

This data comes from the metasnfn package.

mock_data_list *Mock example of a data_list metasn object*

Description

Mock example of a data_list metasn object

Usage

```
mock_data_list
```

Format

```
mock_data_list:
```

A data list containing 4 data frames with 100 observations each: - subcortical volume (30 features) - cortical surface area (151 features) - household income (1 feature) - pubertal status (1 feature)
Used as an example of an data_list metasn object.

Source

This data comes from the metasn package.

mock_dist_fns_list *Mock example of a dist_fns_list metasn object*

Description

Mock example of a dist_fns_list metasn object

Usage

```
mock_dist_fns_list
```

Format

```
mock_dist_fns_list:
```

A dist_fns_list object containing a variety of distance metrics. Extracted from mock_snf_config.

Source

This data comes from the metasn package.

mock_ext_solutions_df *Mock example of a ext_solutions_df metasnf object*

Description

An `ext_solutions_df` class object generated by extending the `mock_rep_solutions_df` object against `mock_data_list` as the target data list.

Usage

```
mock_ext_solutions_df
```

Format

```
mock_ext_solutions_df:  
Contains 20 cluster solutions.
```

Source

This data comes from the `metasnf` package.

mock_mc_solutions_df *Mock example of a mc_solutions_df metasnf object*

Description

Mock example of a `mc_solutions_df` metasnf object

Usage

```
mock_mc_solutions_df
```

Format

```
mock_mc_solutions_df:  
A meta cluster labeled solutions data frame derived from mock_solutions_df. Contains 20  
cluster solutions.
```

Source

This data comes from the `metasnf` package.

mock_rep_solutions_df *Mock example of a rep_solutions_df metasnf object*

Description

A solutions_df class object derived by filtering the mock_mc_solutions_df to its representative solutions.

Usage

```
mock_rep_solutions_df
```

Format

```
mock_rep_solutions_df:  
Contains 4 cluster solutions.
```

Source

This data comes from the metasnf package.

mock_settings_df *Mock example of a settings_df metasnf object*

Description

Mock example of a settings_df metasnf object

Usage

```
mock_settings_df
```

Format

```
mock_settings_df:  
Settings for 20 cluster solutions.
```

Source

This data comes from the metasnf package.

mock_snf_config	<i>Mock example of a snf_config metasnf object</i>
-----------------	--

Description

Mock example of a snf_config metasnf object

Usage

```
mock_snf_config
```

Format

```
mock_snf_config:
```

An SNF config containing hyperparameters and functions defined for generating 20 cluster solutions from a data list. The config has been specified to: - limit the k hyperparameter to 40 - make use of uniformly distributed random weights - randomly select between using spectral clustering where the number of clusters can be 2, 5, decided by the eigen-gap heuristic, or decided by the rotation cost heuristic - use Gower distance for categorical and mixed data, Euclidean distance for ordinal data, and randomly select from Euclidean distance or standard/normalized Euclidean distance for continuous and discrete data The config was built using the mock_data_list loaded into the namespace after calling library("metasnf"). Used as an example of an snf_config metasnf object.

Source

This data comes from the metasnf package.

mock_solutions_df	<i>Mock example of a solutions_df metasnf object</i>
-------------------	--

Description

Mock example of a solutions_df metasnf object

Usage

```
mock_solutions_df
```

Format

```
mock_solutions_df:
```

A solutions data frame containing 20 cluster solutions generated from mock_snf_config and mock_data_list. Used as an example of an solutions_df metasnf object.

Source

This data comes from the metasnf package.

mock_t_solutions_df *Mock example of a t_solutions_df metasnf object*

Description

Mock example of a t_solutions_df metasnf object

Usage

```
mock_t_solutions_df
```

Format

```
mock_t_solutions_df:
```

A transposed solutions data frame containing 20 cluster solutions generated from mock_solutions_df. Used as an example of a t_solutions_df metasnf object.

Source

This data comes from the metasnf package.

mock_weights_matrix *Mock example of a weights_matrix metasnf object*

Description

Mock example of a weights_matrix metasnf object

Usage

```
mock_weights_matrix
```

Format

```
mock_weights_matrix:
```

A weights_matrix class object containing 20 sets of weights for 183 features.

Source

This data comes from the metasnf package.

new_solutions_df *Constructor for solutions_df class object*

Description

Constructor for solutions_df class object

Usage

```
new_solutions_df(sol_df1)
```

Arguments

sol_df1 A solutions data frame-like object to be validated and converted into a solutions data frame.

Value

A solutions_df class object.

plot.ari_matrix *Heatmap of pairwise adjusted rand indices between solutions*

Description

Heatmap of pairwise adjusted rand indices between solutions

Usage

```
## S3 method for class 'ari_matrix'
plot(
  x,
  order = NULL,
  cluster_rows = FALSE,
  cluster_columns = FALSE,
  log_graph = FALSE,
  scale_diag = "none",
  min_colour = "#282828",
  max_colour = "firebrick2",
  col = circlize::colorRamp2(c(min(x), max(x)), c(min_colour, max_colour)),
  ...
)

meta_cluster_heatmap(
  x,
  order = NULL,
```

```

cluster_rows = FALSE,
cluster_columns = FALSE,
log_graph = FALSE,
scale_diag = "none",
min_colour = "#282828",
max_colour = "firebrick2",
col = circlize::colorRamp2(c(min(x), max(x)), c(min_colour, max_colour)),
...
)

```

Arguments

x	Matrix of adjusted rand indices from calc_aris()
order	Numeric vector containing row order of the heatmap.
cluster_rows	Whether rows should be clustered.
cluster_columns	Whether columns should be clustered.
log_graph	If TRUE, log transforms the graph.
scale_diag	Method of rescaling matrix diagonals. Can be "none" (don't change diagonals), "mean" (replace diagonals with average value of off-diagonals), or "zero" (replace diagonals with 0).
min_colour	Colour used for the lowest value in the heatmap.
max_colour	Colour used for the highest value in the heatmap.
col	Colour ramp to use for the heatmap.
...	Additional parameters passed to similarity_matrix_heatmap(), the function that this function wraps.

Value

Returns a heatmap (class "Heatmap" from package ComplexHeatmap) that displays the pairwise adjusted Rand indices (similarities) between the cluster solutions of the provided solutions data frame.

Examples

```

#dl <- data_list(
#   list(cort_sa, "cortical_surface_area", "neuroimaging", "continuous"),
#   list(subc_v, "subcortical_volume", "neuroimaging", "continuous"),
#   list(income, "household_income", "demographics", "continuous"),
#   list(pubertal, "pubertal_status", "demographics", "continuous"),
#   uid = "unique_id"
#)
#
#set.seed(42)
#my_sc <- snf_config(
#   dl = dl,
#   n_solutions = 20,
#   min_k = 20,

```



```

#   max_k = 50
#)
#
#sol_df <- batch_snf(dl, my_sc)
#
#sol_df
#
#sol_aris <- calc_aris(sol_df)
#
#meta_cluster_order <- get_matrix_order(sol_aris)
#
## `split_vec` found by iteratively plotting ari_hm or by ?shiny_annotator()
#split_vec <- c(6, 10, 16)
#ari_hm <- plot(
#   sol_aris,
#   order = meta_cluster_order,
#   split_vector = split_vec
#)

```

plot.data_list

Plot of feature values in a data list

Description

This plot, built on `ComplexHeatmap::Heatmap()`, visualizes the feature values in a data list as a continuous heatmap with observations along the columns and features along the rows.

Usage

```

## S3 method for class 'data_list'
plot(
  x,
  y = NULL,
  cluster_rows = TRUE,
  cluster_columns = TRUE,
  heatmap_legend_param = NULL,
  row_title = "Observation",
  column_title = "Feature",
  show_row_names = FALSE,
  ...
)

```

Arguments

x	A data_list object.
y	Optional argument to plot, not used in this method.
cluster_rows	Logical indicating whether to cluster the rows (observations).

cluster_columns Logical indicating whether to cluster the columns (features).
heatmap_legend_param A list of parameters for the heatmap legend.
row_title Title for the rows (observations).
column_title Title for the columns (features).
show_row_names Logical indicating whether to show row names.
... Additional arguments passed to `ComplexHeatmap::Heatmap()`.

Value

A heatmap visualization of feature values.

`plot.ext_solutions_df` *Plot of cluster assignments in an extended solutions data frame*

Description

This plot, built on `ComplexHeatmap::Heatmap()`, visualizes the cluster assignments in a solutions data frame as a categorical heatmap with observations along the columns and clusters along the rows.

Usage

```

## S3 method for class 'ext_solutions_df'
plot(
  x,
  y = NULL,
  cluster_rows = TRUE,
  cluster_columns = TRUE,
  show_row_names = TRUE,
  show_column_names = TRUE,
  heatmap_legend_param = NULL,
  row_title = "Solution",
  column_title = "Observation",
  ...
)

## S3 method for class 't_ext_solutions_df'
plot(x, ...)

```

Arguments

x An `ext_solutions_df` object.
y Optional argument to plot, not used in this method.

cluster_rows	If the value is a logical, it controls whether to make cluster on rows. The value can also be a <code>hclust</code> or a <code>dendrogram</code> which already contains clustering. Check https://jokergoo.github.io/ComplexHeatmap-reference/book/a-single-heatmap.html#clustering .
cluster_columns	Whether make cluster on columns? Same settings as <code>cluster_rows</code> .
show_row_names	Whether show row names.
show_column_names	Whether show column names.
heatmap_legend_param	A list contains parameters for the heatmap legends. See color_mapping_legend , ColorMapping-method for all available parameters.
row_title	Title on the row.
column_title	Title on the column.
...	Additional arguments passed to <code>ComplexHeatmap::Heatmap()</code> .

Value

A `ComplexHeatmap::Heatmap()` object visualization of cluster assignments.

plot.snf_config	<i>Heatmap for visualizing an SNF config</i>
-----------------	--

Description

Create a heatmap where each row corresponds to a different set of hyperparameters in an SNF config object. Numeric parameters are scaled normalized and non-numeric parameters are added as heatmap annotations. Rows can be reordered to match prior meta clustering results.

Usage

```
## S3 method for class 'snf_config'
plot(
  x,
  order = NULL,
  hide_fixed = FALSE,
  show_column_names = TRUE,
  show_row_names = TRUE,
  rect_gp = grid::gpar(col = "black"),
  colour_breaks = c(0, 1),
  colours = c("black", "darkseagreen"),
  column_split_vector = NULL,
  row_split_vector = NULL,
  column_split = NULL,
  row_split = NULL,
  column_title = NULL,
```

```
    include_weights = TRUE,
    include_settings = TRUE,
    ...
  )

  config_heatmap(
    x,
    order = NULL,
    hide_fixed = FALSE,
    show_column_names = TRUE,
    show_row_names = TRUE,
    rect_gp = grid::gpar(col = "black"),
    colour_breaks = c(0, 1),
    colours = c("black", "darkseagreen"),
    column_split_vector = NULL,
    row_split_vector = NULL,
    column_split = NULL,
    row_split = NULL,
    column_title = NULL,
    include_weights = TRUE,
    include_settings = TRUE,
    ...
  )

  ## S3 method for class 'settings_df'
  plot(
    x,
    order = NULL,
    hide_fixed = FALSE,
    show_column_names = TRUE,
    show_row_names = TRUE,
    rect_gp = grid::gpar(col = "black"),
    colour_breaks = c(0, 1),
    colours = c("black", "darkseagreen"),
    column_split_vector = NULL,
    row_split_vector = NULL,
    column_split = NULL,
    row_split = NULL,
    column_title = NULL,
    include_weights = TRUE,
    include_settings = TRUE,
    ...
  )

  ## S3 method for class 'weights_matrix'
  plot(
    x,
    order = NULL,
```

```

hide_fixed = FALSE,
show_column_names = TRUE,
show_row_names = TRUE,
rect_gp = grid::gpar(col = "black"),
colour_breaks = c(0, 1),
colours = c("black", "darkseagreen"),
column_split_vector = NULL,
row_split_vector = NULL,
column_split = NULL,
row_split = NULL,
column_title = NULL,
include_weights = TRUE,
include_settings = TRUE,
...
)

```

Arguments

x	An snf_config class object.
order	Numeric vector indicating row ordering of SNF config.
hide_fixed	Whether fixed parameters should be removed.
show_column_names	Whether show column names.
show_row_names	Whether show row names.
rect_gp	Graphic parameters for drawing rectangles (for heatmap body). The value should be specified by <code>gpar</code> and fill parameter is ignored.
colour_breaks	Numeric vector of breaks for the legend.
colours	Vector of colours to use for the heatmap. Should match the length of colour_breaks.
column_split_vector	Vector of indices to split columns by.
row_split_vector	Vector of indices to split rows by.
column_split	Split on columns. For heatmap splitting, please refer to https://jokergoo.github.io/ComplexHeatmap-reference/book/a-single-heatmap.html#heatmap-split .
row_split	Same as split.
column_title	Title on the column.
include_weights	If TRUE, includes feature weights of the weights matrix into the config heatmap.
include_settings	If TRUE, includes columns from the settings data frame into the config heatmap.
...	Additional parameters passed to <code>ComplexHeatmap::Heatmap</code> .

Value

Returns a heatmap (class "Heatmap" from package ComplexHeatmap) that displays the scaled values of the provided SNF config.

Examples

```

dl <- data_list(
  list(income, "household_income", "demographics", "ordinal"),
  list(pubertal, "pubertal_status", "demographics", "continuous"),
  list(fav_colour, "favourite_colour", "demographics", "categorical"),
  list(anxiety, "anxiety", "behaviour", "ordinal"),
  list(depress, "depressed", "behaviour", "ordinal"),
  uid = "unique_id"
)

sc <- snf_config(
  dl,
  n_solutions = 10,
  dropout_dist = "uniform"
)

plot(sc)

```

plot.solutions_df *Plot of cluster assignments in a solutions data frame*

Description

This plot, built on `ComplexHeatmap::Heatmap()`, visualizes the cluster assignments in a solutions data frame as a categorical heatmap with observations along the columns and clusters along the rows.

Usage

```

## S3 method for class 'solutions_df'
plot(
  x,
  y = NULL,
  cluster_rows = FALSE,
  cluster_columns = TRUE,
  heatmap_legend_param = NULL,
  row_title = "Solution",
  column_title = "Observation",
  ...
)

## S3 method for class 't_solutions_df'
plot(x, ...)

```

Arguments

`x` A `solutions_df` object.

`y` Optional argument to plot, not used in this method.

cluster_rows	If the value is a logical, it controls whether to make cluster on rows. The value can also be a hclust or a dendrogram which already contains clustering. Check https://jokergoo.github.io/ComplexHeatmap-reference/book/a-single-heatmap.html#clustering .
cluster_columns	Whether make cluster on columns? Same settings as cluster_rows.
heatmap_legend_param	A list contains parameters for the heatmap legends. See color_mapping_legend , ColorMapping-method for all available parameters.
row_title	Title on the row.
column_title	Title on the column.
...	Additional arguments passed to <code>ComplexHeatmap::Heatmap()</code> .

Value

A `ComplexHeatmap::Heatmap()` object visualization of cluster assignments.

print.ari_matrix *Print method for class ari_matrix*

Description

Custom formatted print for weights matrices that outputs information about feature weights functions to the console.

Usage

```
## S3 method for class 'ari_matrix'
print(x, ...)
```

Arguments

x	A <code>ari_matrix</code> class object.
...	Other arguments passed to <code>print</code> (not used in this function)

Value

Function prints to console but does not return any value.

print.clust_fns_list *Print method for class clust_fns_list*

Description

Custom formatted print for clustering functions list objects that outputs information about the contained clustering functions to the console.

Usage

```
## S3 method for class 'clust_fns_list'  
print(x, ...)
```

Arguments

x	A clust_fns_list class object.
...	Other arguments passed to print (not used in this function)

Value

Function prints to console but does not return any value.

print.data_list *Print method for class data_list*

Description

Custom formatted print for data list objects that outputs information about the contained observations and components to the console.

Usage

```
## S3 method for class 'data_list'  
print(x, ...)
```

Arguments

x	A data_list class object.
...	Other arguments passed to print (not used in this function)

Value

Function prints to console but does not return any value.

```
print.dist_fns_list
```

Print method for class dist_fns_list

Description

Custom formatted print for distance metrics list objects that outputs information about the contained distance metrics to the console.

Usage

```
## S3 method for class 'dist_fns_list'  
print(x, ...)
```

Arguments

x	A dist_fns_list class object.
...	Other arguments passed to print (not used in this function)

Value

Function prints to console but does not return any value.

```
print.ext_solutions_df
```

Print method for class ext_solutions_df

Description

Custom formatted print for extended solutions data frame class objects.

Usage

```
## S3 method for class 'ext_solutions_df'  
print(x, n = NULL, ...)
```

Arguments

x	A ext_solutions_df class object.
n	Number of rows to print, passed into tibble::print.tbl_df().
...	Other arguments passed to print (not used in this function).

Value

Function prints to console but does not return any value.

print.settings_df *Print method for class settings_df*

Description

Custom formatted print for settings data frame that outputs information about SNF hyperparameters to the console.

Usage

```
## S3 method for class 'settings_df'  
print(x, ...)
```

Arguments

x A settings_df class object.
... Other arguments passed to print (not used in this function)

Value

Function prints to console but does not return any value.

print.sim_mats_list *Print method for class sim_mats_list*

Description

Custom formatted print for similarity matrix list

Usage

```
## S3 method for class 'sim_mats_list'  
print(x, ...)
```

Arguments

x A sim_mats_list class object.
... Other arguments passed to print (not used in this function).

print.snf_config *Print method for class snf_config*

Description

Custom formatted print for SNF config

Usage

```
## S3 method for class 'snf_config'  
print(x, ...)
```

Arguments

x A snf_config class object.
... Other arguments passed to print (not used in this function)

Value

Function prints to console but does not return any value.

print.solutions_df *Print method for class solutions_df*

Description

Custom formatted print for weights matrices that outputs information about feature weights functions to the console.

Usage

```
## S3 method for class 'solutions_df'  
print(x, n = NULL, tips = TRUE, ...)
```

Arguments

x A weights_matrix class object.
n Number of rows to print, passed into tibble::print.tbl_df().
tips If TRUE, include lines on how to print more rows / transposed.
... Other arguments passed to print (not used in this function).

Value

Function prints to console but does not return any value.

```
print.t_ext_solutions_df
    Print method for class t_ext_solutions_df
```

Description

Custom formatted print for transposed solutions data frame class objects.

Usage

```
## S3 method for class 't_ext_solutions_df'
print(x, ...)
```

Arguments

x	A t_solutions_df class object.
...	Other arguments passed to print (not used in this function)

Value

Function prints to console but does not return any value.

```
print.t_solutions_df    Print method for class t_solutions_df
```

Description

Custom formatted print for transposed solutions data frame class objects.

Usage

```
## S3 method for class 't_solutions_df'
print(x, ...)
```

Arguments

x	A t_solutions_df class object.
...	Other arguments passed to print (not used in this function)

Value

Function prints to console but does not return any value.

print.weights_matrix *Print method for class weights_matrix*

Description

Custom formatted print for weights matrices that outputs information about feature weights functions to the console.

Usage

```
## S3 method for class 'weights_matrix'  
print(x, ...)
```

Arguments

x A weights_matrix class object.
... Other arguments passed to print (not used in this function)

Value

Function prints to console but does not return any value.

pubertal *Mock ABCD pubertal status data*

Description

Like the mock data frame "abcd_pubertal", but with "unique_id" as the "uid".

Usage

```
pubertal
```

Format

pubertal:

A data frame with 275 rows and 2 columns:

unique_id The unique identifier of the ABCD dataset

pubertal_status Average reported pubertal status between child and parent (1-5 categorical scale)

Source

Though this data is no longer "real" ABCD data, the reference for using ABCD as a data source is below:

Data used in the preparation of this article were obtained from the Adolescent Brain Cognitive Development (ABCD) Study (<https://abcdstudy.org>), held in the NIMH Data Archive (NDA). This is a multisite, longitudinal study designed to recruit more than 10,000 children age 9-10 and follow them over 10 years into early adulthood. The ABCD Study® is supported by the National Institutes of Health and additional federal partners under award numbers U01DA041048, U01DA050989, U01DA051016, U01DA041022, U01DA051018, U01DA051037, U01DA050987, U01DA041174, U01DA041106, U01DA041117, U01DA041028, U01DA041134, U01DA050988, U01DA051039, U01DA041156, U01DA041025, U01DA041120, U01DA051038, U01DA041148, U01DA041093, U01DA041089, U24DA041123, U24DA041147. A full list of supporters is available at <https://abcdstudy.org/federal-partners.html>. A listing of participating sites and a complete listing of the study investigators can be found at https://abcdstudy.org/consortium_members/. ABCD consortium investigators designed and implemented the study and/or provided data but did not necessarily participate in the analysis or writing of this report. This manuscript reflects the views of the authors and may not reflect the opinions or views of the NIH or ABCD consortium investigators.

pval_heatmap

Heatmap of p-values

Description

Heatmap of p-values

Usage

```
pval_heatmap(
  ext_sol_df,
  order = NULL,
  cluster_columns = TRUE,
  cluster_rows = FALSE,
  show_row_names = FALSE,
  show_column_names = TRUE,
  min_colour = "red2",
  max_colour = "white",
  legend_breaks = c(0, 1),
  col = circlize::colorRamp2(legend_breaks, c(min_colour, max_colour)),
  heatmap_legend_param = list(color_bar = "continuous", title = "p-value", at = c(0, 1)),
  rect_gp = grid::gpar(col = "black"),
  column_split_vector = NULL,
  row_split_vector = NULL,
  column_split = NULL,
  row_split = NULL,
  ...
)
```

Arguments

ext_sol_df	An ext_solutions_df class object (produced from the function extend_solutions).
order	Numeric vector containing row order of the heatmap.
cluster_columns	Whether columns should be sorted by hierarchical clustering.
cluster_rows	Whether rows should be sorted by hierarchical clustering.
show_row_names	Whether row names should be shown.
show_column_names	Whether column names should be shown.
min_colour	Colour used for the lowest value in the heatmap.
max_colour	Colour used for the highest value in the heatmap.
legend_breaks	Numeric vector of breaks for the legend.
col	Colour function for ComplexHeatmap::Heatmap()
heatmap_legend_param	Legend function for ComplexHeatmap::Heatmap()
rect_gp	Cell border function for ComplexHeatmap::Heatmap()
column_split_vector	Vector of indices to split columns by.
row_split_vector	Vector of indices to split rows by.
column_split	Standard parameter of ComplexHeatmap::Heatmap.
row_split	Standard parameter of ComplexHeatmap::Heatmap.
...	Additional parameters passed to ComplexHeatmap::Heatmap.

Value

Returns a heatmap (class "Heatmap" from package ComplexHeatmap) that displays the provided p-values.

Examples

```
#dl <- data_list(
#   list(income, "household_income", "demographics", "ordinal"),
#   list(pubertal, "pubertal_status", "demographics", "continuous"),
#   list(fav_colour, "favourite_colour", "demographics", "categorical"),
#   list(anxiety, "anxiety", "behaviour", "ordinal"),
#   list(depress, "depressed", "behaviour", "ordinal"),
#   uid = "unique_id"
#)
#
#sc <- snf_config(
#   dl,
#   n_solutions = 4,
#   dropout_dist = "uniform",
#   max_k = 50
```

```

#)
#
#sol_df <- batch_snf(dl, sc)
#
#ext_sol_df <- extend_solutions(sol_df, dl)
#
#pval_heatmap(ext_sol_df)

```

quality_measures	<i>Quality metrics</i>
------------------	------------------------

Description

These functions calculate conventional metrics of cluster solution quality.

Usage

```
calculate_silhouettes(sol_df)
```

```
calculate_dunn_indices(sol_df)
```

```
calculate_db_indices(sol_df)
```

Arguments

`sol_df` A `solutions_df` class object created by `batch_snf()` with the parameter `return_sim_mats = TRUE`.

Details

`calculate_silhouettes`: A wrapper for `cluster::silhouette` that calculates silhouette scores for all cluster solutions in a provided solutions data frame. Silhouette values range from -1 to +1 and indicate an overall ratio of how close together observations within a cluster are to how far apart observations across clusters are. You can learn more about interpreting the results of this function by calling `?cluster::silhouette`.

`calculate_dunn_indices`: A wrapper for `clv::clv.Dunn` that calculates Dunn indices for all cluster solutions in a provided solutions data frame. Dunn indices, like silhouette scores, similarly reflect similarity within clusters and separation across clusters. You can learn more about interpreting the results of this function by calling `?clv::clv.Dunn`.

`calculate_db_indices`: A wrapper for `clv::clv.Davies.Bouldin` that calculates Davies-Bouldin indices for all cluster solutions in a provided solutions data frame. These values can be interpreted similarly as those above. You can learn more about interpreting the results of this function by calling `?clv::clv.Davies.Bouldin`.

Value

A list of `silhouette` class objects, a vector of Dunn indices, or a vector of Davies-Bouldin indices depending on which function was used.

Examples

```
## Not run:
input_dl <- data_list(
  list(gender_df, "gender", "demographics", "categorical"),
  list(diagnosis_df, "diagnosis", "clinical", "categorical"),
  uid = "patient_id"
)

sc <- snf_config(input_dl, n_solutions = 5)

sol_df <- batch_snf(input_dl, sc, return_sim_mats = TRUE)

# calculate Davies-Bouldin indices
davies_bouldin_indices <- calculate_db_indices(sol_df)

# calculate Dunn indices
dunn_indices <- calculate_dunn_indices(sol_df)

# calculate silhouette scores
silhouette_scores <- calculate_silhouettes(sol_df)

## End(Not run)
```

rbind.ext_solutions_df

Row-binding of solutions data frame class objects

Description

Row-binding of solutions data frame class objects

Usage

```
## S3 method for class 'ext_solutions_df'
rbind(..., reset_indices = FALSE)
```

Arguments

... An arbitrary number of ext_solutions_df class objects.

reset_indices If TRUE, re-labels the "solutions" indices in the solutions data frame from 1 to the number of defined settings.

Value

An ext_solutions_df class object.

`rbind.solutions_df` *Row-binding of solutions data frame class objects*

Description

Row-binding of solutions data frame class objects

Usage

```
## S3 method for class 'solutions_df'  
rbind(..., reset_indices = FALSE)
```

Arguments

`...` An arbitrary number of `solutions_df` class objects.
`reset_indices` If TRUE, re-labels the "solutions" indices in the solutions data frame from 1 to the number of defined settings.

Value

A `solutions_df` class object.

`rbind.t_solutions_df` *Row-binding of t_solutions_df class objects*

Description

Vertically stack two or more `t_solutions_df` class objects.

Usage

```
## S3 method for class 't_solutions_df'  
rbind(...)
```

Arguments

`...` An arbitrary number of `t_solutions_df` class objects.

Value

A `t_solutions_df` class object.

rbind.weights_matrix *Row-bind weights matrices*

Description

Vertically stack two or more weights_matrix class objects.

Usage

```
## S3 method for class 'weights_matrix'  
rbind(...)
```

Arguments

... An arbitrary number of weights_matrix class objects.

Value

A weights_matrix class object.

rename_dl *Rename features in a data list*

Description

Rename features in a data list

Usage

```
rename_dl(dl, name_mapping)
```

Arguments

dl A nested list of input data from data_list().
name_mapping A named vector where the values are the features to be renamed and the names are the new names for those features.

Value

A data list ("list"-class object) with adjusted feature names.

Examples

```
d1 <- data_list(  
  list(pubertal, "pubertal_status", "demographics", "continuous"),  
  list(anxiety, "anxiety", "behaviour", "ordinal"),  
  list(depress, "depressed", "behaviour", "ordinal"),  
  uid = "unique_id"  
)  
  
summary(d1, "feature")  
  
name_changes <- c(  
  "anxiety_score" = "cbcl_anxiety_r",  
  "depression_score" = "cbcl_depress_r"  
)  
  
d1 <- rename_dl(d1, name_changes)  
  
summary(d1, "feature")
```

resample

Helper resampling function found in ?sample

Description

Like `sample`, but when given a single value `x`, returns back that single value instead of a random value from 1 to `x`.

Usage

```
resample(x, ...)
```

Arguments

`x` Vector or single value to sample from
`...` Remaining arguments for `base::sample` function

Value

Numeric vector result of running `base::sample`.

save_heatmap	<i>Save a heatmap object to a file</i>
--------------	--

Description

Save a heatmap object to a file

Usage

```
save_heatmap(heatmap, path, width = 480, height = 480, res = 100)
```

Arguments

heatmap	The heatmap object to save.
path	The path to save the heatmap to.
width	The width of the heatmap.
height	The height of the heatmap.
res	The resolution of the heatmap.

Value

Does not return any value. Saves heatmap to file.

settings_df	<i>Build a settings data frame</i>
-------------	------------------------------------

Description

The settings_df is a data frame whose rows completely specify the hyperparameters and decisions required to transform individual input data frames (found in a data list, see ?data_list) into a single similarity matrix through SNF. The format of the settings data frame is as follows:

- A column named "solution": This column is used to keep track of the rows and should have integer values only.
- A column named "alpha": This column contains the value of the alpha hyperparameter that will be used on that run of the SNF pipeline.
- A column named "k": Like above, but for the K (nearest neighbours) hyperparameter.
- A column named "t": Like above, but for the t (number of iterations) hyperparameter.
- A column named "snf_scheme": Which of 3 pre-defined schemes will be used to integrate the data frames of the data list into a final fused network. The purpose of varying these schemes is primarily to increase the diversity of the generated cluster solutions.

- A value of 1 corresponds to the "individual" scheme, in which all data frames are directly merged by SNF into the final fused network. This scheme corresponds to the approach shown in the original SNF paper.
 - A value of 2 corresponds to the "two-step" scheme, in which all data frames within a domain are first merged into a domain-specific fused network. Next, domain-specific networks are fused once more by SNF into the final fused network. This scheme is useful for fairly re-weighting SNF pipelines with unequal numbers of data frames across domains.
 - A value of 3 corresponds to the "domain" scheme, in which all data frames within a domain are first concatenated into a single domain-specific data frame before being merged by SNF into the final fused network. This approach serves as an alternative way to re-weight SNF pipelines with unequal numbers of data frames across domains. You can learn more about this parameter here: https://branchlab.github.io/metasnf/articles/snf_schemes.html.
- A column named "clust_alg": Specification of which clustering algorithm will be applied to the final similarity matrix. By default, this column can take on the integer values 1 or 2, which correspond to spectral clustering where the number of clusters is determined by the eigen-gap or rotation cost heuristic respectively. You can learn more about this parameter here: https://branchlab.github.io/metasnf/articles/clustering_algorithms.html.
 - A column named "cnt_dist": Specification of which distance metric will be used for data frames of purely continuous data. You can learn about this metric and its defaults here: https://branchlab.github.io/metasnf/articles/distance_metrics.html
 - A column named "dsc_dist": Like above, but for discrete data frames.
 - A column named "ord_dist": Like above, but for ordinal data frames.
 - A column named "cat_dist": Like above, but for categorical data frames.
 - A column named "mix_dist": Like above, but for mixed-type (e.g., both categorical and discrete) data frames.
 - One column for every input data frame in the corresponding data list which can either have the value of 0 or 1. The name of the column should be formatted as "inc_[]" where the square brackets are replaced with the name (as found in `dl_summary(dl)$"name"`) of each data frame. When 0, that data frame will be excluded from that run of the SNF pipeline. When 1, that data frame will be included.

Usage

```
settings_df(
  dl,
  n_solutions = 0,
  min_removed_inputs = 0,
  max_removed_inputs = length(dl) - 1,
  dropout_dist = "exponential",
  min_alpha = NULL,
  max_alpha = NULL,
  min_k = NULL,
  max_k = NULL,
  min_t = NULL,
  max_t = NULL,
```

```

alpha_values = NULL,
k_values = NULL,
t_values = NULL,
possible_snf_schemes = c(1, 2, 3),
clustering_algorithms = NULL,
continuous_distances = NULL,
discrete_distances = NULL,
ordinal_distances = NULL,
categorical_distances = NULL,
mixed_distances = NULL,
dfl = NULL,
snf_input_weights = NULL,
snf_domain_weights = NULL,
retry_limit = 10,
allow_duplicates = FALSE
)

```

Arguments

d1	A nested list of input data from <code>data_list()</code> .
n_solutions	Number of rows to generate for the settings data frame.
min_removed_inputs	The smallest number of input data frames that may be randomly removed. By default, 0.
max_removed_inputs	The largest number of input data frames that may be randomly removed. By default, this is 1 less than all the provided input data frames in the data list.
dropout_dist	Parameter controlling how the random removal of input data frames should occur. Can be "none" (no input data frames are randomly removed), "uniform" (uniformly sample between <code>min_removed_inputs</code> and <code>max_removed_inputs</code> to determine number of input data frames to remove), or "exponential" (pick number of input data frames to remove by sampling from <code>min_removed_inputs</code> to <code>max_removed_inputs</code> with an exponential distribution; the default).
min_alpha	The minimum value that the alpha hyperparameter can have. Random assigned value of alpha for each row will be obtained by uniformly sampling numbers between <code>min_alpha</code> and <code>max_alpha</code> at intervals of 0.1. Cannot be used in conjunction with the <code>alpha_values</code> parameter.
max_alpha	The maximum value that the alpha hyperparameter can have. See <code>min_alpha</code> parameter. Cannot be used in conjunction with the <code>alpha_values</code> parameter.
min_k	The minimum value that the k hyperparameter can have. Random assigned value of k for each row will be obtained by uniformly sampling numbers between <code>min_k</code> and <code>max_k</code> at intervals of 1. Cannot be used in conjunction with the <code>k_values</code> parameter.
max_k	The maximum value that the k hyperparameter can have. See <code>min_k</code> parameter. Cannot be used in conjunction with the <code>k_values</code> parameter.
min_t	The minimum value that the t hyperparameter can have. Random assigned value of t for each row will be obtained by uniformly sampling numbers between

	min_t and max_t at intervals of 1. Cannot be used in conjunction with the t_values parameter.
max_t	The maximum value that the t hyperparameter can have. See min_t parameter. Cannot be used in conjunction with the t_values parameter.
alpha_values	A number or numeric vector of a set of possible values that alpha can take on. Value will be obtained by uniformly sampling the vector. Cannot be used in conjunction with the min_alpha or max_alpha parameters.
k_values	A number or numeric vector of a set of possible values that k can take on. Value will be obtained by uniformly sampling the vector. Cannot be used in conjunction with the min_k or max_k parameters.
t_values	A number or numeric vector of a set of possible values that t can take on. Value will be obtained by uniformly sampling the vector. Cannot be used in conjunction with the min_t or max_t parameters.
possible_snf_schemes	A vector containing the possible snf_schemes to uniformly randomly select from. By default, the vector contains all 3 possible schemes: c(1, 2, 3). 1 corresponds to the "individual" scheme, 2 corresponds to the "domain" scheme, and 3 corresponds to the "two-step" scheme.
clustering_algorithms	A list of clustering algorithms to uniformly randomly pick from when clustering. When not specified, randomly select between spectral clustering using the eigen-gap heuristic and spectral clustering using the rotation cost heuristic. See ?clust_fns_list for more details on running custom clustering algorithms.
continuous_distances	A vector of continuous distance metrics to use when a custom dist_fns_list is provided.
discrete_distances	A vector of categorical distance metrics to use when a custom dist_fns_list is provided.
ordinal_distances	A vector of categorical distance metrics to use when a custom dist_fns_list is provided.
categorical_distances	A vector of categorical distance metrics to use when a custom dist_fns_list is provided.
mixed_distances	A vector of mixed distance metrics to use when a custom dist_fns_list is provided.
df1	List containing distance metrics to vary over. See ?generate_dist_fns_list.
snf_input_weights	Nested list containing weights for when SNF is used to merge individual input measures (see ?generate_snf_weights)
snf_domain_weights	Nested list containing weights for when SNF is used to merge domains (see ?generate_snf_weights)

`retry_limit` The maximum number of attempts to generate a novel row. This function does not return matrices with identical rows. As the range of requested possible settings tightens and the number of requested rows increases, the risk of randomly generating a row that already exists increases. If a new random row has matched an existing row `retry_limit` number of times, the function will terminate.

`allow_duplicates` If TRUE, enables creation of a settings data frame with duplicate non-feature weighting related hyperparameters. This function should only be used when paired with a custom weights matrix that has non-duplicate rows.

Value

A settings data frame

<code>shiny_annotator</code>	<i>Launch a shiny app to identify meta cluster boundaries</i>
------------------------------	---

Description

This function calls the `htShiny()` function from the package `InteractiveComplexHeatmap` to assist users in identifying the indices of the boundaries between meta clusters in a meta cluster heatmap. By providing a heatmap of inter-solution similarities (obtained through `meta_cluster_heatmap()`), users can click on positions within the heatmap that appear to meaningfully separate major sets of similar cluster solutions by visual inspection. The corresponding indices of the clicked positions are printed to the console and also shown within the app. This function can only run from an interactive session of R.

Usage

```
shiny_annotator(ari_heatmap)
```

Arguments

`ari_heatmap` Heatmap of ARIs to divide into meta clusters.

Value

Does not return any value. Launches interactive shiny applet.

Examples

```
#dl <- data_list(
#   list(cort_sa, "cortical_surface_area", "neuroimaging", "continuous"),
#   list(subc_v, "subcortical_volume", "neuroimaging", "continuous"),
#   list(income, "household_income", "demographics", "continuous"),
#   list(pubertal, "pubertal_status", "demographics", "continuous"),
#   uid = "unique_id"
#)
```

```
#
#set.seed(42)
#my_sc <- snf_config(
#  dl = dl,
#  n_solutions = 20,
#  min_k = 20,
#  max_k = 50
#)
#
#sol_df <- batch_snf(dl, my_sc)
#
#sol_aris <- calc_aris(sol_df)
#
#meta_cluster_order <- get_matrix_order(sol_aris)
#
#ari_hm <- meta_cluster_heatmap(sol_aris, order = meta_cluster_order)
#
## Click on meta cluster boundaries to obtain `split_vec` values
#shiny_annotator(ari_hm)
#
#split_vec <- c(6, 10, 16)
#
#ari_hm <- meta_cluster_heatmap(
#  sol_aris,
#  order = meta_cluster_order,
#  split_vector = split_vec
#)
```

similarity_matrix_heatmap

Plot heatmap of similarity matrix

Description

Plot heatmap of similarity matrix

Usage

```
similarity_matrix_heatmap(
  similarity_matrix,
  order = NULL,
  cluster_solution = NULL,
  scale_diag = "mean",
  log_graph = TRUE,
  cluster_rows = FALSE,
  cluster_columns = FALSE,
  show_row_names = FALSE,
  show_column_names = FALSE,
  data = NULL,
```

```

    left_bar = NULL,
    right_bar = NULL,
    top_bar = NULL,
    bottom_bar = NULL,
    left_hm = NULL,
    right_hm = NULL,
    top_hm = NULL,
    bottom_hm = NULL,
    annotation_colours = NULL,
    min_colour = NULL,
    max_colour = NULL,
    split_vector = NULL,
    row_split = NULL,
    column_split = NULL,
    ...
)

```

Arguments

similarity_matrix	A similarity matrix
order	Vector of numbers to reorder the similarity matrix (and data if provided). Overwrites ordering specified by cluster_solution param.
cluster_solution	Row of a solutions data frame or column of a transposed solutions data frame.
scale_diag	Method of rescaling matrix diagonals. Can be "none" (don't change diagonals), "mean" (replace diagonals with average value of off-diagonals), or "zero" (replace diagonals with 0).
log_graph	If TRUE, log transforms the graph.
cluster_rows	Parameter for ComplexHeatmap::Heatmap.
cluster_columns	Parameter for ComplexHeatmap::Heatmap.
show_row_names	Parameter for ComplexHeatmap::Heatmap.
show_column_names	Parameter for ComplexHeatmap::Heatmap.
data	A data frame containing elements requested for annotation.
left_bar	Named list of strings, where the strings are features in df that should be used for a barplot annotation on the left of the plot and the names are the names that will be used to caption the plots and their legends.
right_bar	See left_bar.
top_bar	See left_bar.
bottom_bar	See left_bar.
left_hm	Like left_bar, but with a heatmap annotation instead of a barplot annotation.
right_hm	See left_hm.

top_hm	See left_hm.
bottom_hm	See left_hm.
annotation_colours	Named list of heatmap annotations and their colours.
min_colour	Colour used for the lowest value in the heatmap.
max_colour	Colour used for the highest value in the heatmap.
split_vector	A vector of partition indices.
row_split	Standard parameter of ComplexHeatmap: :Heatmap.
column_split	Standard parameter of ComplexHeatmap: :Heatmap.
...	Additional parameters passed into ComplexHeatmap::Heatmap.

Value

Returns a heatmap (class "Heatmap" from package ComplexHeatmap) that displays the similarities between observations in the provided matrix.

Examples

```
#my_dl <- data_list(
#   list(
#     data = expression_df,
#     name = "expression_data",
#     domain = "gene_expression",
#     type = "continuous"
#   ),
#   list(
#     data = methylation_df,
#     name = "methylation_data",
#     domain = "gene_methylation",
#     type = "continuous"
#   ),
#   uid = "patient_id"
#)
#
#sc <- snf_config(my_dl, n_solutions = 10)
#
#sol_df <- batch_snf(my_dl, sc, return_sim_mats = TRUE)
#
#sim_mats <- sim_mats_list(sol_df)
#
#similarity_matrix_heatmap(
#   sim_mats[[1]],
#   cluster_solution = sol_df[1, ]
#)
```

sim_mats_list	<i>Create or extract a sim_mats_list class object</i>
---------------	---

Description

Create or extract a sim_mats_list class object

Usage

```
sim_mats_list(x)
```

Arguments

x The object to create or extract a sim_mats_list from.

Value

A sim_mats_list class object.

siw_euclidean_distance	<i>Squared (including weights) Euclidean distance</i>
------------------------	---

Description

Squared (including weights) Euclidean distance

Usage

```
siw_euclidean_distance(df, weights_row)
```

Arguments

df data frame containing at least 1 data column.
weights_row Single-row data frame where the column names contain the column names in df and the row contains the corresponding weights.

Value

distance_matrix A distance matrix.

snf_config	<i>Define configuration for generating a set of SNF-based cluster solutions</i>
------------	---

Description

snf_config() constructs an SNF config object which inherits from classes snf_config and list. This object is used to store all settings required to transform data stored in a data_list class object into a space of cluster solutions by SNF. The SNF config object contains the following components: 1. A settings data frame (inherits from settings_df and data.frame). Data frame that stores SNF-specific hyperparameters and information about feature selection and weighting, SNF schemes, clustering algorithms, and distance metrics. Each row of the settings data frame corresponds to a distinct cluster solution. 2. A clustering algorithms list (inherits from clust_fns_list and list), which stores all clustering algorithms that the settings data frame can point to. 3. A distance metrics list (inherits from dist_metrics_list and list), which stores all distance metrics that the settings data frame can point to. 4. A weights matrix (inherits from weights_matrix, matrix, and array'), which stores the feature weights to use prior to distance calculations. Each column of the weights matrix corresponds to a different feature in the data list and each row corresponds to a different row in the settings data frame.

Usage

```
snf_config(
  dl = NULL,
  sdf = NULL,
  dfl = NULL,
  cfl = NULL,
  wm = NULL,
  n_solutions = 0,
  min_removed_inputs = 0,
  max_removed_inputs = length(dl) - 1,
  dropout_dist = "exponential",
  min_alpha = NULL,
  max_alpha = NULL,
  min_k = NULL,
  max_k = NULL,
  min_t = NULL,
  max_t = NULL,
  alpha_values = NULL,
  k_values = NULL,
  t_values = NULL,
  possible_snf_schemes = c(1, 2, 3),
  clustering_algorithms = NULL,
  continuous_distances = NULL,
  discrete_distances = NULL,
  ordinal_distances = NULL,
  categorical_distances = NULL,
```

```

mixed_distances = NULL,
snf_input_weights = NULL,
snf_domain_weights = NULL,
retry_limit = 10,
cnt_dist_fns = NULL,
dsc_dist_fns = NULL,
ord_dist_fns = NULL,
cat_dist_fns = NULL,
mix_dist_fns = NULL,
automatic_standard_normalize = FALSE,
use_default_dist_fns = FALSE,
clust_fns = NULL,
use_default_clust_fns = FALSE,
weights_fill = "ones"
)

```

Arguments

<code>dl</code>	A nested list of input data from <code>data_list()</code> .
<code>sdf</code>	A <code>settings_df</code> class object. Overrides settings data frame related parameters.
<code>df1</code>	A <code>dist_fns_list</code> class object. Overrides distance functions list related parameters.
<code>cfl</code>	A <code>clust_fns_list</code> class object. Overrides clustering functions list related parameters.
<code>wm</code>	A <code>weights_matrix</code> class object. Overrides weights matrix related parameters.
<code>n_solutions</code>	Number of rows to generate for the settings data frame.
<code>min_removed_inputs</code>	The smallest number of input data frames that may be randomly removed. By default, 0.
<code>max_removed_inputs</code>	The largest number of input data frames that may be randomly removed. By default, this is 1 less than all the provided input data frames in the data list.
<code>dropout_dist</code>	Parameter controlling how the random removal of input data frames should occur. Can be "none" (no input data frames are randomly removed), "uniform" (uniformly sample between <code>min_removed_inputs</code> and <code>max_removed_inputs</code> to determine number of input data frames to remove), or "exponential" (pick number of input data frames to remove by sampling from <code>min_removed_inputs</code> to <code>max_removed_inputs</code> with an exponential distribution; the default).
<code>min_alpha</code>	The minimum value that the alpha hyperparameter can have. Random assigned value of alpha for each row will be obtained by uniformly sampling numbers between <code>min_alpha</code> and <code>max_alpha</code> at intervals of 0.1. Cannot be used in conjunction with the <code>alpha_values</code> parameter.
<code>max_alpha</code>	The maximum value that the alpha hyperparameter can have. See <code>min_alpha</code> parameter. Cannot be used in conjunction with the <code>alpha_values</code> parameter.
<code>min_k</code>	The minimum value that the k hyperparameter can have. Random assigned value of k for each row will be obtained by uniformly sampling numbers between

	min_k and max_k at intervals of 1. Cannot be used in conjunction with the k_values parameter.
max_k	The maximum value that the k hyperparameter can have. See min_k parameter. Cannot be used in conjunction with the k_values parameter.
min_t	The minimum value that the t hyperparameter can have. Random assigned value of t for each row will be obtained by uniformly sampling numbers between min_t and max_t at intervals of 1. Cannot be used in conjunction with the t_values parameter.
max_t	The maximum value that the t hyperparameter can have. See min_t parameter. Cannot be used in conjunction with the t_values parameter.
alpha_values	A number or numeric vector of a set of possible values that alpha can take on. Value will be obtained by uniformly sampling the vector. Cannot be used in conjunction with the min_alpha or max_alpha parameters.
k_values	A number or numeric vector of a set of possible values that k can take on. Value will be obtained by uniformly sampling the vector. Cannot be used in conjunction with the min_k or max_k parameters.
t_values	A number or numeric vector of a set of possible values that t can take on. Value will be obtained by uniformly sampling the vector. Cannot be used in conjunction with the min_t or max_t parameters.
possible_snf_schemes	A vector containing the possible snf_schemes to uniformly randomly select from. By default, the vector contains all 3 possible schemes: c(1, 2, 3). 1 corresponds to the "individual" scheme, 2 corresponds to the "domain" scheme, and 3 corresponds to the "two-step" scheme.
clustering_algorithms	A list of clustering algorithms to uniformly randomly pick from when clustering. When not specified, randomly select between spectral clustering using the eigen-gap heuristic and spectral clustering using the rotation cost heuristic. See ?clust_fns_list for more details on running custom clustering algorithms.
continuous_distances	A vector of continuous distance metrics to use when a custom dist_fns_list is provided.
discrete_distances	A vector of categorical distance metrics to use when a custom dist_fns_list is provided.
ordinal_distances	A vector of categorical distance metrics to use when a custom dist_fns_list is provided.
categorical_distances	A vector of categorical distance metrics to use when a custom dist_fns_list is provided.
mixed_distances	A vector of mixed distance metrics to use when a custom dist_fns_list is provided.

snf_input_weights	Nested list containing weights for when SNF is used to merge individual input measures (see ?generate_snf_weights)
snf_domain_weights	Nested list containing weights for when SNF is used to merge domains (see ?generate_snf_weights)
retry_limit	The maximum number of attempts to generate a novel row. This function does not return matrices with identical rows. As the range of requested possible settings tightens and the number of requested rows increases, the risk of randomly generating a row that already exists increases. If a new random row has matched an existing row <code>retry_limit</code> number of times, the function will terminate.
cnt_dist_fns	A named list of continuous distance metric functions.
dsc_dist_fns	A named list of discrete distance metric functions.
ord_dist_fns	A named list of ordinal distance metric functions.
cat_dist_fns	A named list of categorical distance metric functions.
mix_dist_fns	A named list of mixed distance metric functions.
automatic_standard_normalize	If TRUE, will automatically use standard normalization prior to calculation of any numeric distances. This parameter overrides all other distance functions list-related parameters.
use_default_dist_fns	If TRUE, prepend the base distance metrics (euclidean distance for continuous, discrete, and ordinal data and gower distance for categorical and mixed data) to the resulting distance metrics list.
clust_fns	A list of named clustering functions
use_default_clust_fns	If TRUE, prepend the base clustering algorithms (<code>spectral_eigen</code> and <code>spectral_rot</code> , which apply spectral clustering and use the eigen-gap and rotation cost heuristics respectively for determining the number of clusters in the graph) to <code>clust_fns</code> .
weights_fill	String indicating what to populate generate rows with. Can be "ones" (default; fill matrix with 1), "uniform" (fill matrix with uniformly distributed random values), or "exponential" (fill matrix with exponentially distributed random values).

Value

An `snf_config` class object.

Examples

```
# Simple random config for 5 cluster solutions
input_dl <- data_list(
  list(anxiety, "anxiety", "behaviour", "ordinal"),
  list(depress, "depressed", "behaviour", "ordinal"),
  uid = "unique_id"
)
my_sc <- snf_config(
```

```

    dl = input_dl,
    n_solutions = 5
  )

# specifying possible K range
my_sc <- snf_config(
  dl = input_dl,
  n_solutions = 5,
  min_k = 20,
  max_k = 40
)

# Random feature weights across from uniform distribution
my_sc <- snf_config(
  dl = input_dl,
  n_solutions = 5,
  min_k = 20,
  max_k = 40,
  weights_fill = "uniform"
)

# Specifying custom pre-built clustering and distance functions
# - Random alternation between 2-cluster and 5-cluster solutions
# - When continuous or discrete data frames are being processed,
#   randomly alternate between standardized/normalized Euclidean
#   distance and regular Euclidean distance
my_sc <- snf_config(
  dl = input_dl,
  n_solutions = 5,
  min_k = 20,
  max_k = 40,
  weights_fill = "uniform",
  clust_fns = list(
    "two_cluster_spectral" = spectral_two,
    "five_cluster_spectral" = spectral_five
  ),
  cnt_dist_fns = list(
    "euclidean" = euclidean_distance,
    "std_nrm_euc" = sn_euclidean_distance
  ),
  dsc_dist_fns = list(
    "euclidean" = euclidean_distance,
    "std_nrm_euc" = sn_euclidean_distance
  )
)

```

split_parser

Helper function to determine which row and columns to split on

Description

Helper function to determine which row and columns to split on

Usage

```
split_parser(
  row_split_vector = NULL,
  column_split_vector = NULL,
  row_split = NULL,
  column_split = NULL,
  n_rows,
  n_columns
)
```

Arguments

`row_split_vector` A vector of row indices to split on.

`column_split_vector` A vector of column indices to split on.

`row_split` Standard parameter of `ComplexHeatmap::Heatmap`.

`column_split` Standard parameter of `ComplexHeatmap::Heatmap`.

`n_rows` The number of rows in the data.

`n_columns` The number of columns in the data.

Value

"list"-class object containing `row_split` and `column_split` character vectors to pass into `ComplexHeatmap::Heatmap`.

`str.ari_matrix` *Structure of a ari_matrix object*

Description

Structure of a `ari_matrix` object

Usage

```
## S3 method for class 'ari_matrix'
str(object, ...)
```

Arguments

`object` A `ari_matrix` class object.

`...` Additional arguments (not used).

Value

Does not return an object; outputs object structure to console.

str.clust_fns_list *Structure of a clust_fns_list object*

Description

Structure of a clust_fns_list object

Usage

```
## S3 method for class 'clust_fns_list'  
str(object, ...)
```

Arguments

object A clust_fns_list class object.
... Additional arguments (not used).

Value

Does not return an object; outputs object structure to console.

str.data_list *Structure of a data_list object*

Description

Structure of a data_list object

Usage

```
## S3 method for class 'data_list'  
str(object, ...)
```

Arguments

object A data_list class object.
... Additional arguments (not used).

Value

Does not return an object; outputs object structure to console.

str.dist_fns_list *Structure of a dist_fns_list object*

Description

Structure of a dist_fns_list object

Usage

```
## S3 method for class 'dist_fns_list'  
str(object, ...)
```

Arguments

object A dist_fns_list class object.
... Additional arguments (not used).

Value

Does not return an object; outputs object structure to console.

str.ext_solutions_df *Structure of a ext_solutions_df object*

Description

Structure of a ext_solutions_df object

Usage

```
## S3 method for class 'ext_solutions_df'  
str(object, ...)
```

Arguments

object A ext_solutions_df class object.
... Additional arguments (not used).

Value

Does not return an object; outputs object structure to console.

str.settings_df *Structure of a settings_df object*

Description

Structure of a settings_df object

Usage

```
## S3 method for class 'settings_df'  
str(object, ...)
```

Arguments

object A settings_df class object.
... Additional arguments (not used).

Value

Does not return an object; outputs object structure to console.

str.sim_mats_list *Structure of a sim_mats_list object*

Description

Structure of a sim_mats_list object

Usage

```
## S3 method for class 'sim_mats_list'  
str(object, ...)
```

Arguments

object A sim_mats_list class object.
... Additional arguments (not used).

Value

Does not return an object; outputs object structure to console.

str.snf_config	<i>Structure of a snf_config object</i>
----------------	---

Description

Structure of a snf_config object

Usage

```
## S3 method for class 'snf_config'  
str(object, ...)
```

Arguments

object	A snf_config class object.
...	Additional arguments (not used).

Value

Does not return an object; outputs object structure to console.

str.solutions_df	<i>Structure of a solutions_df object</i>
------------------	---

Description

Structure of a solutions_df object

Usage

```
## S3 method for class 'solutions_df'  
str(object, ...)
```

Arguments

object	A solutions_df class object.
...	Additional arguments (not used).

Value

Does not return an object; outputs object structure to console.

str.t_ext_solutions_df

Structure of a t_ext_solutions_df object

Description

Structure of a t_ext_solutions_df object

Usage

```
## S3 method for class 't_ext_solutions_df'  
str(object, ...)
```

Arguments

object A t_ext_solutions_df class object.
... Additional arguments (not used).

Value

Does not return an object; outputs object structure to console.

str.t_solutions_df

Structure of a t_solutions_df object

Description

Structure of a t_solutions_df object

Usage

```
## S3 method for class 't_solutions_df'  
str(object, ...)
```

Arguments

object A t_solutions_df class object.
... Additional arguments (not used).

Value

Does not return an object; outputs object structure to console.

str.weights_matrix *Structure of a weights_matrix object*

Description

Structure of a weights_matrix object

Usage

```
## S3 method for class 'weights_matrix'  
str(object, ...)
```

Arguments

object A weights_matrix class object.
... Additional arguments (not used).

Value

Does not return an object; outputs object structure to console.

subc_v *Mock ABCD subcortical volumes data*

Description

Like the mock data frame "abcd_subc_v", but with "unique_id" as the "uid".

Usage

```
subc_v
```

Format

subc_v:
A data frame with 174 rows and 31 columns:
unique_id The unique identifier of the ABCD dataset
... Subcortical volumes of various ROIs (mm³, I think)

Source

Though this data is no longer "real" ABCD data, the reference for using ABCD as a data source is below:

Data used in the preparation of this article were obtained from the Adolescent Brain Cognitive Development (ABCD) Study (<https://abcdstudy.org>), held in the NIMH Data Archive (NDA). This is a multisite, longitudinal study designed to recruit more than 10,000 children age 9-10 and follow them over 10 years into early adulthood. The ABCD Study® is supported by the National Institutes of Health and additional federal partners under award numbers U01DA041048, U01DA050989, U01DA051016, U01DA041022, U01DA051018, U01DA051037, U01DA050987, U01DA041174, U01DA041106, U01DA041117, U01DA041028, U01DA041134, U01DA050988, U01DA051039, U01DA041156, U01DA041025, U01DA041120, U01DA051038, U01DA041148, U01DA041093, U01DA041089, U24DA041123, U24DA041147. A full list of supporters is available at <https://abcdstudy.org/federal-partners.html>. A listing of participating sites and a complete listing of the study investigators can be found at https://abcdstudy.org/consortium_members/. ABCD consortium investigators designed and implemented the study and/or provided data but did not necessarily participate in the analysis or writing of this report. This manuscript reflects the views of the authors and may not reflect the opinions or views of the NIH or ABCD consortium investigators.

subsample_dl

Create subsamples of a data list

Description

Given a data list, return a list of smaller data lists that are generated through random sampling (without replacement). The results of this function can be passed into `batch_snf_subsamples()` to obtain a list of resampled solutions data frames.

Usage

```
subsample_dl(
  dl,
  n_subsamples,
  subsample_fraction = NULL,
  n_observations = NULL
)
```

Arguments

`dl` A nested list of input data from `data_list()`.

`n_subsamples` Number of subsamples to create.

`subsample_fraction` Percentage of patients to include per subsample.

`n_observations` Number of patients to include per subsample.

Value

A "list" class object containing `n_subsamples` number of data lists. Each of those data lists contains a random `subsample_fraction` fraction of the observations of the provided data list.

Examples

```
my_dl <- data_list(
  list(subc_v, "subcortical_volume", "neuroimaging", "continuous"),
  list(income, "household_income", "demographics", "continuous"),
  list(pubertal, "pubertal_status", "demographics", "continuous"),
  uid = "unique_id"
)

my_dl_subsamples <- subsample_dl(
  my_dl,
  n_subsamples = 20,
  subsample_fraction = 0.85
)
```

subsample_pairwise_aris

Calculate pairwise adjusted Rand indices across subsamples of data

Description

Given a list of subsampled solutions data frames from `batch_snf_subsamples()`, this function calculates the adjusted Rand indices across all the subsamples of each solution. ARI calculation between two subsamples only factors in observations that were present in both subsamples.

Usage

```
subsample_pairwise_aris(subsample_solutions, verbose = FALSE)
```

Arguments

`subsample_solutions`

A list of solutions data frames from subsamples of the data. This object is generated by the function `batch_snf_subsamples()`.

`verbose`

If TRUE, output progress to console.

Value

A two-item list: "raw_aris", a list of inter-subsample pairwise ARI matrices (one for each full cluster solution) and "ari_summary", a data frame containing the mean and SD of the inter-subsample ARIs for each original cluster solution.

Examples

```

my_dl <- data_list(
  list(subc_v, "subcortical_volume", "neuroimaging", "continuous"),
  list(income, "household_income", "demographics", "continuous"),
  list(pubertal, "pubertal_status", "demographics", "continuous"),
  uid = "unique_id"
)

sc <- snf_config(my_dl, n_solutions = 5, max_k = 40)

my_dl_subsamples <- subsample_dl(
  my_dl,
  n_subsamples = 20,
  subsample_fraction = 0.85
)

batch_subsample_results <- batch_snf_subsamples(
  my_dl_subsamples,
  sc
)

pairwise_aris <- subsample_pairwise_aris(
  batch_subsample_results,
  verbose = TRUE
)

# Visualize ARIs
ComplexHeatmap::Heatmap(
  pairwise_aris$"raw_aris"[[1]],
  heatmap_legend_param = list(
    color_bar = "continuous",
    title = "Inter-Subsample\nARI",
    at = c(0, 0.5, 1)
  ),
  show_column_names = FALSE,
  show_row_names = FALSE
)

```

summary.ari_matrix *Summary method for class ari_matrix*

Description

Provides a summary of the `ari_matrix` class object, including the distribution of the adjusted Rand index (ARI) values and the number of solutions.

Usage

```

## S3 method for class 'ari_matrix'
summary(object, ...)

```

Arguments

object A `ari_matrix` class object.
... Other arguments passed to `summary` (not used in this function).

Value

A named list containing the number of solutions and the distribution of ARI values.

summary.clust_fns_list *Summary method for class clust_fns_list*

Description

This summary function simply returns to the console the number of functions contained in the `clust_fns_list` object.

Usage

```
## S3 method for class 'clust_fns_list'  
summary(object, ...)
```

Arguments

object A `clust_fns_list` class object.
... Other arguments passed to `summary` (not used in this function).

Value

Returns no value. Outputs a message to the console.

summary.data_list *Summary method for class data_list*

Description

Returns a data list summary (`data.frame` class object) containing information on components, features, variable types, domains, and component dimensions.

Usage

```
## S3 method for class 'data_list'  
summary(object, scope = "component", ...)
```

Arguments

object	A data_list class object.
scope	The level of detail for the summary. By default, this is set to "component", which returns a summary of the data list at the component level. Can also be set to "feature", resulting in a summary at the feature level.
...	Other arguments passed to summary (not used in this function)

Value

A data.frame class object. If scope is "component", each row shows the name, variable type, domain, and dimensions of each component. If scope is "feature", each row shows the name, variable type, and domain of each feature.

summary.dist_fns_list *Summary method for class dist_fns_list*

Description

This summary function simply returns to the console the number of functions contained in the dist_fns_list object.

Usage

```
## S3 method for class 'dist_fns_list'  
summary(object, ...)
```

Arguments

object	A dist_fns_list class object.
...	Other arguments passed to summary (not used in this function).

Value

Returns no value. Outputs a message to the console.

`summary.ext_solutions_df`*Summary method for class ext_solutions_df*

Description

This summary function provides a summary of the `ext_solutions_df` class object, including the number of solutions, the distribution of the number of clusters, the number of features, the number of observations, and the distribution of p-values.

Usage

```
## S3 method for class 'ext_solutions_df'  
summary(object, ...)
```

Arguments

`object` A `ext_solutions_df` class object.
`...` Other arguments passed to `summary` (not used in this function).

Value

A named list containing the number of solutions, the distribution of the number of clusters, the number of features, the number of observations, and the distribution of p-values.

`summary.settings_df`*Summary method for class settings_df*

Description

This summary function provides a summary of the `settings_df` class object, including the number of settings, the distribution of alpha values, the distribution of k values, and the distribution of clustering functions.

Usage

```
## S3 method for class 'settings_df'  
summary(object, ...)
```

Arguments

`object` A `settings_df` class object.
`...` Other arguments passed to `summary` (not used in this function).

Value

A named list containing summary information of the settings data frame.

summary.sim_mats_list *Summary method for class sim_mats_list*

Description

This summary function simply returns to the console the number of functions contained in the sim_mats_list object.

Usage

```
## S3 method for class 'sim_mats_list'  
summary(object, ...)
```

Arguments

object	A sim_mats_list class object.
...	Other arguments passed to summary (not used in this function).

Value

Returns no value. Outputs a message to the console.

summary.snf_config *Summary method for class snf_config*

Description

This summary function provides a summary of the snf_config class object, including the settings data frame, clustering functions list, distance functions list, and weights matrix.

Usage

```
## S3 method for class 'snf_config'  
summary(object, ...)
```

Arguments

object	A snf_config class object.
...	Other arguments passed to summary (not used in this function).

Value

A named list containing the summaries of objects within the config.

summary.solutions_df *Summary method for class solutions_df*

Description

This summary function provides a summary of the solutions_df class object, including the number of solutions, the distribution of the number of clusters, and the number of observations.

Usage

```
## S3 method for class 'solutions_df'  
summary(object, ...)
```

Arguments

object A t_ext_solutions_df class object.
... Other arguments passed to summary (not used in this function).

Value

A named list containing the number of solutions, the distribution of the number of clusters, and the number of observations.

summary.t_ext_solutions_df
Summary method for class t_ext_solutions_df

Description

This summary function provides a summary of the t_ext_solutions_df class object, including the number of solutions, the distribution of the number of clusters, the number of features, the number of observations, and the distribution of p-values.

Usage

```
## S3 method for class 't_ext_solutions_df'  
summary(object, ...)
```

Arguments

object A t_ext_solutions_df class object.
... Other arguments passed to summary (not used in this function).

Value

A named list containing the number of solutions, the distribution of the number of clusters, the number of features, the number of observations, and the distribution of p-values.

summary.t_solutions_df

Summary method for class t_solutions_df

Description

This summary function provides a summary of the `t_solutions_df` class object, including the number of solutions, the distribution of the number of clusters, the number of features, the number of observations, and the distribution of p-values.

Usage

```
## S3 method for class 't_solutions_df'  
summary(object, ...)
```

Arguments

`object` A `t_solutions_df` class object.
`...` Other arguments passed to `summary` (not used in this function).

Value

A named list containing the number of solutions, the distribution of the number of clusters, the number of features, the number of observations, and the distribution of p-values.

summary.weights_matrix

Summary method for class weights_matrix

Description

This summary function provides a summary of the `weights_matrix` class object, including the minimum, maximum, mean, and standard deviation of the feature weights.

Usage

```
## S3 method for class 'weights_matrix'  
summary(object, ...)
```

Arguments

`object` A `weights_matrix` class object.
`...` Other arguments passed to `summary` (not used in this function).

Value

A named list containing the summary statistics of the weights matrix, the number of solutions, and the number of features.

train_test_assign	<i>Training and testing split</i>
-------------------	-----------------------------------

Description

Given a vector of uid_id and a threshold, returns a list of which members should be in the training set and which should be in the testing set. The function relies on whether or not the absolute value of the Jenkins's one_at_a_time hash function exceeds the maximum possible value (2147483647) multiplied by the threshold.

Usage

```
train_test_assign(train_frac, uids, seed = 42)
```

Arguments

train_frac	The fraction (0 to 1) of observations for training
uids	A character vector of UIDs to be distributed into training and test sets.
seed	Seed used for Jenkins's one_at_a_time hash function.

Value

A named list containing the training and testing uid_ids.

uids	<i>Pull UIDs from an object</i>
------	---------------------------------

Description

Pull UIDs from an object

Usage

```
uids(x)
```

Arguments

x	The object to extract UIDs from.
---	----------------------------------

Value

A character vector of UIDs.

validate_solutions_df *Validator for solutions_df class object*

Description

Validator for solutions_df class object

Usage

```
validate_solutions_df(sol_df1)
```

Arguments

sol_df1 A solutions data frame-like object to be validated and converted into a solutions data frame.

Value

If sol_df1 has a valid structure for a solutions_df class object, returns the input unchanged. Otherwise, raises an error.

var_manhattan_plot *Manhattan plot of feature-feature association p-values*

Description

Manhattan plot of feature-feature association p-values

Usage

```
var_manhattan_plot(  
  dl,  
  key_var,  
  neg_log_pval_thresh = 5,  
  threshold = NULL,  
  point_size = 5,  
  text_size = 20,  
  plot_title = NULL,  
  hide_x_labels = FALSE,  
  bonferroni_line = FALSE  
)
```

Arguments

dl	List of data frames containing data information.
key_var	Feature for which the association p-values of all other features are plotted.
neg_log_pval_thresh	Threshold for negative log p-values.
threshold	p-value threshold to plot dashed line at.
point_size	Size of points in the plot.
text_size	Size of text in the plot.
plot_title	Title of the plot.
hide_x_labels	If TRUE, hides x-axis labels.
bonferroni_line	If TRUE, plots a dashed black line at the Bonferroni-corrected equivalent of the p-value threshold.

Value

A Manhattan plot (class "gg", "ggplot") showing the association p-values of features against one key feature in a data list.

Examples

```
dl <- data_list(
  list(subc_v, "subcortical_volume", "neuroimaging", "continuous"),
  list(income, "household_income", "demographics", "continuous"),
  list(pubertal, "pubertal_status", "demographics", "continuous"),
  list(anxiety, "anxiety", "behaviour", "ordinal"),
  list(depress, "depressed", "behaviour", "ordinal"),
  uid = "unique_id"
)

var_manhattan <- var_manhattan_plot(
  dl,
  key_var = "household_income",
  plot_title = "Correlation of Features with Household Income",
  text_size = 16,
  neg_log_pval_thresh = 3,
  threshold = 0.05
)
```

Description

Function for building a weights matrix independently of an SNF config. The weights matrix contains one row corresponding to each row of the settings data frame in an SNF config (one row for each resulting cluster solution) and one column for each feature in the data list used for clustering. Values of the weights matrix are passed to distance metrics functions during the conversion of input data frames to distance matrices. Typically, there is no need to use this function directly. Instead, users should provide weights matrix-building parameters to the `snf_config()` function.

Usage

```
weights_matrix(dl = NULL, n_solutions = 1, weights_fill = "ones")
```

Arguments

<code>dl</code>	A nested list of input data from <code>data_list()</code> .
<code>n_solutions</code>	Number of rows to generate the template weights matrix for.
<code>weights_fill</code>	String indicating what to populate generate rows with. Can be "ones" (default; fill matrix with 1), "uniform" (fill matrix with uniformly distributed random values), or "exponential" (fill matrix with exponentially distributed random values).

Value

`wm` A properly formatted matrix containing columns for all the features that require weights and rows.

Examples

```
input_dl <- data_list(
  list(subc_v, "subcortical_volume", "neuroimaging", "continuous"),
  list(income, "income", "demographics", "continuous"),
  list(pubertal, "pubertal_status", "demographics", "continuous"),
  uid = "unique_id"
)

sc <- snf_config(input_dl, n_solutions = 5)

wm <- weights_matrix(input_dl, n_solutions = 5, weights_fill = "uniform")

# updating an SNF config in parts
sc$weights_matrix <- wm
```

Index

* datasets

- abcd_anxiety, 6
- abcd_colour, 7
- abcd_cort_sa, 8
- abcd_cort_t, 9
- abcd_depress, 10
- abcd_h_income, 11
- abcd_income, 11
- abcd_pubertal, 12
- abcd_subc_v, 13
- age_df, 17
- anxiety, 19
- cache_a_complete_example_ext_sol_df, 37
- cache_a_complete_example_lp_ext_sol_df, 38
- cache_a_complete_example_sol_df, 38
- cancer_diagnosis_df, 43
- cort_sa, 53
- cort_t, 53
- depress, 56
- diagnosis_df, 57
- expression_df, 65
- fav_colour, 67
- gender_df, 68
- income, 73
- methylation_df, 88
- mock_ari_matrix, 89
- mock_clust_fns_list, 89
- mock_data_list, 90
- mock_dist_fns_list, 90
- mock_ext_solutions_df, 91
- mock_mc_solutions_df, 91
- mock_rep_solutions_df, 92
- mock_settings_df, 92
- mock_snf_config, 93
- mock_solutions_df, 93
- mock_t_solutions_df, 94
- mock_weights_matrix, 94
- pubertal, 109
- subc_v, 137
- abcd_anxiety, 6
- abcd_colour, 7
- abcd_cort_sa, 8
- abcd_cort_t, 9
- abcd_depress, 10
- abcd_h_income, 11
- abcd_income, 11
- abcd_pubertal, 12
- abcd_subc_v, 13
- add_settings_df_rows, 14
- age_df, 17
- alluvial_cluster_plot, 17
- anxiety, 19
- as.data.frame.data_list, 20
- as.data.frame.ext_solutions_df, 20
- as.data.frame.settings_df, 21
- as.data.frame.snf_config, 22
- as.data.frame.solutions_df, 22
- as.data.frame.t_ext_solutions_df, 23
- as.data.frame.t_solutions_df, 24
- as.data.frame.weights_matrix, 24
- as.list.clust_fns_list, 25
- as.list.data_list, 25
- as.list.dist_fns_list, 26
- as.list.sim_mats_list, 26
- as.list.snf_config, 27
- as.matrix.ari_matrix, 27
- as.matrix.weights_matrix, 28
- as_ari_matrix, 30
- as_data_list, 31
- as_settings_df, 31
- as_sim_mats_list, 32
- as_snf_config, 32
- as_weights_matrix, 33
- assemble_data, 28
- assoc_pval_heatmap, 29

- auto_plot, 33
- bar_plot, 34
- batch_snf, 35
- batch_snf_subsamples, 36
- cache_a_complete_example_ext_sol_df, 37
- cache_a_complete_example_lp_ext_sol_df, 38
- cache_a_complete_example_sol_df, 38
- calc_aris, 40
- calc_assoc_pval_matrix, 41
- calc_nmis, 42
- calculate_coclustering, 39
- calculate_db_indices (quality_measures), 112
- calculate_dunn_indices (quality_measures), 112
- calculate_silhouettes (quality_measures), 112
- cancer_diagnosis_df, 43
- cell_significance_fn, 44
- check_dataless_annotations, 45
- check_hm_dependencies, 45
- check_similarity_matrices, 46
- clust_fns, 46
- clust_fns_list, 47
- cocluster_density, 48
- cocluster_heatmap, 50
- colour_scale, 52
- config_heatmap (plot.snf_config), 99
- cort_sa, 53
- cort_t, 53
- data_list, 54
- dendrogram, 99, 103
- depress, 56
- diagnosis_df, 57
- dist_fns, 58
- dist_fns_list, 59
- dlapply, 60
- dplyr_row_slice.ext_solutions_df, 61
- dplyr_row_slice.solutions_df, 62
- esm_manhattan_plot, 62
- estimate_nclust_given_graph, 64
- euclidean_distance (dist_fns), 58
- expression_df, 65
- extend_solutions, 66
- fav_colour, 67
- gender_df, 68
- get_complete_uids, 68
- get_heatmap_order, 69
- get_matrix_order, 70
- get_pvals, 71
- get_representative_solutions, 72
- gower_distance (dist_fns), 58
- gpar, 101
- hamming_distance (dist_fns), 58
- hclust, 99, 103
- income, 73
- is_data_list, 74
- jitter_plot, 75
- label_meta_clusters, 75
- label_propagate, 77
- linear_adjust, 79
- mc_manhattan_plot, 80
- merge.clust_fns_list, 82
- merge.data_list, 83
- merge.dist_fns_list, 83
- merge.ext_solutions_df, 84
- merge.settings_df, 84
- merge.sim_mats_list, 85
- merge.snf_config, 85
- merge.solutions_df, 86
- merge.t_ext_solutions_df, 86
- merge.t_solutions_df, 87
- merge.weights_matrix, 87
- merge_df_list, 88
- meta_cluster_heatmap (plot.ari_matrix), 95
- methylation_df, 88
- mock_ari_matrix, 89
- mock_clust_fns_list, 89
- mock_data_list, 90
- mock_dist_fns_list, 90
- mock_ext_solutions_df, 91
- mock_mc_solutions_df, 91
- mock_rep_solutions_df, 92
- mock_settings_df, 92
- mock_snf_config, 93

mock_solutions_df, 93
mock_t_solutions_df, 94
mock_weights_matrix, 94

new_solutions_df, 95

plot.ari_matrix, 95
plot.data_list, 97
plot.ext_solutions_df, 98
plot.settings_df (plot.snf_config), 99
plot.snf_config, 99
plot.solutions_df, 102
plot.t_ext_solutions_df
 (plot.ext_solutions_df), 98
plot.t_solutions_df
 (plot.solutions_df), 102
plot.weights_matrix (plot.snf_config),
 99

print.ari_matrix, 103
print.clust_fns_list, 104
print.data_list, 104
print.dist_fns_list, 105
print.ext_solutions_df, 105
print.settings_df, 106
print.sim_mats_list, 106
print.snf_config, 107
print.solutions_df, 107
print.t_ext_solutions_df, 108
print.t_solutions_df, 108
print.weights_matrix, 109
pubertal, 109
pval_heatmap, 110

quality_measures, 112

rbind.ext_solutions_df, 113
rbind.solutions_df, 114
rbind.t_solutions_df, 114
rbind.weights_matrix, 115
rename_dl, 115
resample, 116

save_heatmap, 117
settings_df, 117
sew_euclidean_distance (dist_fns), 58
shiny_annotator, 121
sim_mats_list, 125
similarity_matrix_heatmap, 122
siw_euclidean_distance, 125

sn_euclidean_distance (dist_fns), 58
snf_config, 126
spectral_eigen (clust_fns), 46
spectral_eigen_classic (clust_fns), 46
spectral_eight (clust_fns), 46
spectral_five (clust_fns), 46
spectral_four (clust_fns), 46
spectral_nine (clust_fns), 46
spectral_rot (clust_fns), 46
spectral_rot_classic (clust_fns), 46
spectral_seven (clust_fns), 46
spectral_six (clust_fns), 46
spectral_ten (clust_fns), 46
spectral_three (clust_fns), 46
spectral_two (clust_fns), 46
split_parser, 130
str.ari_matrix, 131
str.clust_fns_list, 132
str.data_list, 132
str.dist_fns_list, 133
str.ext_solutions_df, 133
str.settings_df, 134
str.sim_mats_list, 134
str.snf_config, 135
str.solutions_df, 135
str.t_ext_solutions_df, 136
str.t_solutions_df, 136
str.weights_matrix, 137
subc_v, 137
subsample_dl, 138
subsample_pairwise_aris, 139
summary.ari_matrix, 140
summary.clust_fns_list, 141
summary.data_list, 141
summary.dist_fns_list, 142
summary.ext_solutions_df, 143
summary.settings_df, 143
summary.sim_mats_list, 144
summary.snf_config, 144
summary.solutions_df, 145
summary.t_ext_solutions_df, 145
summary.t_solutions_df, 146
summary.weights_matrix, 146

train_test_assign, 147

uids, 147

validate_solutions_df, 148

`var_manhattan_plot`, [148](#)

`weights_matrix`, [149](#)