

# Package ‘metasnf’

November 9, 2024

**Title** Meta Clustering with Similarity Network Fusion

**Version** 1.1.2

**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** cluster, digest, dplyr, ggplot2, grDevices, MASS, mclust, methods, progressr, purrr, rlang, SNFtool, stats, tidyr, utils

**Suggests** circlize, ComplexHeatmap, InteractiveComplexHeatmap, clv, future, future.apply, knitr, rmarkdown, testthat (>= 3.0.0), ggalluvial, 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** 2024-11-08 23:20:02 UTC

## Contents

abcd_anxiety . . . . .	5
abcd_colour . . . . .	6
abcd_cort_sa . . . . .	7
abcd_cort_t . . . . .	8
abcd_depress . . . . .	9
abcd_h_income . . . . .	10
abcd_income . . . . .	10
abcd_pubertal . . . . .	11
abcd_subc_v . . . . .	12
add_columns . . . . .	13
add_settings_matrix_rows . . . . .	13
adjusted_rand_index_heatmap . . . . .	16
age_df . . . . .	17
alluvial_cluster_plot . . . . .	18
anxiety . . . . .	19
arrange_dl . . . . .	20
assemble_data . . . . .	20
assoc_pval_heatmap . . . . .	21

auto_plot . . . . .	22
bar_plot . . . . .	23
batch_nmi . . . . .	24
batch_row_closure . . . . .	25
batch_snf . . . . .	26
batch_snf_subsamples . . . . .	27
calculate_coclustering . . . . .	29
calculate_db_indices . . . . .	30
calculate_dunn_indices . . . . .	30
calculate_silhouettes . . . . .	31
calc_aris . . . . .	31
calc_assoc_pval . . . . .	32
calc_assoc_pval_matrix . . . . .	32
cancer_diagnosis_df . . . . .	33
cell_significance_fn . . . . .	34
char_to_fac . . . . .	34
check_dataless_annotations . . . . .	35
check_hm_dependencies . . . . .	35
check_similarity_matrices . . . . .	36
chi_squared_pval . . . . .	36
coclustering_coverage_check . . . . .	37
cocluster_density . . . . .	37
cocluster_heatmap . . . . .	38
collapse_dl . . . . .	39
colour_scale . . . . .	40
convert_uids . . . . .	40
cort_sa . . . . .	41
cort_t . . . . .	41
depress . . . . .	42
diagnosis_df . . . . .	43
discretisation . . . . .	44
discretisation_evec_data . . . . .	44
dl_has_duplicates . . . . .	45
dl_uid_first_col . . . . .	45
dl_variable_summary . . . . .	46
domains . . . . .	46
domain_merge . . . . .	47
drop_inputs . . . . .	48
esm_manhattan_plot . . . . .	48
estimate_nclust_given_graph . . . . .	49
euclidean_distance . . . . .	50
expression_df . . . . .	50
extend_solutions . . . . .	51
fav_colour . . . . .	52
fisher_exact_pval . . . . .	52
gender_df . . . . .	53
generate_annotations_list . . . . .	54
generate_clust_algs_list . . . . .	55

generate_data_list . . . . .	56
generate_distance_metrics_list . . . . .	58
generate_settings_matrix . . . . .	60
generate_weights_matrix . . . . .	63
get_clusters . . . . .	64
get_cluster_df . . . . .	65
get_cluster_solutions . . . . .	65
get_complete_uids . . . . .	66
get_dist_matrix . . . . .	66
get_dl_subjects . . . . .	67
get_heatmap_order . . . . .	68
get_matrix_order . . . . .	68
get_mean_pval . . . . .	69
get_min_pval . . . . .	69
get_pvals . . . . .	70
get_representative_solutions . . . . .	70
gower_distance . . . . .	71
hamming_distance . . . . .	72
income . . . . .	72
individual . . . . .	73
jitter_plot . . . . .	74
label_prop . . . . .	75
label_splits . . . . .	75
linear_adjust . . . . .	76
linear_model_pval . . . . .	76
list_remove . . . . .	77
lp_solutions_matrix . . . . .	77
mc_manhattan_plot . . . . .	78
merge_data_lists . . . . .	79
merge_df_list . . . . .	80
methylation_df . . . . .	80
no_subs . . . . .	81
numcol_to_numeric . . . . .	81
ord_reg_pval . . . . .	82
parallel_batch_snf . . . . .	82
prefix_dl_sk . . . . .	83
pubertal . . . . .	84
pval_heatmap . . . . .	84
random_removal . . . . .	86
reduce_dl_to_common . . . . .	87
remove_dl_na . . . . .	87
rename_dl . . . . .	88
reorder_dl_subs . . . . .	89
resample . . . . .	89
save_heatmap . . . . .	90
scale_diagonals . . . . .	90
settings_matrix_heatmap . . . . .	91
sew_euclidean_distance . . . . .	92

shiny_annotator . . . . .	92
similarity_matrix_heatmap . . . . .	93
similarity_matrix_path . . . . .	95
siw_euclidean_distance . . . . .	95
snf_step . . . . .	96
sn_euclidean_distance . . . . .	97
spectral_eigen . . . . .	97
spectral_eigen_classic . . . . .	98
spectral_eight . . . . .	98
spectral_five . . . . .	99
spectral_four . . . . .	99
spectral_nine . . . . .	100
spectral_rot . . . . .	100
spectral_rot_classic . . . . .	101
spectral_seven . . . . .	101
spectral_six . . . . .	102
spectral_ten . . . . .	102
spectral_three . . . . .	103
spectral_two . . . . .	103
split_parser . . . . .	104
subc_v . . . . .	104
subs . . . . .	105
subsample_data_list . . . . .	106
subsample_pairwise_aris . . . . .	106
summarize_clust_algs_list . . . . .	107
summarize_dl . . . . .	108
summarize_dml . . . . .	108
summarize_pvals . . . . .	109
train_test_assign . . . . .	109
two_step_merge . . . . .	110
var_manhattan_plot . . . . .	111

**Index****112**

---

`abcd_anxiety`*Mock ABCD anxiety data*

---

**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 DevelopmentSM (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/](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_cbcl_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 `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 DevelopmentSM (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/](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 (mm<sup>2</sup>, 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 DevelopmentSM (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/](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 (mm<sup>3</sup>, 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 DevelopmentSM (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/](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_cbc1_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 DevelopmentSM (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/](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 DevelopmentSM (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/](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 DevelopmentSM (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/](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 DevelopmentSM (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/](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 (mm<sup>3</sup>, 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 DevelopmentSM (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/](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\_columns

*Add columns to a dataframe*


---

### Description

Add new columns to a dataframe by providing a character vector of column names (param newcols) and a value to occupy each row of the new columns (param fill, NA by default).

### Usage

```
add_columns(df, newcols, fill = NA)
```

### Arguments

df	The dataframe to extend
newcols	The vector containing new column names
fill	The values of the elements of the newly added columns. NA by default.

### Value

extended\_df The dataframe containing the added columns

---

 add\_settings\_matrix\_rows

*Add settings matrix rows*


---

### Description

Add settings matrix rows

**Usage**

```

add_settings_matrix_rows(
  settings_matrix,
  seed = NULL,
  nrows = 0,
  min_removed_inputs = 0,
  max_removed_inputs = sum(startsWith(colnames(settings_matrix), "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,
  distance_metrics_list = NULL,
  snf_input_weights = NULL,
  snf_domain_weights = NULL,
  retry_limit = 10
)

```

**Arguments**

settings_matrix	The existing settings matrix
seed	set a seed for the random matrix generation. Setting this value will change the seed of the global environment.
nrows	Number of rows to generate for the settings matrix.
min_removed_inputs	The smallest number of input dataframes that may be randomly removed. By default, 0.
max_removed_inputs	The largest number of input dataframes that may be randomly removed. By default, this is 1 less than all the provided input dataframes in the data_list.
dropout_dist	Parameter controlling how the random removal of input dataframes should occur. Can be "none" (no input dataframes are randomly removed), "uniform" (uniformly sample between min_removed_inputs and max_removed_inputs to

	determine number of input dataframes to remove), or "exponential" (pick number of input dataframes to remove by sampling from min_removed_inputs to max_removed_inputs with an exponential distribution; 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 "twostep" 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 ?generate_clust_algs_list for more details on running custom clustering algorithms.
continuous_distances	A vector of continuous distance metrics to use when a custom distance_metrics_list is provided.

discrete_distances	A vector of categorical distance metrics to use when a custom distance_metrics_list is provided.
ordinal_distances	A vector of categorical distance metrics to use when a custom distance_metrics_list is provided.
categorical_distances	A vector of categorical distance metrics to use when a custom distance_metrics_list is provided.
mixed_distances	A vector of mixed distance metrics to use when a custom distance_metrics_list is provided.
distance_metrics_list	List containing distance metrics to vary over. See ?generate_distance_metrics_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.

**Value**

settings\_matrix A settings matrix

---

adjusted\_rand\_index\_heatmap

*Heatmap of pairwise adjusted rand indices between solutions*

---

**Description**

Heatmap of pairwise adjusted rand indices between solutions

**Usage**

```
adjusted_rand_index_heatmap(
  aris,
  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(aris), max(aris)), c(min_colour, max_colour)),
    ...
)

```

### Arguments

aris	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 matrix.

---

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 *Alluvial plot of patients across cluster counts and important features*

---

**Description**

Alluvial plot of patients across cluster counts and important features

**Usage**

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

**Arguments**

cluster_sequence	A list of clustering algorithms (typically, the same algorithm varied over different numbers of clusters).
similarity_matrix	A similarity matrix.
data_list	A nested list of input data from generate_data_list().
data	A dataframe that contains features to include in the plot.
key_outcome	The name of the feature that determines how each patient stream is coloured in the alluvial plot.
key_label	Name of key outcome to be used for the plot legend.
extra_outcomes	Names of additional features to add to the plot.
title	Title of the plot.

**Value**

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

---

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

---

### Description

Like the mock dataframe "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 DevelopmentSM (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/](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.

---

arrange_dl	<i>Given a data_list object, sort data elements by subjectkey</i>
------------	---

---

**Description**

Given a data\_list object, sort data elements by subjectkey

**Usage**

```
arrange_dl(data_list)
```

**Arguments**

data\_list      A nested list of input data from generate\_data\_list().

**Value**

arranged\_data\_list The arranged data\_list object

---

assemble_data	<i>Collapse a dataframe and/or a data_list into a single dataframe</i>
---------------	--

---

**Description**

Collapse a dataframe and/or a data\_list into a single dataframe

**Usage**

```
assemble_data(data, data_list)
```

**Arguments**

data            A dataframe.

data\_list      A nested list of input data from generate\_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,
  data_list = 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	The results of dl_var_summar - a dataframe that has the domain of every feature in the plotted data. columns of the correlation_matrix. Will be used to "slice" the heatmap into visually separated sections.
data_list	A nested list of input data from generate_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.

---

auto\_plot

*Automatically plot features across clusters*

---

**Description**

Given a single row of a solutions matrix and data provided through data\_list and/or target\_list arguments, this function will return a series of bar and/or jitter plots based on feature types.

**Usage**

```
auto_plot(
  solutions_matrix_row = NULL,
  data_list = NULL,
  cluster_df = NULL,
  target_list = NULL,
  return_plots = TRUE,
  save = NULL,
  jitter_width = 6,
  jitter_height = 6,
  bar_width = 6,
  bar_height = 6,
  verbose = FALSE
)
```

**Arguments**

solutions_matrix_row	A single row of a solutions matrix.
data_list	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.
target_list	A target_list containing data to plot.
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, print 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_nmi	<i>Calculate feature NMIs for a data_list and a derived solutions_matrix</i>
-----------	--

---

### Description

Calculate feature NMIs for a data\_list and a derived solutions\_matrix

### Usage

```
batch_nmi(
  data_list,
  solutions_matrix,
  clust_algs_list = NULL,
  distance_metrics_list = NULL,
  automatic_standard_normalize = FALSE,
  transpose = TRUE,
  ignore_inclusions = TRUE,
  verbose = FALSE
)
```

### Arguments

data_list	A nested list of input data from generate_data_list(). Use the same value as was used in the original call to batch_snf().
solutions_matrix	Result of batch_snf storing cluster solutions and the settings that were used to generate them. Use the same value as was used in the original call to batch_snf().
clust_algs_list	List of custom clustering algorithms to apply to the final fused network. See ?generate_clust_algs_list. Use the same value as was used in the original call to batch_snf().
distance_metrics_list	An optional nested list containing which distance metric function should be used for the various feature types (continuous, discrete, ordinal, categorical, and mixed). Use the same value as was used in the original call to batch_snf().
automatic_standard_normalize	If TRUE, will automatically apply standard normalization prior to calculation of any distance matrices. Use the same value as was used in the original call to batch_snf().
transpose	If TRUE, will transpose the output dataframe.
ignore_inclusions	If TRUE, will ignore the inclusion columns in the solutions matrix and calculate NMIs for all features. If FALSE, will give NAs for features that were dropped on a given settings_matrix row.
verbose	If TRUE, print progress to console.



**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 matrix. Populated values show the calculated NMI score for each feature-solution combination.

---

batch_row_closure	<i>Generate closure function to run batch_snf in an apply-friendly format</i>
-------------------	---

---

**Description**

Generate closure function to run batch\_snf in an apply-friendly format

**Usage**

```
batch_row_closure(
  data_list,
  distance_metrics_list,
  clust_algs_list,
  settings_matrix,
  weights_matrix,
  similarity_matrix_dir,
  return_similarity_matrices,
  prog
)
```

**Arguments**

data_list	A nested list of input data from generate_data_list().
distance_metrics_list	An optional nested list containing which distance metric function should be used for the various feature types (continuous, discrete, ordinal, categorical, and mixed). See ?generate_distance_metrics_list for details on how to build this.
clust_algs_list	List of custom clustering algorithms to apply to the final fused network. See ?generate_clust_algs_list.
settings_matrix	matrix indicating parameters to iterate SNF through.
weights_matrix	A matrix containing feature weights to use during distance matrix calculation. See ?generate_weights_matrix for details on how to build this.
similarity_matrix_dir	If specified, this directory will be used to save all generated similarity matrices.
return_similarity_matrices	If TRUE, function will return a list where the first element is the solutions matrix and the second element is a list of similarity matrices for each row in the solutions_matrix. Default FALSE.
prog	Progressr function to update parallel processing progress

**Value**

A "function" class object used to run batch\_snf in lapply-form for parallel processing.

---

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\_matrix (see ?generate\_settings\_matrix) and a data\_list (see ?generate\_data\_list), run repeated complete SNF pipelines to generate a broad space of post-SNF cluster solutions.

**Usage**

```
batch_snf(
  data_list,
  settings_matrix,
  processes = 1,
  return_similarity_matrices = FALSE,
  similarity_matrix_dir = NULL,
  clust_algs_list = NULL,
  suppress_clustering = FALSE,
  distance_metrics_list = NULL,
  weights_matrix = NULL,
  automatic_standard_normalize = FALSE,
  verbose = FALSE
)
```

**Arguments**

data_list	A nested list of input data from generate_data_list().
settings_matrix	A data.frame where each row completely defines an SNF pipeline transforming individual input dataframes into a final cluster solution. See ?generate_settings_matrix or <a href="https://branchlab.github.io/metasn timer/articles/settings_matrix.html">https://branchlab.github.io/metasn timer/articles/settings_matrix.html</a> for more details.
processes	Specify number of processes used to complete SNF iterations <ul style="list-style-type: none"> <li>• 1 (default) Sequential processing: function will iterate through the settings_matrix 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.</li> <li>• 2 or higher: Parallel processing will use the future.apply::future_apply to distribute the SNF iterations across the specified number of CPU cores. If higher than the number of available cores, a warning will be printed and the maximum number of cores will be used.</li> <li>• max: All available cores will be used.</li> </ul>

return_similarity_matrices	If TRUE, function will return a list where the first element is the solutions matrix and the second element is a list of similarity matrices for each row in the solutions_matrix. Default FALSE.
similarity_matrix_dir	If specified, this directory will be used to save all generated similarity matrices.
clust_algs_list	List of custom clustering algorithms to apply to the final fused network. See <code>?generate_clust_algs_list</code> .
suppress_clustering	If FALSE (default), will apply default or custom clustering algorithms to provide cluster solutions on every iteration of SNF. If TRUE, parameter <code>similarity_matrix_dir</code> must be specified.
distance_metrics_list	An optional nested list containing which distance metric function should be used for the various feature types (continuous, discrete, ordinal, categorical, and mixed). See <code>?generate_distance_metrics_list</code> for details on how to build this.
weights_matrix	A matrix containing feature weights to use during distance matrix calculation. See <code>?generate_weights_matrix</code> for details on how to build this.
automatic_standard_normalize	If TRUE, will automatically apply standard normalization prior to calculation of any distance matrices. This parameter cannot be used in conjunction with a custom distance metrics list. If you wish to supply custom distance metrics but also always have standard normalization, simply ensure that the numeric (continuous, discrete, and ordinal) distance metrics are only populated with distance metric functions that apply standard normalization. See <a href="https://branchlab.github.io/metasn timer/articles/distance_">https://branchlab.github.io/metasn timer/articles/distance_</a> to learn more.
verbose	If TRUE, print time remaining estimates to console.

### Value

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

---

`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(
  data_list_subsamples,
  settings_matrix,
  processes = 1,
  return_similarity_matrices = FALSE,
  clust_algs_list = NULL,
  suppress_clustering = FALSE,
  distance_metrics_list = NULL,
  weights_matrix = NULL,
  automatic_standard_normalize = FALSE,
  return_solutions_matrices = FALSE,
  verbose = FALSE
)
```

**Arguments**

`data_list_subsamples` A list of subsampled data lists. This object is generated by the function `batch_snf_subsamples()`.

`settings_matrix` A settings matrix defining the parameters of the SNF pipelines to be applied to the subsampled data lists.

`processes` See `?batch_snf`.

`return_similarity_matrices` See `?batch_snf`.

`clust_algs_list` See `?batch_snf`.

`suppress_clustering` See `?batch_snf`.

`distance_metrics_list` See `?batch_snf`.

`weights_matrix` See `?batch_snf`.

`automatic_standard_normalize` See `?batch_snf`.

`return_solutions_matrices` If TRUE, includes the solutions matrices corresponding to each subsample in the output.

`verbose` If TRUE, print time remaining estimates to console.

**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_similarity_matrices` and `return_solutions_matrices` to TRUE will turn the result of the function to a three-element list containing the corresponding solutions matrices and final fused similarity matrices of those cluster solutions, should you require these objects for your own stability calculations.

---

`calculate_coclustering`*Calculate coclustering data.*

---

**Description**

Calculate coclustering data.

**Usage**

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

**Arguments**

`subsample_solutions`

A list of containing cluster solutions from distinct subsamples of the data. This object is generated by the function `batch_snf_subsamples()`. These solutions should correspond to the ones in the solutions matrix.

`solutions_matrix`

A solutions matrix. This object is generated by the function `batch_snf()`. The solutions in the solutions matrix should correspond to those in the subsample solutions.

`verbose`

If TRUE, print time remaining estimates to console.

**Value**

A list containing the following components:

- `cocluster_dfs`: A list of dataframes, one per cluster solution, that shows the number of times that every pair of subjects 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 subjects occurred in the same subsample, formatted as a pairwise matrix.
- `cocluster_sc_mats`: The number of times every pair of subjects occurred in the same cluster, formatted as a pairwise matrix.
- `cocluster_cf_mats`: The fraction of times every pair of subjects occurred in the same cluster, formatted as a pairwise matrix.
- `cocluster_summary`: Specifically among pairs of subjects 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 dataframe with one row per cluster solution.

---

calculate\_db\_indices *Calculate Davies-Bouldin indices*

---

**Description**

Given a solutions\_matrix and a list of similarity\_matrices (or a single similarity\_matrix if the solutions\_matrix has only 1 row), return a vector of Davies-Bouldin indices

**Usage**

```
calculate_db_indices(solutions_matrix, similarity_matrices)
```

**Arguments**

solutions\_matrix

A solutions\_matrix (see ?batch\_snf)

similarity\_matrices

A list of similarity matrices (see ?batch\_snf)

**Value**

davies\_bouldin\_indices A vector of Davies-Bouldin indices for each cluster solution.

---

calculate\_dunn\_indices

*Calculate Dunn indices*

---

**Description**

Given a solutions\_matrix and a list of similarity\_matrices (or a single similarity\_matrix if the solutions\_matrix has only 1 row), return vector of Dunn indices

**Usage**

```
calculate_dunn_indices(solutions_matrix, similarity_matrices)
```

**Arguments**

solutions\_matrix

A solutions\_matrix (see ?batch\_snf)

similarity\_matrices

A list of similarity matrices (see ?batch\_snf)

**Value**

dunn\_indices A vector of Dunn indices for each cluster solution

---

calculate\_silhouettes *Calculate silhouette scores*

---

### Description

Given a solutions\_matrix and a list of similarity\_matrices (or a single similarity\_matrix if the solutions\_matrix has only 1 row), return a list of 'silhouette' objects from the cluster package

### Usage

```
calculate_silhouettes(solutions_matrix, similarity_matrices)
```

### Arguments

solutions\_matrix  
A solutions\_matrix (see ?batch\_snf)

similarity\_matrices  
A list of similarity matrices (see ?batch\_snf)

### Value

silhouette\_scores A list of "silhouette" objects from the cluster package.

---

calc\_aris *Meta-cluster calculations*

---

### Description

Generate matrix of pairwise cluster-solution similarities by Adjusted Rand index calculations.

### Usage

```
calc_aris(solutions_matrix, processes = 1, verbose = FALSE)
```

### Arguments

solutions\_matrix  
solutions\_matrix containing cluster solutions to calculate pairwise ARIs for.

processes  
Specify number of processes used to complete calculations

- 1 (default) Sequential processing
- 2 or higher: Parallel processing will use the future.apply::future\_apply to distribute the calculations across the specified number of CPU cores. If higher than the number of available cores, a warning will be printed 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, print progress to console.

**Value**

om\_aris ARIs between clustering solutions of an solutions matrix

---

calc_assoc_pval	<i>Calculate p-values based on feature vectors and their types</i>
-----------------	--

---

**Description**

Calculate p-values based on feature vectors and their types

**Usage**

```
calc_assoc_pval(var1, var2, type1, type2, cat_test = "chi_squared")
```

**Arguments**

var1	A single vector containing a feature.
var2	A single vector containing a feature.
type1	The type of var1 (continuous, discrete, ordinal, categorical).
type2	The type of var2 (continuous, discrete, ordinal, categorical).
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**

pval A p-value from a statistical test based on the provided types. Currently, this will either be the F-test p-value from a linear model if at least one feature is non-categorical, or the chi-squared test p-value if both features are categorical.

---

calc_assoc_pval_matrix	<i>Calculate p-values for all pairwise associations of features in a data_list</i>
------------------------	--

---

**Description**

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

**Usage**

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



**Arguments**

<code>data_list</code>	A nested list of input data from <code>generate_data_list()</code> .
<code>verbose</code>	If TRUE, prints new line everytime a p-value is being calculated.
<code>cat_test</code>	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.

---

`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.

---

char\_to\_fac                    *Convert character-type columns of a dataframe to factor-type*

---

**Description**

Convert character-type columns of a dataframe to factor-type

**Usage**

```
char_to_fac(df)
```

**Arguments**

df                    A dataframe

**Value**

df\_converted The dataframe with factor-type columns instead of char-type columns

---

`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 dataframe 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

---

chi\_squared\_pval

*Chi-squared test p-value (generic)*

---

### Description

Return p-value for chi-squared test for any two features

### Usage

```
chi_squared_pval(cat_var1, cat_var2)
```

### Arguments

cat\_var1 A categorical feature.

cat\_var2 A categorical feature.

### Value

pval A p-value (class "numeric").

---

coclustering\_coverage\_check  
*Coclustering coverage check*

---

**Description**

Check if coclustered data has at least one subsample in which every pair of subjects were a part of simultaneously.

**Usage**

```
coclustering_coverage_check(cocluster_df, action = "warn")
```

**Arguments**

`cocluster_df`     Dataframe containing coclustering data.  
`action`            Control if parent function should warn or stop.

**Value**

This function does not return any value. It checks a `cocluster_df` for complete coverage (all pairs occur in the same solution at least once). Will raise a warning or error if coverage is incomplete depending on the value of the `action` parameter.

---

`cocluster_density`     *Density plot coclustering stability across subsampled data.*

---

**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 dataframe containing coclustering 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 coclustering across pairs and subsamples of the data.

---

cocluster\_heatmap      *Heatmap of observation co-clustering across resampled data.*

---

### 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,
  data_list = 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 dataframe containing coclustering 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> .
data_list	See <code>?similarity_matrix_heatmap</code> .
data	See <code>?similarity_matrix_heatmap</code> .
left_bar	See <code>?similarity_matrix_heatmap</code> .
right_bar	See <code>?similarity_matrix_heatmap</code> .

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.

---

collapse_dl	<i>Collapse a data_list into a single dataframe</i>
-------------	---

---

**Description**

Collapse a data\_list into a single dataframe

**Usage**

```
collapse_dl(data_list)
```

**Arguments**

data\_list      A nested list of input data from generate\_data\_list().

**Value**

A "data.frame"-formatted version of the provided data list.

---

colour_scale	<i>Return a colour ramp for a given vector</i>
--------------	--

---

**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.

---

convert_uids	<i>Convert unique identifiers of data_list to 'subjectkey'</i>
--------------	--

---

**Description**

Column name "subjectkey" is reserved for the unique identifier of subjects. This function ensures all dataframes have their UID set as "subjectkey".

**Usage**

```
convert_uids(data_list, uid = NULL)
```

**Arguments**

data_list	A nested list of input data from <code>generate_data_list()</code> .
uid	(string) the name of the uid column currently used data

**Value**

dl\_renamed\_id data list with 'subjectkey' as UID



---

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

---

### Description

Like the mock dataframe "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 (mm<sup>2</sup>, 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 DevelopmentSM (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/](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 dataframe "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 DevelopmentSM (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/](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.

---

```
depress
```

```
Mock ABCD depression data
```

---

**Description**

Like the mock dataframe "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 DevelopmentSM (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/](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.

---

discretisation	<i>Internal function for estimate_nclust_given_graph</i>
----------------	--

---

**Description**

Internal function taken from SNF tool to use for number of cluster estimation.

**Usage**

```
discretisation(eigenvectors)
```

**Arguments**

eigenvectors    Matrix of eigenvectors.

**Value**

"Matrix" class object, intermediate product in spectral clustering.

---

discretisation_evec_data	<i>Internal function for estimate_nclust_given_graph</i>
--------------------------	--

---

**Description**

Internal function taken from SNF tool to use for number of cluster estimation.

**Usage**

```
discretisation_evec_data(eigenvector)
```

**Arguments**

eigenvector    Matrix of eigenvectors

**Value**

"Matrix" class object discretizing provided eigenvector to values 0 or 1.

---

dl_has_duplicates	<i>Check if data list contains any duplicate features</i>
-------------------	---

---

**Description**

Check if data list contains any duplicate features

**Usage**

```
dl_has_duplicates(data_list)
```

**Arguments**

data\_list      A nested list of input data from generate\_data\_list().

**Value**

Doesn't return any value. Raises warning if there are features with duplicate names in a generated data list.

---

dl_uid_first_col	<i>Make the subjectkey UID columns of a data_list first</i>
------------------	---

---

**Description**

Make the subjectkey UID columns of a data\_list first

**Usage**

```
dl_uid_first_col(data_list)
```

**Arguments**

data\_list      A nested list of input data from generate\_data\_list().

**Value**

A data list ("list"-class object) in which each data-subcomponent has "subjectkey" positioned as its first column.

---

<code>dl_variable_summary</code>	<i>Variable-level summary of a data_list</i>
----------------------------------	--

---

**Description**

Variable-level summary of a `data_list`

**Usage**

```
dl_variable_summary(data_list)
```

**Arguments**

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

**Value**

`variable_level_summary` A dataframe containing the name, type, and domain of every variable in a `data_list`.

---

<code>domains</code>	<i>Domains</i>
----------------------	----------------

---

**Description**

Domains

**Usage**

```
domains(data_list)
```

**Arguments**

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

**Value**

`domain_list` list of domains

---

`domain_merge`*SNF scheme: Domain merge*

---

### Description

Given a `data_list`, returns a new `data_list` where all data objects of a particular domain have been concatenated.

### Usage

```
domain_merge(  
  data_list,  
  cont_dist_fn,  
  disc_dist_fn,  
  ord_dist_fn,  
  cat_dist_fn,  
  mix_dist_fn,  
  weights_row,  
  k,  
  alpha,  
  t  
)
```

### Arguments

<code>data_list</code>	A nested list of input data from <code>generate_data_list()</code> .
<code>cont_dist_fn</code>	distance metric function for continuous data.
<code>disc_dist_fn</code>	distance metric function for discrete data.
<code>ord_dist_fn</code>	distance metric function for ordinal data.
<code>cat_dist_fn</code>	distance metric function for categorical data.
<code>mix_dist_fn</code>	distance metric function for mixed data.
<code>weights_row</code>	dataframe row containing feature weights.
<code>k</code>	k hyperparameter.
<code>alpha</code>	alpha/eta/sigma hyperparameter.
<code>t</code>	SNF number of iterations hyperparameter.

### Value

`fused_network` The final fused network (class "matrix", "array") generated by SNF.

---

drop_inputs	<i>Execute inclusion</i>
-------------	--------------------------

---

**Description**

Given a data list and a settings matrix row, returns a data list of selected inputs

**Usage**

```
drop_inputs(settings_matrix_row, data_list)
```

**Arguments**

settings_matrix_row	Row of a settings matrix.
data_list	A nested list of input data from generate_data_list().

**Value**

A data list (class "list") in which any component with a corresponding 0 value in the provided settings matrix row has been removed.

---

esm_manhattan_plot	<i>Manhattan plot of feature-cluster association p-values</i>
--------------------	---

---

**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 matrix 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 matrix.

---

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.

---

euclidean\_distance      *Distance metric: Euclidean distance*

---

### Description

Distance metric: Euclidean distance

### Usage

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

### Arguments

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

### Value

distance\_matrix A distance matrix.

---

expression\_df      *Modification of SNFtool mock dataframe "Data1"*

---

### Description

Modification of SNFtool mock dataframe "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 an solutions matrix to include outcome evaluations</i>
------------------	--

---

### Description

Extend an solutions matrix to include outcome evaluations

### Usage

```
extend_solutions(
  solutions_matrix,
  target_list = NULL,
  data_list = NULL,
  cat_test = "chi_squared",
  calculate_summaries = TRUE,
  min_pval = 1e-10,
  processes = 1,
  verbose = FALSE
)
```

### Arguments

solutions_matrix	Result of batch_snf storing cluster solutions and the settings that were used to generate them.
target_list	A data_list with features to calculate p-values for. Features in the target list will be included during p-value summary measure calculations.
data_list	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.
calculate_summaries	If TRUE, the function will calculate the minimum and mean p-values for each row of the solutions matrix.
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, print progress to console.

### Value

extended\_solutions\_matrix an extended solutions matrix that contains p-value columns for each outcome in the provided target\_list

---

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

---

### Description

Like the mock dataframe "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 DevelopmentSM (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/](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.

---

fisher_exact_pval	<i>Fisher exact test p-value</i>
-------------------	----------------------------------

---

### Description

Return p-value for Fisher exact test for any two features

**Usage**

```
fisher_exact_pval(cat_var1, cat_var2)
```

**Arguments**

cat\_var1        A categorical feature.

cat\_var2        A categorical feature.

**Value**

pval A p-value (class "numeric").

---

gender\_df                    *Mock gender data*

---

**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.

---

`generate_annotations_list`*Generate annotations list*

---

**Description**

Intermediate function that takes in formatted lists of features and the annotations they should be viewed through and returns annotation objects usable by `ComplexHeatmap::Heatmap`.

**Usage**

```
generate_annotations_list(  
  df,  
  left_bar = NULL,  
  right_bar = NULL,  
  top_bar = NULL,  
  bottom_bar = NULL,  
  left_hm = NULL,  
  right_hm = NULL,  
  top_hm = NULL,  
  bottom_hm = NULL,  
  show_legend = TRUE,  
  annotation_colours = NULL  
)
```

**Arguments**

<code>df</code>	Dataframe containing all the data that is specified in the remaining arguments.
<code>left_bar</code>	Named list of strings, where the strings are features in <code>df</code> 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.
<code>right_bar</code>	See <code>left_bar</code> .
<code>top_bar</code>	See <code>left_bar</code> .
<code>bottom_bar</code>	See <code>left_bar</code> .
<code>left_hm</code>	Like <code>left_bar</code> , but with a heatmap annotation instead of a barplot annotation.
<code>right_hm</code>	See <code>left_hm</code> .
<code>top_hm</code>	See <code>left_hm</code> .
<code>bottom_hm</code>	See <code>left_hm</code> .
<code>show_legend</code>	Add legends to the annotations.
<code>annotation_colours</code>	Named list of heatmap annotations and their colours.

**Value**

`annotations_list` A named list of all the annotations.

---

 generate\_clust\_algs\_list

*Generate a list of custom clustering algorithms*


---

### 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

```
generate_clust_algs_list(..., disable_base = FALSE)
```

### Arguments

... An arbitrary number of named clustering functions (see examples below)

disable\_base If TRUE, do not 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.

### 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
clust_algs_list <- generate_clust_algs_list()

# 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
clust_algs_list <- generate_clust_algs_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...
  # solution_data <- list(
  #   "solution" = solution,
  #   "nclust" = number_of_clusters
  # )
}
```

```

    # return(solution_data)
  }

  # Suppress the base algorithms-----
  # This will contain only user-provided clustering algorithms

  clust_algs_list <- generate_clust_algs_list(
    "two_cluster_spectral" = spectral_two,
    "five_cluster_spectral" = spectral_five,
    disable_base = TRUE
  )

```

---

generate\_data\_list      *Generate a data\_list*

---

### Description

This function generates the major data object that will be processed when iterating through the each SNF pipeline defined in the settings\_matrix. The data\_list is a named and nested list containing input dataframes (data), the name of that input dataframe (for the user's reference), the 'domain' of that dataframe (the broader source of information that the input dataframe is capturing, determined by user's domain knowledge), and the type of feature stored in the dataframe (continuous, discrete, ordinal, categorical, or mixed).

### Usage

```

generate_data_list(
  ...,
  uid = NULL,
  test_subjects = NULL,
  train_subjects = NULL,
  sort_subjects = TRUE,
  remove_missing = TRUE,
  return_missing = FALSE
)

```

### Arguments

...	Any number of list formatted as (df, "df_name", "df_domain", "df_type") OR any number of lists of lists formatted as (df, "df_name", "df_domain", "df_type")
uid	(string) the name of the uid column currently used data
test_subjects	character vector of test subjects (useful if building a full data list for label propagation)
train_subjects	character vector of train subjects (useful if building a full data list for label propagation)
sort_subjects	If TRUE, the subjects in the data_list will be sorted



- `remove_missing` If TRUE (default), subjects with incomplete data will be dropped from `data_list` creation. Setting this value to FALSE may lead to unusual and/or unstable results during SNF, clustering, p-value calculations or label propagation.
- `return_missing` If TRUE, function returns a list where the first element is the `data_list` and the second element is a vector of unique IDs of patients who were removed during the complete data filtration step.

### Value

A nested "list" class object. Each list component contains a 4-item list of a data frame, the user-assigned name of the data frame, the user-assigned domain of the data frame, and the user-labeled type of the 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):
data_list <- generate_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
data_list <- generate_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
summarize_dl(data_list)

# 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 <- generate_data_list(
  list_of_lists,
  uid = "patient_id"
)

```

---

```
generate_distance_metrics_list
```

*Generate a list of distance metrics*

---

## Description

This function can be used to specify custom distance metrics

## Usage

```
generate_distance_metrics_list(
  continuous_distances = NULL,
```

```

    discrete_distances = NULL,
    ordinal_distances = NULL,
    categorical_distances = NULL,
    mixed_distances = NULL,
    keep_defaults = TRUE
  )

```

### Arguments

continuous\_distances  
A named list of distance metric functions

discrete\_distances  
A named list of distance metric functions

ordinal\_distances  
A named list of distance metric functions

categorical\_distances  
A named list of distance metric functions

mixed\_distances  
A named list of distance metric functions

keep\_defaults If TRUE (default), prepend the base distance metrics (euclidean and standard normalized euclidean)

### Value

distance\_metrics\_list A well-formatted list of distance metrics

### Examples

```

# Using just the base distance metrics -----
distance_metrics_list <- generate_distance_metrics_list()

# Adding your own metrics -----
# This will contain the base and user-provided clustering algorithms
my_distance_metric <- function(df) {
  # your code that converts a dataframe to a distance metric here...
  # return(distance_metric)
}

distance_metrics_list <- generate_distance_metrics_list(
  continuous_distances = list(
    "my_distance_metric" = my_distance_metric
  )
)

# Suppress the base metrics-----
# This will contain only user-provided clustering algorithms

distance_metrics_list <- generate_distance_metrics_list(
  continuous_distances = list(
    "my_distance_metric" = my_distance_metric
  ),

```

```

discrete_distances = list(
    "my_distance_metric" = my_distance_metric
),
ordinal_distances = list(
    "my_distance_metric" = my_distance_metric
),
categorical_distances = list(
    "my_distance_metric" = my_distance_metric
),
mixed_distances = list(
    "my_distance_metric" = my_distance_metric
),
keep_defaults = FALSE
)

```

---

generate\_settings\_matrix

*Build a settings matrix*

---

## Description

The settings\_matrix is a dataframe whose rows completely specify the hyperparameters and decisions required to transform individual input dataframes (found in a data\_list, see ?generate\_data\_list) into a single similarity matrix through SNF. The format of the settings matrix is as follows:

- A column named "row\_id": 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 "clust\_alg": Specification of which clustering algorithm will be applied to the final similarity matrix to identify patient subtypes. 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/metasnfn/articles/clustering\\_algorithms.html](https://branchlab.github.io/metasnfn/articles/clustering_algorithms.html).
- A column named "cont\_dist": Specification of which distance metric will be used for dataframes of purely continuous data. You can learn about this metric and its defaults here: [https://branchlab.github.io/metasnfn/articles/continuous\\_distances.html](https://branchlab.github.io/metasnfn/articles/continuous_distances.html).
- A column named "disc\_dist": Like above, but for discrete dataframes.
- A column named "ord\_dist": Like above, but for ordinal dataframes.
- A column named "cat\_dist": Like above, but for categorical dataframes.
- A column named "mixed\_dist": Like above, but for mixed-type (e.g., both categorical and discrete) dataframes.

- One column for every input dataframe 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(data\_list)\$"name") of each dataframe. When 0, that dataframe will be excluded from that run of the SNF pipeline. When 1, that dataframe will be included.

## Usage

```
generate_settings_matrix(
  data_list,
  seed = NULL,
  nrows = 0,
  min_removed_inputs = 0,
  max_removed_inputs = length(data_list) - 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,
  distance_metrics_list = NULL,
  snf_input_weights = NULL,
  snf_domain_weights = NULL,
  retry_limit = 10
)
```

## Arguments

data_list	A nested list of input data from generate_data_list().
seed	(DEPRECATED) set the global seed. To ensure reproducible settings matrices are generated, manually call set.seed() prior to settings matrix generation instead of using this parameter.
nrows	Number of rows to generate for the settings matrix.
min_removed_inputs	The smallest number of input dataframes that may be randomly removed. By default, 0.

max_removed_inputs	The largest number of input dataframes that may be randomly removed. By default, this is 1 less than all the provided input dataframes in the data_list.
dropout_dist	Parameter controlling how the random removal of input dataframes should occur. Can be "none" (no input dataframes are randomly removed), "uniform" (uniformly sample between min_removed_inputs and max_removed_inputs to determine number of input dataframes to remove), or "exponential" (pick number of input dataframes 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 "twostep" 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

	<code>?generate_clust_algs_list</code> for more details on running custom clustering algorithms.
<code>continuous_distances</code>	A vector of continuous distance metrics to use when a custom <code>distance_metrics_list</code> is provided.
<code>discrete_distances</code>	A vector of categorical distance metrics to use when a custom <code>distance_metrics_list</code> is provided.
<code>ordinal_distances</code>	A vector of categorical distance metrics to use when a custom <code>distance_metrics_list</code> is provided.
<code>categorical_distances</code>	A vector of categorical distance metrics to use when a custom <code>distance_metrics_list</code> is provided.
<code>mixed_distances</code>	A vector of mixed distance metrics to use when a custom <code>distance_metrics_list</code> is provided.
<code>distance_metrics_list</code>	List containing distance metrics to vary over. See <code>?generate_distance_metrics_list</code> .
<code>snf_input_weights</code>	Nested list containing weights for when SNF is used to merge individual input measures (see <code>?generate_snf_weights</code> )
<code>snf_domain_weights</code>	Nested list containing weights for when SNF is used to merge domains (see <code>?generate_snf_weights</code> )
<code>retry_limit</code>	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.

**Value**

`settings_matrix` A settings matrix

---

`generate_weights_matrix`

*Generate a matrix to store feature weights*

---

**Description**

Generate a matrix to store feature weights

**Usage**

```
generate_weights_matrix(data_list = NULL, nrow = 1, fill = "ones")
```

**Arguments**

data_list	A nested list of input data from generate_data_list().
nrow	Number of rows to generate the template weights matrix for.
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**

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

---

get_clusters	<i>Extract cluster membership vector from one solutions matrix row</i>
--------------	--

---

**Description**

This function takes in a single row of a solutions matrix and returns a vector containing the cluster assignments for each observation. It is similar to get\_cluster\_df(), which takes a solutions matrix with only one row and returns a dataframe with two columns: "cluster" and "subjectkey" (the UID of the observation) and get\_cluster\_solutions(), which takes a solutions matrix with any number of rows and returns a dataframe indicating the cluster assignments for each of those rows.

**Usage**

```
get_clusters(solutions_matrix_row)
```

**Arguments**

solutions_matrix_row	Output matrix row.
----------------------	--------------------

**Value**

clusters Vector of assigned clusters.



---

get_cluster_df	<i>Extract cluster membership information from one solutions matrix row</i>
----------------	---

---

**Description**

This function takes in a single row of a solutions matrix and returns a dataframe containing the cluster assignments for each subjectkey. It is similar to `get_clusters()`, which takes one solutions matrix row and returns a vector of cluster assignments' and `get_cluster_solutions()`, which takes a solutions matrix with any number of rows and returns a dataframe indicating the cluster assignments for each of those rows.

**Usage**

```
get_cluster_df(solutions_matrix_row)
```

**Arguments**

solutions\_matrix\_row  
One row from a solutions matrix.

**Value**

cluster\_df dataframe of cluster and subjectkey.

---

get_cluster_solutions	<i>Extract cluster membership information from a solutions_matrix</i>
-----------------------	---

---

**Description**

This function takes in a solutions matrix and returns a dataframe containing the cluster assignments for each subjectkey. It is similar to `'get_clusters()`, which takes one solutions matrix row and returns a vector of cluster assignments' and `get_cluster_df()`, which takes a solutions matrix with only one row and returns a dataframe with two columns: "cluster" and "subjectkey" (the UID of the observation).

**Usage**

```
get_cluster_solutions(solutions_matrix)
```

**Arguments**

solutions\_matrix  
A solutions\_matrix.

**Value**

cluster\_solutions A "data.frame" object where each row is an observation and each column (apart from the subjectkey column) indicates the cluster that observation was assigned to for the corresponding solutions matrix row.

---

get\_complete\_uids      *Pull complete-data UIDs from a list of dataframes*

---

**Description**

Pull complete-data UIDs from a list of dataframes

**Usage**

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

**Arguments**

list\_of\_dfs      List of dataframes.  
uid                Name of column across dataframes containing UIDs

**Value**

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

---

get\_dist\_matrix      *Calculate distance matrices*

---

**Description**

Given a dataframe of numerical features, return a euclidean distance matrix.

**Usage**

```
get_dist_matrix(  
  df,  
  input_type,  
  cont_dist_fn,  
  disc_dist_fn,  
  ord_dist_fn,  
  cat_dist_fn,  
  mix_dist_fn,  
  weights_row  
)
```

**Arguments**

df	Raw dataframe with subject IDs in column "subjectkey"
input_type	Either "numeric" (resulting in euclidean distances), "categorical" (resulting in binary distances), or "mixed" (resulting in gower distances)
cont_dist_fn	distance metric function for continuous data
disc_dist_fn	distance metric function for discrete data
ord_dist_fn	distance metric function for ordinal data
cat_dist_fn	distance metric function for categorical data
mix_dist_fn	distance metric function for mixed data
weights_row	Single-row dataframe where the column names contain the column names in df and the row contains the corresponding weights_row.

**Value**

dist\_matrix Matrix of inter-observation distances.

---

get\_dl\_subjects      *Extract subjects from a data\_list*

---

**Description**

Extract subjects from a data\_list

**Usage**

```
get_dl_subjects(data_list, prefix = FALSE)
```

**Arguments**

data_list	A nested list of input data from generate_data_list().
prefix	If TRUE, preserves the "subject_" prefix added to UIDs when creating a data_list.

**Value**

A character vector of the UID labels contained in a data list.

---

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      *Return the hierarchical clustering order of a matrix*

---

### 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 apply to the matrix. Argument is directly passed into stats::dist. Options include "euclidean", "maximum", "manhattan", "canberra", "binary", or "minkowski".
hclust_method	Which agglomerative method to be passed into 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.

---

`get_mean_pval`                      *Get mean p-value*

---

**Description**

Given an solutions matrix row containing evaluated p-values, returns mean.

**Usage**

```
get_mean_pval(solutions_matrix_row)
```

**Arguments**

```
solutions_matrix_row  
                            row of solutions_matrix object
```

**Value**

mean\_pval mean p-value

---

`get_min_pval`                      *Get minimum p-value*

---

**Description**

Given an solutions matrix row containing evaluated p-values, returns min.

**Usage**

```
get_min_pval(solutions_matrix_row)
```

**Arguments**

```
solutions_matrix_row  
                            row of solutions_matrix object
```

**Value**

min\_pval minimum p-value

---

get_pvals	<i>Get p-values from an extended solutions matrix</i>
-----------	---

---

### Description

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

### Usage

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

### Arguments

`extended_solutions_matrix` The output of `extend_solutions`. A dataframe 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 `extended_solutions_matrix`.

---

get_representative_solutions	<i>Extract representative solutions from a matrix of ARIs</i>
------------------------------	---

---

### 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 `adjusted_rand_index_heatmap()` 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,
  split_vector,
  order,
  solutions_matrix,
  filter_fn = NULL
)
```

**Arguments**

`aris` Matrix of adjusted rand indices from `calc_aris()`

`split_vector` A vector of partition indices.

`order` Numeric vector indicating row ordering of settings matrix.

`solutions_matrix` Output of `batch_snf` containing cluster solutions.

`filter_fn` 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 `fn <- function(x) x[x$"nclust" == 4, ]`.

**Value**

A "data.frame" class object corresponding to a subset of the provided solutions matrix in which only one row is present per meta cluster.

---

<code>gower_distance</code>	<i>Distance metric: Gower distance</i>
-----------------------------	--

---

**Description**

Distance metric: Gower distance

**Usage**

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

**Arguments**

`df` Dataframe containing at least 1 data column.

`weights_row` For compatibility - function does not accept weights.

**Value**

`distance_matrix` A distance matrix.

---

hamming_distance	<i>Distance metric: Hamming distance</i>
------------------	--

---

**Description**

Distance metric: Hamming distance

**Usage**

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

**Arguments**

df	Dataframe containing one subjectkey column in the first column and at least 1 categorical data column. All feature data should be categorical.
weights_row	Single-row dataframe where the column names contain the column names in df and the row contains the corresponding weights.

**Value**

distance\_matrix A distance matrix.

---

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

---

**Description**

Like the mock dataframe "abcd\_h\_income", but with "unique\_id" as the "uid".  
 Like the mock dataframe "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 DevelopmentSM (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/](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 DevelopmentSM (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/](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.

---

 individual

*SNF Scheme: Individual*


---

**Description**

The "vanilla" scheme - does distance matrix conversions of each input dataframe in a list and

**Usage**

```
individual(
  data_list,
  cont_dist_fn,
```

```

    disc_dist_fn,
    ord_dist_fn,
    cat_dist_fn,
    mix_dist_fn,
    weights_row,
    k,
    alpha,
    t
)

```

### Arguments

<code>data_list</code>	A nested list of input data from <code>generate_data_list()</code> .
<code>cont_dist_fn</code>	distance metric function for continuous data.
<code>disc_dist_fn</code>	distance metric function for discrete data.
<code>ord_dist_fn</code>	distance metric function for ordinal data.
<code>cat_dist_fn</code>	distance metric function for categorical data.
<code>mix_dist_fn</code>	distance metric function for mixed data.
<code>weights_row</code>	dataframe row containing feature weights.
<code>k</code>	k hyperparameter.
<code>alpha</code>	alpha/eta/sigma hyperparameter.
<code>t</code>	SNF number of iterations hyperparameter.

### Value

`fused_network` The final fused network (class "matrix", "array") generated by SNF.

---

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

---

### Description

Jitter plot separating a feature by cluster

### Usage

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

### Arguments

<code>df</code>	A data.frame containing cluster column and the feature to plot.
<code>feature</code>	The feature to plot.

### Value

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

---

label_prop	<i>Label propagation</i>
------------	--------------------------

---

**Description**

Given a full fused network (one containing both pre-labeled training subjects and unlabeled test-subjects) and the clusters of the pre-labeled subjects, return a label propagated list of clusters for all subjects. This function is derived from `SNFtool::groupPredict`. Modifications are made to take a full fused network as input, rather than taking input dataframes and running SNF internally. This ensures that alternative approaches to data normalization and distance matrix calculations can be chosen by the user.

**Usage**

```
label_prop(full_fused_network, clusters)
```

**Arguments**

full_fused_network	Network made by running SNF on training and test subjects together
clusters	a vector of training subject assigned clusters in matching order as they appear in full_fused_network

**Value**

new\_clusters list of cluster labels for all subjects

---

label_splits	<i>Convert a vector of partition indices into meta cluster labels</i>
--------------	---

---

**Description**

Convert a vector of partition indices into meta cluster labels

**Usage**

```
label_splits(split_vector, nrow)
```

**Arguments**

split_vector	A vector of partition indices.
nrow	The number of rows in the data being partitioned.

**Value**

A character vector that expands the `split_vector` into an `nrow`-length sequence of ascending letters of the alphabet. If the `split_vector` is `c(3, 6)` and the number of rows is 8, the result will be a vector of two "A"s (up to the first index, 3), three "B"s (up to the second index, 6), and three "C"s (up to and including the last index, 8).

---

<code>linear_adjust</code>	<i>Linearly correct data_list by features with unwanted signal</i>
----------------------------	--

---

**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(data_list, unwanted_signal_list, sig_digs = NULL)
```

**Arguments**

<code>data_list</code>	A nested list of input data from <code>generate_data_list()</code> .
<code>unwanted_signal_list</code>	A <code>data_list</code> of categorical features that should have their mean differences removed in the first <code>data_list</code> .
<code>sig_digs</code>	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`.

---

<code>linear_model_pval</code>	<i>Linear model p-value (generic)</i>
--------------------------------	---------------------------------------

---

**Description**

Return p-value of F-test for a linear model of any two features

**Usage**

```
linear_model_pval(predictor, response)
```

**Arguments**

predictor      A categorical or numeric feature.  
 response      A numeric feature.

**Value**

pval A p-value (class "numeric").

---

list_remove	<i>Remove items from a data_list</i>
-------------	--------------------------------------

---

**Description**

Removes specified elements from a provided data\_list

**Usage**

```
list_remove(list_object, ...)
```

**Arguments**

list\_object    The data\_list containing components to be removed  
 ...            Any number of components to remove from the list object, passed as strings

**Value**

A "list"-class object in which any specified elements have been removed.

---

lp_solutions_matrix	<i>Label propagate cluster solutions to unclustered subjects</i>
---------------------	--

---

**Description**

Given a solutions\_matrix derived from training subjects and a full\_data\_list containing both training and test subjects, re-run SNF to generate a total affinity matrix of both train and subjects and use the label propagation algorithm to assigned predicted clusters to test subjects.

**Usage**

```
lp_solutions_matrix(  
  train_solutions_matrix,  
  full_data_list,  
  distance_metrics_list = NULL,  
  weights_matrix = NULL,  
  verbose = FALSE  
)
```

**Arguments**

- `train_solutions_matrix` A `solutions_matrix` derived from the training set. The propagation algorithm is slow and should be used for validating a top or top few meaningful chosen clustering solutions. It is advisable to use only a small subset of rows from the original training `solutions_matrix` for label propagation.
- `full_data_list` A `data_list` containing subjects from both the training and testing sets.
- `distance_metrics_list` Like above - the `distance_metrics_list` (if any) that was used for the original `batch_snf` call.
- `weights_matrix` Like above.
- `verbose` If TRUE, print progress to console.

**Value**

`labeled_df` a dataframe containing a column for subjectkeys, a column for whether the subject was in the train (original) or test (held out) set, and one column per row of the solutions matrix indicating the original and propagated clusters.

---

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

---

**Description**

Given a dataframe 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(
  extended_solutions_matrix,
  data_list = NULL,
  target_list = 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**

extended_solutions_matrix	A solutions_matrix that contains "_pval" columns containing the values to be plotted. This object is the output of extend_solutions().
data_list	List of dataframes containing data information.
target_list	List of dataframes 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 matrix, stratified by meta cluster label.

---

merge_data_lists	<i>Horizontally merge compatible data lists</i>
------------------	---

---

**Description**

Join two data\_lists with the same components (dataframes) but separate observations. To instead merge two data\_lists that have the same observations but different components, simply use c().

**Usage**

```
merge_data_lists(data_list1, data_list2)
```

**Arguments**

data_list1	The first data_list to merge.
data_list2	The second data_list to merge.

**Value**

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

---

merge_df_list	<i>Merge list of dataframes</i>
---------------	---------------------------------

---

**Description**

Merge list of dataframes

**Usage**

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

**Arguments**

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

**Value**

merged\_df inner join of all dataframes in list

---

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

---

**Description**

Modification of SNFtool mock dataframe "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.



---

no_subs	<i>Select all columns of a dataframe not starting with the 'subject_' prefix.</i>
---------	---

---

**Description**

Removes the 'subject\_' prefixed columns from a dataframe. Useful for printing solutions\_matrix structures to the console.

**Usage**

```
no_subs(df)
```

**Arguments**

df                    A dataframe

**Value**

df\_no\_subs Dataframe without subjects

---

numcol_to_numeric	<i>Convert dataframe columns to numeric type</i>
-------------------	--

---

**Description**

Converts all columns in a dataframe that can be converted to numeric type to numeric type.

**Usage**

```
numcol_to_numeric(df)
```

**Arguments**

df                    A dataframe

**Value**

df The dataframe with all possible columns converted to type numeric

---

ord_reg_pval	<i>Ordinal regression p-value</i>
--------------	-----------------------------------

---

**Description**

Returns the overall p-value of an ordinal regression on a categorical predictor and response vectors. If the ordinal response

**Usage**

```
ord_reg_pval(predictor, response)
```

**Arguments**

predictor	A categorical or numeric feature.
response	A numeric feature.

**Value**

pval A p-value (class "numeric").

---

parallel_batch_snf	<i>Parallel processing form of batch_snf</i>
--------------------	--

---

**Description**

Parallel processing form of batch\_snf

**Usage**

```
parallel_batch_snf(  
  data_list,  
  distance_metrics_list,  
  clust_algs_list,  
  settings_matrix,  
  weights_matrix,  
  similarity_matrix_dir,  
  return_similarity_matrices,  
  processes  
)
```

**Arguments**

- `data_list` A nested list of input data from `generate_data_list()`.
- `distance_metrics_list`  
An optional nested list containing which distance metric function should be used for the various feature types (continuous, discrete, ordinal, categorical, and mixed). See `?generate_distance_metrics_list` for details on how to build this.
- `clust_algs_list`  
List of custom clustering algorithms to apply to the final fused network. See `?generate_clust_algs_list`.
- `settings_matrix`  
matrix indicating parameters to iterate SNF through.
- `weights_matrix` A matrix containing feature weights to use during distance matrix calculation. See `?generate_weights_matrix` for details on how to build this.
- `similarity_matrix_dir`  
If specified, this directory will be used to save all generated similarity matrices.
- `return_similarity_matrices`  
If TRUE, function will return a list where the first element is the solutions matrix and the second element is a list of similarity matrices for each row in the solutions\_matrix. Default FALSE.
- `processes` Number of parallel processes used when executing SNF.

**Value**

The same values as `?batch_snf()`.

---

<code>prefix_dl_sk</code>	<i>Add "subject_" prefix to all UID values in subjectkey column</i>
---------------------------	---

---

**Description**

Add "subject\_" prefix to all UID values in subjectkey column

**Usage**

```
prefix_dl_sk(data_list)
```

**Arguments**

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

**Value**

`data_list` A `data_list` without NAs

---

pubertal	<i>Mock ABCD pubertal status data</i>
----------	---------------------------------------

---

### Description

Like the mock dataframe "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 DevelopmentSM (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/](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	<i>Heatmap of p-values</i>
--------------	----------------------------

---

### Description

Heatmap of p-values

**Usage**

```

pval_heatmap(
  pvals,
  order = NULL,
  cluster_columns = TRUE,
  cluster_rows = FALSE,
  show_row_names = FALSE,
  show_column_names = TRUE,
  min_colour = "red2",
  mid_colour = "lightyellow",
  max_colour = "slateblue4",
  legend_breaks = c(0, 0.5, 1),
  col = circlize::colorRamp2(legend_breaks, c(min_colour, mid_colour, max_colour)),
  heatmap_legend_param = list(color_bar = "continuous", title = "p-value", at = c(0, 0.5,
    1)),
  rect_gp = grid::gpar(col = "black"),
  column_split_vector = NULL,
  row_split_vector = NULL,
  column_split = NULL,
  row_split = NULL,
  ...
)

```

**Arguments**

pvals	A matrix of p-values.
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.
mid_colour	Colour used for the middle 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.

---

random_removal	<i>Generate random removal sequence</i>
----------------	---

---

**Description**

Helper function to contribute to rows within the settings matrix. Number of columns removed follows a uniform or exponential probability distribution.

**Usage**

```
random_removal(
  columns,
  min_removed_inputs,
  max_removed_inputs,
  dropout_dist = "exponential"
)
```

**Arguments**

columns Columns of the settings\_matrix that are passed in  
 min\_removed\_inputs The smallest number of input dataframes that may be randomly removed.  
 max\_removed\_inputs The largest number of input dataframes that may be randomly removed.  
 dropout\_dist Indication of how input dataframes should be dropped. can be "none" (no dropout), "uniform" (uniformly draw number between min and max removed inputs), or "exponential" (like uniform, but using an exponential distribution; default).

**Value**

inclusions\_df Dataframe that can be rbind'ed to the settings\_matrix

---

reduce\_dl\_to\_common      *Reduce data\_list to common subjects*

---

**Description**

Given a data\_list object, reduce each nested dataframe to contain only the set of subjects that are shared by all nested dataframes

**Usage**

```
reduce_dl_to_common(data_list)
```

**Arguments**

data\_list      A nested list of input data from generate\_data\_list().

**Value**

reduced\_data\_list The data\_list object subsetted only to subjects shared across all nested dataframes

---

remove\_dl\_na      *Remove NAs from a data\_list object*

---

**Description**

Remove NAs from a data\_list object

**Usage**

```
remove_dl_na(data_list, return_missing = FALSE)
```

**Arguments**

data\_list      A nested list of input data from generate\_data\_list().

return\_missing      If TRUE, function returns a list where the first element is the data\_list and the second element is a vector of unique IDs of patients who were removed during the complete data filtration step.

**Value**

data\_list A data\_list without NAs

---

rename_dl	<i>Rename features in a data_list</i>
-----------	---------------------------------------

---

### Description

Rename features in a data\_list

### Usage

```
rename_dl(data_list, name_mapping)
```

### Arguments

data_list	A nested list of input data from generate_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

```
library(metastnf)

data_list <- generate_data_list(
  list(pubertal, "pubertal_status", "demographics", "continuous"),
  list(anxiety, "anxiety", "behaviour", "ordinal"),
  list(depress, "depressed", "behaviour", "ordinal"),
  uid = "unique_id"
)

summarize_dl(data_list, "feature")

name_changes <- c(
  "anxiety_score" = "cbcl_anxiety_r",
  "depression_score" = "cbcl_depress_r"
)

data_list <- rename_dl(data_list, name_changes)

summarize_dl(data_list, "feature")
```



---

reorder_dl_subs	<i>Reorder the subjects in a data_list</i>
-----------------	--

---

**Description**

Reorder the subjects in a data\_list

**Usage**

```
reorder_dl_subs(data_list, ordered_subjects)
```

**Arguments**

`data_list`        A nested list of input data from `generate_data_list()`.  
`ordered_subjects`        A vector of the subjectkey values in the data\_list in the desired order of the sorted data\_list.

**Value**

A data list ("list"-class object) with reordered observations.

---

resample	<i>Helper resample function found in ?sample</i>
----------	--

---

**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.

---

scale_diagonals	<i>Adjust the diagonals of a matrix</i>
-----------------	---

---

**Description**

Adjust the diagonals of a matrix to reduce contrast with off-diagonals during plotting.

**Usage**

```
scale_diagonals(matrix, method = "mean")
```

**Arguments**

matrix	Matrix to rescale.
method	Method of rescaling. Can be: <ul style="list-style-type: none"> <li>• "mean" (replace diagonals with average value of off-diagonals)</li> <li>• "zero" (replace diagonals with 0)</li> <li>• "min" (replace diagonals with min value of off-diagonals)</li> <li>• "max" (replace diagonals with max value of off-diagonals)</li> </ul>

**Value**

A "matrix" class object with rescaled diagonals.

---

 settings\_matrix\_heatmap

*Heatmap for visualizing a settings matrix*


---

### Description

Scales settings matrix values between 0 and 1 and plots as a heatmap. Rows can be reordered to match prior meta clustering results.

### Usage

```
settings_matrix_heatmap(
  settings_matrix,
  order = NULL,
  remove_fixed_columns = TRUE,
  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,
  ...
)
```

### Arguments

settings_matrix	Matrix indicating parameters to iterate SNF through.
order	Numeric vector indicating row ordering of settings matrix.
remove_fixed_columns	Whether columns that have no variation should be removed.
show_column_names	Whether column names should be shown.
show_row_names	Whether row names should be shown.
rect_gp	Cell border function for ComplexHeatmap::Heatmap.
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	Standard parameter of ComplexHeatmap::Heatmap.
row_split	Standard parameter of ComplexHeatmap::Heatmap.
column_title	Standard parameter of ComplexHeatmap::Heatmap.
...	Additional parameters passed to ComplexHeatmap::Heatmap.

**Value**

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

---

sew\_euclidean\_distance  
*Squared (excluding weights) Euclidean distance*

---

**Description**

Squared (excluding weights) Euclidean distance

**Usage**

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

**Arguments**

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

**Value**

distance\_matrix A distance matrix.

---

shiny\_annotator      *Launch shiny app to identify meta cluster boundaries*

---

**Description**

Launch shiny app to identify meta cluster boundaries

**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.

---

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_list = 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,  
  split_vector = NULL,  
  row_split = NULL,  
  column_split = NULL,  
  ...  
)
```

**Arguments**

similarity\_matrix  
A similarity matrix

<code>order</code>	Vector of numbers to reorder the similarity matrix (and data if provided). Overwrites ordering specified by <code>cluster_solution</code> param.
<code>cluster_solution</code>	Vector containing cluster assignments.
<code>scale_diag</code>	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).
<code>log_graph</code>	If TRUE, log transforms the graph.
<code>cluster_rows</code>	Parameter for <code>ComplexHeatmap::Heatmap</code> .
<code>cluster_columns</code>	Parameter for <code>ComplexHeatmap::Heatmap</code> .
<code>show_row_names</code>	Parameter for <code>ComplexHeatmap::Heatmap</code> .
<code>show_column_names</code>	Parameter for <code>ComplexHeatmap::Heatmap</code> .
<code>data_list</code>	A nested list of input data from <code>generate_data_list()</code> .
<code>data</code>	A dataframe containing elements requested for annotation.
<code>left_bar</code>	Named list of strings, where the strings are features in <code>df</code> 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.
<code>right_bar</code>	See <code>left_bar</code> .
<code>top_bar</code>	See <code>left_bar</code> .
<code>bottom_bar</code>	See <code>left_bar</code> .
<code>left_hm</code>	Like <code>left_bar</code> , but with a heatmap annotation instead of a barplot annotation.
<code>right_hm</code>	See <code>left_hm</code> .
<code>top_hm</code>	See <code>left_hm</code> .
<code>bottom_hm</code>	See <code>left_hm</code> .
<code>annotation_colours</code>	Named list of heatmap annotations and their colours.
<code>min_colour</code>	Colour used for the lowest value in the heatmap.
<code>max_colour</code>	Colour used for the highest value in the heatmap.
<code>split_vector</code>	A vector of partition indices.
<code>row_split</code>	Standard parameter of <code>ComplexHeatmap::Heatmap</code> .
<code>column_split</code>	Standard parameter of <code>ComplexHeatmap::Heatmap</code> .
<code>...</code>	Additional parameters passed into <code>ComplexHeatmap::Heatmap</code> .

**Value**

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

---

`similarity_matrix_path`*Generate a complete path and filename to store an similarity matrix*

---

**Description**

Generate a complete path and filename to store an similarity matrix

**Usage**

```
similarity_matrix_path(similarity_matrix_dir, i)
```

**Arguments**

<code>similarity_matrix_dir</code>	Directory to store similarity matrices
<code>i</code>	Corresponding settings matrix row

**Value**

path Complete path and filename to store an similarity matrix

---

`siw_euclidean_distance`*Squared (including weights) Euclidean distance*

---

**Description**

Squared (including weights) Euclidean distance

**Usage**

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

**Arguments**

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

**Value**

distance\_matrix A distance matrix.

---

snf_step	<i>Convert a data list to a similarity matrix through a variety of SNF schemes</i>
----------	--

---

### Description

Convert a data list to a similarity matrix through a variety of SNF schemes

### Usage

```
snf_step(
  data_list,
  scheme,
  k = 20,
  alpha = 0.5,
  t = 20,
  cont_dist_fn,
  disc_dist_fn,
  ord_dist_fn,
  cat_dist_fn,
  mix_dist_fn,
  weights_row
)
```

### Arguments

<code>data_list</code>	A nested list of input data from <code>generate_data_list()</code> .
<code>scheme</code>	Which SNF system to use to achieve the final fused network.
<code>k</code>	<code>k</code> hyperparameter.
<code>alpha</code>	alpha/eta/sigma hyperparameter.
<code>t</code>	SNF number of iterations hyperparameter.
<code>cont_dist_fn</code>	distance metric function for continuous data.
<code>disc_dist_fn</code>	distance metric function for discrete data.
<code>ord_dist_fn</code>	distance metric function for ordinal data.
<code>cat_dist_fn</code>	distance metric function for categorical data.
<code>mix_dist_fn</code>	distance metric function for mixed data.
<code>weights_row</code>	dataframe row containing feature weights.

### Value

`fused_network` The final fused network (class "matrix", "array") generated by SNF.



---

sn\_euclidean\_distance *Distance metric: Standard normalization then Euclidean*

---

**Description**

Distance metric: Standard normalization then Euclidean

**Usage**

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

**Arguments**

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

**Value**

distance\_matrix A distance matrix.

---

spectral\_eigen *Clustering algorithm: Spectral clustering with eigen-gap heuristic*

---

**Description**

Applies spectral clustering to similarity matrix. Number of clusters is based on the eigen-gap heuristic.

**Usage**

```
spectral_eigen(similarity_matrix)
```

**Arguments**

similarity_matrix	A similarity matrix.
-------------------	----------------------

**Value**

solution\_data A list storing cluster assignments ("solution") and the number of clusters ("nclust").

---

spectral\_eigen\_classic

*Clustering algorithm: Spectral clustering with eigen-gap heuristic*

---

### Description

Applies spectral clustering to similarity matrix. Number of clusters is based on the eigen-gap heuristic. Range of possible cluster solutions is fixed between 2 and 5 inclusive.

### Usage

```
spectral_eigen_classic(similarity_matrix)
```

### Arguments

similarity\_matrix  
A similarity matrix.

### Value

solution\_data A list storing cluster assignments ("solution") and the number of clusters ("nclust").

---

spectral\_eight

*Clustering algorithm: Spectral clustering for a eight cluster solution*

---

### Description

Applies spectral clustering to similarity matrix. Seeks eight clusters.

### Usage

```
spectral_eight(similarity_matrix)
```

### Arguments

similarity\_matrix  
A similarity matrix.

### Value

solution\_data A list storing cluster assignments ("solution") and the number of clusters ("nclust").

---

spectral_five	<i>Clustering algorithm: Spectral clustering for a five cluster solution</i>
---------------	--

---

**Description**

Applies spectral clustering to similarity matrix. Seeks five clusters.

**Usage**

```
spectral_five(similarity_matrix)
```

**Arguments**

```
similarity_matrix  
    A similarity matrix.
```

**Value**

```
solution_data A list storing cluster assignments ("solution") and the number of clusters ("nclust").
```

---

spectral_four	<i>Clustering algorithm: Spectral clustering for a four cluster solution</i>
---------------	--

---

**Description**

Applies spectral clustering to similarity matrix. Seeks four clusters.

**Usage**

```
spectral_four(similarity_matrix)
```

**Arguments**

```
similarity_matrix  
    A similarity matrix.
```

**Value**

```
solution_data A list storing cluster assignments ("solution") and the number of clusters ("nclust").
```

---

spectral_nine	<i>Clustering algorithm: Spectral clustering for a nine cluster solution</i>
---------------	--

---

**Description**

Applies spectral clustering to similarity matrix. Seeks nine clusters.

**Usage**

```
spectral_nine(similarity_matrix)
```

**Arguments**

similarity\_matrix  
A similarity matrix.

**Value**

solution\_data A list storing cluster assignments ("solution") and the number of clusters ("nclust").

---

spectral_rot	<i>Clustering algorithm: Spectral clustering with rotation cost heuristic</i>
--------------	---

---

**Description**

Applies spectral clustering to similarity matrix. Number of clusters is based on the rotation cost heuristic.

**Usage**

```
spectral_rot(similarity_matrix)
```

**Arguments**

similarity\_matrix  
A similarity matrix.

**Value**

solution\_data A list storing cluster assignments ("solution") and the number of clusters ("nclust").

---

spectral\_rot\_classic    *Clustering algorithm: Spectral clustering with rotation cost heuristic*

---

**Description**

Applies spectral clustering to similarity matrix. Number of clusters is based on the rotation cost heuristic.

**Usage**

```
spectral_rot_classic(similarity_matrix)
```

**Arguments**

similarity\_matrix  
    A similarity matrix.

**Value**

solution\_data A list storing cluster assignments ("solution") and the number of clusters ("nclust").

---

spectral\_seven            *Clustering algorithm: Spectral clustering for a seven cluster solution*

---

**Description**

Applies spectral clustering to similarity matrix. Seeks seven clusters.

**Usage**

```
spectral_seven(similarity_matrix)
```

**Arguments**

similarity\_matrix  
    A similarity matrix.

**Value**

solution\_data A list storing cluster assignments ("solution") and the number of clusters ("nclust").

---

`spectral_six`*Clustering algorithm: Spectral clustering for a six cluster solution*

---

**Description**

Applies spectral clustering to similarity matrix. Seeks six clusters.

**Usage**

```
spectral_six(similarity_matrix)
```

**Arguments**

```
similarity_matrix  
    A similarity matrix.
```

**Value**

```
solution_data A list storing cluster assignments ("solution") and the number of clusters ("nclust").
```

---

`spectral_ten`*Clustering algorithm: Spectral clustering for a ten cluster solution*

---

**Description**

Applies spectral clustering to similarity matrix. Seeks ten clusters.

**Usage**

```
spectral_ten(similarity_matrix)
```

**Arguments**

```
similarity_matrix  
    A similarity matrix.
```

**Value**

```
solution_data A list storing cluster assignments ("solution") and the number of clusters ("nclust").
```

---

spectral_three	<i>Clustering algorithm: Spectral clustering for a three cluster solution</i>
----------------	---

---

**Description**

Applies spectral clustering to similarity matrix. Seeks three clusters.

**Usage**

```
spectral_three(similarity_matrix)
```

**Arguments**

similarity\_matrix  
A similarity matrix.

**Value**

solution\_data A list storing cluster assignments ("solution") and the number of clusters ("nclust").

---

spectral_two	<i>Clustering algorithm: Spectral clustering for a two cluster solution</i>
--------------	---

---

**Description**

Applies spectral clustering to similarity matrix. Seeks two clusters.

**Usage**

```
spectral_two(similarity_matrix)
```

**Arguments**

similarity\_matrix  
A similarity matrix.

**Value**

solution\_data A list storing cluster assignments ("solution") and the number of clusters ("nclust").

---

split_parser	<i>Helper function to determine which row and columns to split on</i>
--------------	---

---

**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.

---

subc_v	<i>Mock ABCD subcortical volumes data</i>
--------	---

---

**Description**

Like the mock dataframe "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 (mm<sup>3</sup>, 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 DevelopmentSM (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/](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.

---

subs

*Select all columns of a dataframe starting with a given string prefix.*

---

**Description**

Removes the columns that are not prefixed with 'subject\_' prefixed columns from a dataframe. Useful intermediate step for extracting subject UIDs from an solutions\_matrix structure.

**Usage**

subs(df)

**Arguments**

df                      Dataframe

**Value**

df\_subs Dataframe with only 'subject\_' prefixed columns

---

subsample\_data\_list     *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).

### Usage

```
subsample_data_list(
  data_list,
  n_subsamples,
  subsample_fraction = NULL,
  n_subjects = NULL
)
```

### Arguments

`data_list`     A nested list of input data from `generate_data_list()`.  
`n_subsamples`     Number of subsamples to create.  
`subsample_fraction`     Percentage of patients to include per subsample.  
`n_subjects`     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.

---

subsample\_pairwise\_aris     *Calculate pairwise adjusted Rand indices across subsamples of data*

---

### Description

Calculate pairwise adjusted Rand indices across subsamples of data

### Usage

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

**Arguments**

- `subsample_solutions` A list of containing cluster solutions from distinct subsamples of the data. This object is generated by the function `batch_snf_subsamples()`.
- `return_raw_aris` Whether the ARI matrix used to calculate the average ARI across subsamples should be returned.
- `verbose` If TRUE, print time remaining estimates to console.

**Value**

If `return_raw_aris` is FALSE, this function will return

---

`summarize_clust_algs_list`

*Summarize a clust\_algs\_list object*

---

**Description**

Summarize a `clust_algs_list` object

**Usage**

```
summarize_clust_algs_list(clust_algs_list)
```

**Arguments**

- `clust_algs_list`  
A `clust_algs_list` object

**Value**

summary\_df "data.frame" class object containing the name and index of each clustering algorithm in te provided `clust_algs_list`.

---

summarize_dl	<i>Summarize a data list</i>
--------------	------------------------------

---

**Description**

Summarize a data list

**Usage**

```
summarize_dl(data_list, scope = "component")
```

**Arguments**

data_list	A nested list of input data from generate_data_list().
scope	The level of detail for the summary. Options are: <ul style="list-style-type: none"><li>• "component" (default): One row per component (dataframe) in the data_list.</li><li>• "feature": One row for each feature in the data_list.</li></ul>

**Value**

"data.frame"-class object summarizing all components (or features if scope == "component").

---

summarize_dml	<i>Summarize metrics contained in a distance_metrics_list</i>
---------------	---

---

**Description**

Summarize metrics contained in a distance\_metrics\_list

**Usage**

```
summarize_dml(distance_metrics_list)
```

**Arguments**

distance_metrics_list	A distance_metrics_list.
-----------------------	--------------------------

**Value**

"data.frame"-class object summarizing items in a distance metrics list.

---

summarize_pvals	<i>Summarize p-value columns of an extended solutions matrix</i>
-----------------	--

---

**Description**

Summarize p-value columns of an extended solutions matrix

**Usage**

```
summarize_pvals(extended_solutions_matrix)
```

**Arguments**

extended\_solutions\_matrix  
Result of extend\_solutions

**Value**

The provided extended solutions matrix along with columns for the min, mean, and maximum across p-values for each row.

---

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

---

**Description**

Given a vector of subject\_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, subjects, seed = 42)
```

**Arguments**

train\_frac      The fraction (0 to 1) of subjects for training  
subjects        The available subjects for distribution  
seed            Seed used for Jenkins's one\_at\_a\_time hash function

**Value**

split a named list containing the training and testing subject\_ids

---

two_step_merge	<i>Two step SNF</i>
----------------	---------------------

---

### Description

Individual dataframes into individual similarity matrices into one fused network per domain into one final fused network.

### Usage

```
two_step_merge(  
  data_list,  
  k = 20,  
  alpha = 0.5,  
  t = 20,  
  cont_dist_fn,  
  disc_dist_fn,  
  ord_dist_fn,  
  cat_dist_fn,  
  mix_dist_fn,  
  weights_row  
)
```

### Arguments

<code>data_list</code>	A nested list of input data from <code>generate_data_list()</code> .
<code>k</code>	<code>k</code> hyperparameter.
<code>alpha</code>	alpha/eta/sigma hyperparameter.
<code>t</code>	SNF number of iterations hyperparameter.
<code>cont_dist_fn</code>	distance metric function for continuous data.
<code>disc_dist_fn</code>	distance metric function for discrete data.
<code>ord_dist_fn</code>	distance metric function for ordinal data.
<code>cat_dist_fn</code>	distance metric function for categorical data.
<code>mix_dist_fn</code>	distance metric function for mixed data.
<code>weights_row</code>	dataframe row containing feature weights.

### Value

`fused_network` The final fused network (class "matrix", "array") generated by SNF.

---

var\_manhattan\_plot      *Manhattan plot of feature-feature associaiton p-values*

---

### Description

Manhattan plot of feature-feature associaiton p-values

### Usage

```
var_manhattan_plot(  
  data_list,  
  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

data_list	List of dataframes 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.

# Index

## \* datasets

- abcd\_anxiety, 5
  - abcd\_colour, 6
  - abcd\_cort\_sa, 7
  - abcd\_cort\_t, 8
  - abcd\_depress, 9
  - abcd\_h\_income, 10
  - abcd\_income, 10
  - abcd\_pubertal, 11
  - abcd\_subc\_v, 12
  - age\_df, 17
  - anxiety, 19
  - cancer\_diagnosis\_df, 33
  - cort\_sa, 41
  - cort\_t, 41
  - depress, 42
  - diagnosis\_df, 43
  - expression\_df, 50
  - fav\_colour, 52
  - gender\_df, 53
  - income, 72
  - methylation\_df, 80
  - pubertal, 84
  - subc\_v, 104
- 
- abcd\_anxiety, 5
  - abcd\_colour, 6
  - abcd\_cort\_sa, 7
  - abcd\_cort\_t, 8
  - abcd\_depress, 9
  - abcd\_h\_income, 10
  - abcd\_income, 10
  - abcd\_pubertal, 11
  - abcd\_subc\_v, 12
  - add\_columns, 13
  - add\_settings\_matrix\_rows, 13
  - adjusted\_rand\_index\_heatmap, 16
  - age\_df, 17
  - alluvial\_cluster\_plot, 18
  - anxiety, 19
  - arrange\_dl, 20
  - assemble\_data, 20
  - assoc\_pval\_heatmap, 21
  - auto\_plot, 22
  - bar\_plot, 23
  - batch\_nmi, 24
  - batch\_row\_closure, 25
  - batch\_snf, 26
  - batch\_snf\_subsamples, 27
  - calc\_aris, 31
  - calc\_assoc\_pval, 32
  - calc\_assoc\_pval\_matrix, 32
  - calculate\_coclustering, 29
  - calculate\_db\_indices, 30
  - calculate\_dunn\_indices, 30
  - calculate\_silhouettes, 31
  - cancer\_diagnosis\_df, 33
  - cell\_significance\_fn, 34
  - char\_to\_fac, 34
  - check\_dataless\_annotations, 35
  - check\_hm\_dependencies, 35
  - check\_similarity\_matrices, 36
  - chi\_squared\_pval, 36
  - cocluster\_density, 37
  - cocluster\_heatmap, 38
  - coclustering\_coverage\_check, 37
  - collapse\_dl, 39
  - colour\_scale, 40
  - convert\_uids, 40
  - cort\_sa, 41
  - cort\_t, 41
  - depress, 42
  - diagnosis\_df, 43
  - discretisation, 44
  - discretisation\_evec\_data, 44
  - dl\_has\_duplicates, 45
  - dl\_uid\_first\_col, 45



dl\_variable\_summary, 46  
domain\_merge, 47  
domains, 46  
drop\_inputs, 48  
  
esm\_manhattan\_plot, 48  
estimate\_nclust\_given\_graph, 49  
euclidean\_distance, 50  
expression\_df, 50  
extend\_solutions, 51  
  
fav\_colour, 52  
fisher\_exact\_pval, 52  
  
gender\_df, 53  
generate\_annotations\_list, 54  
generate\_clust\_algs\_list, 55  
generate\_data\_list, 56  
generate\_distance\_metrics\_list, 58  
generate\_settings\_matrix, 60  
generate\_weights\_matrix, 63  
get\_cluster\_df, 65  
get\_cluster\_solutions, 65  
get\_clusters, 64  
get\_complete\_uids, 66  
get\_dist\_matrix, 66  
get\_dl\_subjects, 67  
get\_heatmap\_order, 68  
get\_matrix\_order, 68  
get\_mean\_pval, 69  
get\_min\_pval, 69  
get\_pvals, 70  
get\_representative\_solutions, 70  
gower\_distance, 71  
  
hamming\_distance, 72  
  
income, 72  
individual, 73  
  
jitter\_plot, 74  
  
label\_prop, 75  
label\_splits, 75  
linear\_adjust, 76  
linear\_model\_pval, 76  
list\_remove, 77  
lp\_solutions\_matrix, 77  
  
mc\_manhattan\_plot, 78  
  
merge\_data\_lists, 79  
merge\_df\_list, 80  
methylation\_df, 80  
  
no\_subs, 81  
numcol\_to\_numeric, 81  
  
ord\_reg\_pval, 82  
  
parallel\_batch\_snf, 82  
prefix\_dl\_sk, 83  
pubertal, 84  
pval\_heatmap, 84  
  
random\_removal, 86  
reduce\_dl\_to\_common, 87  
remove\_dl\_na, 87  
rename\_dl, 88  
reorder\_dl\_subs, 89  
resample, 89  
  
save\_heatmap, 90  
scale\_diagonals, 90  
settings\_matrix\_heatmap, 91  
sew\_euclidean\_distance, 92  
shiny\_annotator, 92  
similarity\_matrix\_heatmap, 93  
similarity\_matrix\_path, 95  
siw\_euclidean\_distance, 95  
sn\_euclidean\_distance, 97  
snf\_step, 96  
spectral\_eigen, 97  
spectral\_eigen\_classic, 98  
spectral\_eight, 98  
spectral\_five, 99  
spectral\_four, 99  
spectral\_nine, 100  
spectral\_rot, 100  
spectral\_rot\_classic, 101  
spectral\_seven, 101  
spectral\_six, 102  
spectral\_ten, 102  
spectral\_three, 103  
spectral\_two, 103  
split\_parser, 104  
subc\_v, 104  
subs, 105  
subsample\_data\_list, 106  
subsample\_pairwise\_aris, 106

[summarize\\_clust\\_algs\\_list](#), 107  
[summarize\\_dl](#), 108  
[summarize\\_dml](#), 108  
[summarize\\_pvals](#), 109  
  
[train\\_test\\_assign](#), 109  
[two\\_step\\_merge](#), 110  
  
[var\\_manhattan\\_plot](#), 111