# Package 'maaslin3'

October 24, 2025

**Title** ``Refining and extending generalized multivariate linear models for meta-omic association discovery"

Year 2025

Version 1.1.2

**Depends** R (>= 4.4)

**Description** MaAsLin 3 refines and extends generalized multivariate linear models for meta-omicron association discovery. It finds abundance and prevalence associations between microbiome meta-omics features and complex metadata in population-scale epidemiological studies. The software includes multiple analysis methods (including support for multiple covariates, repeated measures, and ordered predictors), filtering, normalization, and transform options to customize analysis for your specific study.

License MIT + file LICENSE

Imports dplyr, plyr, pbapply, lmerTest, parallel, lme4, optparse, logging, multcomp, ggplot2, RColorBrewer, patchwork, scales, rlang, tibble, ggnewscale, survival, methods, BiocGenerics, SummarizedExperiment, TreeSummarizedExperiment

Suggests knitr, testthat (>= 2.1.0), rmarkdown, markdown, kableExtra

VignetteBuilder knitr

Collate fit.R utility\_scripts.R viz.R maaslin3.R

URL http://huttenhower.sph.harvard.edu/maaslin3

**biocViews** Metagenomics, Software, Microbiome, Normalization, MultipleComparison

BugReports https://github.com/biobakery/maaslin3/issues

NeedsCompilation no

git\_url https://git.bioconductor.org/packages/maaslin3

git\_branch devel

git\_last\_commit e7106d5

git\_last\_commit\_date 2025-10-12

Repository Bioconductor 3.23

Date/Publication 2025-10-24

2 maaslin3

Author William Nickols [aut, cre] (ORCID: <a href="https://orcid.org/0000-0001-8214-9746">https://orcid.org/0000-0001-8214-9746</a>), Jacob Nearing [aut]

Maintainer William Nickols <willnickols@g.harvard.edu>

# **Contents**

	maaslin3	2
	maaslin_check_arguments	8
	maaslin_check_formula	10
	maaslin_compute_formula	12
	maaslin_contrast_test	14
	maaslin_filter	17
	maaslin_fit	19
	maaslin_log_arguments	24
	maaslin_normalize	29
	maaslin_plot_results	31
	maaslin_plot_results_from_output	35
	maaslin_process_metadata	39
	maaslin_read_data	41
	maaslin_reorder_data	43
	maaslin_transform	45
	maaslin_write_results	47
	maaslin_write_results_lefse_format	49
	preprocess_dna_mtx	51
	preprocess_taxa_mtx	53
Index		55

maaslin3

MaAsLin 3: A multivariable statistical framework for finding abundance and prevalence associations between metadata and high-dimensional microbial multi-omics data.

# Description

This wrapper for all MaAsLin 3 steps finds abundance and prevalence associations between microbiome meta-omics features and complex metadata in population-scale epidemiological studies. The software includes multiple analysis methods (including support for multiple covariates, repeated measures, and ordered predictors), filtering, normalization, and transform options to customize analysis for your specific study.

maaslin3 3

## Usage

```
maaslin3(input_data,
    input_metadata = NULL,
    output,
    formula = NULL,
    fixed_effects = NULL,
    reference = NULL,
    random_effects = NULL,
    group_effects = NULL,
    ordered_effects = NULL,
    strata_effects = NULL,
    feature_specific_covariate = NULL,
    feature_specific_covariate_name = NULL,
    feature_specific_covariate_record = NULL,
    min_abundance = 0,
    min_prevalence = 0,
    max_prevalence = 1.01,
    zero_threshold = 0,
    min_variance = 0,
    max_significance = 0.1,
    normalization = 'TSS',
    transform = 'LOG',
    correction = 'BH',
    standardize = TRUE,
    unscaled_abundance = NULL,
    median_comparison_abundance = TRUE,
    median_comparison_prevalence = FALSE,
    median_comparison_abundance_threshold = 0,
    median_comparison_prevalence_threshold = 0,
    subtract_median = FALSE,
    warn_prevalence = TRUE,
    small_random_effects = FALSE,
    augment = TRUE,
    evaluate_only = NULL,
    plot_summary_plot = TRUE,
    summary_plot_first_n = 25,
    coef_plot_vars = NULL,
    heatmap_vars = NULL,
    plot_associations = TRUE,
    max_pngs = 30,
    cores = 1,
    save_models = FALSE,
    save_plots_rds = FALSE,
    verbosity = 'FINEST',
    summary_plot_balanced = FALSE,
    assay.type = 1)
```

4 maaslin3

#### **Arguments**

input\_data

A data frame of feature abundances or read counts, a filepath to a tab-delimited file with abundances, or a SummarizedExperiment or TreeSummarizedExperiment object with the taxa table in 'assays' and metadata in 'colData'. If a data frame or a filepath is supplied, the table should be formatted with features as columns and samples as rows (or the transpose). The column and row names should be the feature names and sample names respectively.

input\_metadata A data frame of per-sample metadata or a filepath to a tab-delimited file with metadata. It should be formatted with variables as columns and samples as rows (or the transpose). The column and row names should be the variable names and sample names respectively.

output

The output folder to write results.

formula

A formula in 1me4 format. Random effects, interactions, and functions of the metadata can be included (note that these functions will be applied after standardization if standardize=TRUE). Group, ordered, and strata variables can be specified as: group(grouping\_variable), ordered(ordered\_variable) and strata(strata\_variable). The other variable options below will not be considered if a formula is set.

fixed effects

A vector of variable names to be included as fixed effects.

reference

For a variable with more than two levels supplied with fixed\_effects, the factor to use as a reference provided as a string of 'variable, reference' semicolon delimited for multiple variables.

random\_effects A vector of variable names to be included as random intercepts.

group\_effects

A factored categorical variable to be included for group testing. An ANOVAstyle test will be performed to assess whether any of the variable's levels are significant, and no coefficients or individual p-values will be returned.

ordered\_effects

A factored categorical variable to be included. Consecutive levels will be tested for significance against each other, and the resulting associations will correspond to effect sizes, standard errors, and significances of each level versus the previous.

strata\_effects A vector with one variable name to be included as the strata variable in casecontrol studies. Strata cannot be combined with random effects.

feature\_specific\_covariate

A table of feature-specific covariates or a filepath to a tab-delimited file with feature-specific covariates. It should be formatted with features as columns and samples as rows (or the transpose). The row names and column names should be the same as those of the input\_data: the column and row names should be the feature names and sample names respectively. Typically, this table should be generated by 'preprocess\_mgx\_mtx' or 'preprocess\_taxa\_mtx' first.

feature\_specific\_covariate\_name

The name for the feature-specific covariates when fitting the models. This string must be parse-able in a formula (e.g., no spaces).

feature\_specific\_covariate\_record

Whether to keep the feature-specific covariates in the outputs when calculating p-values, writing results, and displaying plots.

maaslin3 5

min\_abundance Features with abundances more than min\_abundance in more than min\_prevalence

of the samples will be included for analysis. The threshold is applied after nor-

malization and before transformation.

min\_prevalence See min\_abundance.

max\_prevalence Features with abundances more than min\_abundance in fewer than max\_prevalence

of the samples will be included for analysis. The threshold is applied after nor-

malization and before transformation.

zero\_threshold Abundances less than or equal to zero\_threshold will be treated as zeros. This

is primarily to be used when the abundance table has likely low-abundance false

positives.

min\_variance Features with abundance variances less than or equal to min\_variance will be

dropped. This is primarily used for dropping features that are entirely zero.

max\_significance

The FDR corrected q-value threshold for significance used in selecting which

associations to write as significant and to plot.

normalization The normalization to apply to the features before transformation and analysis.

The option TSS (total sum scaling) is recommended, but CLR (centered log ratio)

and NONE can also be used.

transform The transformation to apply to the features after normalization and before anal-

ysis. The option LOG (base 2) is recommended, but PLOG (pseudo-log) and NONE

can also be used.

correction The correction to obtain FDR-corrected q-values from raw p-values. Any valid

options for p. adjust can be used.

standardize Whether to apply z-scores to continuous metadata variables so they are on the

same scale. This is recommended in order to compare coefficients across metadata variables, but note that functions of the metadata specified in the formula

will apply after standardization.

unscaled\_abundance

A data frame with a single column of absolute abundances or a filepath to such a tab-delimited file. The row names should match the names of the samples in input\_data and input\_metadata. When using spike-ins, the single column should have the same name as one of the features in input\_data, and the unscaled\_abundance should correspond to the absolute quantity of the spike-in. When using total abundance scaling, the single column should have the name 'total', and the unscaled\_abundance should correspond to the total abundance

of each sample.

median\_comparison\_abundance

Test abundance coefficients against a null value corresponding to the median coefficient for a metadata variable across the features. This is recommended for relative abundance data but should not be used for absolute abundance data.

median\_comparison\_prevalence

Test prevalence coefficients against a null value corresponding to the median coefficient for a metadata variable across the features. This is only recommended if the analyst is interested in how feature prevalence associations compare to each other or if there is likely strong compositionality-induced sparsity.

6 maaslin3

median\_comparison\_abundance\_threshold

Coefficients within median\_comparison\_abundance\_threshold of the median association will automatically be counted as insignificant (p-value set to 1) since they likely represent compositionality-induced associations. This threshold will be divided by the metadata variable's standard deviation if the metadatum is continuous to ensure the threshold applies to the right scale.

median\_comparison\_prevalence\_threshold

Same as median\_comparison\_abundance\_threshold but applied to the prevalence associations.

subtract\_median

Subtract the median from the coefficients.

warn\_prevalence

Warn when prevalence associations are likely induced by abundance associations. This requires re-fitting the linear models on the TSS log-transformed data.

small\_random\_effects

Automatically replace random effects with fixed effects in the logistic prevalence model to handle low numbers of observations per group.

augment Add extra lowly-weighted 0s and 1s to avoid linear separability.

evaluate\_only Whether to evaluate just the abundance ("abundance") or prevalence ("prevalence") models

plot\_summary\_plot

Generate a summary plot of significant associations.

summary\_plot\_first\_n

Include the top summary\_plot\_first\_n features with significant associations.

coef\_plot\_vars Vector of variable names to be used in the coefficient plot section of the summary plot. Continuous variables should match the metadata column name, and categorical variables should be of the form "[variable] [level]".

heatmap\_vars Vector of variable names to be used in the heatmap section of the summary plot. Continuous variables should match the metadata column name, and categorical variables should be of the form "[variable] [level]".

plot\_associations

Whether to generate plots for significant associations.

max\_pngs The top max\_pngs significant associations will be plotted.

cores How many cores to use when fitting models. (Using multiple cores will likely be faster only for large datasets or complex models.

save\_models Whether to return the fit models and save them to an RData file.

save\_plots\_rds Whether to return the fit models and save them to an RData file.

verbosity The level of verbosity for the logging package.

summary\_plot\_balanced

If set to TRUE the summary plot will show the top N features of each variable included in coef\_plot\_vars where N is equal to: ceiling(summary\_plot\_first\_n/length(coef\_plot\_toef

Will error if coef\_plot\_vars = NULL

assay.type A string or index to select the assay when using a SummarizedExperiment ob-

ject

maaslin3 7

#### Value

A list containing the following items:

- (1) data: A dataframe of feature abundances with the retained samples for fitting.
- (2) normalized\_data: A dataframe of normalized feature abundances.
- (3) filtered\_data: A dataframe of feature abundances on the original scale after normalization and filtering.
- (4) transformed\_data: A dataframe of feature abundances after filtering, normalization, and transformation.
- (5) metadata: A dataframe of metadata with the retained samples for fitting.
- (6) standardized\_metadata: A dataframe of metadata after scaling (if selected).
- (7) formula: Checked or constructed formula(s) specifying the model to be fit.
- (8) fit\_data\_abundance: The results from the fit abundance models (see maaslin\_fit).
- (9) fit\_data\_prevalence: The results from the fit prevalence models (see maaslin\_fit).

## Author(s)

William Nickols<br/><willnickols@g.harvard.edu>,<br/>Jacob Nearing<br/><nearing@broadinstitute.org>,<br/>Maintainers: Lauren McIver<lauren.j.mciver@gmail.com>,

```
# Read features table
taxa_table_name <- system.file("extdata", "HMP2_taxonomy.tsv", package =</pre>
"maaslin3")
taxa_table <- read.csv(taxa_table_name, sep = '\t', row.names = 1)</pre>
# Read metadata table
metadata_name <- system.file("extdata", "HMP2_metadata.tsv", package =</pre>
"maaslin3")
metadata <- read.csv(metadata_name, sep = '\t', row.names = 1)</pre>
metadata$diagnosis <-</pre>
    factor(metadata$diagnosis, levels = c('nonIBD', 'UC', 'CD'))
metadata$dysbiosis_state <-</pre>
    factor(metadata$dysbiosis_state, levels = c('none', 'dysbiosis_UC',
    'dysbiosis_CD'))
metadata$antibiotics <-
    factor(metadata$antibiotics, levels = c('No', 'Yes'))
#Run MaAsLin3
fit_out <- maaslin3::maaslin3(input_data = taxa_table,</pre>
                             input_metadata = metadata,
                              output = 'output',
                              formula = '~ diagnosis + dysbiosis_state +
                              antibiotics + age + reads',
```

maaslin\_check\_arguments

Check parameter arguments to ensure a successful MaAsLin 3 run.

## **Description**

Check the arguments provided are valid for further MaAsLin 3 use.

#### Usage

# Arguments

feature\_specific\_covariate

A table of feature-specific covariates or a filepath to a tab-delimited file with feature-specific covariates. It should be formatted with features as columns and samples as rows (or the transpose). The row names and column names should be the same as those of the input\_data: the column and row names should be the feature names and sample names respectively. Typically, this table should be generated by 'preprocess\_mgx\_mtx' or 'preprocess\_taxa\_mtx' first.

feature\_specific\_covariate\_name

The name for the feature-specific covariates when fitting the models.

feature\_specific\_covariate\_record

Whether to keep the feature-specific covariates in the outputs when calculating p-values, writing results, and displaying plots.

zero\_threshold Abundances less than or equal to zero\_threshold will be treated as zeros. This is primarily to be used when the abundance table has likely low-abundance false positives.

normalization The normalization to apply to the features before transformation and analysis.

The option TSS (total sum scaling) is recommended, but CLR (centered log ratio)

and NONE can also be used.

transform The transformation to apply to the features after normalization and before anal-

ysis. The option LOG (base 2) is recommended, but PLOG (pseudo-log) and NONE

can also be used.

correction The correction to obtain FDR-corrected q-values from raw p-values. Any valid

options for p. adjust can be used.

warn\_prevalence

evaluate\_only

Warn when prevalence associations are likely induced by abundance associations. This requires re-fitting the linear models on the TSS log-transformed

u

Whether to evaluate just the abundance ("abundance") or prevalence ("preva-

lence") models

unscaled\_abundance

A data frame with a single column of absolute abundances or a filepath to such a tab-delimited file. The row names should match the names of the samples in input\_data and input\_metadata. When using spike-ins, the single column should have the same name as one of the features in input\_data, and the unscaled\_abundance should correspond to the absolute quantity of the spike-in. When using total abundance scaling, the single column should have the name 'total', and the unscaled\_abundance should correspond to the total abundance of each sample.

median\_comparison\_abundance

Test abundance coefficients against a null value corresponding to the median coefficient for a metadata variable across the features. This is recommended for relative abundance data but should not be used for absolute abundance data.

#### Value

No value is returned, but incompatibile arguments will produce an error.

#### Author(s)

William Nickols<a href="willnickols@g.harvard.edu">willnickols@g.harvard.edu</a>, Jacob Nearing<a href="millionequation-parabeta-based-approximation-parabeta-bas

Maintainers: Lauren McIver<lauren.j.mciver@gmail.com>,

```
# Read features table
taxa_table_name <- system.file("extdata", "HMP2_taxonomy.tsv", package =
"maaslin3")
taxa_table <- read.csv(taxa_table_name, sep = '\t', row.names = 1)
# Read metadata table
metadata_name <- system.file("extdata", "HMP2_metadata.tsv", package =</pre>
```

```
"maaslin3")
metadata <- read.csv(metadata_name, sep = '\t', row.names = 1)</pre>
metadata$diagnosis <-</pre>
    factor(metadata$diagnosis, levels = c('nonIBD', 'UC', 'CD'))
metadata$dysbiosis_state <-</pre>
    factor(metadata$dysbiosis_state, levels = c('none', 'dysbiosis_UC',
    'dysbiosis_CD'))
metadata$antibiotics <-</pre>
    factor(metadata$antibiotics, levels = c('No', 'Yes'))
# Prepare parameter lists
maaslin3::maaslin_check_arguments(zero_threshold = 0,
                                  normalization = 'TSS',
                                  transform = 'LOG',
                                  correction = 'BH',
                                 median_comparison_abundance = TRUE)
unlink('output', recursive=TRUE)
logging::logReset()
```

maaslin\_check\_formula Check a MaAsLin 3 formula to ensure a proper MaAsLin 3 run.

# Description

Ensure that the formula provided is valid. Only one of maaslin\_compute\_formula or maaslin\_check\_formula should be used.

#### Usage

#### **Arguments**

data A data frame of feature abundances. It should be formatted with features as

columns and samples as rows. The column and row names should be the feature

names and sample names respectively.

metadata A data frame of per-sample metadata. It should be formatted with variables as

columns and samples as rows. The column and row names should be the variable

names and sample names respectively.

input\_formula A formula in lme4 format. Random effects, interactions, and functions of the

metadata can be included (note that these functions will be applied after standardization if standardize=TRUE). Group, ordered, and strata variables can be specified as: group(grouping\_variable), ordered(ordered\_variable) and

strata(strata\_variable). The other variable options below will not be considered if a formula is set.

feature\_specific\_covariate\_name

The name for the feature-specific covariates when fitting the models.

#### Value

A list containing the following named items:

- (1) formula: The constructed formula.
- (2) random\_effects\_formula: A formula for the random effects.

#### Author(s)

William Nickols<willnickols@g.harvard.edu>, Jacob Nearing<nearing@broadinstitute.org>, Maintainers: Lauren McIver<lauren.j.mciver@gmail.com>,

```
# Read features table
taxa_table_name <- system.file("extdata", "HMP2_taxonomy.tsv", package =</pre>
"maaslin3")
taxa_table <- read.csv(taxa_table_name, sep = '\t', row.names = 1)</pre>
# Read metadata table
metadata_name <- system.file("extdata", "HMP2_metadata.tsv", package =</pre>
"maaslin3")
metadata <- read.csv(metadata_name, sep = '\t', row.names = 1)</pre>
metadata$diagnosis <-</pre>
    factor(metadata$diagnosis, levels = c('nonIBD', 'UC', 'CD'))
metadata$dysbiosis_state <-</pre>
    factor(metadata$dysbiosis_state, levels = c('none', 'dysbiosis_UC',
    'dysbiosis_CD'))
metadata$antibiotics <-
    factor(metadata$antibiotics, levels = c('No', 'Yes'))
#Run MaAsLin3
maaslin3::maaslin_log_arguments(
    input_data = taxa_table,
    input_metadata = metadata,
    output = 'output',
    formula = '~ diagnosis + dysbiosis_state + antibiotics +
    age + reads',
    plot_summary_plot = FALSE,
    plot_associations = FALSE)
    read_data_list <- maaslin3::maaslin_read_data(</pre>
    taxa_table,
    metadata)
    read_data_list <- maaslin3::maaslin_reorder_data(</pre>
```

```
read_data_list$data,
    read_data_list$metadata)

data <- read_data_list$data
metadata <- read_data_list$metadata

formulas <- maaslin3::maaslin_check_formula(
    data,
    metadata,
    input_formula = '~ diagnosis + dysbiosis_state + antibiotics +
    age + reads')

unlink('output', recursive=TRUE)
logging::logReset()</pre>
```

maaslin\_compute\_formula

Compute a formula for a MaAsLin 3 run based on the specified effects.

# **Description**

Compute a formula using variables provided through fixed\_effects, random\_effects, group\_effects, ordered\_effects, and strata\_effects. Only one of maaslin\_compute\_formula or maaslin\_check\_formula should be used.

# Usage

# Arguments

data A data frame of feature abundances. It should be formatted with features as

columns and samples as rows. The column and row names should be the feature

names and sample names respectively.

metadata A data frame of per-sample metadata. It should be formatted with variables as

columns and samples as rows. The column and row names should be the variable

names and sample names respectively.

fixed\_effects A vector of variable names to be included as fixed effects.

random\_effects A vector of variable names to be included as random intercepts.

group\_effects

A factored categorical variable to be included for group testing. An ANOVAstyle test will be performed to assess whether any of the variable's levels are significant, and no coefficients or individual p-values will be returned.

ordered\_effects

A factored categorical variable to be included. Consecutive levels will be tested for significance against each other, and the resulting associations will correspond to effect sizes, standard errors, and significances of each level versus the previous.

strata\_effects A vector with one variable name to be included as the strata variable in case-control studies. Strata cannot be combined with random effects.

feature\_specific\_covariate\_name

The name for the feature-specific covariates when fitting the models.

#### Value

A list containing the following named items:

- (1) formula: The constructed formula.
- (2) random\_effects\_formula: A formula for the random effects.

## Author(s)

William Nickols<br/><willnickols@g.harvard.edu>,<br/>Jacob Nearing<br/><nearing@broadinstitute.org>,<br/>Maintainers: Lauren McIver<lauren.j.mciver@gmail.com>,

```
# Read features table
taxa_table_name <- system.file("extdata", "HMP2_taxonomy.tsv", package =</pre>
taxa_table <- read.csv(taxa_table_name, sep = '\t', row.names = 1)</pre>
# Read metadata table
metadata_name <- system.file("extdata", "HMP2_metadata.tsv", package =</pre>
"maaslin3")
metadata <- read.csv(metadata_name, sep = '\t', row.names = 1)</pre>
metadata$diagnosis <-</pre>
    factor(metadata$diagnosis, levels = c('nonIBD', 'UC', 'CD'))
metadata$dysbiosis_state <-</pre>
    factor(metadata$dysbiosis_state, levels = c('none', 'dysbiosis_UC',
    'dysbiosis_CD'))
metadata$antibiotics <-</pre>
    factor(metadata$antibiotics, levels = c('No', 'Yes'))
#Run MaAsLin3
maaslin3::maaslin_log_arguments(
    input_data = taxa_table,
```

14 maaslin\_contrast\_test

```
input_metadata = metadata,
    output = 'output',
    fixed_effects = c('diagnosis', 'dysbiosis_state', 'antibiotics',
                     'age', 'reads'),
    random_effects = c('participant_id'),
    plot_summary_plot = FALSE,
    plot_associations = FALSE)
read_data_list <- maaslin3::maaslin_read_data(</pre>
    taxa_table,
    metadata)
read_data_list <- maaslin3::maaslin_reorder_data(</pre>
    read_data_list$data,
    read_data_list$metadata)
data <- read_data_list$data</pre>
metadata <- read_data_list$metadata</pre>
formulas <- maaslin3::maaslin_compute_formula(</pre>
    data,
    metadata,
    fixed_effects = c('diagnosis', 'dysbiosis_state', 'antibiotics',
                     'age', 'reads'),
    random_effects = c('participant_id'))
unlink('output', recursive=TRUE)
logging::logReset()
```

## Description

Perform a contrast test (lmerTest::contest for mixed effects linear; multcomp::glht for all others) using a named contrast matrix and right hand side. One contrast test is applied per row of the matrix.

# Usage

maaslin\_contrast\_test 15

#### **Arguments**

maaslin3\_fit The output of maaslin\_fit with save\_models = TRUE.

contrast\_mat A matrix with one row per contrast test to run. The columns will be matched to

the coefficients of the model by name. Contrast vector coefficients need not be specified if they would be zero. If row names are provided, they will be used to

label the test in the results.

rhs The right hand size of the contrast test. The length should be the same as the

number of rows in the contrast\_mat. This will default to 0 or the median

comparison if median\_comparison=TRUE.

max\_significance

The FDR corrected q-value threshold for significance used in selecting which

associations to write as significant and to plot.

correction The correction to obtain FDR-corrected q-values from raw p-values. Any valid

options for p. adjust can be used.

median\_comparison\_abundance

Test abundance coefficients against a null value corresponding to the median coefficient for a metadata variable across the features. This is recommended for relative abundance data but should not be used for absolute abundance data.

median\_comparison\_prevalence

Test prevalence coefficients against a null value corresponding to the median coefficient for a metadata variable across the features. This is only recommended if the analyst is interested in how feature prevalence associations compare to each other or if there is likely strong compositionality-induced sparsity.

subtract\_median

Subtract the median from the coefficients.

small\_random\_effects

Automatically replace random effects with fixed effects in the logistic preva-

lence model to handle low numbers of observations per group.

Whether to evaluate just the abundance ("abundance") or prevalence ("prevalence") models

evaluate\_only

#### Value

A dataframe with the following columns:

- (1) feature: The feature involved in the association.
- (2) test: The contrast test name.
- (3) coef: The coefficient of the association: the slope coefficient in the abundance model and the change in log odds in the prevalence model.
- (4) null\_hypothesis: The value of the null hypothesis against which the coefficients are tested (zero or the per-metadatum median).
- (5) stderr: The standard error of the coefficient.
- (6) pval\_individual: The (uncorrected) p-value of the association.
- (7) qval\_individual: The FDR corrected q-value of the association. FDR correction is performed over all associations in the abundance and prevalence modeling without errors together.

16 maaslin\_contrast\_test

(8) pval\_joint: The p-value of the overall association (combining abundance and prevalence) by taking the minimum of the abundance and logistic p-values and applying the Beta(1,2) CDF. These will be the same in the abundance and prevalence results for an association.

- (9) qval\_joint: The FDR corrected q-value of the association. FDR correction is performed over all joint p-values without errors.
- (10) error: Any error produced by the model during fitting. NA otherwise.
- (11) model: linear for the abundance models and logistic for the prevalence models.
- (12) N: The number of data points for the association's feature.
- (13) N\_not\_zero: The number of non-zero data points for the association's feature.

#### Author(s)

William Nickols<a href="million:willnickols@g.harvard.edu">willnickols@g.harvard.edu</a>, Jacob Nearing<a href="million:nearing@broadinstitute.org">nearing@broadinstitute.org</a>, Maintainers: Lauren McIver<a href="million:lauren.j.mciver@gmail.com">nearing@broadinstitute.org</a>,

```
# Read features table
taxa_table_name <- system.file("extdata", "HMP2_taxonomy.tsv", package =</pre>
"maaslin3")
taxa_table <- read.csv(taxa_table_name, sep = '\t', row.names = 1)</pre>
# Read metadata table
metadata_name <- system.file("extdata", "HMP2_metadata.tsv", package =</pre>
"maaslin3")
metadata <- read.csv(metadata_name, sep = '\t', row.names = 1)</pre>
metadata$diagnosis <-</pre>
    factor(metadata$diagnosis, levels = c('nonIBD', 'UC', 'CD'))
metadata$dysbiosis_state <-</pre>
    factor(metadata$dysbiosis_state, levels = c('none', 'dysbiosis_UC',
    'dysbiosis_CD'))
metadata$antibiotics <-
    factor(metadata$antibiotics, levels = c('No', 'Yes'))
#Run MaAsLin3
fit_out <- maaslin3::maaslin3(input_data = taxa_table,</pre>
                             input_metadata = metadata,
                             output = 'output',
                             formula = '~ diagnosis + dysbiosis_state +
                             antibiotics + age + reads',
                             plot_summary_plot = FALSE,
                             plot_associations = FALSE)
contrast_mat <- matrix(c(1, 1, 0, 0, 0, 0, 1, 1),
    ncol = 4, nrow = 2, byrow = TRUE)
colnames(contrast_mat) <- c("diagnosisUC",</pre>
```

17 maaslin\_filter

```
"diagnosisCD",
                             "dysbiosis_statedysbiosis_UC",
                             "dysbiosis_statedysbiosis_CD")
rownames(contrast_mat) <- c("diagnosis_test", "dysbiosis_test")</pre>
maaslin_contrast_test(maaslin3_fit = fit_out,
                         contrast_mat = contrast_mat)
unlink('output', recursive=TRUE)
logging::logReset()
```

maaslin\_filter

Filter abundance data before MaAsLin 3 model fitting.

## **Description**

Set abundances below zero\_threshold to zero, remove features without abundances more than min\_abundance in min\_prevalence of the samples, remove features with abundances more than min\_abundance in more than max\_prevalence of the samples, and remove features with variances less than or equal to min\_variance.

## Usage

```
maaslin_filter(normalized_data,
            output,
            min_abundance = 0,
            min_prevalence = 0,
            max_prevalence = 1.01,
            zero_threshold = 0,
            min_variance = 0)
```

## **Arguments**

normalized\_data

A data frame of normalized feature abundances. It should be formatted with features as columns and samples as rows. The column and row names should be the feature names and sample names respectively.

output The output folder to write results.

min\_abundance Features with abundances more than min\_abundance in more than min\_prevalence

of the samples will be included for analysis. The threshold is applied after nor-

malization and before transformation.

min\_prevalence See min\_abundance.

max\_prevalence Features with abundances more than min\_abundance in fewer than max\_prevalence

of the samples will be included for analysis. The threshold is applied after nor-

malization and before transformation.

18 maaslin\_filter

zero\_threshold Abundances less than or equal to zero\_threshold will be treated as zeros. This

is primarily to be used when the abundance table has likely low-abundance false positives.

min\_variance Features with abundance variances less than or equal to min\_variance will be

dropped. This is primarily used for dropping features that are entirely zero.

#### Value

A dataframe of filtered features (features are columns; samples are rows).

#### Author(s)

William Nickols<br/>
willnickols@g.harvard.edu>,
Jacob Nearing<br/>
gbroadinstitute.org>,
Maintainers: Lauren McIver<br/>
lauren.j.mciver@gmail.com>,

```
# Read features table
taxa_table_name <- system.file("extdata", "HMP2_taxonomy.tsv", package =</pre>
"maaslin3")
taxa_table <- read.csv(taxa_table_name, sep = '\t', row.names = 1)</pre>
# Read metadata table
metadata_name <- system.file("extdata", "HMP2_metadata.tsv", package =</pre>
metadata <- read.csv(metadata_name, sep = '\t', row.names = 1)</pre>
metadata$diagnosis <-</pre>
    factor(metadata$diagnosis, levels = c('nonIBD', 'UC', 'CD'))
metadata$dysbiosis_state <-</pre>
    factor(metadata$dysbiosis_state, levels = c('none', 'dysbiosis_UC',
    'dysbiosis_CD'))
metadata$antibiotics <-
    factor(metadata$antibiotics, levels = c('No', 'Yes'))
#Run MaAsLin3
maaslin3::maaslin_log_arguments(
    input_data = taxa_table,
    input_metadata = metadata,
    output = 'output',
    formula = '~ diagnosis + dysbiosis_state + antibiotics +
    age + reads',
    plot_summary_plot = FALSE,
    plot_associations = FALSE)
read_data_list <- maaslin3::maaslin_read_data(</pre>
    taxa_table,
    metadata)
read_data_list <- maaslin3::maaslin_reorder_data(</pre>
    read_data_list$data,
```

maaslin\_fit

Fit MaAsLin 3 models.

## **Description**

Fit the abundance data with abundance and prevalence models to discover feature-metadata associations.

## Usage

```
maaslin_fit(filtered_data,
            transformed_data,
            metadata.
            formula,
            random_effects_formula,
            feature_specific_covariate = NULL,
            feature_specific_covariate_name = NULL,
            feature_specific_covariate_record = NULL,
            zero_threshold = 0,
            max_significance = 0.1,
            correction = 'BH',
            median_comparison_abundance = TRUE,
            median_comparison_prevalence = FALSE,
            median_comparison_abundance_threshold = 0,
            median_comparison_prevalence_threshold = 0,
            subtract_median = FALSE,
            warn_prevalence = TRUE,
            small_random_effects = FALSE,
            augment = TRUE,
```

```
evaluate_only = NULL,
cores = 1,
save_models = FALSE,
data = NULL,
min_abundance = 0,
min_prevalence = 0,
max_prevalence = 1.01,
min_variance = 0)
```

## **Arguments**

filtered data

A data frame of filtered feature abundances. It should be formatted with features as columns and samples as rows. The column and row names should be the feature names and sample names respectively.

transformed\_data

A data frame of transformed feature abundances. It should be formatted with features as columns and samples as rows. The column and row names should be the feature names and sample names respectively.

metadata

A data frame of per-sample metadata. It should be formatted with variables as columns and samples as rows. The column and row names should be the variable names and sample names respectively.

formula random\_effects\_formula

A formula in lme4 format as from maaslin\_check\_formula.

A formula in lme4 format as from maaslin\_check\_formula. feature\_specific\_covariate

> A table of feature-specific covariates or a filepath to a tab-delimited file with feature-specific covariates. It should be formatted with features as columns and samples as rows (or the transpose). The row names and column names should be the same as those of the input\_data: the column and row names should be the feature names and sample names respectively. Typically, this table should be generated by 'preprocess\_mgx\_mtx' or 'preprocess\_taxa\_mtx' first.

feature\_specific\_covariate\_name

The name for the feature-specific covariates when fitting the models.

feature\_specific\_covariate\_record

Whether to keep the feature-specific covariates in the outputs when calculating p-values, writing results, and displaying plots.

zero\_threshold Abundances less than or equal to zero\_threshold will be treated as zeros. This is primarily to be used when the abundance table has likely low-abundance false positives.

max\_significance

The FDR corrected q-value threshold for significance used in selecting which associations to write as significant and to plot.

correction

The correction to obtain FDR-corrected q-values from raw p-values. Any valid options for p. adjust can be used.

median\_comparison\_abundance

Test abundance coefficients against a null value corresponding to the median coefficient for a metadata variable across the features. This is recommended for relative abundance data but should not be used for absolute abundance data.

median\_comparison\_prevalence

Test prevalence coefficients against a null value corresponding to the median coefficient for a metadata variable across the features. This is only recommended if the analyst is interested in how feature prevalence associations compare to each other or if there is likely strong compositionality-induced sparsity.

median\_comparison\_abundance\_threshold

Coefficients within median\_comparison\_abundance\_threshold of the median association will automatically be counted as insignificant (p-value set to 1) since they likely represent compositionality-induced associations. This threshold will be divided by the metadata variable's standard deviation if the metadatum is continuous to ensure the threshold applies to the right scale.

median\_comparison\_prevalence\_threshold

Same as median\_comparison\_abundance\_threshold but applied to the prevalence associations.

subtract\_median

Subtract the median from the coefficients.

warn\_prevalence

Warn when prevalence associations are likely induced by abundance associations. This requires re-fitting the linear models on the TSS log-transformed data.

small\_random\_effects

Automatically replace random effects with fixed effects in the logistic preva-

lence model to handle low numbers of observations per group.

augment Add extra lowly-weighted 0s and 1s to avoid linear separability.

evaluate\_only Whether to evaluate just the abundance ("abundance") or prevalence ("preva-

lence") models

cores How many cores to use when fitting models. (Using multiple cores will likely

be faster only for large datasets or complex models.

save\_models Whether to return the fit models and save them to an RData file.

data The original data (only necessary if warn\_prevalence is TRUE).

min\_abundance The original min\_abundance parameter (only necessary if warn\_prevalence is

TRUE).

min\_prevalence The original min\_prevalence parameter (only necessary if warn\_prevalence is

TRUE).

max\_prevalence The original max\_prevalence parameter (only necessary if warn\_prevalence is

TRUE).

min\_variance The original min\_variance parameter (only necessary if min\_variance is TRUE).

#### Value

A list containing the following named items:

- (1) fit\_data\_abundance: The results from the fit abundance models.
- (2) fit\_data\_prevalence: The results from the fit prevalence models.

The fit\_data\_abundance and fit\_data\_prevalence items have the same structure. They are both lists with the following named items:

- (1) results: A results table with the modeled associations (see below).
- (2) residuals: A features (rows) by samples (columns) dataframe of residuals from the models.
- (3) fitted: A features (rows) by samples (columns) dataframe of fitted values from the models.
- (4) ranef: A features (rows) by random effect (columns) dataframe of random effects from the models. If multiple random effects are specified, this is a dataframe of dataframes.
- (5) fits: If save\_models=TRUE, this is a list of the fit models.

The results tables contain the following columns for each association (row):

- (1) feature: The feature involved in the association.
- (2) metadata: The metadata variable involved in the association.
- (3) value: The value of the metadata variable: the metadata variable itself if continuous or the level if categorical.
- (4) name: The name of the model component involved in the association: the metadata variable itself if continuous or a concatenated version of the metadata variable and level if categorical.
- (5) coef: The coefficient of the association: the slope coefficient in the abundance model and the change in log odds in the prevalence model.
- (6) null\_hypothesis: The value of the null hypothesis against which the coefficients are tested (zero or the per-metadatum median).
- (7) stderr: The standard error of the coefficient.
- (8) pval\_individual: The (uncorrected) p-value of the association.
- (9) qval\_individual: The FDR corrected q-value of the association. FDR correction is performed over all associations in the abundance and prevalence modeling without errors together.
- (10) pval\_joint: The p-value of the overall association (combining abundance and prevalence) by taking the minimum of the abundance and logistic p-values and applying the Beta(1,2) CDF. These will be the same in the abundance and prevalence results for an association.
- (11) qval\_joint: The FDR corrected q-value of the association. FDR correction is performed over all joint p-values without errors.
- (12) error: Any error produced by the model during fitting. NA otherwise.
- (13) model: linear for the abundance models and logistic for the prevalence models.
- (14) N: The number of data points for the association's feature.
- (15) N\_not\_zero: The number of non-zero data points for the association's feature.

# Author(s)

William Nickols<br/><willnickols@g.harvard.edu>,<br/>Jacob Nearing<br/><nearing@broadinstitute.org>,<br/>Maintainers: Lauren McIver<lauren.j.mciver@gmail.com>,

```
# Read features table
taxa_table_name <- system.file("extdata", "HMP2_taxonomy.tsv", package =</pre>
"maaslin3")
taxa_table <- read.csv(taxa_table_name, sep = '\t', row.names = 1)</pre>
# Read metadata table
metadata_name <- system.file("extdata", "HMP2_metadata.tsv", package =</pre>
"maaslin3")
metadata <- read.csv(metadata_name, sep = '\t', row.names = 1)</pre>
metadata$diagnosis <-</pre>
    factor(metadata$diagnosis, levels = c('nonIBD', 'UC', 'CD'))
metadata$dysbiosis_state <-</pre>
    factor(metadata$dysbiosis_state, levels = c('none', 'dysbiosis_UC',
    'dysbiosis_CD'))
metadata$antibiotics <-</pre>
    factor(metadata$antibiotics, levels = c('No', 'Yes'))
#Run MaAsLin3
maaslin3::maaslin_log_arguments(
    input_data = taxa_table,
    input_metadata = metadata,
    output = 'output',
    formula = '~ diagnosis + dysbiosis_state + antibiotics +
    age + reads',
    plot_summary_plot = FALSE,
    plot_associations = FALSE)
read_data_list <- maaslin3::maaslin_read_data(</pre>
    taxa_table,
    metadata)
read_data_list <- maaslin3::maaslin_reorder_data(</pre>
    read_data_list$data,
    read_data_list$metadata)
data <- read_data_list$data</pre>
metadata <- read_data_list$metadata</pre>
formulas <- maaslin3::maaslin_check_formula(</pre>
    data,
    metadata,
    input_formula = '~ diagnosis + dysbiosis_state + antibiotics +
    age + reads')
formula <- formulas$formula</pre>
random_effects_formula <- formulas$random_effects_formula</pre>
normalized_data = maaslin3::maaslin_normalize(data,
                                  output = 'output')
filtered_data = maaslin3::maaslin_filter(normalized_data,
```

maaslin\_log\_arguments Log MaAsLin 3 parameters.

# Description

Check that the parameters provided are valid for further MaAsLin 3 use and open a logger to log the parameters.

## Usage

```
maaslin_log_arguments(input_data,
                    input_metadata,
                    output,
                    formula = NULL,
                    fixed_effects = NULL,
                    reference = NULL,
                    random_effects = NULL,
                    group_effects = NULL,
                    ordered_effects = NULL,
                    strata_effects = NULL,
                    feature_specific_covariate = NULL,
                    feature_specific_covariate_name = NULL,
                    feature_specific_covariate_record = NULL,
                    min_abundance = 0,
                    min_prevalence = 0,
                    max_prevalence = 1.01,
                    zero_threshold = 0,
                    min_variance = 0,
```

```
max_significance = 0.1,
normalization = 'TSS',
transform = 'LOG',
correction = 'BH',
standardize = TRUE,
unscaled_abundance = NULL,
median_comparison_abundance = TRUE,
median_comparison_prevalence = FALSE,
median_comparison_abundance_threshold = 0,
median_comparison_prevalence_threshold = 0,
subtract_median = FALSE,
warn_prevalence = TRUE,
small_random_effects = FALSE,
augment = TRUE,
evaluate_only = NULL,
plot_summary_plot = TRUE,
summary_plot_first_n = 25,
coef_plot_vars = NULL,
heatmap_vars = NULL,
plot_associations = TRUE,
max_pngs = 30,
cores = 1,
save_models = FALSE,
save_plots_rds = FALSE,
verbosity = 'FINEST',
summary_plot_balanced = FALSE)
```

## **Arguments**

input\_data

A data frame of feature abundances or read counts or a filepath to a tab-delimited file with abundances. It should be formatted with features as columns and samples as rows (or the transpose). The column and row names should be the feature names and sample names respectively.

input\_metadata

A data frame of per-sample metadata or a filepath to a tab-delimited file with metadata. It should be formatted with variables as columns and samples as rows (or the transpose). The column and row names should be the variable names and sample names respectively.

output

The output folder to write results.

formula

A formula in lme4 format. Random effects, interactions, and functions of the metadata can be included (note that these functions will be applied after standardization if standardize=TRUE). Group, ordered, and strata variables can be specified as: group(grouping\_variable), ordered(ordered\_variable) and strata(strata\_variable). The other variable options below will not be considered if a formula is set.

 ${\tt fixed\_effects}$ 

A vector of variable names to be included as fixed effects.

reference

For a variable with more than two levels supplied with fixed\_effects, the factor to use as a reference provided as a string of 'variable,reference' semi-colon delimited for multiple variables.

random\_effects A vector of variable names to be included as random intercepts.

group\_effects A factored categorical variable to be included for group testing. An ANOVA-

style test will be performed to assess whether any of the variable's levels are significant, and no coefficients or individual p-values will be returned.

ordered\_effects

A factored categorical variable to be included. Consecutive levels will be tested for significance against each other, and the resulting associations will correspond to effect sizes, standard errors, and significances of each level versus the previous.

strata\_effects A vector with one variable name to be included as the strata variable in case-control studies. Strata cannot be combined with random effects.

feature\_specific\_covariate

A table of feature-specific covariates or a filepath to a tab-delimited file with feature-specific covariates. It should be formatted with features as columns and samples as rows (or the transpose). The row names and column names should be the same as those of the input\_data: the column and row names should be the feature names and sample names respectively. Typically, this table should be generated by 'preprocess\_mgx\_mtx' or 'preprocess\_taxa\_mtx' first.

feature\_specific\_covariate\_name

The name for the feature-specific covariates when fitting the models.

feature\_specific\_covariate\_record

Whether to keep the feature-specific covariates in the outputs when calculating p-values, writing results, and displaying plots.

min\_abundance Features with abundances more than min\_abundance in more than min\_prevalence of the samples will be included for analysis. The threshold is applied after nor-

malization and before transformation.

min\_prevalence See min\_abundance.

max\_prevalence Features with abundances more than min\_abundance in fewer than max\_prevalence

of the samples will be included for analysis. The threshold is applied after nor-

malization and before transformation.

zero\_threshold Abundances less than or equal to zero\_threshold will be treated as zeros. This

is primarily to be used when the abundance table has likely low-abundance false

positives.

min\_variance Features with abundance variances less than or equal to min\_variance will be

dropped. This is primarily used for dropping features that are entirely zero.

max\_significance

The FDR corrected q-value threshold for significance used in selecting which

associations to write as significant and to plot.

normalization The normalization to apply to the features before transformation and analysis.

The option TSS (total sum scaling) is recommended, but CLR (centered log ratio)

and NONE can also be used.

transform The transformation to apply to the features after normalization and before anal-

ysis. The option LOG (base 2) is recommended, but PLOG (pseudo-log) and NONE

can also be used.

correction

The correction to obtain FDR-corrected q-values from raw p-values. Any valid options for p. adjust can be used.

standardize

Whether to apply z-scores to continuous metadata variables so they are on the same scale. This is recommended in order to compare coefficients across metadata variables, but note that functions of the metadata specified in the formula will apply after standardization.

#### unscaled\_abundance

A data frame with a single column of absolute abundances or a filepath to such a tab-delimited file. The row names should match the names of the samples in input\_data and input\_metadata. When using spike-ins, the single column should have the same name as one of the features in input\_data, and the unscaled\_abundance should correspond to the absolute quantity of the spike-in. When using total abundance scaling, the single column should have the name 'total', and the unscaled\_abundance should correspond to the total abundance of each sample.

## median\_comparison\_abundance

Test abundance coefficients against a null value corresponding to the median coefficient for a metadata variable across the features. This is recommended for relative abundance data but should not be used for absolute abundance data.

## median\_comparison\_prevalence

Test prevalence coefficients against a null value corresponding to the median coefficient for a metadata variable across the features. This is only recommended if the analyst is interested in how feature prevalence associations compare to each other or if there is likely strong compositionality-induced sparsity.

# median\_comparison\_abundance\_threshold

Coefficients within median\_comparison\_abundance\_threshold of the median association will automatically be counted as insignificant (p-value set to 1) since they likely represent compositionality-induced associations. This threshold will be divided by the metadata variable's standard deviation if the metadatum is continuous to ensure the threshold applies to the right scale.

# median\_comparison\_prevalence\_threshold

Same as median\_comparison\_abundance\_threshold but applied to the prevalence associations.

#### subtract\_median

Subtract the median from the coefficients.

#### warn\_prevalence

Warn when prevalence associations are likely induced by abundance associations. This requires re-fitting the linear models on the TSS log-transformed data.

# small\_random\_effects

Automatically replace random effects with fixed effects in the logistic prevalence model to handle low numbers of observations per group.

## augment

Add extra lowly-weighted 0s and 1s to avoid linear separability.

#### evaluate\_only

Whether to evaluate just the abundance ("abundance") or prevalence ("prevalence") models

#### plot\_summary\_plot

Generate a summary plot of significant associations.

cluded in coef\_plot\_vars where N is equal to: ceiling(summary\_plot\_first\_n/length(coef\_plot.

```
summary_plot_first_n
                  Include the top summary_plot_first_n features with significant associations.
coef_plot_vars Vector of variable names to be used in the coefficient plot section of the sum-
                  mary plot. Continuous variables should match the metadata column name, and
                  categorical variables should be of the form "[variable] [level]".
                  Vector of variable names to be used in the heatmap section of the summary plot.
heatmap_vars
                  Continuous variables should match the metadata column name, and categorical
                  variables should be of the form "[variable] [level]".
plot_associations
                  Whether to generate plots for significant associations.
                  The top max_pngs significant associations will be plotted.
max_pngs
                  How many cores to use when fitting models. (Using multiple cores will likely
cores
                  be faster only for large datasets or complex models.
                  Whether to return the fit models and save them to an RData file.
save_models
save_plots_rds Whether to return the plots to an RDS file.
verbosity
                  The level of verbosity for the logging package.
summary_plot_balanced
                  If set to TRUE the summary plot will show the top N features of each variable in-
```

## Value

No value is returned, but a logger is opened with the parameters logged.

Will error if coef\_plot\_vars = NULL

## Author(s)

William Nickols<br/>
willnickols@g.harvard.edu>,
Jacob Nearing<br/>
proadinstitute.org>,
Maintainers: Lauren McIver<br/>
lauren.j.mciver@gmail.com>,

```
# Read features table
taxa_table_name <- system.file("extdata", "HMP2_taxonomy.tsv", package =
"maaslin3")
taxa_table <- read.csv(taxa_table_name, sep = '\t', row.names = 1)

# Read metadata table
metadata_name <- system.file("extdata", "HMP2_metadata.tsv", package =
"maaslin3")
metadata <- read.csv(metadata_name, sep = '\t', row.names = 1)

metadata$diagnosis <-
factor(metadata$diagnosis, levels = c('nonIBD', 'UC', 'CD'))
metadata$dysbiosis_state <-
factor(metadata$dysbiosis_state, levels = c('none', 'dysbiosis_UC',</pre>
```

maaslin\_normalize 29

maaslin\_normalize

Normalize abundance data for MaAsLin 3 model fitting.

# **Description**

Normalize the abundance data according to the normalization parameter. If unscaled\_abundance is specified, compute the absolute abundances.

## Usage

## **Arguments**

data

A data frame of feature abundances. It should be formatted with features as columns and samples as rows. The column and row names should be the feature names and sample names respectively.

output

The output folder to write results.

zero\_threshold

Abundances less than or equal to zero\_threshold will be treated as zeros. This is primarily to be used when the abundance table has likely low-abundance false positives.

normalization

The normalization to apply to the features before transformation and analysis. The option TSS (total sum scaling) is recommended, but CLR (centered log ratio) and NONE can also be used.

unscaled\_abundance

A data frame with a single column of absolute abundances. The row names should match the names of the samples in input\_data and input\_metadata. When using spike-ins, the single column should have the same name as one of

30 maaslin\_normalize

the features in input\_data, and the unscaled\_abundance should correspond to the absolute quantity of the spike-in. When using total abundance scaling, the single column should have the name 'total', and the unscaled\_abundance should correspond to the total abundance of each sample.

#### Value

A dataframe of normalized features (features are columns; samples are rows).

## Author(s)

William Nickols<br/>
willnickols@g.harvard.edu>,<br/>
Jacob Nearing<br/>
nearing@broadinstitute.org>,<br/>
Maintainers: Lauren McIver<lauren.j.mciver@gmail.com>,

```
# Read features table
taxa_table_name <- system.file("extdata", "HMP2_taxonomy.tsv", package =</pre>
taxa_table <- read.csv(taxa_table_name, sep = '\t', row.names = 1)</pre>
# Read metadata table
metadata_name <- system.file("extdata", "HMP2_metadata.tsv", package =</pre>
"maaslin3")
metadata <- read.csv(metadata_name, sep = '\t', row.names = 1)</pre>
metadata$diagnosis <-</pre>
    factor(metadata$diagnosis, levels = c('nonIBD', 'UC', 'CD'))
metadata$dysbiosis_state <-</pre>
    factor(metadata$dysbiosis_state, levels = c('none', 'dysbiosis_UC',
     'dysbiosis_CD'))
metadata$antibiotics <-</pre>
    factor(metadata$antibiotics, levels = c('No', 'Yes'))
#Run MaAsLin3
maaslin3::maaslin_log_arguments(
    input_data = taxa_table,
    input_metadata = metadata,
    output = 'output',
    formula = '~ diagnosis + dysbiosis_state + antibiotics +
    age + reads',
    plot_summary_plot = FALSE,
    plot_associations = FALSE)
read_data_list <- maaslin3::maaslin_read_data(</pre>
    taxa_table,
    metadata)
    read_data_list <- maaslin3::maaslin_reorder_data(</pre>
    read_data_list$data,
    read_data_list$metadata)
```

# **Description**

Two types of plots are generated. First, the summary plot contains sorted per-feature coefficients plotted with their standard errors for key variables and a heatmap summarizing the remaining variables. Second, for significant features, association plots (scatterplots, boxplots, or tables depending on the association) are generated to visualize and verify the model fits. The data are shown with their transformed values in the association plots since this is the scale on which the models are fit.

#### Usage

```
maaslin_plot_results(output,
                    transformed_data,
                    unstandardized_metadata,
                    fit_data_abundance,
                    fit_data_prevalence,
                    normalization,
                    transform,
                    feature_specific_covariate = NULL,
                    feature_specific_covariate_name = NULL,
                    feature_specific_covariate_record = NULL,
                    median_comparison_abundance = TRUE,
                    median_comparison_prevalence = FALSE,
                    max\_significance = 0.1,
                    plot_summary_plot = TRUE,
                    summary_plot_first_n = 25,
                    coef_plot_vars = NULL,
                    heatmap_vars = NULL,
                    plot_associations = TRUE,
                    max_pngs = 30,
```

```
balanced = FALSE,
save_plots_rds = FALSE)
```

## **Arguments**

output The output folder to write results.

transformed\_data

A data frame of transformed feature abundances. It should be formatted with features as columns and samples as rows. The column and row names should be the feature names and sample names respectively.

unstandardized\_metadata

A data frame of per-sample metadata. It should be formatted with variables as columns and samples as rows. The column and row names should be the variable names and sample names respectively.

fit\_data\_abundance

The abundance outputs of maaslin\_fit.

fit\_data\_prevalence

The prevalence outputs of maaslin\_fit.

normalization

The normalization to apply to the features before transformation and analysis. The option TSS (total sum scaling) is recommended, but CLR (centered log ratio) and NONE can also be used.

transform

The transformation to apply to the features after normalization and before analysis. The option LOG (base 2) is recommended, but PLOG (pseudo-log) and NONE can also be used.

feature\_specific\_covariate

A table of feature-specific covariates or a filepath to a tab-delimited file with feature-specific covariates. It should be formatted with features as columns and samples as rows (or the transpose). The row names and column names should be the same as those of the input\_data: the column and row names should be the feature names and sample names respectively. Typically, this table should be generated by 'preprocess\_mgx\_mtx' or 'preprocess\_taxa\_mtx' first.

feature\_specific\_covariate\_name

The name for the feature-specific covariates when fitting the models.

feature\_specific\_covariate\_record

Whether to keep the feature-specific covariates in the outputs when calculating p-values, writing results, and displaying plots.

median\_comparison\_abundance

Test abundance coefficients against a null value corresponding to the median coefficient for a metadata variable across the features. This is recommended for relative abundance data but should not be used for absolute abundance data.

median\_comparison\_prevalence

Test prevalence coefficients against a null value corresponding to the median coefficient for a metadata variable across the features. This is only recommended if the analyst is interested in how feature prevalence associations compare to each other or if there is likely strong compositionality-induced sparsity.

max\_significance

The FDR corrected q-value threshold for significance used in selecting which associations to write as significant and to plot.

plot\_summary\_plot

Generate a summary plot of significant associations.

summary\_plot\_first\_n

Include the top summary\_plot\_first\_n features with significant associations.

coef\_plot\_vars Vector of variable names to be used in the coefficient plot section of the sum-

mary plot. Continuous variables should match the metadata column name, and

categorical variables should be of the form "[variable] [level]".

heatmap\_vars Vector of variable names to be used in the heatmap section of the summary plot.

Continuous variables should match the metadata column name, and categorical

variables should be of the form "[variable] [level]".

plot\_associations

Whether to generate plots for significant associations.

max\_pngs The top max\_pngs significant associations will be plotted.

balanced If set to TRUE the summary plot will show the top N features of each variable in-

cluded in coef\_plot\_vars where N is equal to: ceiling(summary\_plot\_first\_n/length(coef\_plot\_

Will error if coef\_plot\_vars = NULL

save\_plots\_rds Whether to return the plots to an RDS file.

#### Value

Results will be written to the figures folder within the folder output. The list of individual association plots is returned if plot\_associations=TRUE. In the heatmap of the summary plot, one star corresponds to the user-set max\_significance and two stars corresponds to the user-set max\_signifiance divided by 10.

#### Author(s)

William Nickols<willnickols@g.harvard.edu>,

Jacob Nearing<nearing@broadinstitute.org>,

Maintainers: Lauren McIver<lauren.j.mciver@gmail.com>,

```
# Read features table
taxa_table_name <- system.file("extdata", "HMP2_taxonomy.tsv", package =
"maaslin3")
taxa_table <- read.csv(taxa_table_name, sep = '\t', row.names = 1)
# Read metadata table
metadata_name <- system.file("extdata", "HMP2_metadata.tsv", package =
"maaslin3")
metadata <- read.csv(metadata_name, sep = '\t', row.names = 1)
metadata$diagnosis <-</pre>
```

```
factor(metadata$diagnosis, levels = c('nonIBD', 'UC', 'CD'))
metadata$dysbiosis_state <-</pre>
    factor(metadata$dysbiosis_state, levels = c('none', 'dysbiosis_UC',
    'dysbiosis_CD'))
metadata$antibiotics <-</pre>
    factor(metadata$antibiotics, levels = c('No', 'Yes'))
#Run MaAsLin3
maaslin3::maaslin_log_arguments(
    input_data = taxa_table,
    input_metadata = metadata,
    output = 'output',
    formula = '~ diagnosis + dysbiosis_state + antibiotics +
    age + reads')
read_data_list <- maaslin3::maaslin_read_data(</pre>
    taxa_table,
    metadata)
    read_data_list <- maaslin3::maaslin_reorder_data(</pre>
    read_data_list$data,
    read_data_list$metadata)
data <- read_data_list$data</pre>
metadata <- read_data_list$metadata</pre>
formulas <- maaslin3::maaslin_check_formula(</pre>
    data,
    metadata,
    input_formula = '~ diagnosis + dysbiosis_state + antibiotics +
    age + reads')
formula <- formulas$formula</pre>
random_effects_formula <- formulas$random_effects_formula</pre>
normalized_data = maaslin3::maaslin_normalize(data,
                                 output = 'output')
filtered_data = maaslin3::maaslin_filter(normalized_data,
                                 output = 'output')
transformed_data = maaslin3::maaslin_transform(filtered_data,
                                 output = 'output')
standardized_metadata = maaslin3::maaslin_process_metadata(
    metadata,
    formula = formula)
maaslin_results = maaslin3::maaslin_fit(
    filtered_data,
    transformed_data,
    standardized_metadata,
    formula,
    random_effects_formula,
```

```
warn_prevalence = FALSE)
maaslin3::maaslin_write_results(
    output = 'output',
    maaslin_results$fit_data_abundance,
    maaslin_results$fit_data_prevalence,
    random_effects_formula)
maaslin3::maaslin_plot_results(
    output = 'output',
    transformed_data,
    metadata,
    maaslin_results$fit_data_abundance,
    maaslin_results$fit_data_prevalence,
    normalization = "TSS",
    transform = "LOG")
unlink('output', recursive=TRUE)
logging::logReset()
```

```
maaslin_plot_results_from_output

Plot the results from a MaAsLin 3 run.
```

## Description

Two types of plots are generated. First, the summary plot contains sorted per-feature coefficients plotted with their standard errors for key variables and a heatmap summarizing the remaining variables. Second, for significant features, association plots (scatterplots, boxplots, or tables depending on the association) are generated to visualize and verify the model fits. The data are shown with their transformed values in the association plots since this is the scale on which the models are fit. In comparison to maaslin\_plot\_results that needs the entire maaslin\_fit list, only the parameter list and an outputs directory containing a completed run are needed for maaslin\_plot\_results\_from\_output.

## **Usage**

coef\_plot\_vars = NULL,
heatmap\_vars = NULL,
plot\_associations = TRUE,
max\_pngs = 30,
balanced = FALSE,
save\_plots\_rds = FALSE)

## **Arguments**

output The output folder to write results.

metadata A data frame of per-sample metadata. It should be formatted with variables as

columns and samples as rows. The column and row names should be the variable

names and sample names respectively.

normalization The normalization to apply to the features before transformation and analysis.

The option TSS (total sum scaling) is recommended, but CLR (centered log ratio)

and NONE can also be used.

transform The transformation to apply to the features after normalization and before anal-

ysis. The option LOG (base 2) is recommended, but PLOG (pseudo-log) and NONE

can also be used.

feature\_specific\_covariate

A table of feature-specific covariates or a filepath to a tab-delimited file with feature-specific covariates. It should be formatted with features as columns and samples as rows (or the transpose). The row names and column names should be the same as those of the input\_data: the column and row names should be the feature names and sample names respectively. Typically, this table should be generated by 'preprocess\_mgx\_mtx' or 'preprocess\_taxa\_mtx' first.

feature\_specific\_covariate\_name

The name for the feature-specific covariates when fitting the models.

feature\_specific\_covariate\_record

Whether to keep the feature-specific covariates in the outputs when calculating p-values, writing results, and displaying plots.

median\_comparison\_abundance

Test abundance coefficients against a null value corresponding to the median coefficient for a metadata variable across the features. This is recommended for relative abundance data but should not be used for absolute abundance data.

median\_comparison\_prevalence

Test prevalence coefficients against a null value corresponding to the median coefficient for a metadata variable across the features. This is only recommended if the analyst is interested in how feature prevalence associations compare to each other or if there is likely strong compositionality-induced sparsity.

max\_significance

The FDR corrected q-value threshold for significance used in selecting which associations to write as significant and to plot.

plot\_summary\_plot

Generate a summary plot of significant associations.

summary\_plot\_first\_n

Include the top summary\_plot\_first\_n features with significant associations.

coef\_plot\_vars Vector of variable names to be used in the coefficient plot section of the sum-

mary plot. Continuous variables should match the metadata column name, and

categorical variables should be of the form "[variable] [level]".

heatmap\_vars Vector of variable names to be used in the heatmap section of the summary plot.

Continuous variables should match the metadata column name, and categorical

variables should be of the form "[variable] [level]".

plot\_associations

Whether to generate plots for significant associations.

max\_pngs The top max\_pngs significant associations will be plotted.

balanced If set to TRUE the summary plot will show the top N features of each variable in-

cluded in coef\_plot\_vars where N is equal to: ceiling(summary\_plot\_first\_n/length(coef\_plot\_

Will error if coef\_plot\_vars = NULL

save\_plots\_rds Whether to return the plots to an RDS file.

#### Value

Results will be written to the figures folder within the folder output. The list of individual association plots is returned if plot\_associations=TRUE. In the heatmap of the summary plot, one star corresponds to the user-set max\_significance and two stars corresponds to the user-set max\_signifiance divided by 10.

## Author(s)

William Nickols<willnickols@g.harvard.edu>,
Jacob Nearing<nearing@broadinstitute.org>,

Maintainers: Lauren McIver<lauren.j.mciver@gmail.com>,

```
# Read features table
taxa_table_name <- system.file("extdata", "HMP2_taxonomy.tsv", package =</pre>
"maaslin3")
taxa_table <- read.csv(taxa_table_name, sep = '\t', row.names = 1)</pre>
# Read metadata table
metadata_name <- system.file("extdata", "HMP2_metadata.tsv", package =</pre>
"maaslin3")
metadata <- read.csv(metadata_name, sep = '\t', row.names = 1)</pre>
metadata$diagnosis <-</pre>
    factor(metadata$diagnosis, levels = c('nonIBD', 'UC', 'CD'))
metadata$dysbiosis_state <-</pre>
    factor(metadata$dysbiosis_state, levels = c('none', 'dysbiosis_UC',
    'dysbiosis_CD'))
metadata$antibiotics <-</pre>
    factor(metadata$antibiotics, levels = c('No', 'Yes'))
#Run MaAsLin3
```

```
maaslin3::maaslin_log_arguments(
    input_data = taxa_table,
    input_metadata = metadata,
    output = 'output',
    formula = '~ diagnosis + dysbiosis_state + antibiotics +
    age + reads')
read_data_list <- maaslin3::maaslin_read_data(</pre>
    taxa_table,
    metadata)
read_data_list <- maaslin3::maaslin_reorder_data(</pre>
    read_data_list$data,
    read_data_list$metadata)
data <- read_data_list$data</pre>
metadata <- read_data_list$metadata</pre>
formulas <- maaslin3::maaslin_check_formula(</pre>
    data,
    metadata,
    input_formula = '~ diagnosis + dysbiosis_state + antibiotics +
    age + reads')
formula <- formulas$formula</pre>
random_effects_formula <- formulas$random_effects_formula</pre>
normalized_data = maaslin3::maaslin_normalize(data,
                             output = 'output')
filtered_data = maaslin3::maaslin_filter(normalized_data,
                             output = 'output')
transformed_data = maaslin3::maaslin_transform(filtered_data,
                             output = 'output')
standardized_metadata = maaslin3::maaslin_process_metadata(
    metadata,
    formula = formula)
maaslin_results = maaslin3::maaslin_fit(
    filtered_data,
    transformed_data,
    standardized_metadata,
    formula,
    {\tt random\_effects\_formula,}
    warn_prevalence = FALSE)
maaslin3::maaslin_write_results(
    output = 'output',
    maaslin_results$fit_data_abundance,
    maaslin_results$fit_data_prevalence,
    random_effects_formula)
```

```
maaslin3::maaslin_plot_results_from_output(
   output = 'output',
   metadata,
   normalization = "TSS",
   transform = "LOG")

unlink('output', recursive=TRUE)
logging::logReset()
```

maaslin\_process\_metadata

Process metadata before MaAsLin 3 model fitting.

# Description

Check that references are set properly if the metadata variables are categorical and provided through fixed\_effects. Standardize the continuous metadata variables as a z-score (subtract the mean, divide by the standard deviation) if standardize is set.

## Usage

## **Arguments**

metadata A data frame of per-sample metadata. It should be formatted with variables as

columns and samples as rows. The column and row names should be the variable

names and sample names respectively.

formula A formula in lme4 format. Random effects, interactions, and functions of the

metadata can be included (note that these functions will be applied after standardization if standardize=TRUE). Group, ordered, and strata variables can be specified as: group(grouping\_variable), ordered(ordered\_variable) and strata(strata\_variable). The other variable options below will not be con-

sidered if a formula is set.

fixed\_effects A vector of variable names to be included as fixed effects.

reference For a variable with more than two levels supplied with fixed\_effects, the

factor to use as a reference provided as a string of 'variable,reference' semi-

colon delimited for multiple variables.

feature\_specific\_covariate\_name

The name for the feature-specific covariates when fitting the models.

standardize

Whether to apply z-scores to continuous metadata variables so they are on the same scale. This is recommended in order to compare coefficients across metadata variables, but note that functions of the metadata specified in the formula will apply after standardization.

#### Value

The processed metadata.

# Author(s)

William Nickols<willnickols@g.harvard.edu>,
Jacob Nearing<nearing@broadinstitute.org>,
Maintainers: Lauren McIver<lauren.j.mciver@gmail.com>,

```
# Read features table
taxa_table_name <- system.file("extdata", "HMP2_taxonomy.tsv", package =</pre>
"maaslin3")
taxa_table <- read.csv(taxa_table_name, sep = '\t', row.names = 1)</pre>
# Read metadata table
metadata_name <- system.file("extdata", "HMP2_metadata.tsv", package =</pre>
"maaslin3")
metadata <- read.csv(metadata_name, sep = '\t', row.names = 1)</pre>
metadata$diagnosis <-</pre>
    factor(metadata$diagnosis, levels = c('nonIBD', 'UC', 'CD'))
metadata$dysbiosis_state <-</pre>
    factor(metadata$dysbiosis_state, levels = c('none', 'dysbiosis_UC',
     'dysbiosis_CD'))
metadata$antibiotics <-</pre>
    factor(metadata$antibiotics, levels = c('No', 'Yes'))
#Run MaAsLin3
maaslin3::maaslin_log_arguments(
    input_data = taxa_table,
    input_metadata = metadata,
    output = 'output',
    formula = '~ diagnosis + dysbiosis_state + antibiotics +
    age + reads',
    plot_summary_plot = FALSE,
    plot_associations = FALSE)
read_data_list <- maaslin3::maaslin_read_data(</pre>
    taxa_table,
    metadata)
    read_data_list <- maaslin3::maaslin_reorder_data(</pre>
    read_data_list$data,
    read_data_list$metadata)
```

41 maaslin\_read\_data

```
data <- read_data_list$data
metadata <- read_data_list$metadata</pre>
formulas <- maaslin3::maaslin_check_formula(</pre>
    data,
    metadata,
    input_formula = '~ diagnosis + dysbiosis_state + antibiotics +
    age + reads')
formula <- formulas$formula</pre>
random_effects_formula <- formulas$random_effects_formula</pre>
normalized_data = maaslin3::maaslin_normalize(data,
                                 output = 'output')
filtered_data = maaslin3::maaslin_filter(normalized_data,
                                 output = 'output')
standardized_metadata = maaslin3::maaslin_process_metadata(
    metadata,
    formula = formula)
unlink('output', recursive=TRUE)
logging::logReset()
```

maaslin\_read\_data

Read in the abundance data and metadata.

## **Description**

Read in the abundance data and metadata from files if necessary.

## **Usage**

```
maaslin_read_data(input_data,
                input_metadata,
                feature_specific_covariate = NULL,
                unscaled_abundance = NULL)
```

# Arguments

input\_data

A data frame of feature abundances or read counts or a filepath to a tab-delimited file with abundances. It should be formatted with features as columns and samples as rows (or the transpose). The column and row names should be the feature names and sample names respectively.

input\_metadata A data frame of per-sample metadata or a filepath to a tab-delimited file with metadata. It should be formatted with variables as columns and samples as rows (or the transpose). The column and row names should be the variable names and sample names respectively.

42 maaslin\_read\_data

feature\_specific\_covariate

A table of feature-specific covariates or a filepath to a tab-delimited file with feature-specific covariates. It should be formatted with features as columns and samples as rows (or the transpose). The row names and column names should be the same as those of the input\_data: the column and row names should be the feature names and sample names respectively. Typically, this table should be generated by 'preprocess\_mgx\_mtx' or 'preprocess\_taxa\_mtx' first.

## unscaled\_abundance

A data frame with a single column of absolute abundances or a filepath to such a tab-delimited file. The row names should match the names of the samples in input\_data and input\_metadata. When using spike-ins, the single column should have the same name as one of the features in input\_data, and the unscaled\_abundance should correspond to the absolute quantity of the spike-in. When using total abundance scaling, the single column should have the name 'total', and the unscaled\_abundance should correspond to the total abundance of each sample.

#### Value

A list containing the following items:

- (1) data: A data frame of feature abundances.
- (2) metadata: A data frame of metadata.
- (3) feature\_specific\_covariate: A data frame of feature specific covariates.
- (4) unscaled\_abundance: A data frame of unscaled abundances.

#### Author(s)

William Nickols<br/><willnickols@g.harvard.edu>,<br/>Jacob Nearing<br/>@broadinstitute.org>,<br/>Maintainers: Lauren McIver<lauren.j.mciver@gmail.com>,

maaslin\_reorder\_data 43

```
metadata$antibiotics <-</pre>
    factor(metadata$antibiotics, levels = c('No', 'Yes'))
#Run MaAsLin3
maaslin3::maaslin_log_arguments(
    input_data = taxa_table,
    input_metadata = metadata,
    output = 'output',
    formula = '~ diagnosis + dysbiosis_state + antibiotics +
    age + reads',
    plot_summary_plot = FALSE,
    plot_associations = FALSE)
read_data_list <- maaslin3::maaslin_read_data(</pre>
    taxa_table,
    metadata)
unlink('output', recursive=TRUE)
logging::logReset()
```

maaslin\_reorder\_data Reorder the abundance data and metadata.

## Description

Reorder the abundance data and metadata to ensure samples are rows and remove any samples without abundances or metadata.

# Usage

## **Arguments**

data

A data frame of feature abundances or read counts. It should be formatted with features as columns and samples as rows (or the transpose). The column and row names should be the feature names and sample names respectively.

metadata

A data frame of per-sample metadata. It should be formatted with variables as columns and samples as rows (or the transpose). The column and row names should be the variable names and sample names respectively.

feature\_specific\_covariate

A table of feature-specific covariates or a filepath to a tab-delimited file with feature-specific covariates. It should be formatted with features as columns and samples as rows (or the transpose). The row names and column names should be the same as those of the input\_data: the column and row names should be the feature names and sample names respectively. Typically, this table should be generated by 'preprocess\_mgx\_mtx' or 'preprocess\_taxa\_mtx' first.

44 maaslin\_reorder\_data

unscaled\_abundance

A data frame with a single column of absolute abundances. The row names should match the names of the samples in input\_data and input\_metadata. When using spike-ins, the single column should have the same name as one of the features in input\_data, and the unscaled\_abundance should correspond to the absolute quantity of the spike-in. When using total abundance scaling, the single column should have the name 'total', and the unscaled\_abundance should correspond to the total abundance of each sample.

#### Value

A list containing the following items:

- (1) data: A data frame of feature abundances.
- (2) metadata: A data frame of metadata.
- (3) feature\_specific\_covariate: A data frame of feature specific covariates.
- (4) unscaled\_abundance: A data frame of unscaled abundances.

#### Author(s)

William Nickols<br/><willnickols@g.harvard.edu>,<br/>Jacob Nearing<br/><nearing@broadinstitute.org>,<br/>Maintainers: Lauren McIver<lauren.j.mciver@gmail.com>,

```
# Read features table
taxa_table_name <- system.file("extdata", "HMP2_taxonomy.tsv", package =</pre>
"maaslin3")
taxa_table <- read.csv(taxa_table_name, sep = '\t', row.names = 1)</pre>
# Read metadata table
metadata_name <- system.file("extdata", "HMP2_metadata.tsv", package =</pre>
"maaslin3")
metadata <- read.csv(metadata_name, sep = '\t', row.names = 1)</pre>
metadata$diagnosis <-</pre>
    factor(metadata$diagnosis, levels = c('nonIBD', 'UC', 'CD'))
metadata$dysbiosis_state <-</pre>
    factor(metadata$dysbiosis_state, levels = c('none', 'dysbiosis_UC',
    'dysbiosis_CD'))
metadata$antibiotics <-</pre>
    factor(metadata$antibiotics, levels = c('No', 'Yes'))
#Run MaAsLin3
maaslin3::maaslin_log_arguments(
    input_data = taxa_table,
    input_metadata = metadata,
    output = 'output',
    formula = '~ diagnosis + dysbiosis_state + antibiotics +
```

maaslin\_transform 45

```
age + reads',
plot_summary_plot = FALSE,
plot_associations = FALSE)
read_data_list <- maaslin3::maaslin_read_data(
    taxa_table,
    metadata)
read_data_list <- maaslin3::maaslin_reorder_data(
    taxa_table,
    metadata)
unlink('output', recursive=TRUE)
logging::logReset()</pre>
```

maaslin\_transform

Transform abundance data for MaAsLin 3 modeling.

# **Description**

Transform the abundance data according to the transform parameter.

## Usage

# Arguments

filtered\_data A data frame of filtered feature abundances. It should be formatted with features

as columns and samples as rows. The column and row names should be the

feature names and sample names respectively.

output The output folder to write results.

transform The transformation to apply to the features after normalization and before anal-

ysis. The option LOG (base 2) is recommended, but PLOG (pseudo-log) and NONE

can also be used.

## Value

A dataframe of transformed features (features are columns; samples are rows).

## Author(s)

William Nickols<br/>
willnickols@g.harvard.edu>,<br/>
Jacob Nearing<nearing@broadinstitute.org>,<br/>
Maintainers: Lauren McIver<lauren.j.mciver@gmail.com>,

46 maaslin\_transform

```
# Read features table
taxa_table_name <- system.file("extdata", "HMP2_taxonomy.tsv", package =</pre>
"maaslin3")
taxa_table <- read.csv(taxa_table_name, sep = '\t', row.names = 1)</pre>
# Read metadata table
metadata_name <- system.file("extdata", "HMP2_metadata.tsv", package =</pre>
"maaslin3")
metadata <- read.csv(metadata_name, sep = '\t', row.names = 1)</pre>
metadata$diagnosis <-</pre>
    factor(metadata$diagnosis, levels = c('nonIBD', 'UC', 'CD'))
metadata$dysbiosis_state <-</pre>
    factor(metadata$dysbiosis_state, levels = c('none', 'dysbiosis_UC',
    'dysbiosis_CD'))
metadata$antibiotics <-</pre>
    factor(metadata$antibiotics, levels = c('No', 'Yes'))
#Run MaAsLin3
maaslin3::maaslin_log_arguments(
    input_data = taxa_table,
    input_metadata = metadata,
    output = 'output',
    formula = '~ diagnosis + dysbiosis_state + antibiotics +
    age + reads',
    plot_summary_plot = FALSE,
    plot_associations = FALSE)
read_data_list <- maaslin3::maaslin_read_data(</pre>
    taxa_table,
    metadata)
read_data_list <- maaslin3::maaslin_reorder_data(</pre>
    read_data_list$data,
    read_data_list$metadata)
data <- read_data_list$data</pre>
metadata <- read_data_list$metadata</pre>
formulas <- maaslin3::maaslin_check_formula(</pre>
    data,
    metadata,
    input_formula = '~ diagnosis + dysbiosis_state + antibiotics +
    age + reads')
normalized_data = maaslin3::maaslin_normalize(data,
                                  output = 'output')
filtered_data = maaslin3::maaslin_filter(normalized_data,
                                  output = 'output')
transformed_data = maaslin3::maaslin_transform(filtered_data,
```

maaslin\_write\_results 47

```
output = 'output')
unlink('output', recursive=TRUE)
logging::logReset()
```

maaslin\_write\_results Write the results from a MaAsLin 3 run.

#### **Description**

Write the results from a MaAsLin 3 run to the output folder as a TSV.

### Usage

## **Arguments**

#### Value

Results will be written to the all\_results.tsv and significant\_results.tsv files in the folder output. The file all\_results.tsv will contain all results in the fit\_data\_abundance and fit\_data\_prevalence items of the input list (with 'linear' and 'logistic' replaced by 'abundance' and 'prevalence' in the model column). The file significant\_results.tsv will contain all results with joint or individual q-values below the 'max\_significance' parameter. No value is returned.

#### Author(s)

```
William Nickols<willnickols@g.harvard.edu>,
Jacob Nearing<nearing@broadinstitute.org>,
Maintainers: Lauren McIver<lauren.j.mciver@gmail.com>,
```

# **Examples**

48

```
# Read features table
    taxa_table_name <- system.file("extdata", "HMP2_taxonomy.tsv", package =</pre>
    "maaslin3")
    taxa_table <- read.csv(taxa_table_name, sep = '\t', row.names = 1)</pre>
    # Read metadata table
   metadata_name <- system.file("extdata", "HMP2_metadata.tsv", package =</pre>
    "maaslin3")
   metadata <- read.csv(metadata_name, sep = '\t', row.names = 1)</pre>
   metadata$diagnosis <-</pre>
        factor(metadata$diagnosis, levels = c('nonIBD', 'UC', 'CD'))
   metadata$dysbiosis_state <-</pre>
        factor(metadata$dysbiosis_state, levels = c('none', 'dysbiosis_UC',
        'dysbiosis_CD'))
   metadata$antibiotics <-</pre>
        factor(metadata$antibiotics, levels = c('No', 'Yes'))
    #Run MaAsLin3
   maaslin3::maaslin_log_arguments(
        input_data = taxa_table,
        input_metadata = metadata,
        output = 'output',
        formula = '~ diagnosis + dysbiosis_state + antibiotics +
        age + reads',
        plot_summary_plot = FALSE,
        plot_associations = FALSE)
    read_data_list <- maaslin3::maaslin_read_data(</pre>
        taxa_table,
        metadata)
        read_data_list <- maaslin3::maaslin_reorder_data(</pre>
        read_data_list$data,
        read_data_list$metadata)
    data <- read_data_list$data</pre>
   metadata <- read_data_list$metadata</pre>
    formulas <- maaslin3::maaslin_check_formula(</pre>
        data,
        metadata,
        input_formula = '~ diagnosis + dysbiosis_state + antibiotics +
        age + reads')
    formula <- formulas$formula</pre>
    random_effects_formula <- formulas$random_effects_formula</pre>
   normalized_data = maaslin3::maaslin_normalize(data,
                                                   output = 'output')
    filtered_data = maaslin3::maaslin_filter(normalized_data,
```

```
output = 'output')
transformed_data = maaslin3::maaslin_transform(filtered_data,
                                            output = 'output')
standardized_metadata = maaslin3::maaslin_process_metadata(
    metadata,
    formula = formula)
maaslin_results = maaslin3::maaslin_fit(
    filtered_data,
    transformed_data,
    standardized_metadata,
    formula,
    random_effects_formula,
    warn_prevalence = FALSE)
maaslin3::maaslin_write_results(
   output = 'output',
    maaslin_results$fit_data_abundance,
    maaslin_results$fit_data_prevalence,
    random_effects_formula)
unlink('output', recursive=TRUE)
logging::logReset()
```

```
maaslin_write_results_lefse_format

Write the results from a MaAsLin 3 run in LEfSe format.
```

# Description

Write the results from a MaAsLin 3 run to the output folder in LEfSe format.

## Usage

# **Arguments**

#### Value

Results will be written to the lefse\_style\_results\_abundance.res file in the folder output. No value is returned.

#### Author(s)

William Nickols<willnickols@g.harvard.edu>, Jacob Nearing<nearing@broadinstitute.org>, Maintainers: Lauren McIver<lauren.j.mciver@gmail.com>,

```
# Read features table
taxa_table_name <- system.file("extdata", "HMP2_taxonomy.tsv", package =</pre>
"maaslin3")
taxa_table <- read.csv(taxa_table_name, sep = '\t', row.names = 1)</pre>
# Read metadata table
metadata_name <- system.file("extdata", "HMP2_metadata.tsv", package =</pre>
"maaslin3")
metadata <- read.csv(metadata_name, sep = '\t', row.names = 1)</pre>
metadata$diagnosis <-</pre>
    factor(metadata$diagnosis, levels = c('nonIBD', 'UC', 'CD'))
metadata$dysbiosis_state <-</pre>
    factor(metadata$dysbiosis_state, levels = c('none', 'dysbiosis_UC',
    'dysbiosis_CD'))
metadata$antibiotics <-
    factor(metadata$antibiotics, levels = c('No', 'Yes'))
#Run MaAsLin3
maaslin3::maaslin_log_arguments(
    input_data = taxa_table,
    input_metadata = metadata,
    output = 'output',
    formula = '~ diagnosis + dysbiosis_state + antibiotics +
    age + reads',
    plot_summary_plot = FALSE,
    plot_associations = FALSE)
read_data_list <- maaslin3::maaslin_read_data(</pre>
    taxa_table,
    metadata)
    read_data_list <- maaslin3::maaslin_reorder_data(</pre>
    read_data_list$data,
    read_data_list$metadata)
data <- read_data_list$data</pre>
metadata <- read_data_list$metadata</pre>
formulas <- maaslin3::maaslin_check_formula(</pre>
```

preprocess\_dna\_mtx 51

```
data,
    metadata,
    input_formula = '~ diagnosis + dysbiosis_state + antibiotics +
    age + reads')
formula <- formulas$formula</pre>
random_effects_formula <- formulas$random_effects_formula</pre>
normalized_data = maaslin3::maaslin_normalize(data,
                                 output = 'output')
filtered_data = maaslin3::maaslin_filter(normalized_data,
                                 output = 'output')
transformed_data = maaslin3::maaslin_transform(filtered_data,
                                 output = 'output')
standardized_metadata = maaslin3::maaslin_process_metadata(
    metadata,
    formula = formula)
maaslin_results = maaslin3::maaslin_fit(
    filtered_data,
    transformed_data,
    standardized_metadata,
    formula,
    random_effects_formula,
    warn_prevalence = FALSE)
maaslin3::maaslin_write_results_lefse_format(
    output = 'output',
    maaslin_results$fit_data_abundance,
    maaslin_results$fit_data_prevalence)
unlink('output', recursive=TRUE)
logging::logReset()
```

preprocess\_dna\_mtx

Pre-process the DNA covariates for metatranscriptomics

## **Description**

Pre-process the DNA covariates for metatranscriptomics by total-sum-scaling DNA abundances per sample and then, for each sample in each feature:

- 1. Log 2 transforming the DNA abundance if the DNA abundance is >=0
- 2. Setting the DNA abundance to log2([minimum non-zero relative abundance in the dataset] / 2) if the corresponding RNA abundance is non-zero but the DNA abundance is zero
- 3. Setting the DNA abundance to NA if both are zero

52 preprocess\_dna\_mtx

When the DNA is present, the RNA data can be modeled as usual in MaAsLin 3 with log2(DNA) as a covariate. When the DNA is not present, if the RNA is present, we assume the DNA was missed due to finite read depth, so the DNA abundance is imputed with a small pseudo-count. When neither the DNA nor RNA is present, we assume the gene/microbe was not in the sample and therefore no information about the transcription level can be obtained. Setting the DNA covariate to NA has the effect of dropping the sample when fitting the relevant feature's model in MaAsLin 3. Unlike most MaAsLin functions that will infer the samples from the row names and column names, the rna\_table must be formated as samples (rows) by features (columns).

## Usage

```
preprocess_dna_mtx(dna_table, rna_table)
```

## **Arguments**

dna\_table The samples (rows) by features (columns) data frame of DNA abundances to

preprocess. These can be relative abundances or counts.

rna\_table The samples (rows) by features (columns) data frame of RNA to preprocess.

These can be relative abundances or counts.

#### Value

A list containing the following named items:

- 1. dna\_table: The table of log2 transformed DNA relative abundances with NAs for any feature-sample pairs for which both the DNA and RNA abundances were 0.
- 2. rna\_table: The table of total sum scaled RNA abundances. These are not log2 transformed.

# Author(s)

William Nickols<willnickols@g.harvard.edu>,

Jacob Nearing<nearing@broadinstitute.org>,

Maintainers: Lauren McIver<lauren.j.mciver@gmail.com>,

preprocess\_taxa\_mtx 53

preprocess\_taxa\_mtx

Pre-process the taxa covariates for metatranscriptomics

# Description

Pre-process the taxa covariates for metatranscriptomics by total-sum-scaling the taxa, matching the taxa to the RNAs coming from those taxa, and then, for each sample in each feature:

- 1. Log 2 transforming the taxon abundance if the taxon abundance is  $\geq 0$
- 2. Setting the taxon abundance to log2([minimum non-zero relative abundance in the dataset] / 2) if any of the corresponding RNA abundances are non-zero but the taxon abundance is zero
- 3. Setting the taxon abundance to NA if both are zero

When the taxon is present, the RNA data can be modeled as usual in MaAsLin 3 with log2(taxon) as a covariate. When the taxon is not present, if any of its RNA is present, we assume the taxon was missed due to finite read depth, so the taxon abundance is imputed with a small pseudo-count. When neither the taxon nor RNA is present, we assume the gene/microbe was not in the sample and therefore no information about the transcription level can be obtained. Setting the taxon covariate to NA has the effect of dropping the sample when fitting the relevant feature's model in MaAsLin 3. Unlike most MaAsLin functions that will infer the samples from the row names and column names, the rna\_table must be formated as samples (rows) by features (columns).

## Usage

```
preprocess_taxa_mtx(taxa_table, rna_table, rna_per_taxon)
```

## **Arguments**

taxa_table	The samples (rows) by features (columns) data frame of taxon abundances to preprocess. These can be relative abundances or counts.
rna_table	The samples (rows) by features (columns) data frame of RNA to preprocess. These can be relative abundances or counts.
rna_per_taxon	A dataframe with the columns 'RNA' and 'taxon' with one row per 'RNA' column found in 'rna_table' giving both the 'RNA' column and which 'taxon' column it corresponds to in 'taxa_table'.

#### Value

A list containing the following named items:

- 1. dna\_table: The table of log2 transformed taxon relative abundances with NAs for any feature-sample pairs for which both the taxon and RNA abundances were 0.
- 2. rna\_table: The table of total sum scaled RNA abundances. These are not log2 transformed.

54 preprocess\_taxa\_mtx

# Author(s)

William Nickols<br/><willnickols@g.harvard.edu>,<br/>Jacob Nearing<nearing@broadinstitute.org>,<br/>Maintainers: Lauren McIver<lauren.j.mciver@gmail.com>,

# **Index**

```
maaslin3, 2
maaslin_check_arguments, 8
maaslin_check_formula, 10, 20, 47
{\tt maaslin\_compute\_formula, 12}
maaslin_contrast_test, 14
maaslin_filter, 17
maaslin_fit, 7, 15, 19, 32, 47, 49
{\tt maaslin\_log\_arguments}, 24
maaslin_normalize, 29
maaslin_plot_results, 31
maaslin_plot_results_from_output, 35
maaslin_process_metadata, 39
maaslin_read_data, 41
maaslin_reorder_data, 43
maaslin_transform, 45
maaslin_write_results, 47
maaslin_write_results_lefse_format, 49
preprocess_dna_mtx, 51
preprocess_taxa_mtx, 53
```