

Package ‘MicrobiomeStat’

March 3, 2026

Type Package

Title Statistical Methods for Microbiome Compositional Data

Version 1.4

Date 2026-03-03

Author Xianyang Zhang [aut],
Jun Chen [aut, cre],
Huijuan Zhou [ctb],
Linsui Deng [ctb]

Maintainer Jun Chen <chen.jun2@mayo.edu>

Description A suite of methods for powerful and robust microbiome data analysis addressing zero-inflation, phylogenetic structure and compositional effects. Includes the LinDA method for differential abundance analysis (Zhou et al. (2022)<[doi:10.1186/s13059-022-02655-5](https://doi.org/10.1186/s13059-022-02655-5)>), the BMDD (Bimodal Dirichlet Distribution) method for accurate modeling and imputation of zero-inflated microbiome sequencing data (Zhou et al. (2025)<[doi:10.1371/journal.pcbi.1013124](https://doi.org/10.1371/journal.pcbi.1013124)>) and compositional sparse CCA methods for microbiome multi-omics data integration (Deng et al. (2024) <[doi:10.3389/fgene.2024.1489694](https://doi.org/10.3389/fgene.2024.1489694)>).

Depends R (>= 3.5.0)

Imports ggplot2, matrixStats, parallel, stats, utils, Matrix, statmod,
MASS, ggrepel, lmerTest, foreach, modeest, dplyr, Rcpp, mlr3,
mlr3mbo, bbotk, paradox

LinkingTo Rcpp, RcppArmadillo

Suggests DiceKriging, randomForest

NeedsCompilation yes

SystemRequirements NLOpt library (optional, for high-performance BMDD mode)

License GPL-3

Encoding UTF-8

LazyData true

RoxygenNote 7.3.3

Repository CRAN

Date/Publication 2026-03-03 21:40:02 UTC

Contents

bmdd	2
bmdd.nlopt	5
escca	6
escca.CV	8
DGP_OC	10
linda	11
linda.plot	15
linda.wald.test	17
smokers	19
Index	20

 bmdd

Bimodal Dirichlet Distribution Estimation

Description

Estimates parameters of the bimodal Dirichlet distribution using a variational EM algorithm. Automatically selects the optimal implementation: high-performance C++ (NLOpt) when possible, or R when covariates are needed.

Usage

```
bmdd(W, type = c("count", "proportion"),
     Z = NULL, formula.Z = NULL, U = NULL, formula.U = NULL,
     Z.standardizing = TRUE, U.standardizing = TRUE,
     alp.eta = FALSE, alp.kap = FALSE,
     pi.xi = FALSE, pi.zeta = FALSE,
     para.alp.init = NULL, para.pi.init = NULL, gam.init = NULL,
     iterlim = 500, tol = 1e-6, trace = FALSE,
     method = c("auto", "nlopt", "R"),
     inner.loop = TRUE, inner.iterlim = 20, inner.tol = 1e-6,
     alp.iterlim = 100, alp.tol = 1e-6,
     alp.min = 1e-3, alp.max = 1e3, ...)
```

Arguments

W	Matrix of observed data (m taxa × n samples)
type	Data type: "count" or "proportion"
Z	Optional matrix of covariates for alpha (forces R implementation)
formula.Z	Optional formula for Z covariates
U	Optional matrix of covariates for pi (forces R implementation)
formula.U	Optional formula for U covariates
Z.standardizing	Standardize Z covariates (default TRUE)

U.standardizing	Standardize U covariates (default TRUE)
alp.eta	Model alpha0 as function of Z (default FALSE)
alp.kap	Model alpha1 as function of Z (default FALSE)
pi.xi	Model pi as function of U (default FALSE)
pi.zeta	Model pi variance as function of U (default FALSE)
para.alp.init	Initial alpha values
para.pi.init	Initial pi values
gam.init	Initial gamma values
iterlim	Maximum iterations (default 500)
tol	Convergence tolerance (default 1e-6)
trace	Print progress (default FALSE)
method	Force method: "auto", "nlopt", or "R" (default "auto")
inner.loop	Use inner loop for NLOpt (default TRUE)
inner.iterlim	Inner loop max iterations (default 20)
inner.tol	Inner loop tolerance (default 1e-6)
alp.iterlim	Alpha optimization iterations (default 100)
alp.tol	Alpha tolerance (default 1e-6)
alp.min	Minimum alpha (default 1e-3)
alp.max	Maximum alpha (default 1e3)
...	Additional arguments (ignored)

Details

Automatically chooses best implementation:

- **NLOpt C++**: When no covariates. 50-90x faster using L-BFGS-B with analytical gradients. Scientifically equivalent to R ($r > 0.999$).
- **R**: When covariates needed or explicitly requested. Supports full covariate modeling.

Requires NLOpt library for high-performance mode:

- macOS: `brew install nlopt`
- Linux: `sudo apt-get install libnlopt-dev`

Refer to <https://github.com/zhouhj1994/BMDD> for detailed examples about zero imputation and posterior sample generation.

Value

A list containing:

gamma	Estimated gamma parameters (bimodality indicators)
pi	Estimated pi parameters (mixing proportions)
beta	Estimated posterior Dirichlet parameters
alpha	a list with first element being the estimate of alpha0 and second the estimate of alpha1.
para.alpha	if no covariates are incorporated, then para.alpha is (alpha0, alpha1) itself. if there are covariates, then para.alpha is the estimate of the regression coefficients. The length of para.alpha is the same as para.alp.init
para.pi	if no covariates are incorporated, then para.pi is pi itself. if there are covariates, then para.pi is the estimate of the regression coefficients. The length of para.pi the is same as para.pi.init
iter	Number of iterations
method	Method used: "nlopt" or "R"

References

Zhou, H., Chen, J., & Zhang, X. (2025). BMDD: A probabilistic framework for accurate imputation of zero-inflated microbiome sequencing data. *PLoS Computational Biology*, 21(10), e1013124.

Examples

```
## Not run:
# Simulate data
m <- 100 # taxa
n <- 50 # samples
W <- matrix(rpois(m*n, 100), m, n)

# Auto-select method (uses NLOpt for speed)
result <- bmdd(W, type = "count")

# Access results
head(result$beta) # Posterior parameters
head(result$gamma) # Bimodality indicators
result$method # Shows which method was used

# Force specific method
result_nlopt <- bmdd(W, method = "nlopt") # Force NLOpt
result_r <- bmdd(W, method = "R") # Force R

# With covariates (automatically uses R)
Z <- matrix(rnorm(m * 2), m, 2)
result_cov <- bmdd(W, Z = Z, alp.eta = TRUE)

## End(Not run)
```

Description

High-performance implementation using NLopt L-BFGS-B optimizer in C++. Provides 50-90x speedup over R implementation for cases without covariates. This is a low-level function; most users should use `bmdd()` instead.

Usage

```
bmdd.nlopt(W, type = c("count", "proportion"),
           para.alp.init = NULL, para.pi.init = NULL, gam.init = NULL,
           iterlim = 500, tol = 1e-6, trace = FALSE,
           inner.loop = TRUE, inner.iterlim = 20, inner.tol = 1e-6,
           alp.iterlim = 100, alp.tol = 1e-6,
           alp.min = 1e-3, alp.max = 1e3)
```

Arguments

<code>W</code>	Matrix of observed data (m taxa × n samples)
<code>type</code>	Data type: "count" or "proportion"
<code>para.alp.init</code>	Initial alpha values
<code>para.pi.init</code>	Initial pi values
<code>gam.init</code>	Initial gamma values
<code>iterlim</code>	Maximum iterations (default 500)
<code>tol</code>	Convergence tolerance (default 1e-6)
<code>trace</code>	Print progress (default FALSE)
<code>inner.loop</code>	Use inner loop optimization (default TRUE)
<code>inner.iterlim</code>	Inner loop max iterations (default 20)
<code>inner.tol</code>	Inner loop tolerance (default 1e-6)
<code>alp.iterlim</code>	Alpha optimization iterations (default 100)
<code>alp.tol</code>	Alpha tolerance (default 1e-6)
<code>alp.min</code>	Minimum alpha (default 1e-3)
<code>alp.max</code>	Maximum alpha (default 1e3)

Details

This function provides direct access to the NLopt C++ implementation. Most users should use `bmdd()` instead, which automatically selects the optimal implementation.

Does not support covariates. For covariate modeling, use `bmdd()` with `method="R"`.

Value

A list containing estimated parameters (same as `bmdd`)

References

Zhou, H., Chen, J., & Zhang, X. (2025). BMDD: A probabilistic framework for accurate imputation of zero-inflated microbiome sequencing data. *PLoS Computational Biology*, 21(10), e1013124.

See Also

[bmdd](#) for the main user interface

Examples

```
## Not run:
# Simulate data
m <- 100
n <- 50
W <- matrix(rpois(m*n, 100), m, n)

# Direct NLOpt usage (advanced)
result <- bmdd.nlopt(W, type = "count")

# Recommended: use bmdd() instead
result <- bmdd(W, type = "count") # auto-selects NLOpt

## End(Not run)
```

cscca

Compositional Sparse Canonical Correlation Analysis

Description

A compositional sparse canonical correlation analysis (csCCA) framework for integrating microbiome data with other high-dimensional omics data.

Usage

```
cscca(
  Y,
  View.ind,
  lambda.seq,
  a.old = NULL,
  View.type = NULL,
  eps.stop = 1e-04,
  max.step = 30,
  eps = 1e-04,
  T.step = 1
)
```

Arguments

<code>Y</code>	a $n \times (K \times p)$ matrix representing the observations.
<code>View.ind</code>	a $(K \times p)$ integer vector indicating the classes of features. The features with the same <code>View.ind</code> is in the same class.
<code>lambda.seq</code>	a K vector consisting of hyper-parameters.
<code>a.old</code>	Optional initial value for the coefficient vector <code>a.new</code> .
<code>View.type</code>	a K vector encoding the structure type of each feature class. There are two choices: "O" (Omics Data), "C" (Compositional Data).
<code>eps.stop</code>	a numerical value controlling the convergence.
<code>max.step</code>	an integer controlling the maximum step for interaction.
<code>eps</code>	a numerical value controlling the convergence.
<code>T.step</code>	an integer controlling the maximum step for interaction.

Value

`a.new` the estimated coefficient vector.

References

- Deng, L., Tang, Y., Zhang, X., et al. (2024). Structure-adaptive canonical correlation analysis for microbiome multi-omics data. *Frontiers in Genetics*, 15, 1489694.
- Chen, J., Bushman, F. D., Lewis, J. D., et al. (2013). Structure-constrained sparse canonical correlation analysis with an application to microbiome data analysis. *Biostatistics*, 14(2), 244–258.

Examples

```
## Not run:
library(dplyr)

n <- 200
p <- q <- 100
sigma.nu <- 5
sigma.eps <- 1
omega_X <- 0.85*c(rep(1/10,9),-9/10,rep(0,p-10))
omega_Y <- 0.85*c(seq(0.08,0.12,length = 10),rep(0,q-10))
Data1 <- DGP_OC(seed=10,n,p,q,sigma.nu,sigma.eps,omega_X,omega_Y)

library(mIrmBO)
Res.sCCA.CV <- cscca.CV(Y=Data1$Y,View.ind=Data1$View.ind,
                      View.type=c("O","O"),
                      show.info = TRUE)

Res.CsCCA.CV <- cscca.CV(Y=Data1$Y,View.ind=Data1$View.ind,
                        View.type=c("O","C"),
                        show.info = TRUE)

Res.sCCA <- cscca(Y=Data1$Y,View.ind=Data1$View.ind,
```

```

lambda.seq=Res.sCCA.CV$lam.opt.trgt,
View.type=c("0","0"))
Res.CsCCA <- cscca(Y=Data1$Y,View.ind=Data1$View.ind,
lambda.seq=Res.CsCCA.CV$lam.opt.trgt,
View.type=c("0","C"))

## End(Not run)

```

cscca.CV

Compositional Sparse Canonical Correlation Analysis (Cross Valication Version)

Description

The cross validation version of a compositional sparse canonical correlation analysis (sCCA) framework for integrating microbiome data with other high-dimensional omics data.

Usage

```

cscca.CV(
  Y,
  View.ind,
  View.type = NULL,
  eps.stop = 1e-04,
  max.step = 30,
  eps = 1e-04,
  T.step = 10,
  n_fold = 5,
  seed.sam.ind = NULL,
  hp.lower = NULL,
  hp.upper = NULL,
  hp.eta.lower = NULL,
  hp.eta.upper = NULL,
  eta.warm.stat.mat = NULL,
  opt_n_design = 30,
  opt_n_iter = 20,
  Criterion = "cov",
  des.init = NULL,
  is.refit = F,
  is.refix.eta = T,
  opt_n_design.eta_warm = 30,
  opt_n_iter.eta_warm = 20,
  is.opt.hyper = F,
  hyper_n_grid = 20,
  ...
)

```

Arguments

<code>Y</code>	a $n \times (K \times p)$ matrix representing the observations.
<code>View.ind</code>	a $(K \times p)$ integer vector indicating the classes of features. The features with the same <code>View.ind</code> is in the same class.
<code>View.type</code>	a K vector encoding the structure type of each feature class. There are two choices: "O" (Omics Data), "C" (Compositional Data).
<code>eps.stop</code>	a numerical value controlling the convergence.
<code>max.step</code>	an integer controlling the maximum step for interaction.
<code>eps</code>	a numerical value controlling the convergence.
<code>T.step</code>	an integer controlling the maximum step for interaction.
<code>n_fold</code>	an integer representing the number of folds for cross validation.
<code>seed.sam.ind</code>	a vector of the seeds for sampling.
<code>hp.lower</code>	a numerical value or K vector specifying the lower bound of the hyper-parameter.
<code>hp.upper</code>	a numerical value or K vector specifying the upper bound of the hyper-parameter.
<code>hp.eta.lower</code>	a numerical value or K vector specifying the lower bound of the hyper-parameter for eta.
<code>hp.eta.upper</code>	a numerical value or K vector specifying the upper bound of the hyper-parameter for eta.
<code>eta.warm.stat.mat</code>	a matrix providing statistics for warm start of eta.
<code>opt_n_design</code>	an integer controlling the number of design points in the hyperparameter optimization.
<code>opt_n_iter</code>	an integer controlling the number of iterations in the hyperparameter optimization.
<code>Criterion</code>	a character indicating the criterion we choose for cross validation.
<code>des.init</code>	an initial design for hyperparameter optimization.
<code>is.refit</code>	a bool suggesting whether to refit the model using the optimal hyper-parameters.
<code>is.refix.eta</code>	a bool suggesting whether eta is fixed during refitting.
<code>opt_n_design.eta_warm</code>	an integer controlling the number of design points for eta warm-start optimization.
<code>opt_n_iter.eta_warm</code>	an integer controlling the number of iterations for eta warm-start optimization.
<code>is.opt.hyper</code>	a bool suggesting whether to optimize the hyper-parameters.
<code>hyper_n_grid</code>	an integer controlling the grid size for hyperparameter search.
<code>...</code>	additional arguments passed to the internal optimization procedures.

Value

A list containing the following elements: (1) `a.hat.opt.trgt`: The coefficient vector estimated with the optimal hyper-parameter vector; (2) `lam.opt.trgt`: The optimal hyper-parameter vector.

References

1. Deng, L., Tang, Y., Zhang, X., et al. (2024). Structure-adaptive canonical correlation analysis for microbiome multi-omics data. *Frontiers in Genetics*, 15, 1489694.
2. Chen, J., Bushman, F. D., Lewis, J. D., et al. (2013). Structure-constrained sparse canonical correlation analysis with an application to microbiome data analysis. *Biostatistics*, 14(2), 244–258.

Examples

```
## Not run:
library(dplyr)

n <- 200
p <- q <- 100
sigma.nu <- 5
sigma.eps <- 1
omega_X <- 0.85*c(rep(1/10,9),-9/10,rep(0,p-10))
omega_Y <- 0.85*c(seq(0.08,0.12,length = 10),rep(0,q-10))
Data1 <- DGP_OC(seed=10,n,p,q,sigma.nu,sigma.eps,omega_X,omega_Y)

library(mIrmBO)
Res.sCCA.CV <- cscca.CV(Y=Data1$Y,View.ind=Data1$View.ind,
                      View.type=c("0","0"),
                      show.info = TRUE)

Res.CsCCA.CV <- cscca.CV(Y=Data1$Y,View.ind=Data1$View.ind,
                       View.type=c("0","C"),
                       show.info = TRUE)

Res.sCCA <- cscca(Y=Data1$Y,View.ind=Data1$View.ind,
                 lambda.seq=Res.sCCA.CV$lam.opt.trgt,
                 View.type=c("0","0"))
Res.CsCCA <- cscca(Y=Data1$Y,View.ind=Data1$View.ind,
                  lambda.seq=Res.CsCCA.CV$lam.opt.trgt,
                  View.type=c("0","C"))
print(Res.sCCA.CV$Cri.opt.trgt)
print(Res.CsCCA.CV$Cri.opt.trgt)

## End(Not run)
```

DGP_OC

Data Generating Process (Omics Data versus Compositional data)

Description

Data Generating Process (Omics Data versus Compositional data)

Usage

DGP_OC(seed = 10, n, p, q, sigma.nu, sigma.eps, omega_X, omega_Y)

Arguments

seed	an integer for the initial seed.
n	an integer representing the sample size.
p	an integer representing the feature size of the omics dataset.
q	an integer representing the feature size of the compositional dataset.
sigma.nu	a numerical value representing the strength of correlation.
sigma.eps	a numerical value representing the strength of noise.
omega_X	a p vector representing the coefficient for the omics data.
omega_Y	a q vector representing the coefficient for the compositional data.

Value

A list containing the following elements: (a) Y: a $n \times (2p)$ matrix representing the full observations; (b) View.ind: a $2p$ integer vector indicating the classes of features. The features with the same View.ind is in the same class; (c) omega a $2p$ vector representing the true coefficients.

Examples

```
library(dplyr)
n <- 200
p <- q <- 100
sigma.nu <- 5
sigma.eps <- 1
omega_X <- 0.85*c(rep(1/10,9),-9/10,rep(0,p-10))
omega_Y <- 0.85*c(seq(0.08,0.12,length = 10),rep(0,q-10))
Data1 <- DGP_OC(seed=10,n,p,q,sigma.nu,sigma.eps,omega_X,omega_Y)
```

linda	<i>Linear (Lin) Model for Differential Abundance (DA) Analysis of High-dimensional Compositional Data</i>
-------	---

Description

The function implements a simple, robust and highly scalable approach to tackle the compositional effects in differential abundance analysis of high-dimensional compositional data. It fits linear regression models on the centered log₂-ratio transformed data, identifies a bias term due to the transformation and compositional effect, and corrects the bias using the mode of the regression coefficients. It could fit mixed-effect models for analysis of correlated data.

Usage

```
linda(
  feature.dat,
  meta.dat,
  formula,
```

```

feature.dat.type = c('count', 'proportion'),
prev.filter = 0,
mean.abund.filter = 0,
max.abund.filter = 0,
is.winsor = TRUE,
outlier.pct = 0.03,
adaptive = TRUE,
zero.handling = c('pseudo-count', 'imputation'),
pseudo.cnt = 0.5,
corr.cut = 0.1,
p.adj.method = "BH",
alpha = 0.05,
n.cores = 1,
verbose = TRUE
)

```

Arguments

<code>feature.dat</code>	a matrix of counts/proportions, row - features (OTUs, genes, etc) , column - samples.
<code>meta.dat</code>	a data frame containing the sample meta data. If there are NAs, the corresponding samples will be removed in the analysis.
<code>formula</code>	a character string for the formula. The formula should conform to that used by <code>lm</code> (independent data) or <code>lmer</code> (correlated data). For example: <code>formula = '~x1*x2+x3+(1 id)'</code> . At least one fixed effect is required.
<code>feature.dat.type</code>	the type of the feature data. It could be "count" or "proportion".
<code>prev.filter</code>	the prevalence (percentage of non-zeros) cutoff, under which the features will be filtered. The default is 0.
<code>mean.abund.filter</code>	the mean relative abundance cutoff, under which the features will be filtered. The default is 0.
<code>max.abund.filter</code>	the max relative abundance cutoff, under which the features will be filtered. The default is 0.
<code>is.winsor</code>	a logical value indicating whether winsorization should be performed to replace outliers (high values). The default is TRUE.
<code>outlier.pct</code>	the expected percentage of outliers. These outliers will be winsorized. The default is 0.03.
<code>adaptive</code>	a logical value indicating whether the approach to handle zeros (pseudo-count or imputation) will be determined based on the correlations between the log(sequencing depth) and the explanatory variables in <code>formula</code> when <code>feature.dat</code> is 'count'. If TRUE and the correlation p-value for any explanatory variable is smaller than or equal to <code>corr.cut</code> , the imputation approach will be used; otherwise, the pseudo-count approach will be used.

<code>zero.handling</code>	a character string of 'pseudo-count' or 'imputation' indicating the zero handling method used when <code>feature.dat</code> is 'count'. If 'pseudo-count', <code>apseudo.cnt</code> will be added to each value in <code>feature.dat</code> . If 'imputation', then we use the imputation approach using the formula in the referenced paper. Basically, zeros are imputed with values proportional to the sequencing depth. When <code>feature.dat</code> is 'proportion', this parameter will be ignored and zeros will be imputed by half of the minimum for each feature.
<code>pseudo.cnt</code>	a positive numeric value for the pseudo-count to be added if <code>zero.handling</code> is 'pseudo-count'. Default is 0.5.
<code>corr.cut</code>	a numerical value between 0 and 1, indicating the significance level used for determining the zero-handling approach when <code>adaptive</code> is TRUE. Default is 0.1.
<code>p.adj.method</code>	a character string indicating the p-value adjustment approach for addressing multiple testing. See R function <code>p.adjust</code> . Default is 'BH'.
<code>alpha</code>	a numerical value between 0 and 1 indicating the significance level for declaring differential features. Default is 0.05.
<code>n.cores</code>	a positive integer. If <code>n.cores</code> > 1 and <code>formula</code> is in a form of mixed-effect model, <code>n.cores</code> parallels will be conducted. Default is 1.
<code>verbose</code>	a logical value indicating whether the trace information should be printed out.

Value

A list with the elements

<code>variables</code>	A vector of variable names of all fixed effects in <code>formula</code> . For example: <code>formula = '~x1*x2+x3+(1 id)'</code> . Suppose <code>x1</code> and <code>x2</code> are numerical, and <code>x3</code> is a categorical variable of three levels: a, b and c. Then the elements of <code>variables</code> would be ('x1', 'x2', 'x3b', 'x3c', 'x1:x2').
<code>bias</code>	numeric vector; each element corresponds to one variable in <code>variables</code> ; the estimated bias of the regression coefficients due to the compositional effect.
<code>output</code>	a list of data frames with columns 'baseMean', 'log2FoldChange', 'lfcSE', 'stat', 'pvalue', 'padj', 'reject', 'df'; <code>names(output)</code> is equal to <code>variables</code> ; the rows of the data frame corresponds to <code>taxa</code> . Note: if there are <code>taxa</code> being excluded due to <code>prev.cut</code> , the number of the rows of the output data frame will be not equal to the number of the rows of <code>otu.tab</code> . <code>Taxa</code> are identified by the <code>rownames</code> . If the <code>rownames</code> of <code>otu.tab</code> are NULL, then <code>1 : nrow(otu.tab)</code> is set as the <code>rownames</code> of <code>otu.tab</code> . baseMean: 2 to the power of the intercept coefficients (normalized by one million) log2FoldChange: bias-corrected coefficients lfcSE: standard errors of the coefficients stat: <code>log2FoldChange / lfcSE</code> pvalue: <code>2 * pt(-abs(stat), df)</code> padj: <code>p.adjust(pvalue, method = p.adj.method)</code> reject: <code>padj <= alpha</code>

	df: degrees of freedom. The number of samples minus the number of explanatory variables (intercept included) for fixed-effect models; estimates from R package <code>lmerTest</code> with Satterthwaite method of approximation for mixed-effect models.
covariance	a list of data frames; the data frame records the covariances between a regression coefficient with other coefficients; <code>names(covariance)</code> is equal to <code>variables</code> ; the rows of the data frame corresponds to taxa. If the length of <code>variables</code> is equal to 1, then the covariance is <code>NULL</code> .
otu.tab.use	the OTU table used in the abundance analysis (the <code>otu.tab</code> after the preprocessing: samples that have NAs in the variables in <code>formula</code> or have less than <code>lib.cut</code> read counts are removed; taxa with prevalence less than <code>prev.cut</code> are removed and data is winsorized if <code>!is.null(winsor.quan)</code> ; and zeros are treated, i.e., imputed or pseudo-count added).
meta.use	the meta data used in the abundance analysis (only variables in <code>formula</code> are stored; samples that have NAs or have less than <code>lib.cut</code> read counts are removed; numerical variables are scaled).
wald	<p>a list for use in Wald test. If the fitting model is a linear model, then it includes</p> <p>beta: a matrix of the biased regression coefficients including intercept and all fixed effects; the columns correspond to taxa</p> <p>sig: the standard errors; the elements corresponding to taxa</p> <p>X: the design matrix of the fitting model</p> <p>bias: the estimated biases of the regression coefficients including intercept and all fixed effects</p> <p>If the fitting model is a linear mixed-effect model, then it includes</p> <p>beta: a matrix of the biased regression coefficients including intercept and all fixed effects; the columns correspond to taxa</p> <p>beta.cov: a list of covariance matrices of <code>beta</code>; the elements corresponding to taxa</p> <p>rand.cov: a list with covariance matrices of variance parameters of random effects; the elements corresponding to taxa; see more details in the paper of 'lmerTest'</p> <p>Joc.beta.cov.rand: a list of a list of Jacobian matrices of <code>beta.cov</code> with respect to the variance parameters; the elements corresponding to taxa</p> <p>bias: the estimated biases of the regression coefficients including intercept and all fixed effects</p>

Author(s)

Huijuan Zhou, Jun Chen, Xianyang Zhang

References

Zhou, H., He, K., Chen, J., Zhang, X. (2022). LinDA: linear models for differential abundance analysis of microbiome compositional data. *Genome biology*, 23(1), 95.

Examples

```

data(smokers)

ind <- smokers$meta$AIRWAYSITE == 'Throat'
otu.tab <- as.data.frame(smokers$otu[, ind])
depth <- colSums(otu.tab)
meta <- cbind.data.frame(Smoke = factor(smokers$meta$SMOKER[ind]),
                        Sex = factor(smokers$meta$SEX[ind]),
                        Site = factor(smokers$meta$SIDE OF BODY[ind]),
                        SubjectID = factor(smokers$meta$HOST_SUBJECT_ID[ind]))

# Differential abundance analysis using the left throat data
ind1 <- meta$Site == 'Left' & depth >= 1000
linda.obj <- linda(otu.tab[, ind1], meta[ind1, ], formula = '~Smoke+Sex',
                 feature.dat.type = 'count',
                 prev.filter = 0.1, is.winsor = TRUE, outlier.pct = 0.03,
                 p.adj.method = "BH", alpha = 0.1)

linda.plot(linda.obj, c('Smokey', 'Sexmale'),
          titles = c('Smoke: n v.s. y', 'Sex: female v.s. male'),
          alpha = 0.1, lfc.cut = 1, legend = TRUE, directory = NULL,
          width = 11, height = 8)

rownames(linda.obj $output[[1]])[which(linda.obj $output[[1]]$reject)]

# Differential abundance analysis pooling both the left and right throat data
# Mixed effects model is used

ind <- depth >= 1000
linda.obj <- linda(otu.tab[, ind], meta[ind, ], formula = '~Smoke+Sex+(1|SubjectID)',
                 feature.dat.type = 'count',
                 prev.filter = 0.1, is.winsor = TRUE, outlier.pct = 0.03,
                 p.adj.method = "BH", alpha = 0.1)

# For proportion data
otu.tab.p <- t(t(otu.tab) / colSums(otu.tab))
ind1 <- meta$Site == 'Left' & depth >= 1000
linda.obj <- linda(otu.tab.p[, ind1], meta[ind1, ], formula = '~Smoke+Sex',
                 feature.dat.type = 'proportion',
                 prev.filter = 0.1, is.winsor = TRUE, outlier.pct = 0.03,
                 p.adj.method = "BH", alpha = 0.1)

```

Description

The function produces the effect size plot of the differential features and volcano plot based on the output from `linda`.

Usage

```
linda.plot(
  linda.obj,
  variables.plot,
  titles = NULL,
  alpha = 0.05,
  lfc.cut = 1,
  legend = FALSE,
  directory = NULL,
  width = 11,
  height = 8
)
```

Arguments

<code>linda.obj</code>	return from function <code>linda</code> .
<code>variables.plot</code>	vector; variables whose results are to be plotted. For example, suppose the return value <code>variables</code> is equal to <code>('x1', 'x2', 'x3b', 'x3c', 'x1:x2')</code> , then one could set <code>variables.plot = c('x3b', 'x1:x2')</code> .
<code>titles</code>	vector; titles of the effect size plot and volcano plot for each variable in <code>variables.plot</code> . Default is <code>NULL</code> . If <code>NULL</code> , the titles will be set as <code>variables.plot</code> .
<code>alpha</code>	a numerical value between 0 and 1; cutoff for <code>padj</code> .
<code>lfc.cut</code>	a positive numerical value; cutoff for <code>log2FoldChange</code> .
<code>legend</code>	<code>TRUE</code> or <code>FALSE</code> ; whether to show the legends of the effect size plot and volcano plot.
<code>directory</code>	character; the directory to save the figures, e.g., <code>getwd()</code> . Default is <code>NULL</code> . If <code>NULL</code> , figures will not be saved.
<code>width</code>	the width of the graphics region in inches. See R function <code>pdf</code> .
<code>height</code>	the height of the graphics region in inches. See R function <code>pdf</code> .

Value

A list of `ggplot2` objects.

<code>plot.lfc</code>	a list of effect size plots. Each plot corresponds to one variable in <code>variables.plot</code> .
<code>plot.volcano</code>	a list of volcano plots. Each plot corresponds to one variable in <code>variables.plot</code> .

Author(s)

Huijuan Zhou, Jun Chen, Xianyang Zhang

References

Zhou, H., He, K., Chen, J., Zhang, X. (2022). LinDA: linear models for differential abundance analysis of microbiome compositional data. *Genome biology*, 23(1), 95.

Examples

```
data(smokers)
ind <- smokers$meta$AIRWAYSITE == 'Throat' & smokers$meta$SIDEOFBODY == 'Left'
otu.tab <- as.data.frame(smokers$otu[, ind])
depth <- colSums(otu.tab)
meta <- cbind.data.frame(Smoke = factor(smokers$meta$SMOKER[ind]),
                        Sex = factor(smokers$meta$SEX[ind]))

ind <- depth >= 1000
linda.obj <- linda(otu.tab[, ind], meta[ind, ], formula = '~Smoke+Sex',
                 feature.dat.type = 'count',
                 prev.filter = 0.1, is.winsor = TRUE, outlier.pct = 0.03,
                 p.adj.method = "BH", alpha = 0.1)

linda.plot(linda.obj, c('Smokey', 'Sexmale'),
           titles = c('Smoke: n v.s. y', 'Sex: female v.s. male'),
           alpha = 0.1, lfc.cut = 1, legend = TRUE, directory = NULL,
           width = 11, height = 8)
```

linda.wald.test	<i>Wald test for bias-corrected regression coefficients</i>
-----------------	---

Description

The function implements Wald test for bias-corrected regression coefficients generated from the `linda` function. One can utilize the function to perform ANOVA-type analyses. For example, if you have a categorical variable with three levels, you can test whether all levels have the same effect.

Usage

```
linda.wald.test(
  linda.obj,
  L,
  model = c("LM", "LMM"),
  alpha = 0.05,
  p.adj.method = "BH"
)
```

Arguments

`linda.obj` return from the `linda` function.

L	A matrix for testing $Lb = 0$, where b includes the intercept and all fixed effects from running linda. Thus the number of columns of L must be equal to $\text{length}(\text{variables})+1$, where variables is from linda.obj, which does not include the intercept.
model	'LM' or 'LMM' indicating the model fitted in {linda} is linear model or linear mixed-effect model.
alpha	significance level for testing $Lb = 0$.
p.adj.method	P-value adjustment approach. See R function p.adjust. The default is 'BH'.

Value

A data frame with columns

Fstat	Wald statistics for each taxon
df1	The numerator degrees of freedom
df2	The denominator degrees of freedom
pvalue	$1 - \text{pf}(\text{Fstat}, \text{df1}, \text{df2})$
padj	$\text{p.adjust}(\text{pvalue}, \text{method} = \text{p.adj.method})$
reject	$\text{padj} \leq \alpha$

Author(s)

Huijuan Zhou <huijuanzhou2019@gmail.com> Jun Chen <Chen.Jun2@mayo.edu> Xianyang Zhang <zhangxiany@stat.tamu.edu>

References

Zhou, H., He, K., Chen, J., Zhang, X. (2022). LinDA: linear models for differential abundance analysis of microbiome compositional data. *Genome biology*, 23(1), 95.

Examples

```
data(smokers)

ind <- smokers$meta$AIRWAYSITE == 'Throat'
otu.tab <- as.data.frame(smokers$otu[, ind])
depth <- colSums(otu.tab)
meta <- cbind.data.frame(Smoke = factor(smokers$meta$SMOKER[ind]),
                        Sex = factor(smokers$meta$SEX[ind]),
                        Site = factor(smokers$meta$SIDE OF BODY[ind]),
                        SubjectID = factor(smokers$meta$HOST_SUBJECT_ID[ind]))

ind <- depth >= 1000
linda.obj <- linda(otu.tab[, ind], meta[ind, ], formula = '~Smoke+Sex+(1|SubjectID)',
                  feature.dat.type = 'count',
                  prev.filter = 0.1, is.winsor = TRUE, outlier.pct = 0.03,
                  p.adj.method = "BH", alpha = 0.1)
# L matrix (2x3) is designed to test the second (Smoke) and the third (Sex) coefficient to be 0.
# For a categorical variable > two levels, similar L can be designed to do ANOVA-type test.
```

```
L <- matrix(c(0, 1, 0, 0, 0, 1), nrow = 2, byrow = TRUE)
result <- linda.wald.test(linda.obj, L, 'LMM', alpha = 0.1)
```

smokers	<i>Microbiome data from the human upper respiratory tract (left and right throat)</i>
---------	---

Description

A dataset containing "otu", "tax", "meta", "genus", "family"

Usage

```
data(smokers)
```

Format

A list with elements

otu otu table, 2156 taxa by 290 samples

tax taxonomy table, 2156 taxa by 7 taxonomic ranks

meta meta table, 290 samples by 53 sample variables

genus 304 by 290

family 113 by 290

Source

<https://qiita.ucsd.edu/> study ID:524 Reference: Charlson ES, Chen J, Custers-Allen R, Bittinger K, Li H, et al. (2010) Disordered Microbial Communities in the Upper Respiratory Tract of Cigarette Smokers. PLoS ONE 5(12): e15216.

Index

* datasets

smokers, [19](#)

bmdd, [2](#), [6](#)

bmdd.nlopt, [5](#)

cscca, [6](#)

cscca.CV, [8](#)

DGP_OC, [10](#)

linda, [11](#)

linda.plot, [15](#)

linda.wald.test, [17](#)

smokers, [19](#)