

# Package ‘zinck’

November 28, 2023

**Title** Zero-Inflated Compositional Knockoff Filter

**Version** 0.0.0.9000

**Description** Zinck exploits a zero-inflated variant of the Latent Dirichlet Allocation (LDA) model to generate valid knockoffs that capture the key characteristics of microbiome data - mainly its compositional nature and high sparsity. It exhibits the properties of simultaneous variable selection and FDR control to identify microbial biomarkers. This package provides an implementation of zinck, which is trained either using the Automatic Differentiation Variational Inference (ADVI) algorithm or using a collapsed Gibbs sampler, facilitating variable selection for both continuous as well as binary outcomes.

**License** MIT

**Encoding** UTF-8

**Roxygen** list(markdown = TRUE)

**RoxygenNote** 7.2.3

**Imports** dplyr,  
reshape2,  
knockoff,  
glmnet,  
randomForest,  
caret,  
rstan (== 2.21.8),  
stats,  
fitdistrplus,  
ggplot2,  
MLmetrics,  
phyloseq,  
GUniFrac,  
kosel,  
gridExtra,  
zinLDA

**Suggests** knitr,  
rmarkdown

**VignetteBuilder** knitr

**LazyData** true

## R topics documented:

combo . . . . . 2

|                            |   |
|----------------------------|---|
| count . . . . .            | 2 |
| count.genus . . . . .      | 3 |
| draw_heatmap . . . . .     | 3 |
| fit.zinck . . . . .        | 4 |
| generateKnockoff . . . . . | 5 |
| log_normalize . . . . .    | 6 |
| optimal_k . . . . .        | 6 |
| zinck.filter . . . . .     | 7 |

|       |   |
|-------|---|
| Index | 9 |
|-------|---|

---

|       |                 |
|-------|-----------------|
| combo | <i>BMI data</i> |
|-------|-----------------|

---

**Description**

The COMBO study is a cross-sectional study designed to examine the association between micronutrients and gut microbiome composition, where the fecal samples of 98 healthy individuals from the University of Pennsylvania were collected along with demographic data like Body Mass Index (BMI), age, and sex. The data are available at the Qiita repository under the accession ID: 101135.

**Usage**

data(combo)

**Format**

A list containing 17 elements

**References**

Wu GD, Chen J, Hoffmann C, Bittinger K, Chen YY, Keilbaugh SA, Bewtra M, Knights D, Walters WA, Knight R, Sinha R, Gilroy E, Gupta K, Baldassano R, Nessel L, Li H, Bushman FD, Lewis JD. Linking long-term dietary patterns with gut microbial enterotypes. *Science*. 2011 Oct 7;334(6052):105-8. doi: 10.1126/science.1208344. Epub 2011 Sep 1. PMID: 21885731; PMCID: PMC3368382.

---

|       |                                 |
|-------|---------------------------------|
| count | <i>CRC data (species level)</i> |
|-------|---------------------------------|

---

**Description**

This dataset contains data from a meta-analysis of five geographically and technically diverse fecal shotgun metagenomic studies of colorectal cancer (CRC, n = 574). All raw sequencing data across studies were reprocessed using the same bioinformatics pipeline for taxonomic profiling. The 5 studies involve subjects from 5 different countries. The sample size of the studies is 109, 127, 120, 114, and 104; and the number of cases and controls are roughly balanced in each study. This data-set contains p = 849 species found among all the investigated metagenomic CRC studies.

**Usage**

```
data(count)
```

**Format**

An integer matrix with 574 rows and 849 columns.

**Source**

Zeller Lab

---

|             |                               |
|-------------|-------------------------------|
| count.genus | <i>CRC data (genus level)</i> |
|-------------|-------------------------------|

---

**Description**

This dataset contains data from a meta-analysis of five geographically and technically diverse fecal shotgun metagenomic studies of colorectal cancer (CRC, n = 574). All raw sequencing data across studies were reprocessed using the same bioinformatics pipeline for taxonomic profiling. The 5 studies involve subjects from 5 different countries. The sample size of the studies is 109, 127, 120, 114, and 104; and the number of cases and controls are roughly balanced in each study. This data-set contains p = 133 distinct genera found among all the investigated metagenomic CRC studies.

**Usage**

```
data(count.genus)
```

**Format**

An integer matrix with 574 rows and 133 columns.

**Source**

Zeller Lab

---

|              |  |
|--------------|--|
| draw_heatmap | <i>Draw Heatmap for Microbial Sample Taxa Matrix</i> |
|--------------|--|

---

**Description**

This function creates a heatmap for a given microbial sample taxa matrix. It is specifically designed for visualizing abundance patterns across different samples and taxa in microbiome studies. The function applies an arcsinh transformation to the data for normalization and better visualization of abundance patterns, especially useful in handling highly skewed microbiome data.

**Usage**

```
draw_heatmap(X, title = "")
```

**Arguments**

|                    |  |
|--------------------|--|
| <code>X</code>     | A numeric matrix representing microbial sample taxa data, where rows represent samples and columns represent taxa. |
| <code>title</code> | An optional title for the heatmap  |

**Details**

The heatmap is generated using `ggplot2` and `reshape2` packages, with taxa on the x-axis and samples on the y-axis. The color intensity in the heatmap represents the arcsinh-transformed abundance of each taxa in each sample.

**Value**

A `ggplot` object representing the heatmap. This can be further customized or directly plotted.

**Examples**

```
# Create a sample matrix representing microbial sample taxa data
mat <- matrix(runif(20), nrow = 5)

# Draw heatmap
draw_heatmap(mat)
```

---

fit.zinck

*Fit the zinck model to the data.*


---

**Description**

This function fits the zinck hierarchical model to the observed microbiome sample taxa matrix. The posterior is approximated using either ADVI or collapsed Gibbs sampling.

**Usage**

```
fit.zinck(
  X,
  num_clusters,
  method = c("ADVI", "Gibbs"),
  tuned = FALSE,
  seed = NULL
)
```

**Arguments**

|                           |  |
|---------------------------|--|
| <code>X</code>            | An OTU matrix with dimensions $D \times p$ .   |
| <code>num_clusters</code> | An integer specifying the number of clusters.  |
| <code>method</code>       | A character string, either "ADVI" or "Gibbs", specifying the method to fit the model.                              |
| <code>tuned</code>        | A logical value. If TRUE and method is "ADVI", the model will be tuned; otherwise, it will use default parameters. |
| <code>seed</code>         | An integer used to set the seed for reproducibility.   |

**Value**

A list containing the posterior estimates of beta and theta.

**References**

Kucukelbir, A., Tran, D., Ranganath, R., Gelman, A., and Blei, D. (2017). Automatic Differentiation Variational Inference. *Journal of Machine Learning Research* 18(14).

Deek, R., and Li, H. (2021). A Zero-Inflated Latent Dirichlet Allocation Model for Microbiome Studies. *Frontiers in Genetics* 11.

---

generateKnockoff

*Generate Knockoff Copy of Microbial Sample Data*


---

**Description**

This function generates a knockoff copy of microbial sample data given a matrix  $X$  and matrices  $\Theta$  and  $B$ . The function adjusts the column structure of  $B$  to match  $X$ , generates samples based on  $\Theta$  and  $B$ , and then compiles these into a knockoff count matrix.

**Usage**

```
generateKnockoff(X, Theta, Beta, seed = NULL)
```

**Arguments**

|          |   |
|----------|---|
| $X$      | A numeric matrix representing original microbial sample data.       |
| $\Theta$ | A numeric matrix representing cluster mixing probabilities.         |
| $B$      | A numeric matrix representing feature proportions for each cluster. |
| seed     | An optional integer seed for reproducibility of random generation.  |

**Value**

A numeric matrix representing the knockoff copy of the microbial sample data.

**Examples**

```
X <- matrix(runif(40), nrow = 10)
Theta <- matrix(runif(30), nrow = 10)
Beta <- matrix(runif(20), nrow = 5)
knockoff_data <- generateKnockoff(X, Theta, Beta)
```

---

log\_normalize

*Log-Normalization for Microbiome Compositional Data*


---

### Description

This function performs log-normalization on a given matrix, typically used in microbiome data analysis. In microbiome studies, data are often compositional and contain many zero counts. Log-normalization, with the addition of a pseudo-count, is a standard approach to handle zeros and maintain the compositional nature of the data. This function takes a numeric matrix, adds a pseudo-count to zero values, and then applies a log transformation, preserving the relative proportions in the data.

### Usage

```
log_normalize(X)
```

### Arguments

X                      A numeric matrix to log-normalize.

### Details

The function first adds a pseudo-count of 0.5 to zero entries to handle zeros in the data. It then divides each count by the total counts in its row (sample) to make the data compositional. Finally, it applies a natural logarithm transformation to the normalized data. #'

### Value

A matrix of the same dimensions as X, with each element being the log-normalized value of the corresponding element in X. The base of the logarithm used for normalization is e (natural logarithm).

### Examples

```
# Create a sample matrix
mat <- matrix(1:4, nrow = 2)

# Perform log-normalization
log_normalized_mat <- log_normalize(mat)
```

---

optimal\_k

*Optimal Number of Clusters based on Jensen-Shannon Divergence*


---

### Description

This function identifies the optimal number of clusters for fitting a zinck model using Jensen-Shannon Divergence. The function fits the model for various values of clusters and calculates the average Jensen-Shannon Divergence for each, returning the number of clusters that minimizes this value.

**Usage**

```
optimal_k(X, kmin, kmax, seed_list = NULL)
```

**Arguments**

|           |  |
|-----------|--|
| X         | A OTU matrix with dimensions $D \times p$ .  |
| kmin      | Numeric; the minimum number of clusters to be considered.                            |
| kmax      | Numeric; the maximum number of clusters to be considered.                            |
| seed_list | List of numeric values; seeds for reproducibility for each k value, default is NULL. |

**Value**

The optimal number of clusters, K, that minimizes the average Jensen-Shannon Divergence.

**Examples**

```
## Not run:
data("combo")
X <- combo[["abund_list"]][["Genus"]]
result <- optimal_k(X, kmin = 2, kmax = 10, seed_list=list(1,1,2,4,123,6,123,123,12))
print(result)

## End(Not run)
```

---

zinck.filter

*Zinck Filter for Variable Selection*


---

**Description**

Performs variable selection by fitting the augmented set of features to the response, using a glmnet or Random Forest model

**Usage**

```
zinck.filter(
  X,
  X_tilde,
  Y,
  model,
  fdr = 0.1,
  offset = 1,
  seed = NULL,
  ntrees = 1000,
  tune_mtry = FALSE
)
```

**Arguments**

|                        |   |
|------------------------|---|
| <code>X</code>         | An OTU matrix with dimensions $D \times p$ .                                |
| <code>X_tilde</code>   | A knockoff matrix corresponding to <code>X</code> with the same dimensions. |
| <code>Y</code>         | The response variable, either continuous or binary.                         |
| <code>model</code>     | A string; the model to be used ("glmnet" or "Random Forest").               |
| <code>fdr</code>       | Numeric; the false discovery rate. Default is 0.1                           |
| <code>offset</code>    | Numeric; either 0 or 1. Default is 1  |
| <code>seed</code>      | Numeric; the seed for reproducibility.                                      |
| <code>ntrees</code>    | Numeric; the number of trees for Random Forest. Default is 1000.            |
| <code>tune_mtry</code> | Logical; whether to tune mtry in Random Forest. Default is FALSE.           |

**Value**

A vector of selected variables at a target false discovery rate



# Index

- \* **datasets**
  - combo, [2](#)
  - count, [2](#)
  - count.genus, [3](#)
- combo, [2](#)
- count, [2](#)
- count.genus, [3](#)
- draw\_heatmap, [3](#)
- fit.zinck, [4](#)
- generateKnockoff, [5](#)
- log\_normalize, [6](#)
- optimal\_k, [6](#)
- zinck.filter, [7](#)