

# Package ‘MAI’

December 20, 2024

**Type** Package

**Title** Mechanism-Aware Imputation

**Version** 1.12.0

**Description** A two-step approach to imputing missing data in metabolomics.

Step 1 uses a random forest classifier to classify missing values as either Missing Completely at Random/Missing At Random (MCAR/MAR) or Missing Not At Random (MNAR). MCAR/MAR are combined because it is often difficult to distinguish these two missing types in metabolomics data. Step 2 imputes the missing values based on the classified missing mechanisms, using the appropriate imputation algorithms. Imputation algorithms tested and available for MCAR/MAR include Bayesian Principal Component Analysis (BPCA), Multiple Imputation No-Skip K-Nearest Neighbors (Multi\_nsKNN), and Random Forest. Imputation algorithms tested and available for MNAR include nsKNN and a single imputation approach for imputation of metabolites where left-censoring is present.

**License** GPL-3

**Encoding** UTF-8

**Imports** caret, parallel, doParallel, foreach, e1071, future.apply, future, missForest, pcaMethods, tidyverse, stats, utils, methods, SummarizedExperiment, S4Vectors

**biocViews** Software, Metabolomics, StatisticalMethod, Classification

**Suggests** knitr, rmarkdown, BiocStyle, testthat (>= 3.0.0)

**VignetteBuilder** knitr

**Config/testthat/edition** 3

**URL** <https://github.com/KechrisLab/MAI>

**BugReports** <https://github.com/KechrisLab/MAI/issues>

**git\_url** <https://git.bioconductor.org/packages/MAI>

**git\_branch** RELEASE\_3\_20

**git\_last\_commit** da4c5d7

**git\_last\_commit\_date** 2024-10-29

**Repository** Bioconductor 3.20

**Date/Publication** 2024-12-19

**Author** Jonathan Dekermanjian [aut, cre],  
 Elin Shaddox [aut],  
 Debmalya Nandy [aut],  
 Debashis Ghosh [aut],  
 Katerina Kechris [aut]

**Maintainer** Jonathan Dekermanjian <Jonathan.Dekermanjian@CUAnschutz.edu>

## Contents

MAI . . . . .	2
untargeted_LCMS_data . . . . .	4
<b>Index</b>	<b>5</b>

---

MAI	<i>Mechanism-Aware Imputation</i>
-----	-----------------------------------

---

## Description

A two-step approach to imputing missing data in metabolomics. Step 1 uses a random forest classifier to classify missing values as either Missing Completely at Random/Missing At Random (MCAR/MAR) or Missing Not At Random (MNAR). MCAR/MAR are combined because it is often difficult to distinguish these two missing types in metabolomics data. Step 2 imputes the missing values based on the classified missing mechanisms, using the appropriate imputation algorithms. Imputation algorithms tested and available for MCAR/MAR include Bayesian Principal Component Analysis (BPCA), Multiple Imputation No-Skip K-Nearest Neighbors (Multi\_nsKNN), and Random Forest. Imputation algorithms tested and available for MNAR include nsKNN and a single imputation approach for imputation of metabolites where left-censoring is present.

## Usage

```
MAI(data_miss,
     MCAR_algorithm = c("BPCA", "Multi_nsKNN", "random_forest"),
     MNAR_algorithm = c("nsKNN", "Single"),
     n_cores = 1,
     assay_ix = 1,
     forest_list_args = list(
       ntree = 300,
       proximity = FALSE
     ),
     verbose = TRUE
  )
```

## Arguments

<code>data_miss</code>	A matrix or dataframe, or a SummarizedExperiment containing missing values designated by "NA" to impute
<code>MCAR_algorithm</code>	The imputation algorithm you wish to use to impute MCAR predicted missing values. possible algorithms c("BPCA", "Multi_nsKNN", "random_forest")
<code>MNAR_algorithm</code>	The imputation algorithm you wish to use to impute MNAR predicted missing values. possible algorithms c("Single", "nsKNN")

n_cores	The number of cores you want to utilize. Default is 1 core. To use all cores specify n_cores = -1.
assay_ix	If data is a Summarized Experiment then this argument defines the index of the assay to impute. Default is set to the first assay.
forest_list_args	Random forest named arguments to pass to the random forest training process. Default args are ntree = 300 and proximity = FALSE
verbose	A toggle to suppress console output. Default is TRUE

## Value

When matrix or dataframe returns a list containing the following:

Imputed Data Returns dataframes of MAI imputation

Estimated Parameters

Returns the estimated  $\alpha$ ,  $\beta$ , and  $\gamma$  parameters that define the missingness pattern in the data set. These parameters estimate the ratio of MCAR/MAR to MNAR in the data. The parameters  $\alpha$  and  $\beta$  separate high, medium, and low average abundance metabolites, while the parameter  $\gamma$  is used to impose missingness in the medium and low abundance metabolites. A smaller  $\alpha$  corresponds to more MCAR/MAR being present, while larger  $\beta$  and  $\gamma$  values imply more MNAR values being present. The returned estimated parameters are then used to impose known missingness in the complete subset of the input data. Subsequently, a random forest classifier is trained to classify the known missingness in the complete subset of the input data. Once the classifier is established it is applied to the unknown missingness of the full input data to predict the missingness. Finally, the missing values are imputed using a specific algorithm, chosen by the user, according to the predicted missingness mechanism.

When a Summarized Experiment returns:

Imputed Assay Returns the imputed data in the specified assay based on the assay\_ix assigned

Estimated Parameters

Returns estimated parameters in the metadata of the Summarized Experiment as a list

## Examples

```
data(untargeted_LCMS_data)
MAI(data_miss=untargeted_LCMS_data,
    MCAR_algorithm = "BPCA",
    MNAR_algorithm="Single",
    n_cores = 1,
    assay_ix = 1,
    forest_list_args = list(
      ntree = 300,
      proximity = FALSE
    ),
    verbose = TRUE)
```

---

untargeted\_LCMS\_data *Example data set containing missing values*

---

**Description**

This data set is randomly generated. We impose 30 percent missing values using the Mixed missingness algorithm developed by Styczynski et al. Where the parameters alpha, beta, and gamma were chosen to be 30, 70, and 40 percent, respectively.

**References**

Lee JY, Styczynski MP. NS-kNN: a modified k-nearest neighbors approach for imputing metabolomics data. *Metabolomics*. 2018;14(12):153.

# Index

MAI, [2](#)

untargeted\_LCMS\_data, [4](#)