

# Package ‘miceRanger’

January 21, 2020

**Title** Multiple Imputation by Chained Equations with Random Forests

**Version** 1.2.0

**Maintainer** Sam Wilson <samwilson303@gmail.com>

**Description** Multiple Imputation has been shown to be a flexible method to impute missing values by Van Buuren (2007) <doi:10.1177/0962280206074463>. Expanding on this, random forests have been shown to be an accurate model by Stekhoven and Buhlmann <arXiv:1105.0828> to impute missing values in datasets. They have the added benefits of returning out of bag error and variable importance estimates, as well as being simple to run in parallel.

**URL** <https://github.com/FarrellDay/miceRanger>

**BugReports** <https://github.com/FarrellDay/miceRanger/issues>

**Encoding** UTF-8

**LazyData** true

**License** MIT + file LICENSE

**Depends** R (>= 3.5.0)

**Imports** ranger, data.table, stats, FNN, ggplot2, crayon, corplot, ggpubr, DescTools, foreach

**Suggests** knitr, rmarkdown, doParallel, testthat (>= 2.1.0)

**RoxygenNote** 7.0.2

**VignetteBuilder** knitr

**NeedsCompilation** no

**Author** Sam Wilson [aut, cre]

**Repository** CRAN

**Date/Publication** 2020-01-20 23:40:02 UTC

## R topics documented:

addDatasets . . . . .	2
addIterations . . . . .	3
amputeData . . . . .	4
completeData . . . . .	4
getVarImps . . . . .	5
miceRanger . . . . .	6
plotCorrelations . . . . .	8
plotDistributions . . . . .	9
plotImputationVariance . . . . .	10
plotModelError . . . . .	11
plotVarConvergence . . . . .	11
plotVarImportance . . . . .	12
print.miceDefs . . . . .	13
sampleMiceDefs . . . . .	13
<b>Index</b>	<b>15</b>

---

addDatasets	<i>addDatasets</i>
-------------	--------------------

---

### Description

Add datasets to a current miceDefs object. Adds the same number of iterations as other datasets.

### Usage

```
addDatasets(miceObj, datasets = 3, parallel = FALSE, verbose = TRUE, ...)
```

### Arguments

miceObj	A miceDefs object created by miceRanger.
datasets	The number of datasets to add.
parallel	Should the process run in parallel? This process will take advantage of any cluster set up when miceRanger is called.
verbose	should progress be printed?
...	other parameters passed to ranger() to control model building.

### Value

an updated miceDefs object with additional datasets.

**Examples**

```
data("sampleMiceDefs")
miceObj <- addIterations(
  sampleMiceDefs
  , iters = 2
  , verbose = FALSE
  , num.threads = 1
  , num.trees=5
)
```

---

addIterations	<i>addIterations</i>
---------------	----------------------

---

**Description**

Add iterations to a current miceDefs object. Adds iterations for all datasets.

**Usage**

```
addIterations(miceObj, iters = 5, parallel = FALSE, verbose = TRUE, ...)
```

**Arguments**

miceObj	A miceDefs object created by miceRanger.
iters	The number of iterations to add to each dataset.
parallel	Should the process run in parallel? This process will take advantage of any cluster set up when miceRanger is called.
verbose	should progress be printed?
...	other parameters passed to ranger() to control model building.

**Value**

an updated miceDefs object with additional iterations.

**Examples**

```
data("sampleMiceDefs")
miceObj <- addIterations(
  sampleMiceDefs
  , iters=2
  , verbose=FALSE
  , num.threads = 1
  , num.trees=5
)
```

---

 amputeData

*amputeData*


---

### Description

randomly imputes data.

### Usage

```
amputeData(data, perc = 0.1, cols = names(data))
```

### Arguments

data	The data to be amputed
perc	A scalar. The percentage (0-1) to be amputed.
cols	The columns to ampute.

### Value

The same dataset with random values set to NA in vars.

### Examples

```
data(iris)
head(iris,10)

ampIris <- amputeData(iris)
head(ampIris,10)
```

---

 completeData

*completeData*


---

### Description

Return the completed datasets.

### Usage

```
completeData(miceObj, datasets = 1:miceObj$callParams$m, verbose = TRUE)
```

### Arguments

miceObj	an object of class miceDefs.
datasets	a vector of the datasets you want to return.
verbose	a warning is thrown if integers are converted to doubles. To suppress this warning, set to FALSE.

**Value**

A list of imputed datasets.

**Examples**

```
data("sampleMiceDefs")  
imputedList <- completeData(sampleMiceDefs)
```

---

*getVarImps*

*getVarImps*

---

**Description**

Returns all of the imputations for the specified datasets and variable.

**Usage**

```
getVarImps(miceObj, datasets = 1:miceObj$callParams$m, var)
```

**Arguments**

- miceObj*            A *miceDefs* object created by *miceRanger*.
- datasets*           The datasets to return. Can be a number, or a numeric vector.
- var*                The variable to return the imputations for.

**Value**

A matrix of imputations for a single variable. Each column represents a different dataset.

**Examples**

```
data("sampleMiceDefs")  
getVarImps(sampleMiceDefs, "Sepal.Width")
```

---

miceRanger

*miceRanger*


---

### Description

Performs multiple imputation by chained random forests. Returns a miceDefs object, which contains information about the imputation process.

### Usage

```
miceRanger(
  data,
  m = 5,
  maxiter = 5,
  vars,
  valueSelector = c("meanMatch", "value"),
  meanMatchCandidates = pmax(round(nrow(data) * 0.01), 5),
  parallel = FALSE,
  verbose = TRUE,
  ...
)
```

### Arguments

data	The data to be imputed. Can contain variables that are not going to be imputed, which you want to use as features.
m	The number of datasets to produce
maxiter	The number of iterations to run for each dataset.
vars	Specifies which and how variables should be imputed. Can be specified in 3 different ways: <ul style="list-style-type: none"> <li>• <code>&lt;missing&gt;</code> If not provided, all columns will be imputed using all columns. If a column contains no missing values, it will still be used as a feature to impute missing columns.</li> <li>• <code>&lt;Character Vector&gt;</code> If a character vector of column names is passed, these columns will be imputed using all available columns in the dataset. The order of this vector will determine the order in which the variables are imputed.</li> <li>• <code>&lt;Named List of Character Vectors&gt;</code> Predictors can be specified for each variable with named list. List names are the variables to impute. Elements in the vectors should be features used to impute that variable. The order of this list will determine the order in which the variables are imputed.</li> </ul>
valueSelector	How to select the value to be imputed from the model predictions. Can be "meanMatching", "value", or a named vector containing a mixture of those values. If a named vector is passed, the names must equal the variables to be imputed specified in vars.

meanMatchCandidates	Used for regression. Specifies the number of candidate values. If a variable is being imputed using mean matching, this
parallel	Should the process run in parallel? Usually not necessary. This process will take advantage of any cluster set up when miceRanger is called.
verbose	should progress be printed?
...	other parameters passed to ranger() to control forest growth.

**Value**

a miceDefs object, containing the following:

callParams	The parameters of the object.
data	The original data provided by the user.
naWhere	Logical index of missing data, having the same dimensions as data.
missingCounts	The number of missing values for each variable
rawClasses	The original classes provided in data
newClasses	The new classes of the returned dataset. Classes can be changed if necessary.
allImps	The imputations of all variables at each iteration, for each dataset.
allImport	The variable importance metrics at each iteration, for each dataset.
allError	The OOB model error for all variables at each iteration, for each dataset.
finalImps	The final imputations for each dataset.
finalImport	The final variable importance metrics for each dataset.
finalError	The final model error for each variable in every dataset.
imputationTime	The total time in seconds taken to create the imputations for the specified datasets and iterations. Does not include any setup time.

**Examples**

```
# Using Mice to create 5 imputed datasets
data(iris)

ampIris <- amputeData(iris)

miceObj <- miceRanger(
  ampIris
  , m = 2
  , maxiter = 2
  , verbose=FALSE
  , num.threads = 1
  , num.trees=5
)

# Run in parallel
data(iris)
```

```

ampIris <- amputeData(iris)

library(doParallel)
cl <- makeCluster(2)
registerDoParallel(cl)

# Perform mice
parTime <- system.time(
  miceObjPar <- miceRanger(
    ampIris
    , m=2
    , maxiter = 2
    , parallel = TRUE
    , verbose = FALSE
  )
)
stopCluster(cl)
registerDoSEQ()

```

---

plotCorrelations

*plotCorrelations*


---

### Description

Plot the correlation of imputed values between every combination of datasets for each variable.

### Usage

```

plotCorrelations(
  miceObj,
  vars = names(miceObj$callParams$vars),
  factCorrMetric = "CramerV",
  numbCorrMetric = "pearson",
  ...
)

```

### Arguments

**miceObj** an object of class `miceDefs`, created by the `miceRanger` function.

**vars** the variables you want to plot. Default is to plot all variables. Can be a vector of variable names, or one of `'allNumeric'` or `'allCategorical'`

**factCorrMetric** The correlation metric for categorical variables. Can be one of:

- "CramerV" Cramer's V correlation metric.
- "Chisq" Chi Square test statistic.
- "TschuprowT" Tschuprow's T correlation metric.
- "Phi" (Binary Variables Only) Phi coefficient.



- "YuleY" (Binary Variables Only) Yule's Y, also known as coefficient of colligation
  - "YuleQ" (Binary Variables Only) Yule's Q, related to Yule's Y by  $Q=2Y/(1+Y^2)$
- numbCorrMetric The correlation metric for numeric variables. Can be one of:
- "pearson" Pearson's Correlation Coefficient
  - "spearman" Spearman's Rank Correlation Coefficient
  - "kendall" Kendall's Rank Correlation Coefficient
  - "Rsquared" R-squared
- ... Other arguments to pass to ggarrange()

**Value**

an object of class ggarrange.

**Examples**

```
data("sampleMiceDefs")
plotCorrelations(sampleMiceDefs)
```

---

plotDistributions      *plotDistributions*

---

**Description**

Plots the distribution of the original data beside the imputed data.

**Usage**

```
plotDistributions(
  miceObj,
  vars = names(miceObj$callParams$vars),
  title = NULL,
  dotsize = 0.5,
  ...
)
```

**Arguments**

- miceObj      an object of class miceDefs, created by the miceRanger function.
- vars          the variables you want to plot. Default is to plot all variables. Can be a vector of variable names, or one of 'allNumeric' or 'allCategorical'
- title        The title of the plot. Default is no title.
- dotsize      Passed to geom\_dotplot(). Depending on the number of graphs plotted, you may want to change the dot size for categorical variables.
- ...          additional parameters passed to ggarrange().

**Value**

an object of class ggarrange.

**Examples**

```
data("sampleMiceDefs")
plotDistributions(sampleMiceDefs)
```

---

```
plotImputationVariance
      plotImputationVariance
```

---

**Description**

plots the distribution of the difference between datasets of the imputed values. For categorical variables, the distribution of the number of distinct levels imputed for each sample is shown next to the expected hypergeometric distribution, if the imputation was completely random. For numeric variables, the density of the standard deviation (between datasets) of imputations is plotted. The shaded area represents the samples that had a standard deviation lower than the total nonmissing standard deviation for the original data.

**Usage**

```
plotImputationVariance(
  miceObj,
  vars = names(miceObj$callParams$vars),
  monteCarloSimulations = 10000,
  ...
)
```

**Arguments**

miceObj	an object of class miceDefs, created by the miceRanger function.
vars	the variables you want to plot. Default is to plot all variables. Can be a vector of variable names, or one of 'allNumeric' or 'allCategorical'
monteCarloSimulations	The number of simulations to run to determine the distribution of unique categorical levels drawn if the draws were completely random.
...	additional parameters passed to ggarrange.

**Value**

an object of class ggarrange.

**Examples**

```
data("sampleMiceDefs")
plotImputationVariance(sampleMiceDefs)
```

---

plotModelError	<i>plotModelError</i>
----------------	-----------------------

---

**Description**

Plot the Out Of Bag model error for specified variables over all iterations.

**Usage**

```
plotModelError(
  miceObj,
  vars = names(miceObj$callParams$vars),
  pointSize = 1,
  ...
)
```

**Arguments**

miceObj	an object of class miceDefs, created by the miceRanger function.
vars	the variables you want to plot. Default is to plot all variables. Can be a vector of variable names, or one of 'allNumeric' or 'allCategorical'
pointSize	passed to geom_point, allows user to change dot size.
...	other arguments passed to ggarrange()

**Value**

an object of class ggarrange.

**Examples**

```
data("sampleMiceDefs")
plotModelError(sampleMiceDefs)
```

---

plotVarConvergence	<i>plotVarConvergence</i>
--------------------	---------------------------

---

**Description**

Plot the evolution of the dispersion and center of each variable. For numeric variables, the center is the mean, and the dispersion is the standard deviation. For categorical variables, the center is the mode, and the dispersion is the entropy of the distribution.

**Usage**

```
plotVarConvergence(miceObj, vars = names(miceObj$callParams$vars), ...)
```

**Arguments**

<code>miceObj</code>	an object of class <code>miceDefs</code> , created by the <code>miceRanger</code> function.
<code>vars</code>	the variables you want to plot. Default is to plot all variables. Can be a vector of variable names, or one of 'allNumeric' or 'allCategorical'
<code>...</code>	options passed to <code>ggarrange()</code>

**Value**

an object of class `ggarrange`.

**Examples**

```
data("sampleMiceDefs")
plotVarConvergence(sampleMiceDefs)
```

---

<code>plotVarImportance</code>	<i>plotVarImportance</i>
--------------------------------	--------------------------

---

**Description**

Plot the variable importance for each imputed variable.

**Usage**

```
plotVarImportance(
  miceObj,
  display = c("Relative", "Absolute"),
  dataset = 1,
  ...
)
```

**Arguments**

<code>miceObj</code>	an object of class <code>miceDefs</code> , created by the <code>miceRanger</code> function.
<code>display</code>	How do you want to display variable importance? <ul style="list-style-type: none"> <li>• "Relative" Scales the importance measure between 0-1 for each variable.</li> <li>• "Absolute" Displays the variable importance as is. May be highly skewed.</li> </ul>
<code>dataset</code>	The dataset you want to plot the variable importance of.
<code>...</code>	Other arguments passed to <code>corrplot</code> .

**Value**

nothing.

**Examples**

```
data("sampleMiceDefs")
plotVarImportance(sampleMiceDefs)
```

---

print.miceDefs	<i>Print a miceDefs object</i>
----------------	--------------------------------

---

**Description**

Print a miceDefs object

**Usage**

```
## S3 method for class 'miceDefs'
print(x, ...)
```

**Arguments**

x	Object of class miceDefs
...	required to use S3 method

**Value**

NULL

---

sampleMiceDefs	<i>Sample miceDefs object built off of iris dataset. Included so examples don't run for too long.</i>
----------------	---

---

**Description**

Sample miceDefs object built off of iris dataset. Included so examples don't run for too long.

**Usage**

```
sampleMiceDefs
```

**Format**

A miceDefs object. See `“?miceRanger”` for details.

**Source**

```
set.seed(1991) data(iris) ampIris <- amputeData(iris) sampleMiceDefs <- miceRanger( ampIris
, m=3 , maxiter=3 , vars=list( Petal.Width = c("Sepal.Length", "Sepal.Width", "Petal.Length", "Species")
, Species = c("Sepal.Length", "Sepal.Width", "Petal.Length", "Petal.Width") ) )
```

**Examples**

```
## Not run:  
  sampleMiceDefs  
  
## End(Not run)
```

# Index

## \*Topic **datasets**

sampleMiceDefs, [13](#)

addDatasets, [2](#)

addIterations, [3](#)

amputeData, [4](#)

completeData, [4](#)

getVarImps, [5](#)

miceRanger, [6](#)

plotCorrelations, [8](#)

plotDistributions, [9](#)

plotImputationVariance, [10](#)

plotModelError, [11](#)

plotVarConvergence, [11](#)

plotVarImportance, [12](#)

print.miceDefs, [13](#)

sampleMiceDefs, [13](#)