

Package ‘synMicrodata’

April 7, 2024

Type Package

Title Synthetic Microdata Generator

Version 2.0.0

Date 2024-04-06

Maintainer Hang J. Kim <hangkim0@gmail.com>

Description This tool fits a non-parametric Bayesian model called a ``hierarchically coupled mixture model with local dependence (HCMM-LD)'' to the original microdata in order to generate synthetic microdata for privacy protection. The non-parametric feature of the adopted model is useful for capturing the joint distribution of the original input data in a highly flexible manner, leading to the generation of synthetic data whose distributional features are similar to that of the input data. The package allows the original input data to have missing values and impute them with the posterior predictive distribution, so no missing values exist in the synthetic data output. The method builds on the work of Murray and Reiter (2016) <[doi:10.1080/01621459.2016.1174132](https://doi.org/10.1080/01621459.2016.1174132)>.

License GPL (>= 3)

Imports methods, stats, graphics, utils, Rcpp

LinkingTo Rcpp, RcppArmadillo

RcppModules IO_module

NeedsCompilation yes

Author Hang J. Kim [aut, cre],
Juhee Lee [aut],
Young-Min Kim [aut],
Jared Murray [aut]

Repository CRAN

Date/Publication 2024-04-07 08:00:02 UTC

R topics documented:

| | |
|-----------------------|---|
| createModel | 2 |
| modelobject | 2 |

| | |
|----------------------------------|---|
| multipleSyn | 3 |
| plot.synMicro_object | 4 |
| Rcpp_modelobject-class | 5 |
| readData | 6 |

Index**8**

| | |
|-------------|------------------------------|
| createModel | <i>Create a model object</i> |
|-------------|------------------------------|

Description

Create a model object for `multipleSyn`.

Usage

```
createModel(data_obj, max_R_S_K = c(30, 50, 20))
```

Arguments

| | |
|-----------|---|
| data_obj | data object produced by <code>readData</code> |
| max_R_S_K | maximum value of the number of mixture component index (r, s, k). |

Value

`createModel` returns a [Rcpp_modelobject](#)

See Also

[multipleSyn](#), [readData](#)

| | |
|-------------|---|
| modelobject | <i>RCPP Implementation of the Library</i> |
|-------------|---|

Description

[Rcpp_modelobject-class](#)

Value

No return value

multipleSyn *Generate synthetic micro datasets*

Description

Generate synthetic micro datasets using a hierarchically coupled mixture model with local dependence (HCMM-LC).

Usage

```
multipleSyn(data_obj, model_obj, n_burnin, m, interval_btwn_Syn, show_iter = TRUE)

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

Arguments

| | |
|--------------------------------|--|
| <code>data_obj</code> | data object produced by <code>readData</code> . |
| <code>model_obj</code> | model object produced by <code>createModel</code> . |
| <code>n_burnin</code> | size of burn-in. |
| <code>m</code> | number of synthetic micro datasets to be generated. |
| <code>interval_btwn_Syn</code> | interval between MCMC iterations for generating synthetic micro datasets. |
| <code>show_iter</code> | logical value. If <code>TRUE</code> , <code>multipleSyn</code> will print history of (<code>r</code> , <code>s</code> , <code>k</code>) components on console. |
| <code>x</code> | object of class <code>synMicro_object</code> ; a result of a call to <code>multipleSyn()</code> . |
| <code>...</code> | further arguments passed to or from other methods. |

Value

`multipleSyn` returns a list of the following components:

| | |
|------------------------|--|
| <code>synt_data</code> | list of <code>m</code> synthetic micro datasets. |
| <code>comp_mat</code> | list of matrices of the mixture component indices. |
| <code>orig_data</code> | original dataset. |

References

Murray, J. S. and Reiter, J. P. (2016). Multiple imputation of missing categorical and continuous values via Bayesian mixture models with local dependence. *Journal of the American Statistical Association*, **111**(516), pp.1466-1479.

See Also

`readData`, `createModel`, `plot.synMicro_object`

Examples

```
## preparing to generate synthetic datasets
dat_obj <- readData(Y_input = iris[,1:4],
                      X_input = data.frame(Species = iris[,5]))
mod_obj <- createModel(dat_obj, max_R_S_K=c(30,50,20))

## generating synthetic datasets
res_obj <- multipleSyn(dat_obj, mod_obj, n_burnin = 100, m = 5,
                        interval_btwn_Syn = 50, show_iter = FALSE)

print(res_obj)
```

plot.synMicro_object *Plot Comparing Synthetic Data with Original Input Data*

Description

The `plot` method for `synMicro_object` object. This method compares synthetic datasets with original input data.

Usage

```
## S3 method for class 'synMicro_object'
plot(x, vars, ...)
```

Arguments

| | |
|-------------------|--|
| <code>x</code> | <code>synMicro_object</code> object. |
| <code>vars</code> | vector of names or indices of the variables to compare. |
| <code>...</code> | other parameters to be passed through to plotting functions. |

Details

The `plot` takes input variables and draws the graph. The type of graph produced is contingent upon the number of categories in selected variables.

- Putting a continuous variable produces a *box plot* of the selected variable.
- Putting more than two continuous variables produces *pairwise scatter plots* for each pair of selected variables.
- Putting categorical variables produce *bar plot* of each selected variable.

See Also

[multipleSyn](#)

Examples

```

## preparing to generate synthetic datasets
dat_obj <- readData(Y_input = iris[,1:4],
                      X_input = data.frame(Species = iris[,5]))
mod_obj <- createModel(dat_obj, max_R_S_K=c(30,50,20))

## generating synthetic datasets
res_obj <- multipleSyn(dat_obj, mod_obj, n_burnin = 100, m = 2,
                        interval_btwn_Syn = 50, show_iter = FALSE)

print(res_obj)

## plotting synthesis datasets
### box plot
par(mfrow=c(3,2))
plot(res_obj, vars = "Sepal.Length") ## variable names

### pairwise scatter plot
plot(res_obj, vars = c(1,2)) ## or variable index

### bar plot
plot(res_obj, vars = "Species")

```

Rcpp_modelobject-class

Class "Rcpp_modelobject"

Description

This class implements a joint modeling approach to generate synthetic microdata with continuous and categorical variables with possibly missing values. The method builds on the work of Murray and Reiter (2016)

Details

Rcpp_modelobject should be created with [createModel](#). Please see the example below.

Extends

Class "[C++Object](#)", directly.

Fields

- data_obj input dataset generated from [readData](#).

Methods

- `multipleSyn` generates synthetic micro datasets.

References

Murray, J. S. and Reiter, J. P. (2016). Multiple imputation of missing categorical and continuous values via Bayesian mixture models with local dependence. *Journal of the American Statistical Association*, **111(516)**, pp.1466-1479.

Examples

```
## preparing to generate synthetic datasets
dat_obj <- readData(Y_input = iris[,1:4],
                      X_input = data.frame(Species = iris[,5]))
mod_obj <- createModel(dat_obj, max_R_S_K=c(30,50,20))

## generating synthetic datasets
res_obj <- multipleSyn(dat_obj, mod_obj, n_burnin = 100, m = 5,
                        interval_btwn_Syn = 50, show_iter = FALSE)

print(res_obj)
```

`readData`

Read the original datasets

Description

Read the original input datasets to be learned for synthetic data generation. The package allows the input data to have missing values and impute them with the posterior predictive distribution, so no missing values exist in the synthetic data output.

Usage

```
readData(Y_input, X_input, RandomSeed = 99)
```

Arguments

| | |
|-------------------------|--|
| <code>Y_input</code> | data.frame consisting of continuous variables of the original data. It should consist only of <code>numeric</code> . |
| <code>X_input</code> | data.frame consisting of categorical variables of the original data. It should consist only of <code>factor</code> . |
| <code>RandomSeed</code> | random seed number. |

Value

`readData` returns an object of "readData_passed" class.

An object of class "readData_passed" is a list containing the following components:

| | |
|-----------------------------|---|
| <code>n_sample</code> | number of records in the input dataset. |
| <code>p_Y</code> | number of continuous variables. |
| <code>Y_mat_std</code> | matrix with standardized values of <code>Y_input</code> , with mean 0 and standard deviation 1. |
| <code>mean_Y_input</code> | mean vectors of original <code>Y_input</code> . |
| <code>sd_Y_input</code> | standard deviation vectors of original <code>Y_input</code> . |
| <code>NA_Y_mat</code> | matrix indicating missing values in <code>Y_input</code> . |
| <code>p_X</code> | number of categorical variables. |
| <code>D_l_vec</code> | numbers of levels of each categorical variable. |
| <code>X_mat_std</code> | matrix with the numeric-transformed values of <code>X_input</code> . |
| <code>levels_X_input</code> | list of levels of each categorical variable. |
| <code>NA_X_mat</code> | matrix indicating missing values in <code>X_input</code> . |
| <code>var_names</code> | list containing variable names of <code>X_input</code> and <code>Y_input</code> . |
| <code>orig_data</code> | original dataset. |

See Also

[multipleSyn](#), [createModel](#)

Index

* classes

Rcpp_modelobject-class, [5](#)

C++Object, [5](#)

createModel, [2](#), [3](#), [5](#), [7](#)

modelobject, [2](#)

multipleSyn, [2](#), [3](#), [4](#), [7](#)

plot.synMicro_object, [3](#), [4](#)

print.synMicro_object (multipleSyn), [3](#)

Rcpp_modelobject, [2](#)

Rcpp_modelobject

(Rcpp_modelobject-class), [5](#)

Rcpp_modelobject-class, [5](#)

readData, [2](#), [3](#), [5](#), [6](#)