

# Package ‘fakeR’

May 26, 2016

**Type** Package

**Title** Simulates Data from a Data Frame of Different Variable Types

**Version** 1.0

**Date** 2016-05-25

**Description** Generates fake data from a dataset of different variable types.

The package contains the functions `simulate_dataset` and `simulate_dataset_ts` to simulate time-independent and time-dependent data. It randomly samples character and factor variables from contingency tables and numeric and ordered factors from a multivariate normal distribution. It currently supports the simulation of stationary and zero-inflated count time series.

**License** CC0

**Imports** mvtnorm, polycor, pscl, VGAM, stats

**Suggests** knitr, rmarkdown, testthat

**VignetteBuilder** knitr

**NeedsCompilation** no

**Author** Lily Zhang [aut, cre],  
Dustin Tingley [aut]

**Maintainer** Lily Zhang <lilyhzhang1029@gmail.com>

**Repository** CRAN

**Date/Publication** 2016-05-26 17:54:23

## R topics documented:

fakeR-package . . . . .	2
simulate_dataset . . . . .	3
simulate_dataset_ts . . . . .	5

<b>Index</b>	<b>8</b>
--------------	----------

fakeR-package

*Simulates Data from a Data Frame of Different Variable Types***Description**

Generates fake data from a dataset of different variable types. The package contains the functions `simulate_dataset` and `simulate_dataset_ts` to simulate time-independent and time-dependent data. It randomly samples character and factor variables from contingency tables and numeric and ordered factors from a multivariate normal distribution. It currently supports the simulation of stationary and zero-inflated count time series.

**Details**

The DESCRIPTION file:

```
Package:      fakeR
Type:        Package
Title:       Simulates Data from a Data Frame of Different Variable Types
Version:     1.0
Date:       2016-05-25
Authors@R:  c(person("Lily", "Zhang", email = "lilyhzhang1029@gmail.com", role = c("aut", "cre")), person("Dustin",
Description: Generates fake data from a dataset of different variable types. The package contains the functions simulate
License:    CC0
Imports:    mvtnorm, polycor, pscl, VGAM, stats
Suggests:  knitr, rmarkdown, testthat
VignetteBuilder: knitr
Author:     Lily Zhang [aut, cre], Dustin Tingley [aut]
Maintainer: Lily Zhang <lilyhzhang1029@gmail.com>
```

Index of help topics:

```
fakeR-package      Simulates Data from a Data Frame of Different
                   Variable Types
simulate_dataset   Simulate from a data frame of time-independent
                   data.
simulate_dataset_ts Simulate a dataframe of time series data
```

This package is used to simulate datasets of different variable types. The package contains the functions `simulate_dataset` and `simulate_dataset_ts` to simulate time-independent and time-dependent data.

**Author(s)**

NA

Maintainer: NA

## References

~~ Literature or other references for background information ~~

## Examples

```
## time-independent data frame of multiple types
# single column of an unordered, string factor
state_df <- data.frame(division=state.division)
# character variable
state_df$division <- as.character(state_df$division)
# numeric variable
state_df$area <- state.area
# factor variable
state_df$region <- state.region
state_sim <- simulate_dataset(state_df)

## time-independent data frame with missingness
df <- mtcars
# change one of the variable types to an unordered factor
df$carb <- as.factor(df$carb)
# change another variable type to an ordered factor
df$gear <- as.ordered(as.factor(df$gear))
df[2,] <- NA
sim_df <- simulate_dataset(df, stealth.level=2)

## time series dataframe
tree_ring <- data.frame(treering)
tree_ring$year <- c(1: nrow(tree_ring))
sim_tree_ring <- simulate_dataset_ts(tree_ring,
                                   cluster="treering",
                                   time.variable="year")
par(mfrow = c(2, 1), mar = c(3, 3, 4, 2), mgp = 0.9 * 2:0)
plot (tree_ring$year, tree_ring$treering, type='l',
      main=paste("Original","Normalized ring width"),
      ylab="Ring width", xlab="Year index")
plot (tree_ring$year, tree_ring$treering, type='l',
      main=paste("Simulated","Normalized ring width"),
      ylab="Ring width", xlab="Year index")
```

---

simulate\_dataset

*Simulate from a data frame of time-independent data.*

---

## Description

This function takes as argument an existing dataset in the form of a data frame and outputs a randomized version of all its columns. The function accepts the following types: character variables, numeric variables, and ordered and unordered factor variables.

**Usage**

```
simulate_dataset(dataset, digits=2, n=NA,
                 use.levels=TRUE, use.miss=TRUE,
                 mvt.method="eigen", het.ML=FALSE,
                 het.suppress=TRUE, stealth.level=1,
                 level3.noise=FALSE, ignore=NA)
```

**Arguments**

dataset	the data frame from which to generate a randomized version
digits	the number of digits after the decimal point to include in the new values
n	number of rows in the new data frame. Equal to the number of rows in the original if set to NA, the default.
use.levels	when set to true, gives the simulated factor variables the same number of levels as the original.
use.miss	when set to TRUE, inserts the missing data like is present in the original (i.e. based on the distribution of missingness in the original data).
mvt.method	specifies the matrix decomposition to be used in sampling from the multivariate normal.
het.ML	as per the hetcor function, if TRUE, compute maximum-likelihood estimates; if FALSE, compute quick two-step estimates in computing the heterogeneous correlation matrix.
het.suppress	when set to TRUE, suppresses stops from the het.corr function.
stealth.level	when set to 1 (default), takes into account the covariances between all the unordered factors and the covariances between the numeric and ordered factors. When set to 2, simulates each variable independently. When set to 3, does not take into account any covariances and instead randomly samples from a uniform distribution ranging from the min to the max of the data for each variable.
level3.noise	when set to TRUE, add Gaussian noise to the min and max parameter for the uniform distribution in stealth.level 3. The noise term has a variance of one fourth of the range of the data for any particular variable.
ignore	specifies which columns to ignore (i.e. to leave as is instead of simulate). Takes in a list of column names as input.

**Details**

This function does not account for clustered time series data (see `simulate_dataset_ts`).

This function randomly samples each character and factor variable from the population distribution given in the original dataset. It simulates numeric and ordered factors from a multivariate normal distribution. When both numeric and ordered factors are included, a heterogeneous correlation matrix is used, coercing the means of the ordered factor variables to be 0.

The function only accounts for between-column correlations for numeric and ordered factor variables. Each unordered factor and character column is treated as independent.

The order of the columns in the simulated dataset may differ from the order of the original dataset since the function puts the numeric and ordered factor data in the front and the character and unordered factor data afterwards. The column names stay consistent, however.

**Value**

Returns a data frame with the same number of columns and same type for each.

**Author(s)**

Lily Zhang Dustin Tingley

**References**

Inspired by the fakeR function originally created by Ryne Estabrook.

**Examples**

```
# single column of an unordered, string factor
state_df <- data.frame(division=state.division)
# character variable
state_df$division <- as.character(state_df$division)
# numeric variable
state_df$area <- state.area
# factor variable
state_df$region <- state.region
state_sim <- simulate_dataset(state_df)
```

---

simulate\_dataset\_ts     *Simulate a dataframe of time series data*

---

**Description**

This function simulates clustered numeric time series data from an ARIMA model fit or, if specified, a zero-inflated Poisson regression model fit, with each column variable regressed on the first lag.

**Usage**

```
simulate_dataset_ts(dataset, digits=2, n=NA, cluster=NA, time.variable=NA,
                    date.index=FALSE, complete.panel=FALSE, zero.inflate=FALSE,
                    stealth.level=2, level3.noise=FALSE, use.miss=TRUE, ignore=NA)
```

**Arguments**

dataset	the data frame from which to generate a randomized version
digits	the number of digits after the decimal point to include in the new values
n	number of rows in the new data frame. Equal to the number of rows in the original if set to NA, the default.
cluster	the column names of the time series variables. Argument should be in the form of a list if multiple values.



```
                                time.variable="year")
par(mfrow = c(2, 1), mar = c(3, 3, 4, 2), mgp = 0.9 * 2:0)
plot (tree_ring$year, tree_ring$treering, type='l',
      main=paste("Original", "Normalized ring width"),
      ylab="Ring width", xlab="Year index")
plot (tree_ring$year, tree_ring$treering, type='l',
      main=paste("Simulated", "Normalized ring width"),
      ylab="Ring width", xlab="Year index")
```

# Index

\*Topic **datagen**

simulate\_dataset, 3

simulate\_dataset\_ts, 5

\*Topic **manip**

simulate\_dataset, 3

simulate\_dataset\_ts, 5

\*Topic **package**

fakeR-package, 2

fakeR (fakeR-package), 2

fakeR-package, 2

simulate\_dataset, 3

simulate\_dataset\_ts, 5