

# Package ‘rstatsToolkit’

May 7, 2015

**Title** Bundles up all my most used functions for doing statistical analysis

**Version** 0.1

**Description** This package is mainly my personal collection of code that I commonly use in my research and statistical analysis. Eventually I would like to develop it into a toolkit that other graduate students (I'm a graduate student right now btw) and in the future for my own graduate students to use and develop. Until that point, I am slowly developing this package into toolkit for analyzing and exploring data.

**Imports** gee (>= 4.13.18),  
magrittr,  
data.table (>= 1.9.4),  
ggplot2 (>= 1.0.0),  
grid (>= 3.1.1),  
visreg (>= 2.0.5),  
reshape2 (>= 1.4),  
dplyr (>= 0.4.1)

**Depends** R (>= 3.1.1)

**License** Creative Commons 0

**LazyData** true

## R topics documented:

bivarPlot . . . . .	2
createCI . . . . .	3
createFormulaList . . . . .	4
diagnosticPlots . . . . .	5
extractBetaFromListGEE . . . . .	6
extractBetaGEE . . . . .	7
heatmapCorr . . . . .	7
jitterBoxplot . . . . .	9
loopOutputToListGEE . . . . .	10
multiPlot . . . . .	11
plotForest . . . . .	12
plotSpaghetti . . . . .	13

plotVisreg . . . . .	14
rstatsToolkit . . . . .	15
smoothPlot . . . . .	15
smoothTimePlot . . . . .	16
summarySE . . . . .	17
themeWhite . . . . .	18
unlistAndFilterIndep . . . . .	18
<b>Index</b>	<b>20</b>

---

bivarPlot	<i>Bivariate plot</i>
-----------	-----------------------

---

**Description**

In development ..

**Usage**

bivarPlot(x, y, data, ...)

**Arguments**

- x
- y
- data
- ...

**Details**

.. content for details ..

**Value**

Outputs a plot

**Author(s)**

Luke Johnston

---

createCI	<i>Compute confidence interval for GEE</i>
----------	--

---

## Description

Generate confidence intervals (upper and lower) and p-values from a regression (e.g. GEE or GLM).

## Usage

```
createCI(data, dist = test.distribution, se = standard.error,
  test.distribution = "Naive.z", standard.error = "Naive.S.E.",
  sig.level = 0.95)
```

## Arguments

data	The results output from a GEE analysis from the <a href="#">gee</a> package.
test.distribution, dist	The column that contains the test statistic distribution, e.g. z-score or t-score.
standard.error, se	The column that contains the standard error.
sig.level	The significance level for calculating the confidence interval.

## Value

Columns with the confidence interval and p-value.

## Author(s)

Luke W. Johnston

## Examples

```
## Very simple test example. Merely to show how the function is used.
outcomes <- c('Income', 'Population')
exposures <- c('Frost', 'Illiteracy')
covariates <- c('Murder', 'LifeExp')
interaction <- 'LifeExp'

## This uses the dplyr package.
ds <- cbind(state.region, state.x77) %>%
  as.data.frame() %>%
  rename(LifeExp = `Life Exp`,
    ## Need to rename the id variable to SID (see description
    ## above)
    SID = state.region) %>%
  arrange(SID)

## Without interaction
```

```

loopOutputToListGEE(ds, outcomes, exposures, covariates,
                    corstr = 'exchangeable') %>%
  extractBetaFromListGEE() %>%
  unlistAndFilterIndep(., exposures) %>%
  createCI()

## With interaction
loopOutputToListGEE(ds, outcomes, exposures, covariates,
                    interaction, corstr = 'exchangeable') %>%
  extractBetaFromListGEE() %>%
  unlistAndFilterIndep(., ':', pattern = TRUE) %>%
  ## The ':' represents an interaction
  createCI()

```

---

createFormulaList	<i>Create a list of formulas</i>
-------------------	----------------------------------

---

## Description

createFormulaList returns a list of formulas with all combinations of dependent and independent variables.

## Usage

```
createFormulaList(dependent, independent, covariates, interactions = NULL)
```

## Arguments

dependent	A single or a vector of variables names that the user wishes to use as dependent variables in an analysis; dependent must be as a string/character.
independent	Similar to dependent, except that the variables would be the independent, or exposure, variables of interest.
covariates	Confounders or covariates that would be included in <b>all</b> formulas and hence <b>all</b> models. Also must be a string/character.
interactions	Optional, include <i>one</i> variable from the covariates set that will be assigned as an interaction term with the independent variable. Must be a string/character.

## Details

This function creates a list of formulas for use in a chain of processes. This function's goal is to create all combinations of a set of dependent, or outcome, variables with a set of independent, or exposure variables. Covariates and an interaction term can also be specified and included into the formula.

## Value

Outputs a list of formulas.

**Author(s)**

Luke W. Johnston

**Examples**

```
outcomes <- c('Income', 'Education', 'Job')
exposures <- c('Age', 'Sex', 'Height', 'Race', 'IQ')
covariates <- c('ParentEdu', 'Country', 'City')
interactions <- 'City'

createFormulaList(outcomes, exposures, covariates)
createFormulaList(outcomes, exposures, covariates, interactions)
```

---

diagnosticPlots

*Regression diagnostic plots and tests*


---

**Description**

Generate regression diagnostic plots and tests for linear regression models.

**Usage**

```
diagnosticPlots(data, y, x, covar)
```

**Arguments**

data	The dataset with the variables of interest
y	The dependent or outcome variable (that is, the y in the regression equation)
x	The independent, exposure, or predictor variable (that is, the x in the regression equation)
covar	The variables selected as to condition or adjust for the y and x relationship, also known as the confounding variables

**Details**

This function runs a linear regression on the specified variables and generates diagnostics based on the regression. Basic diagnostics include checking the normality of the residuals, assessing outliers, influence and Cook's D, and multicollinearity. Several tests have been commented out, though they can be uncommented if desired (edit the function to output these if desired). Some of the tests I don't fully understand how to interpret them, but as I learn more I will probably know. This function relies on **MASS** and **gplots**.

**Value**

Outputs multiple plots and textplots with diagnostic information

**Author(s)**

Luke Johnston

---

```
extractBetaFromListGEE
```

*Extract beta coefficients from a GEE list object*

---

## Description

Used as part of a chain, `extractBetaFromListGEE` grabs the beta estimate from a list of GEE objects.

## Usage

```
extractBetaFromListGEE(data)
```

## Arguments

`data` This is variable that contains the list of GEE objects.

## Details

This is used after [loopOutputToListGEE](#) as part of a chain, preferably using `dplyr/magrittr`'s pipe command ( After generating the list of GEE objects, `extractBetaFromListGEE` loops through each GEE object and converts the beta estimate and associated statistics into list of dataframes.

## Value

Outputs the beta coefficients.

## Author(s)

Luke W. Johnston

## Examples

```
## Very simple test example. Merely to show how the function is used.
outcomes <- c('Income', 'Population')
exposures <- c('Frost', 'Illiteracy')
covariates <- c('Murder', 'LifeExp')
interaction <- 'LifeExp'

## This uses the dplyr package.
ds <- cbind(state.region, state.x77) %>%
  as.data.frame() %>%
  rename(LifeExp = `Life Exp`,
    ## Need to rename the id variable to SID (see description
    ## above)
    SID = state.region) %>%
  arrange(SID)

## Without interaction
```

```

loopOutputToListGEE(ds, outcomes, exposures, covariates,
                     corstr = 'exchangeable') %>%
  extractBetaFromListGEE()

## With interaction
loopOutputToListGEE(ds, outcomes, exposures, covariates,
                     interaction, corstr = 'exchangeable') %>%
  extractBetaFromListGEE()

```

---

extractBetaGEE	<i>Extract beta coefficients from a GEE object</i>
----------------	--

---

### Description

Extract the beta object from a GEE object.

### Usage

```
extractBetaGEE(data)
```

### Arguments

data	The variable that contains the GEE object.
------	--

### Author(s)

Luke W. Johnston

---

heatmapCorr	<i>Correlation heatmap</i>
-------------	----------------------------

---

### Description

Generate a matrix or non-matrix style heatmap of correlation coefficients (i.e. the number of columns and rows can be different).

### Usage

```

heatmapCorr(data, x, y, leg.range = c(-1, 1), levels.xlab = NULL,
             levels.ylab = NULL, xlab = "", ylab = "", lo.color = "darkorange2",
             hi.color = "skyblue4", rm.legend = FALSE, matrix.sty = FALSE,
             corr.values = TRUE, leg.title.size = 8, leg.number.size = 7,
             axis.text.size = 10)

```

**Arguments**

<code>data</code>	Dataset that contains the variables of interest
<code>x</code>	Vector of variables that will run along the x-axis
<code>y</code>	Same as the <code>x</code> arg, but for the y-axis
<code>leg.range</code>	Range in values for the legend, between -1 and 1
<code>levels.xlab</code>	Specifies custom variable names for the x-axis (it needs to be a list object)
<code>levels.ylab</code>	Same as the <code>levels.xlab</code> arg, but for the y-axis
<code>xlab</code>	Sets the label for the x-axis
<code>ylab</code>	Same as <code>xlab</code> , but for the y-axis
<code>lo.color</code>	Color of negative correlation coefficients
<code>hi.color</code>	Color of positive correlation coefficients
<code>rm.legend</code>	In development. Goal is it will remove the legend
<code>matrix.sty</code>	Select whether the heatmap will be a matrix style (the same variable on the x- and y- axis) or non-matrix (different variables on both axes)
<code>corr.values</code>	Set whether to have the correlation values on the heatmap, or if it will be blank. Value is either TRUE or FALSE
<code>leg.title.size</code>	Size of the font in the legend title
<code>leg.number.size</code>	Size of the numbers in or near the legend

**Details**

This function takes two arguments, the x variables and the y variables, and generates a heatmap from the variables. A correlation matrix is computed from the data, melted (reshape package), and input into ggplot2 to generate a heatmap. The output is the correlations and the plot object.

Dependencies are: **reshape2** and **ggplot2**

**Value**

Outputs a plot object

**Author(s)**

Luke Johnston

**Examples**

```
xlabs <- list("X Name" = "x1", "X Name Two!" = "x2" ... )
ylabs <- list("Y Name" = "y1", "Y Name Two!" = "y2" ... )
df <- data.frame(replicate(10, sample(0:1, 1000, rep=TRUE)))
plot <- heatmap.corr(data = df, x = 1:3, y = 4:10,
                    levels.xlab=xlabs, levels.ylab=ylabs)
print(plot)
```



---

jitterBoxplot	<i>Univariate jittered boxplot</i>
---------------	------------------------------------

---

## Description

Generates a boxplot of variables on one axis with raw values "jittered" as dots underneath. The variables need to represent a similar concept or have the same units for the plot to make sense.

## Usage

```
jitterBoxplot(subset.ds, dot.size = 2, dot.colour = "grey50",  
  custom.var.names = NULL, xlab = NULL, ylab = NULL)
```

## Arguments

subset.ds	The dataset that only contains the series of variables that will be plotted
dot.size	Size of the dot for <code>geom_jitter</code>
dot.colour	Color of the dots for <code>geom_jitter</code>
custom.var.names	List object that contains the custom (alternative) variable names for the variable (column) names you passed into the function
xlab	The x-axis label
ylab	The y-axis label

## Details

This function is useful for exploring the distribution of a series of variables that share a common unit, such as kilogram. The values for each variable are plotted as jittered dots with a boxplot of the distribution layered on top of the dots. The function takes a subsetting dataset that contains only the series of variables that share a common unit. The output object is the plot. This function depends on **ggplot2** and **reshape2**.

## Value

Outputs the plot object

## Author(s)

Luke Johnston

## Examples

```
varnames <- list("Name of X1" = "X1", "Name of X2" = "X2", ...)  
df <- subset(ds, VN == 0, select=nefa)  
names(df)  
jitter.boxplot(df, xlab="Concentration (nmol/mL)",  
  ylab="Fatty acid", custom.var.names=varnames)
```

---

loopOutputToListGEE     *Loop through GEE analyses.*

---

## Description

Run a GEE analysis on a list of formulas and create a list of GEE objects.

## Usage

```
loopOutputToListGEE(data, dependent, independent, covariates,
  interactions = NULL, corstr = "exchangeable")
```

## Arguments

data	The dataset with the variables of interest
dependent	A single or a vector of variables names that the user wishes to use as dependent variables in an analysis; dependent must be as a string/character.
independent	Similar to dependent, except that the variables would be the independent, or exposure, variables of interest.
covariates	Confounders or covariates that would be included in <b>all</b> formulas and hence <b>all</b> models. Also must be a string/character.
interactions	Optional, include <i>one</i> variable from the covariates set that will be assigned as an interaction term with the independent variable. Must be a string/character.
corstr	The working correlation structure to use in the <a href="#">gee</a> function call. Options can be found in the gee package, but include 'exchangeable', 'unstructured', 'AR-M' (M =

## Details

This function is merely a wrapper around [gee](#), and therefore all inquiries into GEE should start there. As a note, I don't know how (yet) to provide an option to use a custom id variable for the GEE analysis, so **make sure to rename your id variable to 'SID'**. Also, according to the GEE documentations, the **ordering of the ID variables matters!** Make sure to sort the ID variables as per how you want them! For documentation on the function to create the formula list, see [createFormulaList](#).

## Value

Creates a list of GEE objects

## Author(s)

Luke W. Johnston

## Examples

```
## Very simple test example. Merely to show how the function is used.
outcomes <- c('Income', 'Population')
exposures <- c('Frost', 'Illiteracy')
covariates <- c('Murder', 'LifeExp')
interaction <- 'LifeExp'

## This uses the dplyr package.
ds <- cbind(state.region, state.x77) %>%
  as.data.frame() %>%
  rename(LifeExp = `Life Exp`,
         ## Need to rename the id variable to SID (see description
         ## above)
         SID = state.region) %>%
  arrange(SID)

## Without interaction
loopOutputToListGEE(ds, outcomes, exposures, covariates,
                    corstr = 'exchangeable')

## With interaction
loopOutputToListGEE(ds, outcomes, exposures, covariates,
                    interaction, corstr = 'exchangeable')
```

---

multiPlot

---

*Multiple plots on page*


---

## Description

Lay out multiple ggplots on one frame or page.

## Usage

```
multiPlot(..., plotlist = NULL, file, cols = 1, layout = NULL)
```

## Arguments

...	Where the ggplot objects are placed to be laid out on the graph grid
plotlist	Can be used in place of the ... argument by specifying the ggplot objects as a list object
file	Not sure what this is used for
cols	Number of columns for the layout. For example cols=2 provides two columns and with four ggplot objects, the resulting output would be a 2 by 2 graphic
layout	A matrix that indicates the plot grid layout. For example, if layout = matrix(c(1, 2, 3, 3), nrow = 2, byrow = TRUE) the result would have plot 1 in the upper left, plot 2 in the upper right, and plot 3 would be go across the bottom

## Details

This function, which was from <http://www.cookbook-r.com/Graphs>, is used to lay out several ggplot objects onto one frame or pdf page. For instance, you can have 3 plots on a page, one going vertically across the top, the other two in each corner on the bottom. This function makes up for the difficulty **ggplot2** has with outputting multiple plots on one grid. This function depends on **grid**.

## Author(s)

Cookbook R

---

plotForest	<i>Forest plot</i>
------------	--------------------

---

## Description

Generate a forest plot without the traditional side table .

## Usage

```
plotForest(data, coeff = coefficient, yvar = y.variables.column,
  ylab = yaxis.label, xlab = xaxis.label, ci = confid.interval,
  dot.pval = pvalue.factor.column, coefficient = "Estimate",
  y.variables.column = "indep", confid.interval = c("lowerCI", "upperCI"),
  pvalue.factor.column = "NULL", groups = NULL, yaxis.label = "Exposures",
  xaxis.label = "Beta estimate")
```

## Arguments

data	Dataset for the forest plot.
coefficient,coeff	The column that contains the beta estimate/coefficient.
y.variables.column,yvar	The column with the exposure variables that will be places on the y-axis of the forest plot.
confid.interval,ci	A vector that contains the lower and upper confidence interval.
pvalue.factor.column,dot.pval	The column that contains the p-value in the form of a factor variable (ie. with levels such as '>0.05' and '<0.05').
groups	The variable to split the plot up.
yaxis.label,ylab	The y-axis label.
xaxis.label,xlab	The x-axis label.

**Details**

Create a forest plot, with a dot and 95 without the usual side table that contains the raw data values. If the `dot.pval` argument is supplied, the dots and confidence lines increase in size and opacity as significance increases. If `groups` is also supplied, the forest plot will be split up vertically for each grouping. Thus, a large amount of information on the results can be provided in a fairly small amount of space.

**Value**

A forest plot

**Author(s)**

Luke W. Johnston

**Examples**

```
## Very simple test example. Merely to show how the function is used.
outcomes <- c('Income', 'Population')
exposures <- c('Frost', 'Illiteracy')
covariates <- c('Murder', 'LifeExp')
interaction <- 'LifeExp'

## This uses the dplyr package.
ds <- cbind(state.region, state.x77) %>%
  as.data.frame() %>%
  rename(LifeExp = `Life Exp`,
         ## Need to rename the id variable to SID (see description
         ## above)
         SID = state.region) %>%
  arrange(SID)

loopOutputToListGEE(ds, outcomes, exposures, covariates,
                    corstr = 'exchangeable') %>%
  extractBetaFromListGEE() %>%
  unlistAndFilterIndep(., exposures) %>%
  createCI() %>%
  plotForest(., dot.pval = 'f.pvalue', groups = 'dep')
```

---

plotSpaghetti

*Spaghetti plot*


---

**Description**

Plot subjects in a longitudinal dataset, making a 'spaghetti' plot.

**Usage**

```
plotSpaghetti(data, y, x, groups = "SID")
```

**Arguments**

data	The dataset to plot.
y	The variable to go on the y-axis.
x	The variable to go on the x-axis.
groups	The unique ID variable to differentiate subjects in a longitudinal dataset.

**Author(s)**

Luke W. Johnston

**Examples**

```
## Not run:
## A pretend case (not a real example)
plotSpaghetti(dataset, 'Height', 'Year', 'SubjectID')

## End(Not run)
```

---

plotVisreg

---

*Visualizing adjusted linear regression models*


---

**Description**

Generates plots of a linear regression model which includes confounding variables.

**Usage**

```
plotVisreg(data, y, x, covar, ylabel = x, xlabel = y, ...)
```

**Arguments**

data	Dataset with the variables of interest
y	The dependent or outcome variable in the regression equation
x	The independent or exposure variable in the regression equation
covar	The confounding variables, that is the variables being adjusted for
ylabel	The y-axis label
xlabel	The x-axis label
...	Other options. In development

**Details**

This function runs a linear regression on the specified variables and plots the partial residuals. This allows for visualizing the relationship between the outcome and the exposure, after adjusting for confounders. A linear slope is plotted through the partial residuals, with a confidence interval band around it. The output is a plot. This function depends on **visreg**.

**Value**

Outputs a plot of the regression model

**Author(s)**

Luke Johnston

---

rstatsToolkit

*rstatsToolkit.*

---

**Description**

rstatsToolkit.

---

smoothPlot

*Smooth plot*

---

**Description**

In development ..

**Usage**

```
smoothPlot(x, y, data, ...)
```

**Arguments**

x

y

data

...

**Details**

.. content for details ..

**Value**

Outputs a plot

**Author(s)**

Luke Johnston

---

smoothTimePlot	<i>Smooth time plots</i>
----------------	--------------------------

---

### Description

Generates a plot of two variables, particularly for longitudinal data. The x-axis is typically the time variable.

### Usage

```
smoothTimePlot(dsn, x, y, id, smooth.type = "loess", ylab = y, xlab = x)
```

### Arguments

dsn	The dataset name.
x	Specify the variable for the x-axis. Must be either in quotes or as a number.
y	As with x, but for the y-axis.
id	Specify the ID variable for longitudinal or repeated measures type data.
smooth.type	Set the smoothing method, which includes "loess", "lm", "rlm".
ylab	Optionally set the y-axis label.
xlab	Optionally set the x-axis label.

### Details

This function creates a plot with smooth lines, typically using the loess method, along with the means of a discrete time variable (eg. time 0, 1, 2, 3 etc).

Dependencies are: **ggplot2**

### Value

Outputs a ggplot object

### Author(s)

Luke Johnston



---

summarySE	<i>Summarize means and standard errors of the mean</i>
-----------	--

---

### Description

Calculates the sample size, mean, standard deviation, standard error of the mean, and the confidence interval of specified variables.

### Usage

```
summarySE(data = NULL, measurevar, groupvars = NULL, na.rm = FALSE,  
  conf.interval = 0.95, .drop = TRUE)
```

### Arguments

data	A dataset (dataframe) that contains the values to be summarized
measurevar	The name of a column that contains the variable to be summarized
groupvars	A vector containing names of columns that contain grouping variables
na.rm	A binary (boolean) response that indicates whether to ignore missing (NA) data
conf.interval	Percent range of the confidence interval

### Details

I took this function on 2014-01-21 from the website <http://www.cookbook-r.com/Graphs>. It basically summarizes the provided data by giving count, mean, standard deviation, standard error of the mean, and confidence interval (default 95). The dependencies are **plyr**.

### Value

Outputs a dataframe that contains the summarized statistics (means, etc.)

### Author(s)

Cookbook R

---

themeWhite	<i>Custom white ggplot theme</i>
------------	----------------------------------

---

### Description

Creates a white, simple theme for **ggplot2** objects

### Usage

```
themeWhite()
```

### Details

The default **ggplot2** theme is decent for most purposes, but is visually unappealing. This function aims to correct that by setting the theme to something more similar to the default in the base R plot package. The function depends on **ggplot2**.

### Author(s)

Luke Johnston

### Examples

```
## This creates a white theme
themeWhite()
```

---

unlistAndFilterIndep	<i>Unlist a list of dataframes and filter a pattern or variable.</i>
----------------------	--

---

### Description

Unlist a list of dataframes, convert into a single dataframe, and filter out a string or pattern from the 'indep' column.

### Usage

```
unlistAndFilterIndep(data, x, pattern = FALSE)
```

### Arguments

data	The list object with the dataframes
x	The variables of interest (eg. exposures) that are within the 'indep' column.
pattern	Logical: TRUE if x is a pattern rather than an explicit list of variables and FALSE if x is an explicit list of variables (eg. exposures).

## Details

This function is generally used within a chain of other commands that creates a list of dataframes, generally from a regression. The list gets unlisted and converted into a single dataset. Afterward, the dataset gets filtered by the variables of interest (eg. the exposures) that are contained within the 'indep' column.

## Value

Outputs a single dataframe with only the rows with the variables from x.

## Author(s)

Luke W. Johnston

## Examples

```
## Very simple test example. Merely to show how the function is used.
outcomes <- c('Income', 'Population')
exposures <- c('Frost', 'Illiteracy')
covariates <- c('Murder', 'LifeExp')
interaction <- 'LifeExp'

## This uses the dplyr package.
ds <- cbind(state.region, state.x77) %>%
  as.data.frame() %>%
  rename(LifeExp = `Life Exp`,
         ## Need to rename the id variable to SID (see description
         ## above)
         SID = state.region) %>%
  arrange(SID)

## Without interaction
loopOutputToListGEE(ds, outcomes, exposures, covariates,
                    corstr = 'exchangeable') %>%
  extractBetaFromListGEE() %>%
  unlistAndFilterIndep(., exposures)

## With interaction
loopOutputToListGEE(ds, outcomes, exposures, covariates,
                    interaction, corstr = 'exchangeable') %>%
  extractBetaFromListGEE() %>%
  unlistAndFilterIndep(., ':', pattern = TRUE)
  ## The ':' represents an interaction
```

# Index

bivarPlot, [2](#)

createCI, [3](#)  
createFormulaList, [4](#), [10](#)

diagnosticPlots, [5](#)

extractBetaFromListGEE, [6](#)  
extractBetaGEE, [7](#)

gee, [3](#), [10](#)  
geom\_jitter, [9](#)

heatmapCorr, [7](#)

jitterBoxplot, [9](#)

loopOutputToListGEE, [6](#), [10](#)

multiPlot, [11](#)

plotForest, [12](#)  
plotSpaghetti, [13](#)  
plotVisreg, [14](#)

rstatsToolkit, [15](#)  
rstatsToolkit-package (rstatsToolkit),  
[15](#)

smoothPlot, [15](#)  
smoothTimePlot, [16](#)  
summarySE, [17](#)

themeWhite, [18](#)

unlistAndFilterIndep, [18](#)