

Package ‘explore’

August 27, 2019

Type Package

Title Simplifies Exploratory Data Analysis

Version 0.4.4

Author Roland Krasser

Maintainer Roland Krasser <roland.krasser@gmail.com>

Description Interactive data exploration with one line of code or use an easy to remember set of tidy functions for exploratory data analysis. Introduces two main verbs. describe() to describe a variable or table, explore() to graphically explore a variable or table.

License GPL-3

Encoding UTF-8

LazyData true

URL <http://github.com/rolkra/explore>

Imports broom, dplyr, DBI, DT, forcats, ggplot2 (>= 3.0.0), gridExtra, magrittr, MASS, odbc, rlang, rpart, rpart.plot, shiny, rmarkdown

RoxygenNote 6.1.1

Suggests knitr

VignetteBuilder knitr

NeedsCompilation no

Repository CRAN

Date/Publication 2019-08-27 12:00:12 UTC

R topics documented:

balance_target	2
clean_var	3
data_dict_md	4
decrypt	5
describe	5
describe_all	6

describe_cat	7
describe_num	7
describe_tbl	8
dwh_connect	9
dwh_disconnect	9
dwh_fastload	10
dwh_read_data	11
dwh_read_table	11
encrypt	12
explain_logreg	13
explain_tree	13
explore	14
explore_all	15
explore_bar	16
explore_cor	17
explore_density	18
explore_shiny	19
explore_tbl	19
format_num_kMB	20
format_num_space	20
format_target	21
format_type	21
get_nrow	22
get_type	22
guess_cat_num	23
plot_text	23
plot_var_info	24
replace_na_with	24
report	25
target_explore_cat	25
target_explore_num	26

Index	28
--------------	-----------

balance_target	<i>Balance target variable</i>
----------------	--------------------------------

Description

Balances the target variable in your dataset. Target must be 0/1, FALSE/TRUE ore no/yes

Usage

```
balance_target(data, target, min_prop = 0.1)
```

Arguments

data	A dataset
target	Target variable (0/1, TRUE/FALSE, yes/no)
min_prop	Minimum proportion of one of the target categories

Value

Data

Examples

```
iris$is_versicolor <- ifelse(iris$Species == "versicolor", 1, 0)
balanced <- balance_target(iris, target = is_versicolor, min_prop = 0.5)
describe(balanced, is_versicolor)
```

clean_var	<i>Clean variable</i>
-----------	-----------------------

Description

Clean variable (replace NA values, set min_val and max_val)

Usage

```
clean_var(data, var, na = NA, min_val = NA, max_val = NA,
          max_cat = NA, name = NA)
```

Arguments

data	A dataset
var	Name of variable
na	Value that replaces NA
min_val	All values < min_val are converted to min_val (var numeric or character)
max_val	All values > max_val are converted to max_val (var numeric or character)
max_cat	Maximum number of different factor levels for categorical variable (if more, .OTHER is added)
name	New name of variable (as string)

Value

Dataset

Examples

```
clean_var(iris, Sepal.Width, max_val = 3.5, name = "sepal_width")
```

data_dict_md	<i>Create a data dictionary Markdown file</i>
--------------	---

Description

Create a data dictionary Markdown file

Usage

```
data_dict_md(data, title = "", description = NA,  
             output_file = "data_dict.md", output_dir)
```

Arguments

data	A dataframe (data dictionary for all variables)
title	Title of the data dictionary
description	Detailed description of variables in data (dataframe with columns 'variable' and 'description')
output_file	Output filename for Markdown file
output_dir	Directory where the Markdown file is saved

Value

Create Markdown file

Examples

```
# Data dictionary of a dataframe  
data_dict_md(iris,  
             title = "iris flower data set",  
             output_dir = tempdir())  
  
# Data dictionary of a dataframe with additional description of variables  
description <- data.frame(  
  variable = c("Species"),  
  description = c("Species of Iris flower"))  
data_dict_md(iris,  
             title = "iris flower data set",  
             description = description,  
             output_dir = tempdir())
```

decrypt	<i>decrypt text</i>
---------	---------------------

Description

decrypt text

Usage

```
decrypt(text, codeletters = c(toupper(letters), letters, 0:9),
        shift = 18)
```

Arguments

text	A text (character)
codeletters	A string of letters that are used for decryption
shift	Number of elements shifted

Value

Decrypted text

Examples

```
decrypt("zw336 E693v")
```

describe	<i>Describe a dataset or variable</i>
----------	---------------------------------------

Description

Describe a dataset or variable (depending on input parameters)

Usage

```
describe(data, var, target, out = "text", ...)
```

Arguments

data	A dataset
var	A variable of the dataset
target	Target variable (0/1 or FALSE/TRUE)
out	Output format ("text" "list") of variable description
...	Further arguments

Value

Description as table, text or list

Examples

```
# Load package
library(magrittr)

# Describe a dataset
iris %>% describe()

# Describe a variable
iris %>% describe(Species)
iris %>% describe(Sepal.Length)
```

describe_all	<i>Describe all variables of a dataset</i>
--------------	--

Description

Describe all variables of a dataset

Usage

```
describe_all(data = NA, out = "large")
```

Arguments

data	A dataset
out	Output format ("small" "large")

Value

Dataset

Examples

```
describe_all(iris)
```

describe_cat	<i>Describe categorical variable</i>
--------------	--------------------------------------

Description

Describe categorical variable

Usage

```
describe_cat(data, var, max_cat = 10, out = "text", margin = 0)
```

Arguments

data	A dataset
var	Variable or variable name
max_cat	Maximum number of categories displayed
out	Output format ("text" "list")
margin	Left margin for text output (number of spaces)

Value

Description as text or list

Examples

```
describe_cat(iris, Species)
```

describe_num	<i>Describe numerical variable</i>
--------------	------------------------------------

Description

Describe numerical variable

Usage

```
describe_num(data, var, out = "text", margin = 0)
```

Arguments

data	A dataset
var	Variable or variable name
out	Output format ("text" "list")
margin	Left margin for text output (number of spaces)

Value

Description as text or list

Examples

```
describe_num(iris, Sepal.Length)
```

describe_tbl	<i>Describe table</i>
--------------	-----------------------

Description

Describe table (e.g. number of rows and columns of dataset)

Usage

```
describe_tbl(data, target, out = "text")
```

Arguments

data	A dataset
target	Target variable (binary)
out	Output format ("text" "list")

Value

Description as text or list

Examples

```
describe_tbl(iris)

iris$is_virginica <- ifelse(iris$Species == "virginica", 1, 0)
describe_tbl(iris, is_virginica)
```

dwh_connect	<i>connect to DWH</i>
-------------	-----------------------

Description

connect to datawarehouse (DWH) using ODBC

Usage

```
dwh_connect(dsn, user = NA, pwd = NA, pwd_crypt = FALSE, ...)
```

Arguments

dsn	DSN string
user	user name
pwd	password of user
pwd_crypt	is password encryption used?
...	Further arguments to be passed to DBI::dbConnect()

Value

connection

Examples

```
## Not run:  
con <- dwh_connect(dsn = "DWH1", user = "u12345")  
  
## End(Not run)
```

dwh_disconnect	<i>disconnect from DWH</i>
----------------	----------------------------

Description

disconnect from datawarehouse (DWH) using a ODBC connection

Usage

```
dwh_disconnect(connection, ...)
```

Arguments

connection	channel (ODBC connection)
...	Further arguments to be passed to DBI::dbDisconnect()

Examples

```
## Not run:  
dwh_disconnect(con)  
  
## End(Not run)
```

dwh_fastload	<i>write data to a DWH table</i>
--------------	----------------------------------

Description

write data fast to a DWH table using a ODBC connection Function uses packages DBI/odbc to write data faster than RODBC Connects, writes data and disconnects

Usage

```
dwh_fastload(data, dsn, table, ...)
```

Arguments

data	dataframe
dsn	DSN string
table	table name (character string)
...	Further arguments to be passed to DBI::dbConnect()

Value

status

Examples

```
## Not run:  
dwh_fastload(data, "DWH", "database.table_test")  
  
## End(Not run)
```

dwh_read_data	<i>read data from DWH</i>
---------------	---------------------------

Description

read data from DWH using a ODBC connection

Usage

```
dwh_read_data(connection, sql, names_lower = TRUE, ...)
```

Arguments

connection	DWH connection
sql	sql (character string)
names_lower	convert field names to lower (default = TRUE)
...	Further arguments to be passed to DBI::dbGetQuery()

Value

dataframe containing table data

Examples

```
## Not run:
dwh_read_data(con, "select * from database.table_test")

## End(Not run)
```

dwh_read_table	<i>read a table from DWH</i>
----------------	------------------------------

Description

read a table from DWH using a ODBC connection

Usage

```
dwh_read_table(connection, table, names_lower = TRUE, ...)
```

Arguments

connection	DWH connection
table	table name (character string)
names_lower	convert field names to lower (default = TRUE)
...	Further arguments to be passed to DBI::dbGetQuery()

Value

dataframe containing table data

Examples

```
## Not run:  
dwh_read_table(con, "database.table_test")  
  
## End(Not run)
```

encrypt

encrypt text

Description

encrypt text

Usage

```
encrypt(text, codeletters = c(toupper(letters), letters, 0:9),  
        shift = 18)
```

Arguments

text	A text (character)
codeletters	A string of letters that are used for encryption
shift	Number of elements shifted

Value

Encrypted text

Examples

```
encrypt("hello world")
```

explain_logreg	<i>Explain a binary target using logistic regression</i>
----------------	--

Description

Explain a binary target using logistic regression

Usage

```
explain_logreg(data, target, ...)
```

Arguments

data	A dataset
target	Target variable (binary)
...	Further arguments

Value

Dataset with results (term, estimate, std.error, z.value, p.value)

Examples

```
data <- iris
data$is_versicolor <- ifelse(iris$Species == "versicolor", 1, 0)
data$Species <- NULL
explain_logreg(data, target = is_versicolor)
```

explain_tree	<i>Explain a target using a simple decision tree (classification or regression)</i>
--------------	---

Description

Explain a target using a simple decision tree (classification or regression)

Usage

```
explain_tree(data, target, max_cat = 10, maxdepth = 3, minsplit = 20,
  cp = 0, size = 0.7, ...)
```

Arguments

<code>data</code>	A dataset
<code>target</code>	Target variable
<code>max_cat</code>	Drop categorical variables with higher number of levels
<code>maxdepth</code>	Maximal depth of the tree (rpart-parameter)
<code>minsplit</code>	The minimum number of observations that must exist in a node in order for a split to be attempted (rpart-parameter)
<code>cp</code>	Complexity parameter (rpart-parameter)
<code>size</code>	Textsize of plot
<code>...</code>	Further arguments

Value

Plot

Examples

```
data <- iris
data$is_versicolor <- ifelse(iris$Species == "versicolor", 1, 0)
data$Species <- NULL
explain_tree(data, target = is_versicolor)
```

explore

Explore a dataset or variable

Description

Explore a dataset or variable

Usage

```
explore(data, var, var2, target, split, min_val = NA, max_val = NA,
        auto_scale = TRUE, na = NA, ...)
```

Arguments

<code>data</code>	A dataset
<code>var</code>	A variable
<code>var2</code>	A variable for checking correlation
<code>target</code>	Target variable (0/1 or FALSE/TRUE)
<code>split</code>	Split by target variable (FALSE/TRUE)
<code>min_val</code>	All values < min_val are converted to min_val
<code>max_val</code>	All values > max_val are converted to max_val

auto_scale	Use 0.2 and 0.98 quantile for min_val and max_val (if min_val and max_val are not defined)
na	Value to replace NA
...	Further arguments

Value

Plot object

Examples

```
## Launch Shiny app (in interactive R sessions)
if (interactive()) {
  explore(iris)
}

## Explore grafically

# Load library
library(magrittr)

# Explore a variable
iris %>% explore(Species)
iris %>% explore(Sepal.Length)
iris %>% explore(Sepal.Length, min_val = 4, max_val = 7)

# Explore a variable with a target
iris$is_virginica <- ifelse(iris$Species == "virginica", 1, 0)
iris %>% explore(Species, target = is_virginica)
iris %>% explore(Sepal.Length, target = is_virginica)

# Explore correlation between two variables
iris %>% explore(Species, Petal.Length)
iris %>% explore(Sepal.Length, Petal.Length)

# Explore correlation between two variables and split by target
iris %>% explore(Sepal.Length, Petal.Length, target = is_virginica)
```

explore_all

Explore all variables

Description

Explore all variables of a dataset (create plots)

Usage

```
explore_all(data, target, ncol = 2, split = TRUE)
```

Arguments

data	A dataset
target	Target variable (0/1 or FALSE/TRUE)
ncol	Layout of plots (number of columns)
split	Split by target (TRUE FALSE)

Value

Plot

Examples

```
explore_all(iris)

iris$is_virginica <- ifelse(iris$Species == "virginica", 1, 0)
explore_all(iris, target = is_virginica)
```

explore_bar	<i>Explore categorical variable using bar charts</i>
-------------	--

Description

Create a barplot to explore a categorical variable. If a target is selected, the barplot is created for all levels of the target.

Usage

```
explore_bar(data, var, target, flip = TRUE, title = "", max_cat = 30,
  max_target_cat = 5, legend_position = "right", label,
  label_size = 2.7)
```

Arguments

data	A dataset
var	variable
target	target (can have more than 2 levels)
flip	Should plot be flipped? (change of x and y)
title	Title of the plot (if empty var name)
max_cat	Maximum number of categories to be plotted
max_target_cat	Maximum number of categories to be plotted for target (except NA)
legend_position	Position of the legend ("bottom" "top" "none")
label	Show labels? (if empty, automatic)
label_size	Size of labels
...	Further arguments

Value

Plot object (bar chart)

explore_cor

Explore the correlation between two variables

Description

Explore the correlation between two variables

Usage

```
explore_cor(data, x, y, target, bins = 8, min_val = NA, max_val = NA,  
            auto_scale = TRUE, color = "grey", ...)
```

Arguments

data	A dataset
x	Variable on x axis
y	Variable on y axis
target	Target variable (categorical)
bins	Number of bins
min_val	All values < min_val are converted to min_val
max_val	All values > max_val are converted to max_val
auto_scale	Use 0.2 and 0.98 quantile for min_val and max_val (if min_val and max_val are not defined)
color	Color of the plot
...	Further arguments

Value

Plot

Examples

```
explore_cor(iris, x = Sepal.Length, y = Sepal.Width)
```

explore_density	<i>Explore density of variable</i>
-----------------	------------------------------------

Description

Create a density plot to explore numerical variable

Usage

```
explore_density(data, var, target, min_val = NA, max_val = NA,  
  color = "grey", auto_scale = TRUE, max_target_cat = 5, ...)
```

Arguments

data	A dataset
var	Variable
target	Target variable (0/1 or FALSE/TRUE)
min_val	All values < min_val are converted to min_val
max_val	All values > max_val are converted to max_val
color	Color of plot
auto_scale	Use 0.02 and 0.98 percent quantile for min_val and max_val (if min_val and max_val are not defined)
max_target_cat	Maximum number of levels of target shown in the plot (except NA).
...	Further arguments

Value

Plot object (density plot)

Examples

```
explore_density(iris, "Sepal.Length")  
iris$is_virginica <- ifelse(iris$Species == "virginica", 1, 0)  
explore_density(iris, Sepal.Length, target = is_virginica)
```

explore_shiny	<i>Explore dataset interactive</i>
---------------	------------------------------------

Description

Launches a shiny app to explore a dataset

Usage

```
explore_shiny(data, target)
```

Arguments

data	A dataset
target	Target variable (0/1 or FALSE/TRUE)

Examples

```
# Only run examples in interactive R sessions
if (interactive()) {
  explore_shiny(iris)
}
```

explore_tbl	<i>Explore table</i>
-------------	----------------------

Description

Explore a table. Plots variable types, variables with no variance and variables with NA

Usage

```
explore_tbl(data)
```

Arguments

data	A dataset
------	-----------

Examples

```
explore_tbl(iris)
```

format_num_kMB	<i>Format number</i>
----------------	----------------------

Description

Formats a big number as k (1 000), M (1 000 000) or B (1 000 000 000)

Usage

```
format_num_kMB(number = 0, digits = 1)
```

Arguments

number	A number (integer or real)
digits	Number of digits

Value

Formatted number as text

Examples

```
format_num_kMB(5500, digits = 2)
```

format_num_space	<i>Format number</i>
------------------	----------------------

Description

Formats a big number using space as big.mark (1000 = 1 000)

Usage

```
format_num_space(number = 0, digits = 1)
```

Arguments

number	A number (integer or real)
digits	Number of digits

Value

Formatted number as text

Examples

```
format_num_space(5500, digits = 2)
```

format_target	<i>Format target</i>
---------------	----------------------

Description

Formats a target as a 0/1 variable. If target is numeric, 1 = above average.

Usage

```
format_target(target)
```

Arguments

target Variable as vector

Value

Formatted target

Examples

```
iris$is_virginica <- ifelse(iris$Species == "virginica", "yes", "no")
iris$target <- format_target(iris$is_virginica)
table(iris$target)
```

format_type	<i>Format type description</i>
-------------	--------------------------------

Description

Format type description of variable to 3 letters (intdblglchrdat)

Usage

```
format_type(type)
```

Arguments

type Type description ("integer", "double", "logical", "character", "date")

Value

Formatted type description (intdblglchrdat)

Examples

```
format_type(typeof(iris$Species))
```

get_nrow	<i>Get number of rows for a grid plot</i>
----------	---

Description

Get number of rows for a grid plot

Usage

```
get_nrow(varnames, exclude = 0, ncol = 2)
```

Arguments

varnames	List of variables to be plotted
exclude	Number of variables that will be excluded from plot
ncol	Number of columns (default = 2)

Value

Number of rows

Examples

```
get_nrow(names(iris), ncol = 2)
```

get_type	<i>Return type of variable</i>
----------	--------------------------------

Description

Return value of typeof, except if variable contains <hide>, then return "other"

Usage

```
get_type(var)
```

Arguments

var	A vector (dataframe column)
-----	-----------------------------

Value

Value of typeof or "other"

Examples

```
get_type(iris$Species)
```

guess_cat_num	<i>Return if variable is categorial or nomerical</i>
---------------	--

Description

Guess if variable is categorial or numerical based on name, type and values of variable

Usage

```
guess_cat_num(var)
```

Arguments

var	A vector (dataframe column)
-----	-----------------------------

Value

"cat" (categorial), "num" (numerical) or "oth" (other)

Examples

```
guess_cat_num(iris$Species)
```

plot_text	<i>Plot a text</i>
-----------	--------------------

Description

Plots a text (base plot) and let you choose text-size and color

Usage

```
plot_text(text = "hello world", size = 1.2, color = "black")
```

Arguments

text	Text as string
size	Text-size
color	Text-color

Value

Plot

Examples

```
plot_text("hello", size = 2, color = "red")
```

plot_var_info	<i>Plot a variable info</i>
---------------	-----------------------------

Description

Creates a ggplot with the variable-name as title and a text

Usage

```
plot_var_info(data, var, info = "")
```

Arguments

data	A dataset
var	Variable
info	Text to plot

Value

Plot (ggplot)

replace_na_with	<i>Replace NA</i>
-----------------	-------------------

Description

Replace NA values of a variable in a dataframe

Usage

```
replace_na_with(data, var_name, with)
```

Arguments

data	A dataframe
var_name	Name of variable where NAs are replaced
with	Value instead of NA

Value

Updated dataframe

Examples

```
data <- data.frame(nr = c(1,2,3,NA,NA))  
replace_na_with(data, "nr", 0)
```

report	<i>Generate a report of all variables</i>
--------	---

Description

Generate a report of all variables. If target is defined, the relation to the target is reported.

Usage

```
report(data, target, split = TRUE, output_file, output_dir)
```

Arguments

data	A dataset
target	Target variable (0/1 or FALSE/TRUE)
split	Split by target? (TRUE/FALSE)
output_file	Filename of the html report
output_dir	Directory where to save the html report

Examples

```
if (rmarkdown::pandoc_available("1.12.3")) {
  report(iris, output_dir = tempdir())
}
```

target_explore_cat	<i>Explore categorical variable + target</i>
--------------------	--

Description

Create a plot to explore relation between categorical variable and a binary target.

Usage

```
target_explore_cat(data, var, target = "target_ind", min_val = NA,
  max_val = NA, flip = TRUE, num2char = TRUE, title = NA,
  auto_scale = TRUE, na = NA, max_cat = 30,
  legend_position = "bottom")
```

Arguments

data	A dataset
var	Categorical variable
target	Target variable (0/1 or FALSE/TRUE)
min_val	All values < min_val are converted to min_val
max_val	All values > max_val are converted to max_val
flip	Should plot be flipped? (change of x and y)
num2char	If TRUE, numeric values in variable are converted into character
title	Title of plot
auto_scale	Not used, just for compatibility
na	Value to replace NA
max_cat	Maximum numbers of categories to be plotted
legend_position	Position of legend ("right" "bottom" "non")

Value

Plot object

target_explore_num	<i>Explore categorical variable + target</i>
--------------------	--

Description

Create a plot to explore relation between numerical variable and a binary target

Usage

```
target_explore_num(data, var, target = "target_ind", min_val = NA,
  max_val = NA, flip = TRUE, title = NA, auto_scale = TRUE,
  na = NA, legend_position = "bottom")
```

Arguments

data	A dataset
var	Numerical variable
target	Target variable (0/1 or FALSE/TRUE)
min_val	All values < min_val are converted to min_val
max_val	All values > max_val are converted to max_val
flip	Should plot be flipped? (change of x and y)
title	Title of plot

target_explore_num

27

<code>auto_scale</code>	Use 0.02 and 0.98 quantile for <code>min_val</code> and <code>max_val</code> (if <code>min_val</code> and <code>max_val</code> are not defined)
<code>na</code>	Value to replace NA
<code>legend_position</code>	Position of legend ("right" "bottom" "non")

Value

Plot object

Index

balance_target, [2](#)

clean_var, [3](#)

data_dict_md, [4](#)
decrypt, [5](#)
describe, [5](#)
describe_all, [6](#)
describe_cat, [7](#)
describe_num, [7](#)
describe_tbl, [8](#)
dwh_connect, [9](#)
dwh_disconnect, [9](#)
dwh_fastload, [10](#)
dwh_read_data, [11](#)
dwh_read_table, [11](#)

encrypt, [12](#)
explain_logreg, [13](#)
explain_tree, [13](#)
explore, [14](#)
explore_all, [15](#)
explore_bar, [16](#)
explore_cor, [17](#)
explore_density, [18](#)
explore_shiny, [19](#)
explore_tbl, [19](#)

format_num_kMB, [20](#)
format_num_space, [20](#)
format_target, [21](#)
format_type, [21](#)

get_nrow, [22](#)
get_type, [22](#)
guess_cat_num, [23](#)

plot_text, [23](#)
plot_var_info, [24](#)

replace_na_with, [24](#)

report, [25](#)

target_explore_cat, [25](#)
target_explore_num, [26](#)