

Package ‘dlookr’

November 17, 2020

Type Package

Title Tools for Data Diagnosis, Exploration, Transformation

Version 0.3.14

Description A collection of tools that support data diagnosis, exploration, and transformation. Data diagnostics provides information and visualization of missing values and outliers and unique and negative values to help you understand the distribution and quality of your data. Data exploration provides information and visualization of the descriptive statistics of univariate variables, normality tests and outliers, correlation of two variables, and relationship between target variable and predictor. Data transformation supports binning for categorizing continuous variables, imputates missing values and outliers, resolving skewness. And it creates automated reports that support these three tasks.

License GPL-2 | file LICENSE

Encoding UTF-8

LazyData true

Depends R (>= 3.2.0), mice

Imports dplyr, magrittr, tidyr, ggplot2, RcmdrMisc, corrplot, rlang, purrr, tibble, tidyselect, classInt, kableExtra, prettydoc, smbinning, xtable, knitr, rmarkdown, RColorBrewer, gridExtra, tinytex, methods, DMwR, rpart, cli, car, broom, forcats, reshape2

Suggests ISLR, nycflights13, randomForest, dbplyr, DBI, RSQLite, stringr, testthat

Author Choonghyun Ryu [aut, cre]

Maintainer Choonghyun Ryu <choonghyun.ryu@gmail.com>

BugReports <https://github.com/choonghyunryu/dlookr/issues>

VignetteBuilder knitr

RoxygenNote 7.1.1

NeedsCompilation no

Repository CRAN

Date/Publication 2020-11-17 08:20:22 UTC

R topics documented:

dlookr-package	3
binning	4
binning_by	6
compare_category	8
compare_numeric	10
correlate	13
correlate.tbl_dbi	16
describe	18
describe.tbl_dbi	20
diagnose	23
diagnose.tbl_dbi	25
diagnose_category	27
diagnose_category.tbl_dbi	29
diagnose_numeric	31
diagnose_numeric.tbl_dbi	33
diagnose_outlier	35
diagnose_outlier.tbl_dbi	37
diagnose_report	39
diagnose_report.tbl_dbi	41
eda_report	44
eda_report.tbl_dbi	46
find_class	49
find_na	51
find_outliers	52
find_skewness	53
get_class	54
get_column_info	55
get_os	56
imputate_na	57
imputate_outlier	59
kurtosis	60
normality	61
normality.tbl_dbi	63
plot.bins	66
plot.compare_category	67
plot.compare_numeric	68
plot.imputation	69
plot.optimal_bins	70
plot.relate	71
plot.transform	73
plot.univar_category	74
plot.univar_numeric	75
plot_correlate	77
plot_correlate.tbl_dbi	79
plot_na_hclust	81
plot_na_intersect	82

plot_na_pareto	83
plot_normality	85
plot_normality.tbl_dbi	87
plot_outlier	89
plot_outlier.tbl_dbi	90
print.relate	93
relate	94
skewness	97
summary.bins	97
summary.compare_category	99
summary.compare_numeric	101
summary.imputation	104
summary.transform	105
summary.univar_category	106
summary.univar_numeric	108
target_by	110
target_by.tbl_dbi	111
transform	113
transformation_report	115
univar_category	116
univar_numeric	118

Index	122
--------------	------------

dlookr-package	<i>dlookr: Tools for Data Diagnosis, Exploration, Transformation</i>
----------------	--

Description

dlookr provides data diagnosis, data exploration and transformation of variables during data analysis.

Details

It has three main goals:

- When data is acquired, it is possible to judge whether data is erroneous or to select a variable to be corrected or removed through data diagnosis.
- Understand the distribution of data in the EDA process. It can also understand the relationship between target variables and predictor variables for the prediction model.
- Imputes missing value and outlier to standardization and resolving skewness. And, To convert a continuous variable to a categorical variable, bin the continuous variables.

To learn more about dlookr, start with the vignettes: `'browseVignettes(package = "dlookr")'`

Author(s)

Maintainer: Choonghyun Ryu <choonghyun.ryu@gmail.com>

See Also

Useful links:

- Report bugs at <https://github.com/choonghyunryu/dlookr/issues>

 binning

Binning the Numeric Data

Description

The `binning()` converts a numeric variable to a categorization variable.

Usage

```
binning(
  x,
  nbins,
  type = c("quantile", "equal", "pretty", "kmeans", "bclust"),
  ordered = TRUE,
  labels = NULL,
  approxy.lab = TRUE
)
```

Arguments

<code>x</code>	numeric. numeric vector for binning.
<code>nbins</code>	integer. number of intervals(bins). required. if missing, <code>nclass.Sturges</code> is used.
<code>type</code>	character. binning method. Choose from "quantile", "equal", "equal", "pretty", "kmeans" and "bclust". The "quantile" sets breaks with quantiles of the same interval. The "equal" sets breaks at the same interval. The "pretty" chooses a number of breaks not necessarily equal to <code>nbins</code> using <code>base::pretty</code> function. The "kmeans" uses <code>stats::kmeans</code> function to generate the breaks. The "bclust" uses <code>e1071::bclust</code> function to generate the breaks using bagged clustering. "kmeans" and "bclust" was implemented by <code>classInt::classIntervals</code> function.
<code>ordered</code>	logical. whether to build an ordered factor or not.
<code>labels</code>	character. the label names to use for each of the bins.
<code>approxy.lab</code>	logical. If TRUE, large number breaks are approximated to pretty numbers. If FALSE, the original breaks obtained by <code>type</code> are used.

Details

This function is useful when used with the `mutate/transmute` function of the `dplyr` package.

Value

An object of bins class. Attributes of bins class is as follows.

- type : binning type, "quantile", "equal", "pretty", "kmeans", "bclust".
- breaks : the number of intervals into which x is to be cut.
- levels : levels of binned value.
- raw : raw data, numeric vector corresponding to x argument.

"bins" class attributes information

Attributes of the "bins" class that is as follows.

- class : "bins"
- levels : levels of factor or ordered factor
- type : binning method
- breaks : breaks for binning
- raw : raw data before binning

See vignette("transformation") for an introduction to these concepts.

See Also

[binning_by](#), [print.bins](#), [summary.bins](#).

Examples

```
# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA
# Binning the carat variable. default type argument is "quantile"
bin <- binning(carseats$Income)
# Print bins class object
bin
# Summarise bins class object
summary(bin)
# Plot bins class object
plot(bin)
# Using labels argument
bin <- binning(carseats$Income, nbins = 4,
               labels = c("LQ1", "UQ1", "LQ3", "UQ3"))
bin
# Using another type argument
bin <- binning(carseats$Income, nbins = 5, type = "equal")
bin
bin <- binning(carseats$Income, nbins = 5, type = "pretty")
bin
bin <- binning(carseats$Income, nbins = 5, type = "kmeans")
bin
```

```

bin <- binning(carseats$Income, nbins = 5, type = "bclust")
bin

x <- sample(1:1000, size = 50) * 12345679
bin <- binning(x)
bin
bin <- binning(x, approxy.lab = FALSE)
bin

# -----
# Using pipes & dplyr
# -----
library(dplyr)

carseats %>%
  mutate(Income_bin = binning(carseats$Income)) %>%
  group_by(ShelveLoc, Income_bin) %>%
  summarise(freq = n()) %>%
  arrange(desc(freq)) %>%
  head(10)

```

binning_by

Optimal Binning for Scoring Modeling

Description

The `binning_by()` finding intervals for numerical variable using optical binning. Optimal binning categorizes a numeric characteristic into bins for ulterior usage in scoring modeling.

Usage

```
binning_by(df, y, x, p = 0.05, ordered = TRUE, labels = NULL)
```

Arguments

<code>df</code>	a data frame.
<code>y</code>	character. name of binary response variable(0, 1). The variable must contain only the integers 0 and 1 as element. However, in the case of factor having two levels, it is performed while type conversion is performed in the calculation process. It does not support that the variable name is "default" and that the dot is included in the variable name.
<code>x</code>	character. name of continuous characteristic variable. At least 5 different values. Inf is not allowed. It does not support that the variable name that the dot is included in the variable name.
<code>p</code>	numeric. percentage of records per bin. Default 5% (0.05). This parameter only accepts values greater than 0.00 (0%) and lower than 0.50 (50%).
<code>ordered</code>	logical. whether to build an ordered factor or not.
<code>labels</code>	character. the label names to use for each of the bins.

Details

This function is useful when used with the mutate/transmute function of the dplyr package. And this function is implemented using sbinning() function of sbinning package.

Value

an object of "optimal_bins" class. Attributes of "optimal_bins" class is as follows.

- class : "optimal_bins".
- type : binning type, "optimal".
- breaks : numeric. the number of intervals into which x is to be cut.
- levels : character. levels of binned value.
- raw : numeric. raw data, x argument value.
- ivtable : data.frame. information value table
- iv : numeric. information value
- target : integer. binary response variable

attributes of "optimal_bins" class

Attributes of the "optimal_bins" class that is as follows.

- class : "optimal_bins".
- levels : character. factor or ordered factor levels
- type : character. binning method
- breaks : numeric. breaks for binning
- raw : numeric. before the binned the raw data
- ivtable : data.frame. information value table
- iv : numeric. information value
- target : integer. binary response variable

See vignette("transformation") for an introduction to these concepts.

See Also

[binning](#), [sbinning](#).

Examples

```
# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

# optimal binning
bin <- binning_by(carseats, "US", "Advertising")
bin
```

```
# summary optimal_bins class
summary(bin)
# visualize optimal_bins class
plot(bin, sub = "bins of Advertising variable")
```

compare_category *Compare categorical variables*

Description

The `compare_category()` compute information to examine the relationship between categorical variables.

Usage

```
compare_category(.data, ...)

## S3 method for class 'data.frame'
compare_category(.data, ...)
```

Arguments

<code>.data</code>	a <code>data.frame</code> or a <code>tbl_df</code> .
<code>...</code>	one or more unquoted expressions separated by commas. You can treat variable names like they are positions. Positive values select variables; negative values to drop variables. These arguments are automatically quoted and evaluated in a context where column names represent column positions. They support unquoting and splicing.

Details

It is important to understand the relationship between categorical variables in EDA. `compare_category()` compares relations by pair combination of all categorical variables. and return `compare_category` class that based list object.

Value

An object of the class as compare based list. The information to examine the relationship between categorical variables is as follows each components.

- `var1` : factor. The level of the first variable to compare. 'var1' is the name of the first variable to be compared.
- `var2` : factor. The level of the second variable to compare. 'var2' is the name of the second variable to be compared.
- `n` : integer. frequency by var1 and var2.
- `rate` : double. relative frequency.
- `first_rate` : double. relative frequency in first variable.
- `second_rate` : double. relative frequency in second variable.

Attributes of return object

Attributes of compare_category class is as follows.

- variables : character. List of variables selected for comparison.
- combination : matrix. It consists of pairs of variables to compare.

See Also

[summary.compare_category](#), [print.compare_category](#), [plot.compare_category](#).

Examples

```
# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

library(dplyr)
library(stringr)

# Compare the all categorical variables
all_var <- compare_category(carseats)

# Print compare_numeric class object
all_var

# Compare the categorical variables that case of joint the US variable
all_var %>%
  "["(str_detect(names(all_var), "US"))

# Compare the two categorical variables
two_var <- compare_category(carseats, ShelveLoc, Urban)

# Print compare_numeric class object
two_var

# Filtering the case of US included NA
two_var %>%
  "["(1) %>%
  filter(!is.na(Urban))

# Summary the all case : Return a invisible copy of an object.
stat <- summary(all_var)

# Summary by returned object
stat

# component of table
stat$table

# component of chi-square test
stat$chisq
```

```

# component of chi-square test
summary(all_var, "chisq")

# component of chi-square test (first, third case)
summary(all_var, "chisq", pos = c(1, 3))

# component of relative frequency table
summary(all_var, "relative")

# component of table without missing values
summary(all_var, "table", na.rm = TRUE)

# component of table include marginal value
margin <- summary(all_var, "table", marginal = TRUE)
margin

# component of chi-square test
summary(two_var, method = "chisq")

# verbose is FALSE
summary(all_var, "chisq", verbose = FALSE)

#' # Using pipes & dplyr -----
# If you want to use dplyr, set verbose to FALSE
summary(all_var, "chisq", verbose = FALSE) %>%
  filter(p.value < 0.26)

# Extract component from list by index
summary(all_var, "table", na.rm = TRUE, verbose = FALSE) %>%
  "[["(1)

# Extract component from list by name
summary(all_var, "table", na.rm = TRUE, verbose = FALSE) %>%
  "[["("ShelveLoc vs Urban")

# plot all pair of variables
plot(all_var)

# plot a pair of variables
plot(two_var)

# plot all pair of variables by prompt
plot(all_var, prompt = TRUE)

# plot a pair of variables
plot(two_var, las = 1)

```

Description

The `compare_numeric()` compute information to examine the relationship between numerical variables.

Usage

```
compare_numeric(.data, ...)

## S3 method for class 'data.frame'
compare_numeric(.data, ...)
```

Arguments

`.data` a data.frame or a `tbl_df`.

`...` one or more unquoted expressions separated by commas. You can treat variable names like they are positions. Positive values select variables; negative values to drop variables. These arguments are automatically quoted and evaluated in a context where column names represent column positions. They support unquoting and splicing.

Details

It is important to understand the relationship between numerical variables in EDA. `compare_numeric()` compares relations by pair combination of all numerical variables. and return `compare_numeric` class that based list object.

Value

An object of the class as compare based list. The information to examine the relationship between numerical variables is as follows each components. - correlation component : Pearson's correlation coefficient.

- `var1` : factor. The level of the first variable to compare. 'var1' is the name of the first variable to be compared.
- `var2` : factor. The level of the second variable to compare. 'var2' is the name of the second variable to be compared.
- `coef_corr` : double. Pearson's correlation coefficient.

- linear component : linear model summaries

- `var1` : factor. The level of the first variable to compare. 'var1' is the name of the first variable to be compared.
- `var2` : factor. The level of the second variable to compare. 'var2' is the name of the second variable to be compared.
- `r.squared` : double. The percent of variance explained by the model.
- `adj.r.squared` : double. `r.squared` adjusted based on the degrees of freedom.
- `sigma` : double. The square root of the estimated residual variance.
- `statistic` : double. F-statistic.

- `p.value` : double. p-value from the F test, describing whether the full regression is significant.
- `df` : integer degrees of freedom.
- `logLik` : double. the log-likelihood of data under the model.
- `AIC` : double. the Akaike Information Criterion.
- `BIC` : double. the Bayesian Information Criterion.
- `deviance` : double. deviance.
- `df.residual` : integer residual degrees of freedom.

Attributes of return object

Attributes of `compare_numeric` class is as follows.

- `raw` : a data.frame or a `tbl_df`. Data containing variables to be compared. Save it for visualization with `plot.compare_numeric()`.
- `variables` : character. List of variables selected for comparison.
- `combination` : matrix. It consists of pairs of variables to compare.

See Also

[correlate](#), [summary.compare_numeric](#), [print.compare_numeric](#), [plot.compare_numeric](#).

Examples

```
# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

library(dplyr)

# Compare the all numerical variables
all_var <- compare_numeric(carseats)

# Print compare_numeric class object
all_var

# Compare the correlation that case of joint the Price variable
all_var %>%
  "$"(correlation) %>%
  filter(var1 == "Price" | var2 == "Price") %>%
  arrange(desc(abs(coef_corr)))

# Compare the correlation that case of abs(coef_corr) > 0.3
all_var %>%
  "$"(correlation) %>%
  filter(abs(coef_corr) > 0.3)

# Compare the linear model that case of joint the Price variable
all_var %>%
```

```

"$"(linear) %>%
  filter(var1 == "Price" | var2 == "Price") %>%
  arrange(desc(r.squared))

# Compare the two numerical variables
two_var <- compare_numeric(carseats, Price, CompPrice)

# Print compare_numeric class object
two_var

# Summary the all case : Return a invisible copy of an object.
stat <- summary(all_var)

# Just correlation
summary(all_var, method = "correlation")

# Just correlation condition by r > 0.2
summary(all_var, method = "correlation", thres_corr = 0.2)

# linear model summaries condition by R^2 > 0.05
summary(all_var, thres_rs = 0.05)

# verbose is FALSE
summary(all_var, verbose = FALSE)

# plot all pair of variables
plot(all_var)

# plot a pair of variables
plot(two_var)

# plot all pair of variables by prompt
plot(all_var, prompt = TRUE)

```

correlate

Compute the correlation coefficient between two numerical data

Description

The `correlate()` compute Pearson's the correlation coefficient of the numerical data.

Usage

```
correlate(.data, ...)
```

```
## S3 method for class 'data.frame'
```

```
correlate(.data, ..., method = c("pearson", "kendall", "spearman"))
```

Arguments

<code>.data</code>	a <code>data.frame</code> or a <code>tbl_df</code> .
<code>...</code>	one or more unquoted expressions separated by commas. You can treat variable names like they are positions. Positive values select variables; negative values to drop variables. If the first expression is negative, <code>correlate()</code> will automatically start with all variables. These arguments are automatically quoted and evaluated in a context where column names represent column positions. They support unquoting and splicing. See <code>vignette("EDA")</code> for an introduction to these concepts.
<code>method</code>	a character string indicating which correlation coefficient (or covariance) is to be computed. One of "pearson" (default), "kendall", or "spearman": can be abbreviated.

Details

This function is useful when used with the `group_by()` function of the `dplyr` package. If you want to compute by level of the categorical data you are interested in, rather than the whole observation, you can use `grouped_df` as the `group_by()` function. This function is computed `stats::cor()` function by use = "pairwise.complete.obs" option.

Correlation coefficient information

The information derived from the numerical data compute is as follows.

- `var1` : names of numerical variable
- `var2` : name of the corresponding numeric variable
- `coef_corr` : Pearson's correlation coefficient

See Also

`cor`, `correlate.tbl_dbi`.

Examples

```
# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

# Correlation coefficients of all numerical variables
correlate(carseats)

# Select the variable to compute
correlate(carseats, Sales, Price)
correlate(carseats, -Sales, -Price)
correlate(carseats, "Sales", "Price")
correlate(carseats, 1)
# Non-parametric correlation coefficient by kendall method
correlate(carseats, Sales, method = "kendall")
```

```

# Using dplyr::grouped_dt
library(dplyr)

gdata <- group_by(carseats, ShelveLoc, US)
correlate(gdata, "Sales")
correlate(gdata)

# Using pipes -----
# Correlation coefficients of all numerical variables
carseats %>%
  correlate()
# Positive values select variables
carseats %>%
  correlate(Sales, Price)
# Negative values to drop variables
carseats %>%
  correlate(-Sales, -Price)
# Positions values select variables
carseats %>%
  correlate(1)
# Positions values select variables
carseats %>%
  correlate(-1, -2, -3, -5, -6)
# Non-parametric correlation coefficient by spearman method
carseats %>%
  correlate(Sales, Price, method = "spearman")

# -----
# Correlation coefficient
# that eliminates redundant combination of variables
carseats %>%
  correlate() %>%
  filter(as.integer(var1) > as.integer(var2))

carseats %>%
  correlate(Sales, Price) %>%
  filter(as.integer(var1) > as.integer(var2))

# Using pipes & dplyr -----
# Compute the correlation coefficient of Sales variable by 'ShelveLoc'
# and 'US' variables. And extract only those with absolute
# value of correlation coefficient is greater than 0.5
carseats %>%
  group_by(ShelveLoc, US) %>%
  correlate(Sales) %>%
  filter(abs(coef_corr) >= 0.5)

# extract only those with 'ShelveLoc' variable level is "Good",
# and compute the correlation coefficient of 'Sales' variable
# by 'Urban' and 'US' variables.
# And the correlation coefficient is negative and smaller than 0.5
carseats %>%

```

```

filter(ShelveLoc == "Good") %>%
group_by(Urban, US) %>%
correlate(Sales) %>%
filter(coef_corr < 0) %>%
filter(abs(coef_corr) > 0.5)

```

correlate.tbl_dbi	<i>Compute the correlation coefficient between two numerical data</i>
-------------------	---

Description

The `correlate()` compute Pearson's the correlation coefficient of the numerical(INTEGER, NUMBER, etc.) column of the DBMS table through `tbl_dbi`.

Usage

```

## S3 method for class 'tbl_dbi'
correlate(
  .data,
  ...,
  in_database = FALSE,
  collect_size = Inf,
  method = c("pearson", "kendall", "spearman")
)

```

Arguments

<code>.data</code>	a <code>tbl_dbi</code> .
<code>...</code>	one or more unquoted expressions separated by commas. You can treat variable names like they are positions. Positive values select variables; negative values to drop variables. If the first expression is negative, <code>correlate()</code> will automatically start with all variables. These arguments are automatically quoted and evaluated in a context where column names represent column positions. They support unquoting and splicing.
<code>in_database</code>	Specifies whether to perform in-database operations. If <code>TRUE</code> , most operations are performed in the DBMS. if <code>FALSE</code> , table data is taken in R and operated in-memory. Not yet supported <code>in_database = TRUE</code> .
<code>collect_size</code>	a integer. The number of data samples from the DBMS to R. Applies only if <code>in_database = FALSE</code> .
<code>method</code>	a character string indicating which correlation coefficient (or covariance) is to be computed. One of "pearson" (default), "kendall", or "spearman": can be abbreviated. See <code>vignette("EDA")</code> for an introduction to these concepts.

Details

This function is useful when used with the `group_by()` function of the `dplyr` package. If you want to compute by level of the categorical data you are interested in, rather than the whole observation, you can use `grouped_df` as the `group_by()` function. This function is computed `stats::cor()` function by use = "pairwise.complete.obs" option.

Correlation coefficient information

The information derived from the numerical data compute is as follows.

- `var1` : names of numerical variable
- `var2` : name of the corresponding numeric variable
- `coef_corr` : Pearson's correlation coefficient

See Also

[correlate.data.frame](#), [cor](#).

Examples

```
library(dplyr)

# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

# connect DBMS
con_sqlite <- DBI::dbConnect(RSQLite::SQLite(), ":memory:")

# copy carseats to the DBMS with a table named TB_CARSEATS
copy_to(con_sqlite, carseats, name = "TB_CARSEATS", overwrite = TRUE)

# Using pipes -----
# Correlation coefficients of all numerical variables
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  correlate()

# Positive values select variables
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  correlate(Sales, Price)

# Negative values to drop variables, and In-memory mode and collect size is 200
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  correlate(-Sales, -Price, collect_size = 200)

# Positions values select variables
con_sqlite %>%
```

```

tbl("TB_CARSEATS") %>%
  correlate(1)

# Positions values select variables
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  correlate(-1, -2, -3, -5, -6)

# -----
# Correlation coefficient
# that eliminates redundant combination of variables
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  correlate() %>%
  filter(as.integer(var1) > as.integer(var2))

con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  correlate(Sales, Price) %>%
  filter(as.integer(var1) > as.integer(var2))

# Using pipes & dplyr -----
# Compute the correlation coefficient of Sales variable by 'ShelveLoc'
# and 'US' variables. And extract only those with absolute
# value of correlation coefficient is greater than 0.5
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  group_by(ShelveLoc, US) %>%
  correlate(Sales) %>%
  filter(abs(coef_corr) >= 0.5)

# extract only those with 'ShelveLoc' variable level is "Good",
# and compute the correlation coefficient of 'Sales' variable
# by 'Urban' and 'US' variables.
# And the correlation coefficient is negative and smaller than -0.5
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  filter(ShelveLoc == "Good") %>%
  group_by(Urban, US) %>%
  correlate(Sales) %>%
  filter(coef_corr < 0) %>%
  filter(abs(coef_corr) > 0.5)

```

describe

Compute descriptive statistic

Description

The describe() compute descriptive statistic of numeric variable for exploratory data analysis.

Usage

```
describe(.data, ...)

## S3 method for class 'data.frame'
describe(.data, ...)
```

Arguments

`.data` a `data.frame` or a `tbl_df`.

`...` one or more unquoted expressions separated by commas. You can treat variable names like they are positions. Positive values select variables; negative values to drop variables. If the first expression is negative, `describe()` will automatically start with all variables. These arguments are automatically quoted and evaluated in a context where column names represent column positions. They support unquoting and splicing.

See `vignette("EDA")` for an introduction to these concepts.

Details

This function is useful when used with the `group_by` function of the `dplyr` package. If you want to calculate the statistic by level of the categorical data you are interested in, rather than the whole statistic, you can use `grouped_df` as the `group_by()` function.

Value

An object of the same class as `.data`.

Descriptive statistic information

The information derived from the numerical data `describe` is as follows.

- `n` : number of observations excluding missing values
- `na` : number of missing values
- `mean` : arithmetic average
- `sd` : standard deviation
- `se_mean` : standard error mean. sd/\sqrt{n}
- `IQR` : interquartile range (Q3-Q1)
- `skewness` : skewness
- `kurtosis` : kurtosis
- `p25` : Q1. 25% percentile
- `p50` : median. 50% percentile
- `p75` : Q3. 75% percentile
- `p01`, `p05`, `p10`, `p20`, `p30` : 1%, 5%, 20%, 30% percentiles
- `p40`, `p60`, `p70`, `p80` : 40%, 60%, 70%, 80% percentiles
- `p90`, `p95`, `p99`, `p100` : 90%, 95%, 99%, 100% percentiles

See Also

[describe.tbl_dbi, diagnose_numeric.data.frame.](#)

Examples

```
# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

# Describe descriptive statistics of numerical variables
describe(carseats)

# Select the variable to describe
describe(carseats, Sales, Price)
describe(carseats, -Sales, -Price)
describe(carseats, 5)

# Using dplyr::grouped_dt
library(dplyr)

gdata <- group_by(carseats, ShelveLoc, US)
describe(gdata, "Income")

# Using pipes -----
# Positive values select variables
carseats %>%
  describe(Sales, CompPrice, Income)

# Negative values to drop variables
carseats %>%
  describe(-Sales, -CompPrice, -Income)

# Using pipes & dplyr -----
# Find the statistic of all numerical variables by 'ShelveLoc' and 'US',
# and extract only those with 'ShelveLoc' variable level is "Good".
carseats %>%
  group_by(ShelveLoc, US) %>%
  describe() %>%
  filter(ShelveLoc == "Good")

# extract only those with 'Urban' variable level is "Yes",
# and find 'Sales' statistics by 'ShelveLoc' and 'US'
carseats %>%
  filter(Urban == "Yes") %>%
  group_by(ShelveLoc, US) %>%
  describe(Sales)
```

Description

The describe() compute descriptive statistic of numerical(INTEGER, NUMBER, etc.) column of the DBMS table through tbl_dbi for exploratory data analysis.

Usage

```
## S3 method for class 'tbl_dbi'  
describe(.data, ..., in_database = FALSE, collect_size = Inf)
```

Arguments

.data	a tbl_dbi.
...	one or more unquoted expressions separated by commas. You can treat variable names like they are positions. Positive values select variables; negative values to drop variables. If the first expression is negative, describe() will automatically start with all variables. These arguments are automatically quoted and evaluated in a context where column names represent column positions. They support unquoting and splicing.
in_database	Specifies whether to perform in-database operations. If TRUE, most operations are performed in the DBMS. if FALSE, table data is taken in R and operated in-memory. Not yet supported in_database = TRUE.
collect_size	a integer. The number of data samples from the DBMS to R. Applies only if in_database = FALSE. See vignette("EDA") for an introduction to these concepts.

Details

This function is useful when used with the [group_by](#) function of the dplyr package. If you want to calculate the statistic by level of the categorical data you are interested in, rather than the whole statistic, you can use grouped_df as the group_by() function.

Value

An object of the same class as .data.

Descriptive statistic information

The information derived from the numerical data describe is as follows.

- n : number of observations excluding missing values
- na : number of missing values
- mean : arithmetic average
- sd : standard deviation
- se_mean : standard error mean. sd/\sqrt{n}
- IQR : interquartile range (Q3-Q1)
- skewness : skewness

- kurtosis : kurtosis
- p25 : Q1. 25% percentile
- p50 : median. 50% percentile
- p75 : Q3. 75% percentile
- p01, p05, p10, p20, p30 : 1%, 5%, 20%, 30% percentiles
- p40, p60, p70, p80 : 40%, 60%, 70%, 80% percentiles
- p90, p95, p99, p100 : 90%, 95%, 99%, 100% percentiles

See Also

[describe.data.frame](#), [diagnose_numeric.tbl_dbi](#).

Examples

```
library(dplyr)

# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

# connect DBMS
con_sqlite <- DBI::dbConnect(RSQLite::SQLite(), ":memory:")

# copy carseats to the DBMS with a table named TB_CARSEATS
copy_to(con_sqlite, carseats, name = "TB_CARSEATS", overwrite = TRUE)

# Using pipes -----
# Positive values select variables
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  describe(Sales, CompPrice, Income)

# Negative values to drop variables, and In-memory mode and collect size is 200
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  describe(-Sales, -CompPrice, -Income, collect_size = 200)

# Using pipes & dplyr -----
# Find the statistic of all numerical variables by 'ShelveLoc' and 'US',
# and extract only those with 'ShelveLoc' variable level is "Good".
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  group_by(ShelveLoc, US) %>%
  describe() %>%
  filter(ShelveLoc == "Good")

# extract only those with 'Urban' variable level is "Yes",
# and find 'Sales' statistics by 'ShelveLoc' and 'US'
con_sqlite %>%
```

```
tbl("TB_CARSEATS") %>%  
filter(Urban == "Yes") %>%  
group_by(ShelveLoc, US) %>%  
describe(Sales)
```

diagnose

Diagnose data quality of variables

Description

The `diagnose()` produces information for diagnosing the quality of the variables of `data.frame` or `tbl_df`.

Usage

```
diagnose(.data, ...)  
  
## S3 method for class 'data.frame'  
diagnose(.data, ...)
```

Arguments

`.data` a `data.frame` or a `tbl_df`.
`...` one or more unquoted expressions separated by commas. You can treat variable names like they are positions. Positive values select variables; negative values to drop variables. If the first expression is negative, `diagnose()` will automatically start with all variables. These arguments are automatically quoted and evaluated in a context where column names represent column positions. They support unquoting and splicing.

Details

The scope of data quality diagnosis is information on missing values and unique value information. Data quality diagnosis can determine variables that require missing value processing. Also, the unique value information can determine the variable to be removed from the data analysis.

Value

An object of `tbl_df`.

Diagnostic information

The information derived from the data diagnosis is as follows.:

- `variables` : variable names
- `types` : data type of the variable or to select a variable to be corrected or removed through data diagnosis.

- integer, numeric, factor, ordered, character, etc.
- `missing_count` : number of missing values
- `missing_percent` : percentage of missing values
- `unique_count` : number of unique values
- `unique_rate` : ratio of unique values. `unique_count / number of observation`

See vignette("diagnosis") for an introduction to these concepts.

See Also

[diagnose.tbl_dbi](#), [diagnose_category.data.frame](#), [diagnose_numeric.data.frame](#).

Examples

```
# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

# Diagnosis of all variables
diagnose(carseats)

# Select the variable to diagnose
diagnose(carseats, Sales, Income, Age)
diagnose(carseats, -Sales, -Income, -Age)
diagnose(carseats, "Sales", "Income", "Age")
diagnose(carseats, 1, 3, 8)

# Using pipes -----
library(dplyr)

# Diagnosis of all variables
carseats %>%
  diagnose()
# Positive values select variables
carseats %>%
  diagnose(Sales, Income, Age)
# Negative values to drop variables
carseats %>%
  diagnose(-Sales, -Income, -Age)
# Positions values select variables
carseats %>%
  diagnose(1, 3, 8)
# Positions values select variables
carseats %>%
  diagnose(-8, -9, -10)

# Using pipes & dplyr -----
# Diagnosis of missing variables
carseats %>%
  diagnose() %>%
  filter(missing_count > 0)
```

diagnose.tbl_dbi	<i>Diagnose data quality of variables in the DBMS</i>
------------------	---

Description

The `diagnose()` produces information for diagnosing the quality of the column of the DBMS table through `tbl_dbi`.

Usage

```
## S3 method for class 'tbl_dbi'
diagnose(.data, ..., in_database = TRUE, collect_size = Inf)
```

Arguments

<code>.data</code>	a <code>tbl_dbi</code> .
<code>...</code>	one or more unquoted expressions separated by commas. You can treat variable names like they are positions. Positive values select variables; negative values to drop variables. If the first expression is negative, <code>diagnose()</code> will automatically start with all variables. These arguments are automatically quoted and evaluated in a context where column names represent column positions. They support unquoting and splicing.
<code>in_database</code>	a logical. Specifies whether to perform in-database operations. If <code>TRUE</code> , most operations are performed in the DBMS. if <code>FALSE</code> , table data is taken in R and operated in-memory.
<code>collect_size</code>	a integer. The number of data samples from the DBMS to R. Applies only if <code>in_database = FALSE</code> .

Details

The scope of data quality diagnosis is information on missing values and unique value information. Data quality diagnosis can determine variables that require missing value processing. Also, the unique value information can determine the variable to be removed from the data analysis.

Value

An object of `tbl_df`.

Diagnostic information

The information derived from the data diagnosis is as follows.:

- `variables` : column names
- `types` : data type of the variable or to select a variable to be corrected or removed through data diagnosis.
 - integer, numeric, factor, ordered, character, etc.

- `missing_count` : number of missing values
- `missing_percent` : percentage of missing values
- `unique_count` : number of unique values
- `unique_rate` : ratio of unique values. `unique_count / number of observation`

See vignette("diagonosis") for an introduction to these concepts.

See Also

[diagnose.data.frame](#), [diagnose.category.tbl_dbi](#), [diagnose.numeric.tbl_dbi](#).

Examples

```
library(dplyr)

# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

# connect DBMS
con_sqlite <- DBI::dbConnect(RSQLite::SQLite(), ":memory:")

# copy carseats to the DBMS with a table named TB_CARSEATS
copy_to(con_sqlite, carseats, name = "TB_CARSEATS", overwrite = TRUE)

# Using pipes -----
# Diagnosis of all columns
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  diagnose()

# Positive values select columns
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  diagnose(Sales, Income, Age)

# Negative values to drop columns
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  diagnose(-Sales, -Income, -Age)

# Positions values select columns, and In-memory mode
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  diagnose(1, 3, 8, in_database = FALSE)

# Positions values select columns, and In-memory mode and collect size is 200
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  diagnose(-8, -9, -10, in_database = FALSE, collect_size = 200)
```

```
# Using pipes & dplyr -----
# Diagnosis of missing variables
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  diagnose() %>%
  filter(missing_count > 0)
```

diagnose_category *Diagnose data quality of categorical variables*

Description

The `diagnose_category()` produces information for diagnosing the quality of the variables of `data.frame` or `tbl_df`.

Usage

```
diagnose_category(.data, ...)

## S3 method for class 'data.frame'
diagnose_category(.data, ..., top = 10, add_character = TRUE)
```

Arguments

<code>.data</code>	a <code>data.frame</code> or a <code>tbl_df</code> .
<code>...</code>	one or more unquoted expressions separated by commas. You can treat variable names like they are positions. Positive values select variables; negative values to drop variables. If the first expression is negative, <code>diagnose_category()</code> will automatically start with all variables. These arguments are automatically quoted and evaluated in a context where column names represent column positions. They support unquoting and splicing.
<code>top</code>	an integer. Specifies the upper top rank to extract. Default is 10.
<code>add_character</code>	logical. Decide whether to include text variables in the diagnosis of categorical data. The default value is <code>TRUE</code> , which also includes character variables.

Details

The scope of the diagnosis is the occupancy status of the levels in categorical data. If a certain level of occupancy is close to 100 then the removal of this variable in the forecast model will have to be considered. Also, if the occupancy of all levels is close to 0 variable is likely to be an identifier.

Value

an object of `tbl_df`.

Categorical diagnostic information

The information derived from the categorical data diagnosis is as follows.

- variables : variable names
- levels: level names
- N : number of observation
- freq : number of observation at the levels
- ratio : percentage of observation at the levels
- rank : rank of occupancy ratio of levels

See vignette("diagonosis") for an introduction to these concepts.

See Also

[diagnose_category.tbl_dbi](#), [diagnose.data.frame](#), [diagnose_numeric.data.frame](#), [diagnose_outlier.data.frame](#)

Examples

```
# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

# Diagnosis of categorical variables
diagnose_category(carseats)

# Select the variable to diagnose
diagnose_category(carseats, ShelveLoc, Urban)
diagnose_category(carseats, -ShelveLoc, -Urban)
diagnose_category(carseats, "ShelveLoc", "Urban")
diagnose_category(carseats, 7)

# Using pipes -----
library(dplyr)

# Diagnosis of all categorical variables
carseats %>%
  diagnose_category()
# Positive values select variables
carseats %>%
  diagnose_category(Urban, US)
# Negative values to drop variables
carseats %>%
  diagnose_category(-Urban, -US)
# Positions values select variables
carseats %>%
  diagnose_category(7)
# Positions values select variables
carseats %>%
  diagnose_category(-7)
```

```

# Top rank levels with top argument
carseats %>%
  diagnose_category(top = 2)

# Using pipes & dplyr -----
# Extraction of level that is more than 60% of categorical data
carseats %>%
  diagnose_category() %>%
  filter(ratio >= 60)

```

diagnose_category.tbl_dbi

Diagnose data quality of categorical variables in the DBMS

Description

The `diagnose_category()` produces information for diagnosing the quality of the character (CHAR, VARCHAR, VARCHAR2, etc.) column of the DBMS table through `tbl_dbi`.

Usage

```

## S3 method for class 'tbl_dbi'
diagnose_category(.data, ..., top = 10, in_database = TRUE, collect_size = Inf)

```

Arguments

<code>.data</code>	a <code>tbl_dbi</code> .
<code>...</code>	one or more unquoted expressions separated by commas. You can treat variable names like they are positions. Positive values select variables; negative values to drop variables. If the first expression is negative, <code>diagnose_category()</code> will automatically start with all variables. These arguments are automatically quoted and evaluated in a context where column names represent column positions. They support unquoting and splicing.
<code>top</code>	an integer. Specifies the upper top rank to extract. Default is 10.
<code>in_database</code>	Specifies whether to perform in-database operations. If <code>TRUE</code> , most operations are performed in the DBMS. if <code>FALSE</code> , table data is taken in R and operated in-memory.
<code>collect_size</code>	a integer. The number of data samples from the DBMS to R. Applies only if <code>in_database = FALSE</code> .

Details

The scope of the diagnosis is the occupancy status of the levels in categorical data. If a certain level of occupancy is close to 100 then the removal of this variable in the forecast model will have to be considered. Also, if the occupancy of all levels is close to 0 variable is likely to be an identifier.

Value

an object of tbl_df.

Categorical diagnostic information

The information derived from the categorical data diagnosis is as follows.

- variables : variable names
- levels: level names
- N : number of observation
- freq : number of observation at the levels
- ratio : percentage of observation at the levels
- rank : rank of occupancy ratio of levels

See vignette("diagnosis") for an introduction to these concepts.

See Also

[diagnose_category.data.frame](#), [diagnose.tbl_dbi](#), [diagnose_category.tbl_dbi](#), [diagnose_numeric.tbl_dbi](#), [diagnose_outlier.tbl_dbi](#).

Examples

```
library(dplyr)

# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

# connect DBMS
con_sqlite <- DBI::dbConnect(RSQLite::SQLite(), ":memory:")

# copy carseats to the DBMS with a table named TB_CARSEATS
copy_to(con_sqlite, carseats, name = "TB_CARSEATS", overwrite = TRUE)

# Using pipes -----
# Diagnosis of all categorical variables
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  diagnose_category()

# Positive values select variables
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  diagnose_category(Urban, US)

# Negative values to drop variables, and In-memory mode
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
```

```

diagnose_category(-Urban, -US, in_database = FALSE)

# Positions values select variables, and In-memory mode and collect size is 200
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  diagnose_category(7, in_database = FALSE, collect_size = 200)

# Positions values select variables
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  diagnose_category(-7)

# Top rank levels with top argument
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  diagnose_category(top = 2)

# Using pipes & dplyr -----
# Extraction of level that is more than 60% of categorical data
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  diagnose_category() %>%
  filter(ratio >= 60)

```

diagnose_numeric	<i>Diagnose data quality of numerical variables</i>
------------------	---

Description

The `diagnose_numeric()` produces information for diagnosing the quality of the numerical data.

Usage

```

diagnose_numeric(.data, ...)

## S3 method for class 'data.frame'
diagnose_numeric(.data, ...)

```

Arguments

<code>.data</code>	a <code>data.frame</code> or a <code>tbl_df</code> .
<code>...</code>	one or more unquoted expressions separated by commas. You can treat variable names like they are positions. Positive values select variables; negative values to drop variables. If the first expression is negative, <code>diagnose_numeric()</code> will automatically start with all variables. These arguments are automatically quoted and evaluated in a context where column names represent column positions. They support unquoting and splicing.

Details

The scope of the diagnosis is to calculate a statistic that can be used to understand the distribution of numerical data. min, Q1, mean, median, Q3, max can be used to estimate the distribution of data. If the number of zero or minus is large, it is necessary to suspect the error of the data. If the number of outliers is large, a strategy of eliminating or replacing outliers is needed.

Value

an object of `tbl_df`.

Numerical diagnostic information

The information derived from the numerical data diagnosis is as follows.

- variables : variable names
- min : minimum
- Q1 : 25 percentile
- mean : arithmetic average
- median : median. 50 percentile
- Q3 : 75 percentile
- max : maximum
- zero : count of zero values
- minus : count of minus values
- outlier : count of outliers

See vignette("diagnosis") for an introduction to these concepts.

See Also

[diagnose_numeric.tbl_dbi](#), [diagnose.data.frame](#), [diagnose_category.data.frame](#), [diagnose_outlier.data.frame](#)

Examples

```
# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

# Diagnosis of numerical variables
diagnose_numeric(carseats)

# Select the variable to diagnose
diagnose_numeric(carseats, Sales, Income)
diagnose_numeric(carseats, -Sales, -Income)
diagnose_numeric(carseats, "Sales", "Income")
diagnose_numeric(carseats, 5)

# Using pipes -----
```



```

library(dplyr)

# Diagnosis of all numerical variables
carseats %>%
  diagnose_numeric()
# Positive values select variables
carseats %>%
  diagnose_numeric(Sales, Income)
# Negative values to drop variables
carseats %>%
  diagnose_numeric(-Sales, -Income)
# Positions values select variables
carseats %>%
  diagnose_numeric(5)
# Positions values select variables
carseats %>%
  diagnose_numeric(-1, -5)

# Using pipes & dplyr -----
# Information records of zero variable more than 0
carseats %>%
  diagnose_numeric() %>%
  filter(zero > 0)

```

```
diagnose_numeric.tbl_dbi
```

Diagnose data quality of numerical variables in the DBMS

Description

The `diagnose_numeric()` produces information for diagnosing the quality of the numerical(INTEGER, NUMBER, etc.) column of the DBMS table through `tbl_dbi`.

Usage

```
## S3 method for class 'tbl_dbi'
diagnose_numeric(.data, ..., in_database = FALSE, collect_size = Inf)
```

Arguments

<code>.data</code>	a <code>tbl_dbi</code> .
<code>...</code>	one or more unquoted expressions separated by commas. You can treat variable names like they are positions. Positive values select variables; negative values to drop variables. If the first expression is negative, <code>diagnose_numeric()</code> will automatically start with all variables. These arguments are automatically quoted and evaluated in a context where column names represent column positions. They support unquoting and splicing.

<code>in_database</code>	Specifies whether to perform in-database operations. If TRUE, most operations are performed in the DBMS. If FALSE, table data is taken in R and operated in-memory. Not yet supported in_database = TRUE.
<code>collect_size</code>	a integer. The number of data samples from the DBMS to R. Applies only if in_database = FALSE.

Details

The scope of the diagnosis is to calculate a statistic that can be used to understand the distribution of numerical data. min, Q1, mean, median, Q3, max can be used to estimate the distribution of data. If the number of zero or minus is large, it is necessary to suspect the error of the data. If the number of outliers is large, a strategy of eliminating or replacing outliers is needed.

Value

an object of `tbl_df`.

Numerical diagnostic information

The information derived from the numerical data diagnosis is as follows.

- `variables` : variable names
- `min` : minimum
- `Q1` : 25 percentile
- `mean` : arithmetic average
- `median` : median. 50 percentile
- `Q3` : 75 percentile
- `max` : maximum
- `zero` : count of zero values
- `minus` : count of minus values
- `outlier` : count of outliers

See vignette("diagnosis") for an introduction to these concepts.

See Also

[diagnose_numeric.data.frame](#), [diagnose.tbl_dbi](#), [diagnose_category.tbl_dbi](#), [diagnose_outlier.tbl_dbi](#).

Examples

```
library(dplyr)

# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

# connect DBMS
```

```

con_sqlite <- DBI::dbConnect(RSQLite::SQLite(), ":memory:")

# copy carseats to the DBMS with a table named TB_CARSEATS
copy_to(con_sqlite, carseats, name = "TB_CARSEATS", overwrite = TRUE)

# Using pipes -----
# Diagnosis of all numerical variables
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  diagnose_numeric()

# Positive values select variables, and In-memory mode and collect size is 200
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  diagnose_numeric(Sales, Income, collect_size = 200)

# Negative values to drop variables
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  diagnose_numeric(-Sales, -Income)

# Positions values select variables
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  diagnose_numeric(5)

# Positions values select variables
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  diagnose_numeric(-1, -5)

# Using pipes & dplyr -----
# Information records of zero variable more than 0
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  diagnose_numeric() %>%
  filter(zero > 0)

```

diagnose_outlier

Diagnose outlier of numerical variables

Description

The `diagnose_outlier()` produces outlier information for diagnosing the quality of the numerical data.

Usage

```
diagnose_outlier(.data, ...)
```

```
## S3 method for class 'data.frame'
diagnose_outlier(.data, ...)
```

Arguments

<code>.data</code>	a <code>data.frame</code> or a <code>tbl_df</code> .
<code>...</code>	one or more unquoted expressions separated by commas. You can treat variable names like they are positions. Positive values select variables; negative values to drop variables. If the first expression is negative, <code>diagnose_outlier()</code> will automatically start with all variables. These arguments are automatically quoted and evaluated in a context where column names represent column positions. They support unquoting and splicing.

Details

The scope of the diagnosis is to provide an outlier information. If the number of outliers is small and the difference between the averages including outliers and the averages not including them is large, it is necessary to eliminate or replace the outliers.

Value

an object of `tbl_df`.

Outlier Diagnostic information

The information derived from the numerical data diagnosis is as follows.

- `variables` : variable names
- `outliers_cnt` : number of outliers
- `outliers_ratio` : percent of outliers
- `outliers_mean` : arithmetic average of outliers
- `with_mean` : arithmetic average of with outliers
- `without_mean` : arithmetic average of without outliers

See vignette("diagnosis") for an introduction to these concepts.

See Also

[diagnose_outlier.tbl_dbi](#), [diagnose.data.frame](#), [diagnose_category.data.frame](#), [diagnose_numeric.data.frame](#)

Examples

```
# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

# Diagnosis of numerical variables
```

```

diagnose_outlier(carseats)

# Select the variable to diagnose
diagnose_outlier(carseats, Sales, Income)
diagnose_outlier(carseats, -Sales, -Income)
diagnose_outlier(carseats, "Sales", "Income")
diagnose_outlier(carseats, 5)

# Using pipes -----
library(dplyr)

# Diagnosis of all numerical variables
carseats %>%
  diagnose_outlier()
# Positive values select variables
carseats %>%
  diagnose_outlier(Sales, Income)
# Negative values to drop variables
carseats %>%
  diagnose_outlier(-Sales, -Income)
# Positions values select variables
carseats %>%
  diagnose_outlier(5)
# Positions values select variables
carseats %>%
  diagnose_outlier(-1, -5)

# Using pipes & dplyr -----
# outlier_ratio is more than 1%
carseats %>%
  diagnose_outlier() %>%
  filter(outliers_ratio > 1)

```

```
diagnose_outlier.tbl_dbi
```

Diagnose outlier of numerical variables in the DBMS

Description

The `diagnose_outlier()` produces outlier information for diagnosing the quality of the numerical (INTEGER, NUMBER, etc.) column of the DBMS table through `tbl_dbi`.

Usage

```
## S3 method for class 'tbl_dbi'
diagnose_outlier(.data, ..., in_database = FALSE, collect_size = Inf)
```

Arguments

<code>.data</code>	a <code>tbl_dbi</code> .
<code>...</code>	one or more unquoted expressions separated by commas. You can treat variable names like they are positions. Positive values select variables; negative values to drop variables. If the first expression is negative, <code>diagnose_outlier()</code> will automatically start with all variables. These arguments are automatically quoted and evaluated in a context where column names represent column positions. They support unquoting and splicing.
<code>in_database</code>	Specifies whether to perform in-database operations. If <code>TRUE</code> , most operations are performed in the DBMS. if <code>FALSE</code> , table data is taken in R and operated in-memory. Not yet supported <code>in_database = TRUE</code> .
<code>collect_size</code>	a integer. The number of data samples from the DBMS to R. Applies only if <code>in_database = FALSE</code> .

Details

The scope of the diagnosis is the provide a outlier information. If the number of outliers is small and the difference between the averages including outliers and the averages not including them is large, it is necessary to eliminate or replace the outliers.

Value

an object of `tbl_df`.

Outlier Diagnostic information

The information derived from the numerical data diagnosis is as follows.

- `variables` : variable names
- `outliers_cnt` : number of outliers
- `outliers_ratio` : percent of outliers
- `outliers_mean` : arithmetic average of outliers
- `with_mean` : arithmetic average of with outliers
- `without_mean` : arithmetic average of without outliers

See vignette("diagonosis") for an introduction to these concepts.

See Also

[diagnose_outlier.data.frame](#), [diagnose.tbl_dbi](#), [diagnose_category.tbl_dbi](#), [diagnose_numeric.tbl_dbi](#).

Examples

```

library(dplyr)

# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

# connect DBMS
con_sqlite <- DBI::dbConnect(RSQLite::SQLite(), ":memory:")

# copy carseats to the DBMS with a table named TB_CARSEATS
copy_to(con_sqlite, carseats, name = "TB_CARSEATS", overwrite = TRUE)

# Using pipes -----
# Diagnosis of all numerical variables
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  diagnose_outlier()

# Positive values select variables, and In-memory mode and collect size is 200
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  diagnose_outlier(Sales, Income, collect_size = 200)

# Negative values to drop variables
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  diagnose_outlier(-Sales, -Income)

# Positions values select variables
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  diagnose_outlier(5)
# Positions values select variables

con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  diagnose_outlier(-1, -5)

# Using pipes & dplyr -----
# outlier_ratio is more than 1%
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  diagnose_outlier() %>%
  filter(outliers_ratio > 1)

```

Description

The `diagnose_report()` report the information for diagnosing the quality of the data.

Usage

```
diagnose_report(.data, output_format, output_file, output_dir, ...)
```

```
## S3 method for class 'data.frame'
diagnose_report(
  .data,
  output_format = c("pdf", "html"),
  output_file = NULL,
  output_dir = tempdir(),
  font_family = NULL,
  browse = TRUE,
  ...
)
```

Arguments

<code>.data</code>	a <code>data.frame</code> or a <code>tbl_df</code> .
<code>output_format</code>	report output type. Choose either "pdf" and "html". "pdf" create pdf file by <code>knitr::knit()</code> . "html" create html file by <code>rmarkdown::render()</code> .
<code>output_file</code>	name of generated file. default is <code>NULL</code> .
<code>output_dir</code>	name of directory to generate report file. default is <code>tempdir()</code> .
<code>...</code>	arguments to be passed to methods.
<code>font_family</code>	character. font family name for figure in pdf.
<code>browse</code>	logical. choose whether to output the report results to the browser.

Details

Generate generalized data diagnostic reports automatically. You can choose to output to pdf and html files. This is useful for diagnosing a data frame with a large number of variables than data with a small number of variables. For pdf output, Korean Gothic font must be installed in Korean operating system.

Reported information

Reported from the data diagnosis is as follows.

- Diagnose Data
 - Overview of Diagnosis
 - * List of all variables quality
 - * Diagnosis of missing data
 - * Diagnosis of unique data(Text and Category)
 - * Diagnosis of unique data(Numerical)

- Detailed data diagnosis
 - * Diagnosis of categorical variables
 - * Diagnosis of numerical variables
 - * List of numerical diagnosis (zero)
 - * List of numerical diagnosis (minus)
- Diagnose Outliers
 - Overview of Diagnosis
 - * Diagnosis of numerical variable outliers
 - * Detailed outliers diagnosis

See vignette("diagonosis") for an introduction to these concepts.

Examples

```

carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

# reporting the diagnosis information -----
# create pdf file. file name is DataDiagnosis_Report.pdf
diagnose_report(carseats)
# create pdf file. file name is Diagn.pdf
diagnose_report(carseats, output_file = "Diagn.pdf")
# create pdf file. file name is ./Diagn.pdf and not browse
# diagnose_report(carseats, output_dir = ".", output_file = "Diagn.pdf",
#   browse = FALSE)
# create html file. file name is Diagnosis_Report.html
diagnose_report(carseats, output_format = "html")
# create html file. file name is Diagn.html
diagnose_report(carseats, output_format = "html", output_file = "Diagn.html")

```

diagnose_report.tbl_dbi

Reporting the information of data diagnosis for table of the DBMS

Description

The `diagnose_report()` report the information for diagnosing the quality of the DBMS table through `tbl_dbi`

Usage

```
## S3 method for class 'tbl_dbi'
diagnose_report(
  .data,
  output_format = c("pdf", "html"),
  output_file = NULL,
  output_dir = tempdir(),
  font_family = NULL,
  in_database = FALSE,
  collect_size = Inf,
  ...
)
```

Arguments

.data	a tbl_dbi.
output_format	report output type. Choose either "pdf" and "html". "pdf" create pdf file by knitr::knit(). "html" create html file by rmarkdown::render().
output_file	name of generated file. default is NULL.
output_dir	name of directory to generate report file. default is tempdir().
font_family	character. font family name for figure in pdf.
in_database	Specifies whether to perform in-database operations. If TRUE, most operations are performed in the DBMS. if FALSE, table data is taken in R and operated in-memory. Not yet supported in_database = TRUE.
collect_size	a integer. The number of data samples from the DBMS to R. Applies only if in_database = FALSE.
...	arguments to be passed to methods.

Details

Generate generalized data diagnostic reports automatically. You can choose to output to pdf and html files. This is useful for diagnosing a data frame with a large number of variables than data with a small number of variables. For pdf output, Korean Gothic font must be installed in Korean operating system.

Reported information

Reported from the data diagnosis is as follows.

- Diagnose Data
 - Overview of Diagnosis
 - * List of all variables quality
 - * Diagnosis of missing data
 - * Diagnosis of unique data(Text and Category)
 - * Diagnosis of unique data(Numerical)
 - Detailed data diagnosis

- * Diagnosis of categorical variables
- * Diagnosis of numerical variables
- * List of numerical diagnosis (zero)
- * List of numerical diagnosis (minus)
- Diagnose Outliers
 - Overview of Diagnosis
 - * Diagnosis of numerical variable outliers
 - * Detailed outliers diagnosis

See vignette("diagonosis") for an introduction to these concepts.

See Also

[diagnose_report.data.frame.](#)

Examples

```
library(dplyr)

# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

# connect DBMS
con_sqlite <- DBI::dbConnect(RSQLite::SQLite(), ":memory:")

# copy carseats to the DBMS with a table named TB_CARSEATS
copy_to(con_sqlite, carseats, name = "TB_CARSEATS", overwrite = TRUE)

# reporting the diagnosis information -----
# create pdf file. file name is DataDiagnosis_Report.pdf
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  diagnose_report()

# create pdf file. file name is Diagn.pdf, and collect size is 350
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  diagnose_report(collect_size = 350, output_file = "Diagn.pdf")

# create html file. file name is Diagnosis_Report.html
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  diagnose_report(output_format = "html")

# create html file. file name is Diagn.html
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  diagnose_report(output_format = "html", output_file = "Diagn.html")
```

 eda_report

Reporting the information of EDA

Description

The `eda_report()` report the information of exploratory data analysis for object inheriting from `data.frame`.

Usage

```
eda_report(.data, ...)

## S3 method for class 'data.frame'
eda_report(
  .data,
  target = NULL,
  output_format = c("pdf", "html"),
  output_file = NULL,
  output_dir = tempdir(),
  font_family = NULL,
  browse = TRUE,
  ...
)
```

Arguments

<code>.data</code>	a <code>data.frame</code> or a <code>tbl_df</code> .
<code>...</code>	arguments to be passed to methods.
<code>target</code>	target variable.
<code>output_format</code>	character. report output type. Choose either "pdf" and "html". "pdf" create pdf file by <code>knitr::knit()</code> . "html" create html file by <code>rmarkdown::render()</code> .
<code>output_file</code>	character. name of generated file. default is NULL.
<code>output_dir</code>	character. name of directory to generate report file. default is <code>tempdir()</code> .
<code>font_family</code>	character. font family name for figure in pdf.
<code>browse</code>	logical. choose whether to output the report results to the browser.

Details

Generate generalized EDA report automatically. You can choose to output as pdf and html files. This feature is useful for EDA of data with many variables, rather than data with fewer variables. For pdf output, Korean Gothic font must be installed in Korean operating system.

Reported information

The EDA process will report the following information:

- Introduction
 - Information of Dataset
 - Information of Variables
 - About EDA Report
- Univariate Analysis
 - Descriptive Statistics
 - Normality Test of Numerical Variables
 - * Statistics and Visualization of (Sample) Data
- Relationship Between Variables
 - Correlation Coefficient
 - * Correlation Coefficient by Variable Combination
 - * Correlation Plot of Numerical Variables
- Target based Analysis
 - Grouped Descriptive Statistics
 - * Grouped Numerical Variables
 - * Grouped Categorical Variables
 - Grouped Relationship Between Variables
 - * Grouped Correlation Coefficient
 - * Grouped Correlation Plot of Numerical Variables

See vignette("EDA") for an introduction to these concepts.

Examples

```
library(dplyr)

# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

## target variable is categorical variable
# reporting the EDA information
# create pdf file. file name is EDA_Report.pdf
eda_report(carseats, US)

# create pdf file. file name is EDA_carseats.pdf
eda_report(carseats, "US", output_file = "EDA_carseats.pdf")

# create pdf file. file name is EDA_carseats.pdf and not browse
# eda_report(carseats, "US", output_dir = ".", output_file = "EDA_carseats.pdf", browse = FALSE)
```

```

# create html file. file name is EDA_Report.html
eda_report(carseats, "US", output_format = "html")

# create html file. file name is EDA_carseats.html
eda_report(carseats, US, output_format = "html", output_file = "EDA_carseats.html")

## target variable is numerical variable
# reporting the EDA information
eda_report(carseats, Sales)

# create pdf file. file name is EDA2.pdf
eda_report(carseats, "Sales", output_file = "EDA2.pdf")

# create html file. file name is EDA_Report.html
eda_report(carseats, "Sales", output_format = "html")

# create html file. file name is EDA2.html
eda_report(carseats, Sales, output_format = "html", output_file = "EDA2.html")

## target variable is null
# reporting the EDA information
eda_report(carseats)

# create pdf file. file name is EDA2.pdf
eda_report(carseats, output_file = "EDA2.pdf")

# create html file. file name is EDA_Report.html
eda_report(carseats, output_format = "html")

# create html file. file name is EDA2.html
eda_report(carseats, output_format = "html", output_file = "EDA2.html")

```

eda_report.tbl_dbi *Reporting the information of EDA for table of the DBMS*

Description

The `eda_report()` report the information of Exploratory data analysis for object inheriting from the DBMS table through `tbl_dbi`

Usage

```

## S3 method for class 'tbl_dbi'
eda_report(
  .data,
  target = NULL,

```

```

    output_format = c("pdf", "html"),
    output_file = NULL,
    font_family = NULL,
    output_dir = tempdir(),
    in_database = FALSE,
    collect_size = Inf,
    ...
)

```

Arguments

.data	a tbl_dbi.
target	target variable.
output_format	report output type. Choose either "pdf" and "html". "pdf" create pdf file by knitr::knit(). "html" create html file by rmarkdown::render().
output_file	name of generated file. default is NULL.
font_family	character. font family name for figure in pdf.
output_dir	name of directory to generate report file. default is tempdir().
in_database	Specifies whether to perform in-database operations. If TRUE, most operations are performed in the DBMS. if FALSE, table data is taken in R and operated in-memory. Not yet supported in_database = TRUE.
collect_size	a integer. The number of data samples from the DBMS to R. Applies only if in_database = FALSE.
...	arguments to be passed to methods.

Details

Generate generalized data EDA reports automatically. You can choose to output to pdf and html files. This is useful for EDA a data frame with a large number of variables than data with a small number of variables. For pdf output, Korean Gothic font must be installed in Korean operating system.

Reported information

The EDA process will report the following information:

- Introduction
 - Information of Dataset
 - Information of Variables
 - About EDA Report
- Univariate Analysis
 - Descriptive Statistics
 - Normality Test of Numerical Variables
 - * Statistics and Visualization of (Sample) Data
- Relationship Between Variables

- Correlation Coefficient
 - * Correlation Coefficient by Variable Combination
 - * Correlation Plot of Numerical Variables
- Target based Analysis
 - Grouped Descriptive Statistics
 - * Grouped Numerical Variables
 - * Grouped Categorical Variables
 - Grouped Relationship Between Variables
 - * Grouped Correlation Coefficient
 - * Grouped Correlation Plot of Numerical Variables

See vignette("EDA") for an introduction to these concepts.

See Also

[eda_report.data.frame.](#)

Examples

```
library(dplyr)

# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

# connect DBMS
con_sqlite <- DBI::dbConnect(RSQLite::SQLite(), ":memory:")

# copy carseats to the DBMS with a table named TB_CARSEATS
copy_to(con_sqlite, carseats, name = "TB_CARSEATS", overwrite = TRUE)

## target variable is categorical variable
# reporting the EDA information
# create pdf file. file name is EDA_Report.pdf
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  eda_report(US)

# create pdf file. file name is EDA_TB_CARSEATS.pdf
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  eda_report("US", output_file = "EDA_TB_CARSEATS.pdf")

# create html file. file name is EDA_Report.html
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  eda_report("US", output_format = "html")
```



```
# create html file. file name is EDA_TB_CARSEATS.html
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  eda_report(US, output_format = "html", output_file = "EDA_TB_CARSEATS.html")

## target variable is numerical variable
# reporting the EDA information, and collect size is 350
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  eda_report(Sales, collect_size = 350)

# create pdf file. file name is EDA2.pdf
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  eda_report("Sales", output_file = "EDA2.pdf")

# create html file. file name is EDA_Report.html
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  eda_report("Sales", output_format = "html")

# create html file. file name is EDA2.html
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  eda_report(Sales, output_format = "html", output_file = "EDA2.html")

## target variable is null
# reporting the EDA information
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  eda_report()

# create pdf file. file name is EDA2.pdf
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  eda_report(output_file = "EDA2.pdf")

# create html file. file name is EDA_Report.html
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  eda_report(output_format = "html")

# create html file. file name is EDA2.html
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  eda_report(output_format = "html", output_file = "EDA2.html")
```

Description

The `find_class()` extracts variable information having a certain class from an object inheriting `data.frame`.

Usage

```
find_class(  
  df,  
  type = c("numerical", "categorical", "categorical2"),  
  index = TRUE  
)
```

Arguments

<code>df</code>	a <code>data.frame</code> or objects inheriting from <code>data.frame</code>
<code>type</code>	character. Defines a group of classes to be searched. "numerical" searches for "numeric" and "integer" classes, "categorical" searches for "factor" and "ordered" classes. "categorical2" adds "character" class to "categorical".
<code>index</code>	logical. If TRUE is return numeric vector that is variables index. and if FALSE is return character vector that is variables name. default is TRUE.

Value

character vector or numeric vector. The meaning of vector according to data type is as follows.

- character vector : variables name
- numeric vector : variables index

See Also

[get_class](#).

Examples

```
## Not run:  
# data.frame  
find_class(iris, "numerical")  
find_class(iris, "numerical", index = FALSE)  
find_class(iris, "categorical")  
find_class(iris, "categorical", index = FALSE)  
  
# tbl_df  
find_class(ISLR::Carseats, "numerical")  
find_class(ISLR::Carseats, "numerical", index = FALSE)  
find_class(ISLR::Carseats, "categorical")  
find_class(ISLR::Carseats, "categorical", index = FALSE)  
  
# type is "categorical2"  
iris2 <- data.frame(iris, char = "chars",  
                    stringsAsFactors = FALSE)
```

```
find_class(iris2, "categorical", index = FALSE)
find_class(iris2, "categorical2", index = FALSE)

## End(Not run)
```

find_na	<i>Finding variables including missing values</i>
---------	---

Description

Find the variable that contains the missing value in the object that inherits the data.frame or data.frame.

Usage

```
find_na(.data, index = TRUE, rate = FALSE)
```

Arguments

.data	a data.frame or a tbl_df .
index	logical. When representing the information of a variable including missing values, specify whether or not the variable is represented by an index. Returns an index if TRUE or a variable names if FALSE.
rate	logical. If TRUE, returns the percentage of missing values in the individual variable.

Value

Information on variables including missing values.

See Also

[imputate_na](#), [find_na](#).

Examples

```
## Not run:
# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

find_na(carseats)

find_na(carseats, index = FALSE)

find_na(carseats, rate = TRUE)

## using dplyr -----
library(dplyr)
```

```
# Perform simple data quality diagnosis of variables with missing values.
carseats %>%
  select(find_na(.)) %>%
  diagnose()

## End(Not run)
```

find_outliers

Finding variables including outliers

Description

Find the numerical variable that contains outliers in the object that inherits the data.frame or data.frame.

Usage

```
find_outliers(.data, index = TRUE, rate = FALSE)
```

Arguments

.data	a data.frame or a tbl_df .
index	logical. When representing the information of a variable including outliers, specify whether or not the variable is represented by an index. Returns an index if TRUE or a variable names if FALSE.
rate	logical. If TRUE, returns the percentage of outliers in the individual variable.

Value

Information on variables including outliers.

See Also

[find_na](#), [imputate_outlier](#).

Examples

```
## Not run:
# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

find_outliers(carseats)

find_outliers(carseats, index = FALSE)

find_outliers(carseats, rate = TRUE)
```

```
## using dplyr -----
library(dplyr)

# Perform simple data quality diagnosis of variables with outliers.
carseats %>%
  select(find_outliers(.)) %>%
  diagnose()

## End(Not run)
```

find_skewness	<i>Finding skewed variables</i>
---------------	---------------------------------

Description

Find the numerical variable that skewed variable that inherits the data.frame or data.frame.

Usage

```
find_skewness(.data, index = TRUE, value = FALSE, thres = NULL)
```

Arguments

.data	a data.frame or a tbl_df .
index	logical. When representing the information of a skewed variable, specify whether or not the variable is represented by an index. Returns an index if TRUE or a variable names if FALSE.
value	logical. If TRUE, returns the skewness value in the individual variable.
thres	Returns a skewness threshold value that has an absolute skewness greater than thres. The default is NULL to ignore the threshold. but, If value = TRUE, default to 0.5.

Value

Information on variables including skewness.

See Also

[find_na](#), [find_outliers](#).

Examples

```
## Not run:
# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA
```

```
find_skewness(carseats)

find_skewness(carseats, index = FALSE)

find_skewness(carseats, thres = 0.1)

find_skewness(carseats, value = TRUE)

find_skewness(carseats, value = TRUE, thres = 0.1)

## using dplyr -----
library(dplyr)

# Perform simple data quality diagnosis of variables with outliers.
carseats %>%
  select(find_skewness(.)) %>%
  diagnose()

## End(Not run)
```

get_class

Extracting a class of variables

Description

The `get_class()` gets class of variables in `data.frame` or `tbl_df`.

Usage

```
get_class(df)
```

Arguments

`df` a `data.frame` or objects inheriting from `data.frame`

Value

a `data.frame` Variables of `data.frame` is as follows.

- `variable` : variables name
- `class` : class of variables

See Also

[find_class](#).

Examples

```
## Not run:
# data.frame
get_class(iris)

# tbl_df
get_class(ggplot2::diamonds)

library(dplyr)
get_class(ggplot2::diamonds) %>%
  filter(class %in% c("integer", "numeric"))

## End(Not run)
```

get_column_info	<i>Describe column of table in the DBMS</i>
-----------------	---

Description

The `get_column_info()` retrieves the column information of the DBMS table through the `tbl_bdi` object of `dplyr`.

Usage

```
get_column_info(df)
```

Arguments

`df` a `tbl_dbi`.

Value

An object of `data.frame`.

Column information of the DBMS table

- SQLite DBMS connected `RSQLite::SQLite()`:
 - name: column name
 - type: data type in R
- MySQL/MariaDB DBMS connected `RMySQL::MySQL()`:
 - name: column name
 - Sclass: data type in R
 - type: data type of column in the DBMS
 - length: data length in the DBMS
- Oracle DBMS connected `ROracle::dbConnect()`:
 - name: column name

- Sclass: column type in R
- type: data type of column in the DBMS
- len: length of column(CHAR/VARCHAR/VARCHAR2 data type) in the DBMS
- precision: precision of column(NUMBER data type) in the DBMS
- scale: decimal places of column(NUMBER data type) in the DBMS
- nullOK: nullability

Examples

```
library(dplyr)

# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

# connect DBMS
con_sqlite <- DBI::dbConnect(RSQLite::SQLite(), ":memory:")

# copy carseats to the DBMS with a table named TB_CARSEATS
copy_to(con_sqlite, carseats, name = "TB_CARSEATS", overwrite = TRUE)

con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  get_column_info
```

get_os

Finding Users Machine's OS

Description

Get the operating system that users machines.

Usage

```
get_os()
```

Value

OS names. "windows" or "osx" or "linux"

Examples

```
get_os()
```

impute_na	<i>Impute Missing Values</i>
-----------	------------------------------

Description

Missing values are imputed with some representative values and statistical methods.

Usage

```
impute_na(.data, xvar, yvar, method, seed, print_flag, no_attrs)
```

Arguments

.data	a data.frame or a tbl_df .
xvar	variable name to replace missing value.
yvar	target variable.
method	method of missing values imputation.
seed	integer. the random seed used in mice. only used "mice" method.
print_flag	logical. If TRUE, mice will print running log on console. Use print_flag=FALSE for silent computation. Used only when method is "mice".
no_attrs	logical. If TRUE, return numerical variable or categorical variable. else If FALSE, imputation class.

Details

impute_na () creates an imputation class. The 'imputation' class includes missing value position, imputed value, and method of missing value imputation, etc. The 'imputation' class compares the imputed value with the original value to help determine whether the imputed value is used in the analysis.

See vignette("transformation") for an introduction to these concepts.

Value

An object of imputation class. or numerical variable or categorical variable. if no_attrs is FALSE then return imputation class, else no_attrs is TRUE then return numerical vector or factor. Attributes of imputation class is as follows.

- var_type : the data type of predictor to replace missing value.
- method : method of missing value imputation.
 - predictor is numerical variable
 - * "mean" : arithmetic mean
 - * "median" : median
 - * "mode" : mode
 - * "knn" : K-nearest neighbors

- * "rpart" : Recursive Partitioning and Regression Trees
- * "mice" : Multivariate Imputation by Chained Equations
- predictor is categorical variable
 - * "mode" : mode
 - * "rpart" : Recursive Partitioning and Regression Trees
 - * "mice" : Multivariate Imputation by Chained Equations
- na_pos : position of missing value in predictor.
- seed : the random seed used in mice. only used "mice" method.
- type : "missing values". type of imputation.
- message : a message tells you if the result was successful.
- success : Whether the imputation was successful.

See Also

[impute_outlier](#).

Examples

```
# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

# Replace the missing value of the Income variable with median
impute_na(carseats, Income, method = "median")

# Replace the missing value of the Income variable with rpart
# The target variable is US.
impute_na(carseats, Income, US, method = "rpart")

# Replace the missing value of the Urban variable with mode
impute_na(carseats, Urban, method = "mode")

# Replace the missing value of the Urban variable with mice
# The target variable is US.
impute_na(carseats, Urban, US, method = "mice")

## using dplyr -----
library(dplyr)

# The mean before and after the imputation of the Income variable
carseats %>%
  mutate(Income_imp = impute_na(carseats, Income, US, method = "knn", no_attr = TRUE)) %>%
  group_by(US) %>%
  summarise(orig = mean(Income, na.rm = TRUE),
            imputation = mean(Income_imp))

# If the variable of interest is a numerical variable
```

```

income <- impute_na(carseats, Income, US, method = "rpart")
income
summary(income)
plot(income)

# If the variable of interest is a categorical variable
urban <- impute_na(carseats, Urban, US, method = "mice")
urban
summary(urban)
plot(urban)

```

impute_outlier *Impute Outliers*

Description

Outliers are imputed with some representative values and statistical methods.

Usage

```
impute_outlier(.data, xvar, method, no_attr)
```

Arguments

.data	a data.frame or a tbl_df .
xvar	variable name to replace missing value.
method	method of missing values imputation.
no_attr	logical. If TRUE, return numerical variable or categorical variable. else If FALSE, imputation class.

Details

impute_outlier() creates an imputation class. The ‘imputation’ class includes missing value position, imputed value, and method of missing value imputation, etc. The ‘imputation’ class compares the imputed value with the original value to help determine whether the imputed value is used in the analysis.

See vignette("transformation") for an introduction to these concepts.

Value

An object of imputation class. or numerical variable. if no_attr is FALSE then return imputation class, else no_attr is TRUE then return numerical vector. Attributes of imputation class is as follows.

- method : method of missing value imputation.
 - predictor is numerical variable

- * "mean" : arithmetic mean
- * "median" : median
- * "mode" : mode
- * "capping" : Impute the upper outliers with 95 percentile, and Impute the bottom outliers with 5 percentile.

- outlier_pos : position of outliers in predictor.
- outliers : outliers. outliers corresponding to outlier_pos.
- type : "outliers". type of imputation.

See Also

[imputate_na](#).

Examples

```
# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

# Replace the outliers of the Price variable with median
imputate_outlier(carseats, Price, method = "median")

# Replace the outliers of the Price variable with capping
imputate_outlier(carseats, Price, method = "capping")

## using dplyr -----
library(dplyr)

# The mean before and after the imputation of the Price variable
carseats %>%
  mutate(Price_imp = imputate_outlier(carseats, Price, method = "capping", no_attrs = TRUE)) %>%
  group_by(US) %>%
  summarise(orig = mean(Price, na.rm = TRUE),
            imputation = mean(Price_imp, na.rm = TRUE))

# If the variable of interest is a numerical variable
price <- imputate_outlier(carseats, Price)
price
summary(price)
plot(price)
```

kurtosis

Kurtosis of the data

Description

This function calculated kurtosis of given data.

Usage

```
kurtosis(x, na.rm = FALSE)
```

Arguments

x	a numeric vector.
na.rm	logical. Determine whether to remove missing values and calculate them. The default is TRUE.

Value

numeric. calculated kurtosis

See Also

[skewness](#).

Examples

```
set.seed(123)
kurtosis(rnorm(100))
```

normality

Performs the Shapiro-Wilk test of normality

Description

The `normality()` performs Shapiro-Wilk test of normality of numerical values.

Usage

```
normality(.data, ...)

## S3 method for class 'data.frame'
normality(.data, ..., sample = 5000)
```

Arguments

.data	a data.frame or a tbl_df .
...	one or more unquoted expressions separated by commas. You can treat variable names like they are positions. Positive values select variables; negative values to drop variables. If the first expression is negative, <code>normality()</code> will automatically start with all variables. These arguments are automatically quoted and evaluated in a context where column names represent column positions. They support unquoting and splicing.
sample	the number of samples to perform the test. See <code>vignette("EDA")</code> for an introduction to these concepts.

Details

This function is useful when used with the `group_by` function of the `dplyr` package. If you want to test by level of the categorical data you are interested in, rather than the whole observation, you can use `group_tf` as the `group_by` function. This function is computed `shapiro.test` function.

Value

An object of the same class as `.data`.

Normality test information

The information derived from the numerical data test is as follows.

- `statistic` : the value of the Shapiro-Wilk statistic.
- `p_value` : an approximate p-value for the test. This is said in Royston(1995) to be adequate for $p_value < 0.1$.
- `sample` : the number of samples to perform the test. The number of observations supported by the `stats::shapiro.test` function is 3 to 5000.

See Also

[normality.tbl_dbi](#), [diagnose_numeric.data.frame](#), [describe.data.frame](#), [plot_normality.data.frame](#).

Examples

```
# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

# Normality test of numerical variables
normality(carseats)

# Select the variable to describe
normality(carseats, Sales, Price)
normality(carseats, -Sales, -Price)
normality(carseats, 1)
normality(carseats, Sales, Price, sample = 300)

# Using dplyr::grouped_dt
library(dplyr)

gdata <- group_by(carseats, ShelveLoc, US)
normality(gdata, "Sales")
normality(gdata, sample = 250)

# Using pipes -----
# Normality test of all numerical variables
carseats %>%
  normality()
```

```

# Positive values select variables
carseats %>%
  normality(Sales, Price)

# Positions values select variables
carseats %>%
  normality(1)

# Using pipes & dplyr -----
# Test all numerical variables by 'ShelveLoc' and 'US',
# and extract only those with 'ShelveLoc' variable level is "Good".
carseats %>%
  group_by(ShelveLoc, US) %>%
  normality() %>%
  filter(ShelveLoc == "Good")

# extract only those with 'Urban' variable level is "Yes",
# and test 'Sales' by 'ShelveLoc' and 'US'
carseats %>%
  filter(Urban == "Yes") %>%
  group_by(ShelveLoc, US) %>%
  normality(Sales)

# Test log(Income) variables by 'ShelveLoc' and 'US',
# and extract only p.value greater than 0.01.
carseats %>%
  mutate(log_income = log(Income)) %>%
  group_by(ShelveLoc, US) %>%
  normality(log_income) %>%
  filter(p_value > 0.01)

```

normality.tbl_dbi	<i>Performs the Shapiro-Wilk test of normality</i>
-------------------	--

Description

The `normality()` performs Shapiro-Wilk test of normality of numerical(INTEGER, NUMBER, etc.) column of the DBMS table through `tbl_dbi`.

Usage

```

## S3 method for class 'tbl_dbi'
normality(.data, ..., sample = 5000, in_database = FALSE, collect_size = Inf)

```

Arguments

<code>.data</code>	a <code>tbl_dbi</code> .
<code>...</code>	one or more unquoted expressions separated by commas. You can treat variable names like they are positions. Positive values select variables; negative values to drop variables. If the first expression is negative, <code>normality()</code> will automatically start with all variables. These arguments are automatically quoted and evaluated in a context where column names represent column positions. They support unquoting and splicing.
<code>sample</code>	the number of samples to perform the test.
<code>in_database</code>	Specifies whether to perform in-database operations. If <code>TRUE</code> , most operations are performed in the DBMS. if <code>FALSE</code> , table data is taken in R and operated in-memory. Not yet supported in <code>in_database = TRUE</code> .
<code>collect_size</code>	a integer. The number of data samples from the DBMS to R. Applies only if <code>in_database = FALSE</code> . See vignette("EDA") for an introduction to these concepts.

Details

This function is useful when used with the `group_by` function of the `dplyr` package. If you want to test by level of the categorical data you are interested in, rather than the whole observation, you can use `group_tf` as the `group_by` function. This function is computed `shapiro.test` function.

Value

An object of the same class as `.data`.

Normality test information

The information derived from the numerical data test is as follows.

- `statistic` : the value of the Shapiro-Wilk statistic.
- `p_value` : an approximate p-value for the test. This is said in Royston(1995) to be adequate for `p_value < 0.1`.
- `sample` : the numer of samples to perform the test. The number of observations supported by the `stats::shapiro.test` function is 3 to 5000.

See Also

[normality.data.frame](#), [diagnose_numeric.tbl_dbi](#), [describe.tbl_dbi](#), [plot_normality.tbl_dbi](#).

Examples

```
library(dplyr)

# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
```



```

carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

# connect DBMS
con_sqlite <- DBI::dbConnect(RSQLite::SQLite(), ":memory:")

# copy carseats to the DBMS with a table named TB_CARSEATS
copy_to(con_sqlite, carseats, name = "TB_CARSEATS", overwrite = TRUE)

# Using pipes -----
# Normality test of all numerical variables
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  normality()

# Positive values select variables, and In-memory mode and collect size is 200
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  normality(Sales, Price, collect_size = 200)

# Positions values select variables
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  normality(1)

# Using pipes & dplyr -----
# Test all numerical variables by 'ShelveLoc' and 'US',
# and extract only those with 'ShelveLoc' variable level is "Good".
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  group_by(ShelveLoc, US) %>%
  normality() %>%
  filter(ShelveLoc == "Good")

# extract only those with 'Urban' variable level is "Yes",
# and test 'Sales' by 'ShelveLoc' and 'US'
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  filter(Urban == "Yes") %>%
  group_by(ShelveLoc, US) %>%
  normality(Sales)

# Test log(Income) variables by 'ShelveLoc' and 'US',
# and extract only p.value greater than 0.01.

# SQLite extension functions for log
RSQLite::initExtension(con_sqlite)

con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  mutate(log_income = log(Income)) %>%
  group_by(ShelveLoc, US) %>%
  normality(log_income) %>%
  filter(p_value > 0.01)

```

`plot.bins`*Visualize Distribution for a "bins" object*

Description

Visualize two plots on a single screen. The plot at the top is a histogram representing the frequency of the level. The plot at the bottom is a bar chart representing the frequency of the level.

Usage

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

Arguments

`x` an object of class "bins", usually, a result of a call to `binning()`.
`...` arguments to be passed to methods, such as graphical parameters (see `par`).

See Also

[binning](#), [print.bins](#), [summary.bins](#).

Examples

```
# Generate data for the example  
carseats <- ISLR::Carseats  
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA  
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA  
  
# Binning the carat variable. default type argument is "quantile"  
bin <- binning(carseats$Income, nbins = 5)  
plot(bin)  
# Using another type argument  
bin <- binning(carseats$Income, nbins = 5, type = "equal")  
plot(bin)  
bin <- binning(carseats$Income, nbins = 5, type = "pretty")  
plot(bin)  
bin <- binning(carseats$Income, nbins = 5, type = "kmeans")  
plot(bin)  
bin <- binning(carseats$Income, nbins = 5, type = "bclust")  
plot(bin)
```

plot.compare_category *Visualize Information for an "compare_category" Object*

Description

Visualize mosaics plot by attribute of compare_category class.

Usage

```
## S3 method for class 'compare_category'  
plot(x, prompt = FALSE, na.rm = FALSE, ...)
```

Arguments

x	an object of class "compare_category", usually, a result of a call to compare_category().
prompt	logical. The default value is FALSE. If there are multiple visualizations to be output, if this argument value is TRUE, a prompt is output each time.
na.rm	logical. Specifies whether to include NA when plotting mosaics plot. The default is FALSE, so plot NA.
...	arguments to be passed to methods, such as graphical parameters (see par). However, it does not support all parameters.

See Also

[compare_category](#), [print.compare_category](#), [summary.compare_category](#).

Examples

```
# Generate data for the example  
carseats <- ISLR::Carseats  
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA  
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA  
  
# Compare the all categorical variables  
all_var <- compare_category(carseats)  
  
# Print compare class object  
all_var  
  
# Compare the two categorical variables  
two_var <- compare_category(carseats, ShelveLoc, Urban)  
  
# Print compare class object  
two_var  
  
# plot all pair of variables  
plot(all_var)
```

```
# plot a pair of variables
plot(two_var)

# plot all pair of variables by prompt
plot(all_var, prompt = TRUE)

# plot a pair of variables
plot(two_var, las = 1)
```

plot.compare_numeric *Visualize Information for an "compare_numeric" Object*

Description

Visualize scatter plot included box plots by attribute of compare_numeric class.

Usage

```
## S3 method for class 'compare_numeric'
plot(x, prompt = FALSE, ...)
```

Arguments

x	an object of class "compare_numeric", usually, a result of a call to compare_numeric().
prompt	logical. The default value is FALSE. If there are multiple visualizations to be output, if this argument value is TRUE, a prompt is output each time.
...	arguments to be passed to methods, such as graphical parameters (see par). However, it does not support.

See Also

[compare_numeric](#), [print.compare_numeric](#), [summary.compare_numeric](#).

Examples

```
# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

# Compare the all numerical variables
all_var <- compare_numeric(carseats)

# Print compare compare_numeric object
all_var

# Compare the two numerical variables
two_var <- compare_numeric(carseats, CompPrice, Price)
```

```

# Print compare_numeric class object
two_var

# plot all pair of variables
plot(all_var)

# plot a pair of variables
plot(two_var)

# plot all pair of variables by prompt
plot(all_var, prompt = TRUE)

```

plot.imputation	<i>Visualize Information for an "imputation" Object</i>
-----------------	---

Description

Visualize two kinds of plot by attribute of 'imputation' class. The imputation of a numerical variable is a density plot, and the imputation of a categorical variable is a bar plot.

Usage

```

## S3 method for class 'imputation'
plot(x, ...)

```

Arguments

x	an object of class "imputation", usually, a result of a call to <code>impute_na()</code> or <code>impute_outlier()</code> .
...	arguments to be passed to methods, such as graphical parameters (see <code>par</code>). only applies when the model argument is TRUE, and is used for ... of the <code>plot.lm()</code> function.

See Also

[impute_na](#), [impute_outlier](#), [summary.imputation](#).

Examples

```

# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

# Impute missing values -----
# If the variable of interest is a numerical variable

```

```

income <- imputate_na(carseats, Income, US, method = "rpart")
income
summary(income)
plot(income)

# If the variable of interest is a categorical variable
urban <- imputate_na(carseats, Urban, US, method = "mice")
urban
summary(urban)
plot(urban)

# Impute outliers -----
# If the variable of interest is a numerical variable
price <- imputate_outlier(carseats, Price, method = "capping")
price
summary(price)
plot(price)

```

plot.optimal_bins *Visualize Distribution for an "optimal_bins" Object*

Description

It generates plots for understand distribution, bad rate, and weight of evidence after running smbinning and saving its output.

See vignette("transformation") for an introduction to these concepts.

Usage

```

## S3 method for class 'optimal_bins'
plot(x, type = c("dist", "goodrate", "badrate", "WoE"), sub = "", ...)

```

Arguments

x	an object of class "optimal_bins", usually, a result of a call to binning_by().
type	character. options for visualization. Distribution ("dist"), Good Rate ("goodrate"), Bad Rate ("badrate"), and Weight of Evidence ("WoE").
sub	character. sub title for the chart (optional).
...	arguments to be passed to methods, such as graphical parameters (see par). only applies to the first graph that is implemented with the boxplot() function.

See Also

[binning_by](#), [plot.bins](#), [smbinning.plot](#).

Examples

```

# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

# optimal binning
bin <- binning_by(carseats, "US", "Advertising")
bin

# summary optimal_bins class
summary(bin)

# information value
attr(bin, "iv")

# information value table
attr(bin, "ivtable")

# visualize optimal_bins class
plot(bin, sub = "bins of Advertising variable")

```

plot.relate

Visualize Information for an "relate" Object

Description

Visualize four kinds of plot by attribute of relate class.

Usage

```

## S3 method for class 'relate'
plot(
  x,
  model = FALSE,
  hex_thres = 1000,
  pal = RColorBrewer::brewer.pal(7, "YlOrRd"),
  ...
)

```

Arguments

x an object of class "relate", usually, a result of a call to relate().

model logical. This argument selects whether to output the visualization result to the visualization of the object of the lm model to grasp the relationship between the numerical variables.

hex_thres an integer. Use only when the target and predictor are numeric variables. Used when the number of observations is large. Specify the threshold of the observations to draw hexabin plots that are not scatterplots. The default value is 1000.

pal Color palette to paint hexabin. Use only when the target and predictor are numeric variables. Applied only when the number of observations is greater than hex_thres.

... arguments to be passed to methods, such as graphical parameters (see par). only applies when the model argument is TRUE, and is used for ... of the plot.lm() function.

See Also

[relate](#), [print.relate](#).

Examples

```
# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

# If the target variable is a categorical variable
categ <- target_by(carseats, US)

# If the variable of interest is a numerical variable
cat_num <- relate(categ, Sales)
cat_num
summary(cat_num)
plot(cat_num)

# If the variable of interest is a categorical variable
cat_cat <- relate(categ, ShelveLoc)
cat_cat
summary(cat_cat)
plot(cat_cat)

##-----
# If the target variable is a categorical variable
num <- target_by(carseats, Sales)

# If the variable of interest is a numerical variable
num_num <- relate(num, Price)
num_num
summary(num_num)
plot(num_num)
plot(num_num, hex_thres = 400)

# If the variable of interest is a categorical variable
num_cat <- relate(num, ShelveLoc)
num_cat
summary(num_cat)
plot(num_cat)
```

plot.transform	<i>Visualize Information for an "transform" Object</i>
----------------	--

Description

Visualize two kinds of plot by attribute of 'transform' class. The transformation of a numerical variable is a density plot.

Usage

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

Arguments

x an object of class "transform", usually, a result of a call to transform().
... arguments to be passed to methods, such as graphical parameters (see par).

See Also

[transform](#), [summary.transform](#).

Examples

```
# Generate data for the example  
carseats <- ISLR::Carseats  
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA  
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA  
  
# Standardization -----  
advertising_minmax <- transform(carseats$Advertising, method = "minmax")  
advertising_minmax  
summary(advertising_minmax)  
plot(advertising_minmax)  
  
# Resolving Skewness -----  
advertising_log <- transform(carseats$Advertising, method = "log")  
advertising_log  
summary(advertising_log)  
plot(advertising_log)
```

plot.univar_category *Visualize Information for an "univar_category" Object*

Description

Visualize mosaics plot by attribute of univar_category class.

Usage

```
## S3 method for class 'univar_category'  
plot(x, na.rm = TRUE, prompt = FALSE, ...)
```

Arguments

x	an object of class "univar_category", usually, a result of a call to univar_category().
na.rm	logical. Specifies whether to include NA when plotting bar plot. The default is FALSE, so plot NA.
prompt	logical. The default value is FALSE. If there are multiple visualizations to be output, if this argument value is TRUE, a prompt is output each time.
...	arguments to be passed to methods, such as graphical parameters (see par). However, it does not support all parameters.

See Also

[univar_category](#), [print.univar_category](#), [summary.univar_category](#).

Examples

```
# Generate data for the example  
carseats <- ISLR::Carseats  
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA  
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA  
  
library(dplyr)  
  
# Calculates the all categorical variables  
all_var <- univar_category(carseats)  
  
# Print univar_category class object  
all_var  
  
# Calculates the only Urban variable  
urban <- univar_category(carseats, Urban)  
  
# Print univar_category class object  
urban  
  
# plot all variables
```

```

plot(all_var)

# plot urban
plot(urban)

# plot all variables by na.rm = FALSE
plot(all_var, na.rm = FALSE)

# plot all variables by prompt
plot(all_var, prompt = TRUE)

```

plot.univar_numeric *Visualize Information for an "univar_numeric" Object*

Description

Visualize boxplots and histogram by attribute of univar_numeric class.

Usage

```

## S3 method for class 'univar_numeric'
plot(
  x,
  indiv = FALSE,
  viz = c("hist", "boxplot"),
  stand = ifelse(rep(indiv, 4), c("none", "robust", "minmax", "zscore"), c("robust",
    "minmax", "zscore", "none")),
  prompt = FALSE,
  ...
)

```

Arguments

x	an object of class "univar_numeric", usually, a result of a call to univar_numeric().
indiv	logical. Select whether to display information of all variables in one plot when there are multiple selected numeric variables. In case of FALSE, all variable information is displayed in one plot. If TRUE, the information of the individual variables is output to the individual plots. The default is FALSE. If only one variable is selected, TRUE is applied.
viz	character. Describe what to plot visualization. "hist" draws a histogram and "boxplot" draws a boxplot. The default is "hist".
stand	character. Describe how to standardize the original data. "robust" normalizes the raw data through transformation calculated by IQR and median. "minmax" normalizes the original data using minmax transformation. "zscore" standardizes the original data using z-Score transformation. "none" does not perform data transformation. The default is "none" if indiv is TRUE, and "robust" if FALSE.

prompt logical. The default value is FALSE. If there are multiple visualizations to be output, if this argument value is TRUE, a prompt is output each time.

... arguments to be passed to methods, such as graphical parameters (see `par`). However, it does not support.

See Also

[univar_numeric](#), [print.univar_numeric](#), [summary.univar_numeric](#).

Examples

```
# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

# Calculates the all categorical variables
all_var <- univar_numeric(carseats)

# Print univar_numeric class object
all_var

# Summary the all case : Return a invisible copy of an object.
stat <- summary(all_var)

# Summary by returned object
stat

# one plot with all variables
plot(all_var)

# one plot with all normalized variables by Min-Max method
plot(all_var, stand = "minmax")

# one plot with all variables
plot(all_var, stand = "none")

# one plot with all robust standardized variables
plot(all_var, viz = "boxplot")

# one plot with all standardized variables by Z-score method
plot(all_var, viz = "boxplot", stand = "zscore")

# individual boxplot by variables
plot(all_var, indiv = TRUE, "boxplot")

# individual histogram by variables
plot(all_var, indiv = TRUE, "hist")

# individual histogram by robust standardized variable
plot(all_var, indiv = TRUE, "hist", stand = "robust")
```

```
# plot all variables by prompt
plot(all_var, indiv = TRUE, "hist", prompt = TRUE)
```

plot_correlate	<i>Visualize correlation plot of numerical data</i>
----------------	---

Description

The `plot_correlate()` visualize correlation plot for find relationship between two numerical variables.

Usage

```
plot_correlate(.data, ...)
```

```
## S3 method for class 'data.frame'
plot_correlate(.data, ..., method = c("pearson", "kendall", "spearman"))
```

Arguments

<code>.data</code>	a <code>data.frame</code> or a <code>tbl_df</code> .
<code>...</code>	one or more unquoted expressions separated by commas. You can treat variable names like they are positions. Positive values select variables; negative values to drop variables. If the first expression is negative, <code>plot_correlate()</code> will automatically start with all variables. These arguments are automatically quoted and evaluated in a context where column names represent column positions. They support unquoting and splicing. See vignette("EDA") for an introduction to these concepts.
<code>method</code>	a character string indicating which correlation coefficient (or covariance) is to be computed. One of "pearson" (default), "kendall", or "spearman": can be abbreviated.

Details

The scope of the visualization is the provide a correlation information. Since the plot is drawn for each variable, if you specify more than one variable in the `...` argument, the specified number of plots are drawn.

See Also

[plot_correlate.tbl_dbi](#), [plot_outlier.data.frame](#).

Examples

```

# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

# Visualize correlation plot of all numerical variables
plot_correlate(carseats)

# Select the variable to compute
plot_correlate(carseats, Sales, Price)
plot_correlate(carseats, -Sales, -Price)
plot_correlate(carseats, "Sales", "Price")
plot_correlate(carseats, 1)
plot_correlate(carseats, Sales, Price, method = "spearman")

# Using dplyr::grouped_dt
library(dplyr)

gdata <- group_by(carseats, ShelveLoc, US)
plot_correlate(gdata, "Sales")
plot_correlate(gdata)

# Using pipes -----
# Visualize correlation plot of all numerical variables
carseats %>%
  plot_correlate()
# Positive values select variables
carseats %>%
  plot_correlate(Sales, Price)
# Negative values to drop variables
carseats %>%
  plot_correlate(-Sales, -Price)
# Positions values select variables
carseats %>%
  plot_correlate(1)
# Positions values select variables
carseats %>%
  plot_correlate(-1, -2, -3, -5, -6)

# Using pipes & dplyr -----
# Visualize correlation plot of 'Sales' variable by 'ShelveLoc'
# and 'US' variables.
carseats %>%
  group_by(ShelveLoc, US) %>%
  plot_correlate(Sales)

# Extract only those with 'ShelveLoc' variable level is "Good",
# and visualize correlation plot of 'Sales' variable by 'Urban'
# and 'US' variables.
carseats %>%
  filter(ShelveLoc == "Good") %>%

```

```
group_by(Urban, US) %>%
plot_correlate(Sales)
```

```
plot_correlate.tbl_dbi
```

Visualize correlation plot of numerical data

Description

The `plot_correlate()` visualize correlation plot for find relationship between two numerical (INTEGER, NUMBER, etc.) column of the DBMS table through `tbl_dbi`.

Usage

```
## S3 method for class 'tbl_dbi'
plot_correlate(
  .data,
  ...,
  in_database = FALSE,
  collect_size = Inf,
  method = c("pearson", "kendall", "spearman")
)
```

Arguments

<code>.data</code>	a <code>tbl_dbi</code> .
<code>...</code>	one or more unquoted expressions separated by commas. You can treat variable names like they are positions. Positive values select variables; negative values to drop variables. If the first expression is negative, <code>plot_correlate()</code> will automatically start with all variables. These arguments are automatically quoted and evaluated in a context where column names represent column positions. They support unquoting and splicing.
<code>in_database</code>	Specifies whether to perform in-database operations. If <code>TRUE</code> , most operations are performed in the DBMS. if <code>FALSE</code> , table data is taken in R and operated in-memory. Not yet supported <code>in_database = TRUE</code> .
<code>collect_size</code>	a integer. The number of data samples from the DBMS to R. Applies only if <code>in_database = FALSE</code> .
<code>method</code>	a character string indicating which correlation coefficient (or covariance) is to be computed. One of "pearson" (default), "kendall", or "spearman": can be abbreviated. See vignette("EDA") for an introduction to these concepts.

Details

The scope of the visualization is the provide a correlation information. Since the plot is drawn for each variable, if you specify more than one variable in the `...` argument, the specified number of plots are drawn.

See Also

[plot_correlate.data.frame](#), [plot_outlier.tbl_dbi](#).

Examples

```

library(dplyr)

# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

# connect DBMS
con_sqlite <- DBI::dbConnect(RSQLite::SQLite(), ":memory:")

# copy carseats to the DBMS with a table named TB_CARSEATS
copy_to(con_sqlite, carseats, name = "TB_CARSEATS", overwrite = TRUE)

# Using pipes -----
# Visualize correlation plot of all numerical variables
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  plot_correlate()

# Positive values select variables, and In-memory mode and collect size is 200
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  plot_correlate(Sales, Price, collect_size = 200)

# Negative values to drop variables
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  plot_correlate(-Sales, -Price)

# Positions values select variables
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  plot_correlate(1)

# Positions values select variables
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  plot_correlate(-1, -2, -3, -5, -6)

# Using pipes & dplyr -----
# Visualize correlation plot of 'Sales' variable by 'ShelveLoc'
# and 'US' variables.
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  group_by(ShelveLoc, US) %>%
  plot_correlate(Sales)

```



```
# Extract only those with 'ShelveLoc' variable level is "Good",
# and visualize correlation plot of 'Sales' variable by 'Urban'
# and 'US' variables.
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  filter(ShelveLoc == "Good") %>%
  group_by(Urban, US) %>%
  plot_correlate(Sales)
```

plot_na_hclust

Combination chart for missing value

Description

Visualize distribution of missing value by combination of variables.

Usage

```
plot_na_hclust(x, main = NULL, col.left = "#009E73", col.right = "#56B4E9")
```

Arguments

x	data frames, or objects to be coerced to one.
main	character. Main title.
col.left	character. The color of left legend that is frequency of NA. default is "#009E73".
col.right	character. The color of right legend that is percentage of NA. default is "#56B4E9".

Details

Rows are variables containing missing values, and columns are observations. These data structures were grouped into similar groups by applying hclust. So, it was made possible to visually examine how the missing values are distributed for each combination of variables.

Examples

```
# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

# Visualize pareto chart for variables with missing value.
plot_na_hclust(carseats)
plot_na_hclust(airquality)

# Visualize pareto chart for variables with missing value.
plot_na_hclust(mice::boys)
```

```
# Change the main title.
plot_na_hclust(mice::boys, main = "Distribution of missing value")
```

plot_na_intersect *Plot the combination variables that is include missing value*

Description

Visualize the combinations of missing value across cases.

Usage

```
plot_na_intersect(
  x,
  only_na = TRUE,
  n_intersects = NULL,
  n_vars = NULL,
  main = NULL
)
```

Arguments

x	data frames, or objects to be coerced to one.
only_na	logical. The default value is FALSE. If TRUE, only variables containing missing values are selected for visualization. If FALSE, included complete case.
n_intersects	integer. Specifies the number of combinations of variables including missing values. The combination of variables containing many missing values is chosen first.
n_vars	integer. Specifies the number of variables that contain missing values to be visualized. The default value is NULL, which visualizes variables containing all missing values. If this value is greater than the number of variables containing missing values, all variables containing missing values are visualized. Variables containing many missing values are chosen first.
main	character. Main title.

Details

The visualization consists of four parts. The bottom left, which is the most basic, visualizes the case of cross(intersection)-combination. The x-axis is the variable including the missing value, and the y-axis represents the case of a combination of variables. And on the marginal of the two axes, the frequency of the case is expressed as a bar graph. Finally, the visualization at the top right expresses the number of variables including missing values in the data set, and the number of observations including missing values and complete cases .

Examples

```

# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

# Visualize the combination variables that is include missing value.
plot_na_intersect(carseats)

# Diagnose the data with missing_count using diagnose() function
library(dplyr)

mice::boys %>%
  diagnose %>%
  arrange(desc(missing_count))

# Visualize the combination variables that is include missing value
plot_na_intersect(mice::boys)

# Visualize variables containing missing values and complete case
plot_na_intersect(mice::boys, only_na = FALSE)

# Using n_vars argument
plot_na_intersect(mice::boys, n_vars = 5)

# Using n_intersacts argument
plot_na_intersect(mice::boys, only_na = FALSE, n_intersacts = 7)

```

plot_na_pareto	<i>Pareto chart for missing value</i>
----------------	---------------------------------------

Description

Visualize pareto chart for variables with missing value.

Usage

```

plot_na_pareto(
  x,
  only_na = FALSE,
  relative = FALSE,
  main = NULL,
  col = "black",
  grade = list(Good = 0.05, OK = 0.4, Bad = 0.8, Remove = 1),
  plot = TRUE
)

```

Arguments

x	data frames, or objects to be coerced to one.
only_na	logical. The default value is FALSE. If TRUE, only variables containing missing values are selected for visualization. If FALSE, all variables are included.
relative	logical. If this argument is TRUE, it sets the unit of the left y-axis to relative frequency. In case of FALSE, set it to frequency.
main	character. Main title.
col	character. The color of line for display the cumulative percentage.
grade	list. Specifies the cut-off to set the grade of the variable according to the ratio of missing values. The default values are Good: [0, 0.05], OK: (0.05, 0.4], Bad: (0.4, 0.8], Remove: (0.8, 1].
plot	logical. If this value is TRUE then visualize plot. else if FALSE, return aggregate information about missing values.

Examples

```
# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

# Diagnose the data with missing_count using diagnose() function
library(dplyr)
carseats %>%
  diagnose %>%
  arrange(desc(missing_count))

# Visualize pareto chart for variables with missing value.
plot_na_pareto(carseats)
plot_na_pareto(airquality)

# Diagnose the data with missing_count using diagnose() function
mice::boys %>%
  diagnose %>%
  arrange(desc(missing_count))

# Visualize pareto chart for variables with missing value.
plot_na_pareto(mice::boys, col = "darkorange")

# Visualize only variables containing missing values
plot_na_pareto(mice::boys, only_na = TRUE)

# Display the relative frequency
plot_na_pareto(mice::boys, relative = TRUE)

# Change the grade
plot_na_pareto(mice::boys, grade = list(High = 0.1, Middle = 0.6, Low = 1))

# Change the main title.
```

```

plot_na_pareto(mice::boys, relative = TRUE, only_na = TRUE,
main = "Pareto Chart for mice::boys")

# Return the aggregate information about missing values.
plot_na_pareto(mice::boys, only_na = TRUE, plot = FALSE)

```

plot_normality	<i>Plot distribution information of numerical data</i>
----------------	--

Description

The `plot_normality()` visualize distribution information for normality test of the numerical data.

Usage

```

plot_normality(.data, ...)

## S3 method for class 'data.frame'
plot_normality(.data, ...)

```

Arguments

<code>.data</code>	a <code>data.frame</code> or a <code>tbl_df</code> .
<code>...</code>	one or more unquoted expressions separated by commas. You can treat variable names like they are positions. Positive values select variables; negative values to drop variables. If the first expression is negative, <code>plot_normality()</code> will automatically start with all variables. These arguments are automatically quoted and evaluated in a context where column names represent column positions. They support unquoting and splicing. See vignette("EDA") for an introduction to these concepts.

Details

The scope of the visualization is the provide a distribution information. Since the plot is drawn for each variable, if you specify more than one variable in the `...` argument, the specified number of plots are drawn.

Distribution information

The plot derived from the numerical data visualization is as follows.

- histogram by original data
- q-q plot by original data
- histogram by log transfer data
- histogram by square root transfer data

See Also

[plot_normality.tbl_dbi](#), [plot_outlier.data.frame](#).

Examples

```
# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

# Visualization of all numerical variables
plot_normality(carseats)

# Select the variable to plot
plot_normality(carseats, Income, Price)
plot_normality(carseats, -Income, -Price)
plot_normality(carseats, 1)

# Using dplyr::grouped_df
library(dplyr)

gdata <- group_by(carseats, ShelveLoc, US)
plot_normality(carseats)
plot_normality(carseats, "Sales")

# Using pipes -----
# Visualization of all numerical variables
carseats %>%
  plot_normality()

# Positive values select variables
carseats %>%
  plot_normality(Income, Price)

# Positions values select variables
carseats %>%
  plot_normality(1)

# Using pipes & dplyr -----
# Plot 'Sales' variable by 'ShelveLoc' and 'US'
carseats %>%
  group_by(ShelveLoc, US) %>%
  plot_normality(Sales)

# extract only those with 'ShelveLoc' variable level is "Good",
# and plot 'Income' by 'US'
carseats %>%
  filter(ShelveLoc == "Good") %>%
  group_by(US) %>%
  plot_normality(Income)
```

 plot_normality.tbl_dbi

Plot distribution information of numerical data

Description

The `plot_normality()` visualize distribution information for normality test of the numerical (INTEGER, NUMBER, etc.) column of the DBMS table through `tbl_dbi`.

Usage

```
## S3 method for class 'tbl_dbi'
plot_normality(.data, ..., in_database = FALSE, collect_size = Inf)
```

Arguments

<code>.data</code>	a <code>tbl_dbi</code> .
<code>...</code>	one or more unquoted expressions separated by commas. You can treat variable names like they are positions. Positive values select variables; negative values to drop variables. If the first expression is negative, <code>plot_normality()</code> will automatically start with all variables. These arguments are automatically quoted and evaluated in a context where column names represent column positions. They support unquoting and splicing.
<code>in_database</code>	Specifies whether to perform in-database operations. If TRUE, most operations are performed in the DBMS. if FALSE, table data is taken in R and operated in-memory. Not yet supported in <code>in_database = TRUE</code> .
<code>collect_size</code>	a integer. The number of data samples from the DBMS to R. Applies only if <code>in_database = FALSE</code> . See <code>vignette("EDA")</code> for an introduction to these concepts.

Details

The scope of the visualization is the provide a distribution information. Since the plot is drawn for each variable, if you specify more than one variable in the `...` argument, the specified number of plots are drawn.

Distribution information

The plot derived from the numerical data visualization is as follows.

- histogram by original data
- q-q plot by original data
- histogram by log transfer data
- histogram by square root transfer data

See Also

[plot_normality.data.frame](#), [plot_outlier.tbl_dbi](#).

Examples

```
library(dplyr)

# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

# connect DBMS
con_sqlite <- DBI::dbConnect(RSQLite::SQLite(), ":memory:")

# copy carseats to the DBMS with a table named TB_CARSEATS
copy_to(con_sqlite, carseats, name = "TB_CARSEATS", overwrite = TRUE)

# Using pipes -----
# Visualization of all numerical variables
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  plot_normality()

# Positive values select variables, and In-memory mode and collect size is 200
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  plot_normality(Income, Price, collect_size = 200)

# Positions values select variables
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  plot_normality(1)

# Using pipes & dplyr -----
# Plot 'Sales' variable by 'ShelveLoc' and 'US'
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  group_by(ShelveLoc, US) %>%
  plot_normality(Sales)

# extract only those with 'ShelveLoc' variable level is "Good",
# and plot 'Income' by 'US'
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  filter(ShelveLoc == "Good") %>%
  group_by(US) %>%
  plot_normality(Income)
```

plot_outlier	<i>Plot outlier information of numerical data diagnosis</i>
--------------	---

Description

The `plot_outlier()` visualize outlier information for diagnosing the quality of the numerical data.

Usage

```
plot_outlier(.data, ...)  
  
## S3 method for class 'data.frame'  
plot_outlier(.data, ..., col = "lightblue")
```

Arguments

<code>.data</code>	a <code>data.frame</code> or a <code>tbl_df</code> .
<code>...</code>	one or more unquoted expressions separated by commas. You can treat variable names like they are positions. Positive values select variables; negative values to drop variables. If the first expression is negative, <code>plot_outlier()</code> will automatically start with all variables. These arguments are automatically quoted and evaluated in a context where column names represent column positions. They support unquoting and splicing.
<code>col</code>	a color to be used to fill the bars. The default is "lightblue".

Details

The scope of the diagnosis is the provide a outlier information. Since the plot is drawn for each variable, if you specify more than one variable in the `...` argument, the specified number of plots are drawn.

Outlier diagnostic information

The plot derived from the numerical data diagnosis is as follows.

- With outliers box plot
- Without outliers box plot
- With outliers histogram
- Without outliers histogram

See vignette("diagonosis") for an introduction to these concepts.

See Also

[plot_outlier.tbl_dbi](#), [diagnose_outlier.data.frame](#).

Examples

```

# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

# Visualization of all numerical variables
plot_outlier(carseats)

# Select the variable to diagnose
plot_outlier(carseats, Sales, Price)
plot_outlier(carseats, -Sales, -Price)
plot_outlier(carseats, "Sales", "Price")
plot_outlier(carseats, 6)

# Using the col argument
plot_outlier(carseats, Sales, col = "gray")

# Using pipes -----
library(dplyr)

# Visualization of all numerical variables
carseats %>%
  plot_outlier()
# Positive values select variables
carseats %>%
  plot_outlier(Sales, Price)
# Negative values to drop variables
carseats %>%
  plot_outlier(-Sales, -Price)
# Positions values select variables
carseats %>%
  plot_outlier(6)
# Positions values select variables
carseats %>%
  plot_outlier(-1, -5)

# Using pipes & dplyr -----
# Visualization of numerical variables with a ratio of
# outliers greater than 1%
carseats %>%
  plot_outlier(carseats %>%
    diagnose_outlier() %>%
    filter(outliers_ratio > 1) %>%
    select(variables) %>%
    pull())

```

Description

The `plot_outlier()` visualize outlier information for diagnosing the quality of the numerical (INTEGER, NUMBER, etc.) column of the DBMS table through `tbl_dbi`.

Usage

```
## S3 method for class 'tbl_dbi'
plot_outlier(
  .data,
  ...,
  col = "lightblue",
  in_database = FALSE,
  collect_size = Inf
)
```

Arguments

<code>.data</code>	a <code>tbl_dbi</code> .
<code>...</code>	one or more unquoted expressions separated by commas. You can treat variable names like they are positions. Positive values select variables; negative values to drop variables. If the first expression is negative, <code>plot_outlier()</code> will automatically start with all variables. These arguments are automatically quoted and evaluated in a context where column names represent column positions. They support unquoting and splicing.
<code>col</code>	a color to be used to fill the bars. The default is "lightblue".
<code>in_database</code>	Specifies whether to perform in-database operations. If TRUE, most operations are performed in the DBMS. if FALSE, table data is taken in R and operated in-memory. Not yet supported in <code>in_database = TRUE</code> .
<code>collect_size</code>	a integer. The number of data samples from the DBMS to R. Applies only if <code>in_database = FALSE</code> .

Details

The scope of the diagnosis is the provide a outlier information. Since the plot is drawn for each variable, if you specify more than one variable in the `...` argument, the specified number of plots are drawn.

Outlier diagnostic information

The plot derived from the numerical data diagnosis is as follows.

- With outliers box plot
- Without outliers box plot
- With outliers histogram
- Without outliers histogram

See `vignette("diagonosis")` for an introduction to these concepts.

See Also

[plot_outlier.data.frame](#), [diagnose_outlier.tbl_dbi](#).

Examples

```
library(dplyr)

# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

# connect DBMS
con_sqlite <- DBI::dbConnect(RSQLite::SQLite(), ":memory:")

# copy carseats to the DBMS with a table named TB_CARSEATS
copy_to(con_sqlite, carseats, name = "TB_CARSEATS", overwrite = TRUE)

# Using pipes -----
# Visualization of all numerical variables
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  plot_outlier()

# Positive values select variables
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  plot_outlier(Sales, Price)

# Negative values to drop variables, and In-memory mode and collect size is 200
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  plot_outlier(-Sales, -Price, collect_size = 200)

# Positions values select variables
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  plot_outlier(6)

# Positions values select variables
carseats %>%
  plot_outlier(-1, -5)

# Using pipes & dplyr -----
# Visualization of numerical variables with a ratio of
# outliers greater than 1%
con_sqlite %>%
  tbl("TB_CARSEATS") %>%
  plot_outlier(con_sqlite %>%
    tbl("TB_CARSEATS") %>%
    diagnose_outlier() %>%
    filter(outliers_ratio > 1) %>%
```

```
select(variables) %>%
pull()
```

print.relate	<i>Summarizing relate information</i>
--------------	---------------------------------------

Description

print and summary method for "relate" class.

Usage

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

Arguments

x an object of class "relate", usually, a result of a call to relate().
... further arguments passed to or from other methods.

Details

print.relate() tries to be smart about formatting four kinds of relate. summary.relate() tries to be smart about formatting four kinds of relate.

See Also

[plot.relate](#).

Examples

```
## Not run:
# Generate data for the example
diamonds2 <- diamonds
diamonds2[sample(seq(NROW(diamonds2)), 250), "price"] <- NA
diamonds2[sample(seq(NROW(diamonds2)), 20), "clarity"] <- NA

# Binning the carat variable. default type argument is "quantile"
bin <- binning(diamonds2$carat)

# Print bins class object
bin

# Summarize bins class object
summary(bin)

## End(Not run)
```

```

# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

# If the target variable is a categorical variable
categ <- target_by(carseats, US)

# If the variable of interest is a numerical variable
cat_num <- relate(categ, Sales)
cat_num
summary(cat_num)
plot(cat_num)

# If the variable of interest is a categorical variable
cat_cat <- relate(categ, ShelveLoc)
cat_cat
summary(cat_cat)
plot(cat_cat)

##-----
# If the target variable is a categorical variable
num <- target_by(carseats, Sales)

# If the variable of interest is a numerical variable
num_num <- relate(num, Price)
num_num
summary(num_num)
plot(num_num)

# If the variable of interest is a categorical variable
num_cat <- relate(num, ShelveLoc)
num_cat
summary(num_cat)
plot(num_cat)

```

relate

Relationship between target variable and variable of interest

Description

The relationship between the target variable and the variable of interest (predictor) is briefly analyzed.

Usage

```
relate(.data, predictor)
```

Arguments

.data A target_df.
 predictor variable of interest. predictor.
 See vignette("relate") for an introduction to these concepts.

Details

Returns the four types of results that correspond to the combination of the target variable and the data type of the variable of interest.

- target variable: categorical variable
 - predictor: categorical variable
 - * contingency table
 - * c("xtabs", "table") class
 - predictor: numerical variable
 - * descriptive statistic for each levels and total observation.
- target variable: numerical variable
 - predictor: categorical variable
 - * ANOVA test. "lm" class.
 - predictor: numerical variable
 - * simple linear model. "lm" class.

Value

An object of the class as relate. Attributes of relate class is as follows.

- target : name of target variable
- predictor : name of predictor
- model : levels of binned value.
- raw : table_df with two variables target and predictor.

Descriptive statistic information

The information derived from the numerical data describe is as follows.

- mean : arithmetic average
- sd : standard deviation
- se_mean : standrd error mean. sd/\sqrt{n}
- IQR : interqurtle range (Q3-Q1)
- skewness : skewness
- kurtosis : kurtosis
- p25 : Q1. 25% percentile
- p50 : median. 50% percentile

- p75 : Q3. 75% percentile
- p01, p05, p10, p20, p30 : 1%, 5%, 20%, 30% percentiles
- p40, p60, p70, p80 : 40%, 60%, 70%, 80% percentiles
- p90, p95, p99, p100 : 90%, 95%, 99%, 100% percentiles

See Also

[print.relate](#), [plot.relate](#).

Examples

```
# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

# If the target variable is a categorical variable
categ <- target_by(carseats, US)

# If the variable of interest is a numerical variable
cat_num <- relate(categ, Sales)
cat_num
summary(cat_num)
plot(cat_num)

# If the variable of interest is a categorical variable
cat_cat <- relate(categ, ShelveLoc)
cat_cat
summary(cat_cat)
plot(cat_cat)

##-----
# If the target variable is a categorical variable
num <- target_by(carseats, Sales)

# If the variable of interest is a numerical variable
num_num <- relate(num, Price)
num_num
summary(num_num)
plot(num_num)

# If the variable of interest is a categorical variable
num_cat <- relate(num, ShelveLoc)
num_cat
summary(num_cat)
plot(num_cat)
```

skewness	<i>Skewness of the data</i>
----------	-----------------------------

Description

This function calculated skewness of given data.

Usage

```
skewness(x, na.rm = TRUE)
```

Arguments

x	a numeric vector.
na.rm	logical. Determine whether to remove missing values and calculate them. The default is TRUE.

Value

numeric. calculated skewness.

See Also

[kurtosis](#), [find_skewness](#).

Examples

```
set.seed(123)
skewness(rnorm(100))
```

summary.bins	<i>Summarizing Binned Variable</i>
--------------	------------------------------------

Description

summary method for "bins" and "optimal_bins".

Usage

```
## S3 method for class 'bins'
summary(object, ...)

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

Arguments

object	an object of "bins" and "optimal_bins", usually, a result of a call to binning().
...	further arguments passed to or from other methods.
x	an object of class "bins" and "optimal_bins", usually, a result of a call to binning().

Details

print.bins() prints the information of "bins" and "optimal_bins" objects nicely. This includes frequency of bins, binned type, and number of bins. summary.bins() returns data.frame including frequency and relative frequency for each levels(bins).

See vignette("transformation") for an introduction to these concepts.

Value

The function summary.bins() computes and returns a data.frame of summary statistics of the binned given in object. Variables of data frame is as follows.

- levels : levels of factor.
- freq : frequency of levels.
- rate : relative frequency of levels. it is not percentage.

See Also

[binning](#)

Examples

```
# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

# Binning the carat variable. default type argument is "quantile"
bin <- binning(carseats$Income)

# Print bins class object
bin

# Summarise bins class object
summary(bin)
```

summary.compare_category

Summarizing compare_category information

Description

print and summary method for "compare_category" class.

Usage

```
## S3 method for class 'compare_category'
summary(
  object,
  method = c("all", "table", "relative", "chisq"),
  pos = NULL,
  na.rm = TRUE,
  marginal = FALSE,
  verbose = TRUE,
  ...
)

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

Arguments

object	an object of class "compare_category", usually, a result of a call to compare_category().
method	character. Specifies the type of information to be aggregated. "table" create contingency table, "relative" create relative contingency table, and "chisq" create information of chi-square test. and "all" aggregates all information. The default is "all"
pos	integer. Specifies the pair of variables to be summarized by index. The default is NULL, which aggregates all variable pairs.
na.rm	logical. Specifies whether to include NA when counting the contingency tables or performing a chi-square test. The default is TRUE, where NA is removed and aggregated.
marginal	logical. Specifies whether to add marginal values to the contingency table. The default value is FALSE, so no marginal value is added.
verbose	logical. Specifies whether to output additional information during the calculation process. The default is to output information as TRUE. In this case, the function returns the value with invisible(). If FALSE, the value is returned by return().
...	further arguments passed to or from other methods.
x	an object of class "compare_category", usually, a result of a call to compare_category().

Details

print.compare_category() displays only the information compared between the variables included in compare_category. The "type", "variables" and "combination" attributes are not displayed. When using summary.compare_category(), it is advantageous to set the verbose argument to TRUE if the user is only viewing information from the console. It is also advantageous to specify FALSE if you want to manipulate the results.

See Also

[plot.compare_category.](#)

Examples

```
# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

library(dplyr)

# Compare the all categorical variables
all_var <- compare_category(carseats)

# Print compare class object
all_var

# Compare the two categorical variables
two_var <- compare_category(carseats, ShelveLoc, Urban)

# Print compare class object
two_var

# Summary the all case : Return a invisible copy of an object.
stat <- summary(all_var)

# Summary by returned object
stat

# component of table
stat$table

# component of chi-square test
stat$chisq

# component of chi-square test
summary(all_var, "chisq")

# component of chi-square test (first, third case)
summary(all_var, "chisq", pos = c(1, 3))

# component of relative frequency table
summary(all_var, "relative")
```

```

# component of table without missing values
summary(all_var, "table", na.rm = TRUE)

# component of table include marginal value
margin <- summary(all_var, "table", marginal = TRUE)
margin

# component of chi-square test
summary(two_var, method = "chisq")

# verbose is FALSE
summary(all_var, "chisq", verbose = FALSE)

#' # Using pipes & dplyr -----
# If you want to use dplyr, set verbose to FALSE
summary(all_var, "chisq", verbose = FALSE) %>%
  filter(p.value < 0.26)

# Extract component from list by index
summary(all_var, "table", na.rm = TRUE, verbose = FALSE) %>%
  "[ "(1)

# Extract component from list by name
summary(all_var, "table", na.rm = TRUE, verbose = FALSE) %>%
  "[ "(ShelveLoc vs Urban)

```

```
summary.compare_numeric
```

Summarizing compare_numeric information

Description

print and summary method for "compare_numeric" class.

Usage

```

## S3 method for class 'compare_numeric'
summary(
  object,
  method = c("all", "correlation", "linear"),
  thres_corr = 0.3,
  thres_rs = 0.1,
  verbose = TRUE,
  ...
)

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

```

Arguments

object	an object of class "compare_numeric", usually, a result of a call to compare_numeric().
method	character. Select statistics to be aggregated. "correlation" calculates the Pearson's correlation coefficient, and "linear" returns the aggregation of the linear model. "all" returns both information. However, the difference between summary.compare_numeric() and compare_numeric() is that only cases that are greater than the specified threshold are returned. "correlation" returns only cases with a correlation coefficient greater than the thres_corr argument value. "linear" returns only cases with R ² greater than the thres_rs argument.
thres_corr	numeric. This is the correlation coefficient threshold of the correlation coefficient information to be returned. The default is 0.3.
thres_rs	numeric. R ² threshold of linear model summaries information to return. The default is 0.1.
verbose	logical. Specifies whether to output additional information during the calculation process. The default is to output information as TRUE. In this case, the function returns the value with invisible(). If FALSE, the value is returned by return().
...	further arguments passed to or from other methods.
x	an object of class "compare_numeric", usually, a result of a call to compare_numeric().

Details

print.compare_numeric() displays only the information compared between the variables included in compare_numeric. When using summary.compare_numeric(), it is advantageous to set the verbose argument to TRUE if the user is only viewing information from the console. It is also advantageous to specify FALSE if you want to manipulate the results.

Value

An object of the class as compare based list. The information to examine the relationship between numerical variables is as follows each components. - correlation component : Pearson's correlation coefficient.

- var1 : factor. The level of the first variable to compare. 'var1' is the name of the first variable to be compared.
- var2 : factor. The level of the second variable to compare. 'var2' is the name of the second variable to be compared.
- coef_corr : double. Pearson's correlation coefficient.

- linear component : linear model summaries

- var1 : factor. The level of the first variable to compare. 'var1' is the name of the first variable to be compared.
- var2 : factor. The level of the second variable to compare. 'var2' is the name of the second variable to be compared.
- r.squared : double. The percent of variance explained by the model.

- adj.r.squared : double. r.squared adjusted based on the degrees of freedom.
- sigma : double. The square root of the estimated residual variance.
- statistic : double. F-statistic.
- p.value : double. p-value from the F test, describing whether the full regression is significant.
- df : integer degrees of freedom.
- logLik : double. the log-likelihood of data under the model.
- AIC : double. the Akaike Information Criterion.
- BIC : double. the Bayesian Information Criterion.
- deviance : double. deviance.
- df.residual : integer residual degrees of freedom.

See Also

[plot.compare_numeric](#).

Examples

```
# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

library(dplyr)

# Compare the all numerical variables
all_var <- compare_numeric(carseats)

# Print compare class object
all_var

# Compare the correlation that case of joint the Price variable
all_var %>%
  "$"(correlation) %>%
  filter(var1 == "Price" | var2 == "Price") %>%
  arrange(desc(abs(coef_corr)))

# Compare the correlation that case of abs(coef_corr) > 0.3
all_var %>%
  "$"(correlation) %>%
  filter(abs(coef_corr) > 0.3)

# Compare the linear model that case of joint the Price variable
all_var %>%
  "$"(linear) %>%
  filter(var1 == "Price" | var2 == "Price") %>%
  arrange(desc(r.squared))

# Compare the two numerical variables
two_var <- compare_numeric(carseats, Price, CompPrice)
```

```
# Print compare class object
two_var

# Summary the all case : Return a invisible copy of an object.
stat <- summary(all_var)

# Just correlation
summary(all_var, method = "correlation")

# Just correlation condition by r > 0.2
summary(all_var, method = "correlation", thres_corr = 0.2)

# linear model summaries condition by R^2 > 0.05
summary(all_var, thres_rs = 0.05)

# verbose is FALSE
summary(all_var, verbose = FALSE)
```

summary.imputation *Summarizing imputation information*

Description

print and summary method for "imputation" class.

Usage

```
## S3 method for class 'imputation'
summary(object, ...)
```

Arguments

object an object of class "imputation", usually, a result of a call to `imputate_na()` or `imputate_outlier()`.

... further arguments passed to or from other methods.

Details

summary.imputation tries to be smart about formatting two kinds of imputation.

See Also

[imputate_na](#), [imputate_outlier](#), [summary.imputation](#).

Examples

```

# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

# Impute missing values -----
# If the variable of interest is a numerical variable
income <- imputate_na(carseats, Income, US, method = "rpart")
income
summary(income)
plot(income)

# If the variable of interest is a categorical variable
urban <- imputate_na(carseats, Urban, US, method = "mice")
urban
summary(urban)
plot(urban)

# Impute outliers -----
# If the variable of interest is a numerical variable
price <- imputate_outlier(carseats, Price, method = "capping")
price
summary(price)
plot(price)

```

summary.transform

Summarizing transformation information

Description

print and summary method for "transform" class.

Usage

```

## S3 method for class 'transform'
summary(object, ...)

```

Arguments

object an object of class "transform", usually, a result of a call to transform().
... further arguments passed to or from other methods.

Details

summary.transform compares the distribution of data before and after data transformation.

See Also

[transform](#), [plot.transform](#).

Examples

```
# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

# Standardization -----
advertising_minmax <- transform(carseats$Advertising, method = "minmax")
advertising_minmax
summary(advertising_minmax)
plot(advertising_minmax)

# Resolving Skewness -----
advertising_log <- transform(carseats$Advertising, method = "log")
advertising_log
summary(advertising_log)
plot(advertising_log)
```

summary.univar_category

Summarizing univar_category information

Description

print and summary method for "univar_category" class.

Usage

```
## S3 method for class 'univar_category'
summary(object, na.rm = TRUE, ...)

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

Arguments

object	an object of class "univar_category", usually, a result of a call to univar_category().
na.rm	logical. Specifies whether to include NA when performing a chi-square test. The default is TRUE, where NA is removed and aggregated.
...	further arguments passed to or from other methods.
x	an object of class "univar_category", usually, a result of a call to univar_category().

Details

print.univar_category() displays only the information of variables included in univar_category. The "variables" attribute is not displayed.

Value

An object of the class as individual variables based list. The information to examine the relationship between categorical variables is as follows each components.

- variable : factor. The level of the variable. 'variable' is the name of the variable.
- statistic : numeric. the value the chi-squared test statistic.
- p.value : numeric. the p-value for the test.
- df : integer. the degrees of freedom of the chi-squared test.

See Also

[plot.univar_category.](#)

Examples

```
# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

library(dplyr)
library(stringr)

# Calculates the all categorical variavels
all_var <- univar_category(carseats)

# Print univar_category class object
all_var

# Calculates the only Urban variable
all_var %>%
  "["(str_detect(names(all_var), "Urban"))

urban <- univar_category(carseats, Urban)

# Print univar_category class object
urban

# Filtering the case of Urban included NA
urban %>%
  "[["(1) %>%
  filter(!is.na(Urban))

# Summary the all case : Return a invisible copy of an object.
stat <- summary(all_var)
```

```
# Summary by returned object
stat

# Summary chi-squared test by returned object exclude NA
summary(urban)

# Summary chi-squared test by returned object include NA
summary(urban, na.rm = FALSE)
```

```
summary.univar_numeric
    Summarizing univar_numeric information
```

Description

print and summary method for "univar_numeric" class.

Usage

```
## S3 method for class 'univar_numeric'
summary(object, stand = c("robust", "minmax", "zscore"), ...)

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

Arguments

object	an object of class "univar_numeric", usually, a result of a call to univar_numeric().
stand	character Describe how to standardize the original data. "robust" normalizes the raw data through transformation calculated by IQR and median. "minmax" normalizes the original data using minmax transformation. "zscore" standardizes the original data using z-Score transformation. The default is "robust".
...	further arguments passed to or from other methods.
x	an object of class "univar_numeric", usually, a result of a call to univar_numeric().

Details

print.univar_numeric() displays only the information of variables included in univar_numeric The "variables" attribute is not displayed.

Value

An object of the class as individual variables based list. The statistics returned by `summary.univar_numeric()` are different from the statistics returned by `univar_numeric()`. `univar_numeric()` is the statistics for the original data, but `summary.univar_numeric()` is the statistics for the standardized data. A component named "statistics" is a tibble object with the following statistics.:

- `variable` : factor. The level of the variable. 'variable' is the name of the variable.
- `n` : number of observations excluding missing values
- `na` : number of missing values
- `mean` : arithmetic average
- `sd` : standard deviation
- `se_mean` : standard error mean. sd/\sqrt{n}
- `IQR` : interquartile range (Q3-Q1)
- `skewness` : skewness
- `kurtosis` : kurtosis
- `median` : median. 50% percentile

See Also

[plot.univar_numeric.](#)

Examples

```
# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

library(dplyr)

# Calculates the all categorical variables
all_var <- univar_numeric(carseats)

# Print univar_numeric class object
all_var

# Calculates the Price, CompPrice variable
univar_numeric(carseats, Price, CompPrice)

# Summary the all case : Return a invisible copy of an object.
stat <- summary(all_var)

# Summary by returned object
stat

# Statistics of numerical variables normalized by Min-Max method
summary(all_var, stand = "minmax")
```

```
# Statistics of numerical variables standardized by Z-score method
summary(all_var, stand = "zscore")
```

target_by	<i>Target by one variables</i>
-----------	--------------------------------

Description

In the data analysis, a `target_df` class is created to identify the relationship between the target variable and the other variable.

Usage

```
target_by(.data, target, ...)

## S3 method for class 'data.frame'
target_by(.data, target, ...)
```

Arguments

<code>.data</code>	a data.frame or a tbl_df .
<code>target</code>	target variable.
<code>...</code>	arguments to be passed to methods.

Details

Data analysis proceeds with the purpose of predicting target variables that correspond to the facts of interest, or examining associations and relationships with other variables of interest. Therefore, it is a major challenge for EDA to examine the relationship between the target variable and its corresponding variable. Based on the derived relationships, analysts create scenarios for data analysis.

`target_by()` inherits the [grouped_df](#) class and returns a `target_df` class containing information about the target variable and the variable.

See `vignette("EDA")` for an introduction to these concepts.

Value

an object of `target_df` class. Attributes of `target_df` class is as follows.

- `type_y` : the data type of target variable.

See Also

[relate](#).

Examples

```

# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

# If the target variable is a categorical variable
categ <- target_by(carseats, US)

# If the variable of interest is a numerical variable
cat_num <- relate(categ, Sales)
cat_num
summary(cat_num)
plot(cat_num)

# If the variable of interest is a categorical variable
cat_cat <- relate(categ, ShelveLoc)
cat_cat
summary(cat_cat)
plot(cat_cat)

##-----
# If the target variable is a categorical variable
num <- target_by(carseats, Sales)

# If the variable of interest is a numerical variable
num_num <- relate(num, Price)
num_num
summary(num_num)
plot(num_num)

# If the variable of interest is a categorical variable
num_cat <- relate(num, ShelveLoc)
num_cat
summary(num_cat)
plot(num_cat)

```

target_by.tbl_dbi *Target by one column in the DBMS*

Description

In the data analysis, a `target_df` class is created to identify the relationship between the target column and the other column of the DBMS table through `tbl_dbi`

Usage

```

## S3 method for class 'tbl_dbi'
target_by(.data, target, in_database = FALSE, collect_size = Inf, ...)

```

Arguments

<code>.data</code>	a <code>tbl_dbi</code> .
<code>target</code>	target variable.
<code>in_database</code>	Specifies whether to perform in-database operations. If <code>TRUE</code> , most operations are performed in the DBMS. if <code>FALSE</code> , table data is taken in R and operated in-memory. Not yet supported in <code>in_database = TRUE</code> .
<code>collect_size</code>	a integer. The number of data samples from the DBMS to R. Applies only if <code>in_database = FALSE</code> .
<code>...</code>	arguments to be passed to methods.

Details

Data analysis proceeds with the purpose of predicting target variables that correspond to the facts of interest, or examining associations and relationships with other variables of interest. Therefore, it is a major challenge for EDA to examine the relationship between the target variable and its corresponding variable. Based on the derived relationships, analysts create scenarios for data analysis. `target_by()` inherits the `grouped_df` class and returns a `target_df` class containing information about the target variable and the variable.

See `vignette("EDA")` for an introduction to these concepts.

Value

an object of `target_df` class. Attributes of `target_df` class is as follows.

- `type_y` : the data type of target variable.

See Also

[target_by.data.frame](#), [relate](#).

Examples

```
library(dplyr)

# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

# connect DBMS
con_sqlite <- DBI::dbConnect(RSQLite::SQLite(), ":memory:")

# copy carseats to the DBMS with a table named TB_CARSEATS
copy_to(con_sqlite, carseats, name = "TB_CARSEATS", overwrite = TRUE)

# If the target variable is a categorical variable
categ <- target_by(con_sqlite %>% tbl("TB_CARSEATS") , US)

# If the variable of interest is a numerical variable
```



```
cat_num <- relate(categ, Sales)
cat_num
summary(cat_num)
plot(cat_num)

# If the variable of interest is a categorical column
cat_cat <- relate(categ, ShelveLoc)
cat_cat
summary(cat_cat)
plot(cat_cat)

##-----
# If the target variable is a categorical column,
# and In-memory mode and collect size is 350
num <- target_by(con_sqlite %>% tbl("TB_CARSEATS"), Sales, collect_size = 350)

# If the variable of interest is a numerical column
num_num <- relate(num, Price)
num_num
summary(num_num)
plot(num_num)
plot(num_num, hex_thres = 400)

# If the variable of interest is a categorical column
num_cat <- relate(num, ShelveLoc)
num_cat
summary(num_cat)
plot(num_cat)
```

transform

Data Transformations

Description

Performs variable transformation for standardization and resolving skewness of numerical variables.

Usage

```
transform(
  x,
  method = c("zscore", "minmax", "log", "log+1", "sqrt", "1/x", "x^2", "x^3")
)
```

Arguments

x numeric vector for transformation.
method method of transformations.

Details

`transform()` creates an transform class. The ‘transform’ class includes original data, transformed data, and method of transformation.

See `vignette("transformation")` for an introduction to these concepts.

Value

An object of transform class. Attributes of transform class is as follows.

- method : method of transformation data.
 - Standardization
 - * "zscore" : z-score transformation. $(x - \mu) / \sigma$
 - * "minmax" : minmax transformation. $(x - \min) / (\max - \min)$
 - Resolving Skewness
 - * "log" : log transformation. $\log(x)$
 - * "log+1" : log transformation. $\log(x + 1)$. Used for values that contain 0.
 - * "sqrt" : square root transformation.
 - * "1/x" : $1 / x$ transformation
 - * "x^2" : x square transformation
 - * "x^3" : x^3 square transformation

See Also

[summary.transform](#), [plot.transform](#).

Examples

```
# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

# Standardization -----
advertising_minmax <- transform(carseats$Advertising, method = "minmax")
advertising_minmax
summary(advertising_minmax)
plot(advertising_minmax)

# Resolving Skewness -----
advertising_log <- transform(carseats$Advertising, method = "log")
advertising_log
summary(advertising_log)
plot(advertising_log)

# Using dplyr -----
library(dplyr)

carseats %>%
  mutate(Advertising_log = transform(Advertising, method = "log+1")) %>%
  lm(Sales ~ Advertising_log, data = .)
```

transformation_report *Reporting the information of transformation*

Description

The transformation_report() report the information of transform numerical variables for object inheriting from data.frame.

Usage

```
transformation_report(  
  .data,  
  target = NULL,  
  output_format = c("pdf", "html"),  
  output_file = NULL,  
  output_dir = tempdir(),  
  font_family = NULL,  
  browse = TRUE  
)
```

Arguments

.data	a data.frame or a tbl_df .
target	target variable. If the target variable is not specified, the method of using the target variable information is not performed when the missing value is imputed. and Optimal binning is not performed if the target variable is not a binary class.
output_format	report output type. Choose either "pdf" and "html". "pdf" create pdf file by knitr::knit(). "html" create html file by rmarkdown::render().
output_file	name of generated file. default is NULL.
output_dir	name of directory to generate report file. default is tempdir().
font_family	character. font family name for figure in pdf.
browse	logical. choose whether to output the report results to the browser.

Details

Generate transformation reports automatically. You can choose to output to pdf and html files. This is useful for Binning a data frame with a large number of variables than data with a small number of variables. For pdf output, Korean Gothic font must be installed in Korean operating system.

Reported information

The transformation process will report the following information:

- Imputation
 - Missing Values

- * * Variable names including missing value
- Outliers
 - * * Variable names including outliers
- Resolving Skewness
 - Skewed variables information
 - * * Variable names with an absolute value of skewness greater than or equal to 0.5
- Binning
 - Numerical Variables for Binning
 - Binning
 - * Numeric variable names
 - Optimal Binning
 - * Numeric variable names

See vignette("transformation") for an introduction to these concepts.

Examples

```
# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

# reporting the Binning information -----
# create pdf file. file name is Transformation_Report.pdf & No target variable
transformation_report(carseats)
# create pdf file. file name is Transformation_Report.pdf
transformation_report(carseats, US)
# create pdf file. file name is Transformation_carseats.pdf
transformation_report(carseats, "US", output_file = "Transformation_carseats.pdf")
# create html file. file name is Transformation_Report.html
transformation_report(carseats, "US", output_format = "html")
# create html file. file name is Transformation_carseats
transformation_report(carseats, US, output_format = "html",
                      output_file = "Transformation_carseats.html")
```

univar_category

Statistic of univariate categorical variables

Description

The `univar_category()` calculates statistic of categorical variables that is frequency table

Usage

```
univar_category(.data, ...)

## S3 method for class 'data.frame'
univar_category(.data, ...)
```

Arguments

<code>.data</code>	a <code>data.frame</code> or a <code>tbl_df</code> .
<code>...</code>	one or more unquoted expressions separated by commas. You can treat variable names like they are positions. Positive values select variables; negative values to drop variables. These arguments are automatically quoted and evaluated in a context where column names represent column positions. They support unquoting and splicing.

Details

`univar_category()` calculates the frequency table of categorical variables. If a specific variable name is not specified, frequency tables for all categorical variables included in the data are calculated. The `univar_category` class returned by `univar_category()` is useful because it can draw chisquare tests and bar plots as well as frequency tables of individual variables. and return `univar_category` class that based list object.

Value

An object of the class as individual variables based list. The information to examine the relationship between categorical variables is as follows each components.

- `variable` : factor. The level of the variable. 'variable' is the name of the variable.
- `n` : integer. frequency by variable.
- `rate` : double. relative frequency.

Attributes of return object

Attributes of `compare_category` class is as follows.

- `variables` : character. List of variables selected for calculate frequency.

See Also

[summary.univar_category](#), [print.univar_category](#), [plot.univar_category](#).

Examples

```
# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA
```

```
library(dplyr)
library(stringr)

# Calculates the all categorical variables
all_var <- univar_category(carseats)

# Print univar_category class object
all_var

# Calculates the only Urban variable
all_var %>%
  "["(str_detect(names(all_var), "Urban"))

urban <- univar_category(carseats, Urban)

# Print univar_category class object
urban

# Filtering the case of Urban included NA
urban %>%
  "["(1) %>%
  filter(!is.na(Urban))

# Summary the all case : Return a invisible copy of an object.
stat <- summary(all_var)

# Summary by returned object
stat

# plot all variables
plot(all_var)

# plot urban
plot(urban)

# plot all variables by prompt
plot(all_var, prompt = TRUE)
```

univar_numeric

Statistic of univariate numerical variables

Description

The `univar_numeric()` calculates statistic of numerical variables that is frequency table

Usage

```
univar_numeric(.data, ...)
```

```
## S3 method for class 'data.frame'
univar_numeric(.data, ...)
```

Arguments

`.data` a `data.frame` or a `tbl_df`.

`...` one or more unquoted expressions separated by commas. You can treat variable names like they are positions. Positive values select variables; negative values to drop variables. These arguments are automatically quoted and evaluated in a context where column names represent column positions. They support unquoting and splicing.

Details

`univar_numeric()` calculates the popular statistics of numerical variables. If a specific variable name is not specified, statistics for all categorical numerical included in the data are calculated. The statistics obtained by `univar_numeric()` are part of those obtained by `describe()`. Therefore, it is recommended to use `describe()` to simply calculate statistics. However, if you want to visualize the distribution of individual variables, you should use `univar_numeric()`.

Value

An object of the class as individual variables based list. A component named "statistics" is a tibble object with the following statistics.:

- `variable` : factor. The level of the variable. 'variable' is the name of the variable.
- `n` : number of observations excluding missing values
- `na` : number of missing values
- `mean` : arithmetic average
- `sd` : standard deviation
- `se_mean` : standrd error mean. sd/\sqrt{n}
- `IQR` : interquartile range (Q3-Q1)
- `skewness` : skewness
- `kurtosis` : kurtosis
- `median` : median. 50% percentile

Attributes of return object

Attributes of `compare_category` class is as follows.

- `raw` : a `data.frame` or a `tbl_df`. Data containing variables to be compared. Save it for visualization with `plot.univar_numeric()`.
- `variables` : character. List of variables selected for calculate statistics.

See Also

[summary.univar_numeric](#), [print.univar_numeric](#), [plot.univar_numeric](#).

Examples

```
# Generate data for the example
carseats <- ISLR::Carseats
carseats[sample(seq(NROW(carseats)), 20), "Income"] <- NA
carseats[sample(seq(NROW(carseats)), 5), "Urban"] <- NA

# Calculates the all categorical variables
all_var <- univar_numeric(carseats)

# Print univar_numeric class object
all_var

# Calculates the Price, CompPrice variable
univar_numeric(carseats, Price, CompPrice)

# Summary the all case : Return a invisible copy of an object.
stat <- summary(all_var)

# Summary by returned object
stat

# Statistics of numerical variables normalized by Min-Max method
summary(all_var, stand = "minmax")

# Statistics of numerical variables standardized by Z-score method
summary(all_var, stand = "zscore")

# one plot with all variables
plot(all_var)

# one plot with all normalized variables by Min-Max method
plot(all_var, stand = "minmax")

# one plot with all variables
plot(all_var, stand = "none")

# one plot with all robust standardized variables
plot(all_var, viz = "boxplot")

# one plot with all standardized variables by Z-score method
plot(all_var, viz = "boxplot", stand = "zscore")

# individual boxplot by variables
plot(all_var, indiv = TRUE, "boxplot")

# individual histogram by variables
plot(all_var, indiv = TRUE, "hist")

# individual histogram by robust standardized variable
plot(all_var, indiv = TRUE, "hist", stand = "robust")

# plot all variables by prompt
```



```
plot(all_var, indiv = TRUE, "hist", prompt = TRUE)
```

Index

binning, [4](#), [7](#), [66](#), [98](#)
binning_by, [5](#), [6](#), [70](#)

compare_category, [8](#), [67](#)
compare_numeric, [10](#), [68](#)
cor, [14](#), [17](#)
correlate, [12](#), [13](#)
correlate.data.frame, [17](#)
correlate.tbl_dbi, [14](#), [16](#)

describe, [18](#)
describe.data.frame, [22](#), [62](#)
describe.tbl_dbi, [20](#), [20](#), [64](#)
diagnose, [23](#)
diagnose.data.frame, [26](#), [28](#), [32](#), [36](#)
diagnose.tbl_dbi, [24](#), [25](#), [30](#), [34](#), [38](#)
diagnose_category, [27](#)
diagnose_category.data.frame, [24](#), [30](#), [32](#), [36](#)
diagnose_category.tbl_dbi, [26](#), [28](#), [29](#), [30](#), [34](#), [38](#)
diagnose_numeric, [31](#)
diagnose_numeric.data.frame, [20](#), [24](#), [28](#), [34](#), [36](#), [62](#)
diagnose_numeric.tbl_dbi, [22](#), [26](#), [30](#), [32](#), [33](#), [38](#), [64](#)
diagnose_outlier, [35](#)
diagnose_outlier.data.frame, [28](#), [32](#), [38](#), [89](#)
diagnose_outlier.tbl_dbi, [30](#), [34](#), [36](#), [37](#), [92](#)
diagnose_report, [39](#)
diagnose_report.data.frame, [43](#)
diagnose_report.tbl_dbi, [41](#)
dlookr (dlookr-package), [3](#)
dlookr-package, [3](#)

eda_report, [44](#)
eda_report.data.frame, [48](#)
eda_report.tbl_dbi, [46](#)

find_class, [49](#), [54](#)
find_na, [51](#), [51](#), [52](#), [53](#)
find_outliers, [52](#), [53](#)
find_skewness, [53](#), [97](#)

get_class, [50](#), [54](#)
get_column_info, [55](#)
get_os, [56](#)
group_by, [19](#), [21](#), [62](#), [64](#)
grouped_df, [14](#), [17](#), [110](#), [112](#)

imutate_na, [51](#), [57](#), [60](#), [69](#), [104](#)
imutate_outlier, [52](#), [58](#), [59](#), [69](#), [104](#)

kurtosis, [60](#), [97](#)

normality, [61](#)
normality.data.frame, [64](#)
normality.tbl_dbi, [62](#), [63](#)

plot.bins, [66](#), [70](#)
plot.compare_category, [9](#), [67](#), [100](#)
plot.compare_numeric, [12](#), [68](#), [103](#)
plot.imputation, [69](#)
plot.optimal_bins, [70](#)
plot.relate, [71](#), [93](#), [96](#)
plot.transform, [73](#), [106](#), [114](#)
plot.univar_category, [74](#), [107](#), [117](#)
plot.univar_numeric, [75](#), [109](#), [119](#)
plot_correlate, [77](#)
plot_correlate.data.frame, [80](#)
plot_correlate.tbl_dbi, [77](#), [79](#)
plot_na_hclust, [81](#)
plot_na_intersect, [82](#)
plot_na_pareto, [83](#)
plot_normality, [85](#)
plot_normality.data.frame, [62](#), [88](#)
plot_normality.tbl_dbi, [64](#), [86](#), [87](#)
plot_outlier, [89](#)
plot_outlier.data.frame, [77](#), [86](#), [92](#)
plot_outlier.tbl_dbi, [80](#), [88](#), [89](#), [90](#)

`print.bins`, [5](#), [66](#)
`print.bins(summary.bins)`, [97](#)
`print.compare_category`, [9](#), [67](#)
`print.compare_category`
 (`summary.compare_category`), [99](#)
`print.compare_numeric`, [12](#), [68](#)
`print.compare_numeric`
 (`summary.compare_numeric`), [101](#)
`print.relate`, [72](#), [93](#), [96](#)
`print.univar_category`, [74](#), [117](#)
`print.univar_category`
 (`summary.univar_category`), [106](#)
`print.univar_numeric`, [76](#), [119](#)
`print.univar_numeric`
 (`summary.univar_numeric`), [108](#)

`relate`, [72](#), [94](#), [110](#), [112](#)

`shapiro.test`, [62](#), [64](#)
`skewness`, [61](#), [97](#)
`smbinning`, [7](#)
`smbinning.plot`, [70](#)
`summary.bins`, [5](#), [66](#), [97](#)
`summary.compare_category`, [9](#), [67](#), [99](#)
`summary.compare_numeric`, [12](#), [68](#), [101](#)
`summary.imputation`, [69](#), [104](#), [104](#)
`summary.transform`, [73](#), [105](#), [114](#)
`summary.univar_category`, [74](#), [106](#), [117](#)
`summary.univar_numeric`, [76](#), [108](#), [119](#)

`target_by`, [110](#)
`target_by.data.frame`, [112](#)
`target_by.tbl_dbi`, [111](#)
`tbl_df`, [8](#), [11](#), [12](#), [14](#), [19](#), [23](#), [27](#), [31](#), [36](#), [40](#),
 [44](#), [51–53](#), [57](#), [59](#), [61](#), [77](#), [85](#), [89](#),
 [110](#), [115](#), [117](#), [119](#)
`transform`, [73](#), [106](#), [113](#)
`transformation_report`, [115](#)

`univar_category`, [74](#), [116](#)
`univar_numeric`, [76](#), [118](#)