

# Package ‘guideR’

February 14, 2025

**Type** Package

**Title** Miscellaneous Statistical Functions Used in 'guide-R'

**Version** 0.1.0

**Description** Companion package for the manual

'guide-R : Guide pour l'analyse de données d'enquêtes avec R' available at  
<https://larmarange.github.io/guide-R/>. 'guideR' implements miscellaneous functions introduced in 'guide-R' to facilitate statistical analysis and manipulation of survey data.

**License** GPL (>= 3)

**URL** <https://larmarange.github.io/guideR/>,  
<https://github.com/larmarange/guideR>

**BugReports** <https://github.com/larmarange/guideR/issues>

**Depends** R (>= 4.2)

**Imports** cli, dplyr, ggplot2, labelled, lifecycle, pak, patchwork, purrr, renv, rlang, scales, srvyr, tidyverse, tidyselect

**Encoding** UTF-8

**RoxygenNote** 7.3.2

**Suggests** broom.helpers, FactoMineR, nnet, spelling, survey, testthat (>= 3.0.0)

**Config/testthat.edition** 3

**Language** en-US

**NeedsCompilation** no

**Author** Joseph Larmarange [aut, cre] (<<https://orcid.org/0000-0001-7097-700X>>)

**Maintainer** Joseph Larmarange <joseph@larmarange.net>

**Repository** CRAN

**Date/Publication** 2025-02-14 14:30:06 UTC

## Contents

install_dependencies	2
is_different	3
leading_zeros	4
long_to_periods	4
observed_vs_theoretical	5
plot_inertia_from_tree	6
proportion	7
round_preserve_sum	9
step_with_na	10
titanic	12
unrowwise	12

## Index

13

**install\_dependencies** *Install / Update project dependencies*

### Description

This function uses `renv::dependencies()` to identify R package dependencies in a project and then calls `pak::pkg_install()` to install / update these packages.

### Usage

```
install_dependencies(ask = TRUE)
```

### Arguments

`ask` Whether to ask for confirmation when installing a different version of a package that is already installed. Installations that only add new packages never require confirmation.

### Value

(Invisibly) A data frame with information about the installed package(s).

### Examples

```
## Not run:
install_dependencies()

## End(Not run)
```

---

**is\_different***Comparison tests considering NA as values to be compared*

---

## Description

`is_different()` and `is_equal()` performs comparison tests, considering NA values as legitimate values (see examples).

## Usage

```
is_different(x, y)

is_equal(x, y)

cumdifferent(x)

num_cycle(x)
```

## Arguments

`x, y` Vectors to be compared.

## Details

`cum_different()` allows to identify groups of continuous rows that have the same value. `num_cycle()` could be used to identify sub-groups that respect a certain condition (see examples).

`is_equal(x, y)` is equivalent to `(x == y & !is.na(x) & !is.na(y)) | (is.na(x) & is.na(y))`, and `is_different(x, y)` is equivalent to `(x != y & !is.na(x) & !is.na(y)) | xor(is.na(x), is.na(y))`.

## Value

A vector of the same length as `x`.

## Examples

```
v <- c("a", "b", NA)
is_different(v, "a")
is_different(v, NA)
is_equal(v, "a")
is_equal(v, NA)
d <- dplyr::tibble(group = c("a", "a", "b", "b", "a", "b", "c", "a"))
d |>
  dplyr::mutate(
    subgroup = cumdifferent(group),
    sub_a = num_cycle(group == "a")
  )
```

**leading\_zeros** *Add leading zeros*

### Description

Add leading zeros

### Usage

```
leading_zeros(x, left_digits = NULL, digits = 0, prefix = "", suffix = "", ...)
```

### Arguments

x	a numeric vector
left_digits	number of digits before decimal point, automatically computed if not provided
digits	number of digits after decimal point
prefix, suffix	Symbols to display before and after value
...	additional parameters passed to <a href="#">base::formatC()</a> , as big.mark or decimal.mark

### Value

A character vector of the same length as x.

### See Also

[base::formatC\(\)](#), [base::sprintf\(\)](#)

### Examples

```
v <- c(2, 103.24, 1042.147, 12.4566, NA)
leading_zeros(v)
leading_zeros(v, digits = 1)
leading_zeros(v, left_digits = 6, big.mark = " ")
leading_zeros(c(0, 6, 12, 18), prefix = "M")
```

**long\_to\_periods** *Transform a data frame from long format to period format*

### Description

Transform a data frame from long format to period format

### Usage

```
long_to_periods(data, id, start, stop = NULL, by = NULL)
```

**Arguments**

data	A data frame, or a data frame extension (e.g. a tibble).
id	<tidy-select> Column containing individual ids
start	<tidy-select> Time variable indicating the beginning of each row
stop	<tidy-select> Optional time variable indicating the end of each row. If not provided, it will be derived from the dataset, considering that each row ends at the beginning of the next one.
by	<tidy-select> Co-variables to consider (optional)

**Value**

A tibble.

**Examples**

```
d <- dplyr::tibble(
  patient = c(1, 2, 3, 3, 4, 4, 4),
  begin = c(0, 0, 0, 1, 0, 36, 39),
  end = c(50, 6, 1, 16, 36, 39, 45),
  covar = c("no", "no", "no", "yes", "no", "yes", "yes")
)
d

d |> long_to_periods(id = patient, start = begin, stop = end)
d |> long_to_periods(id = patient, start = begin, stop = end, by = covar)

# If stop not provided, it is deduced.
# However, it considers that observation ends at the last start time.
d |> long_to_periods(id = patient, start = begin)
```

**observed\_vs\_theoretical**

*Plot observed vs predicted distribution of a fitted model*

**Description**

Plot observed vs predicted distribution of a fitted model

**Usage**

```
observed_vs_theoretical(model)
```

**Arguments**

model	A statistical model.
-------	----------------------

## Details

Has been tested with `stats::lm()` and `stats::glm()` models. It may work with other types of models, but without any warranty.

## Value

A ggplot2 plot.

## Examples

```
# a linear model
mod <- lm(Sepal.Length ~ Sepal.Width + Species, data = iris)
mod |> observed_vs_theoretical()

# a logistic regression
mod <- glm(
  as.factor(Survived) ~ Class + Sex,
  data = titanic,
  family = binomial()
)
mod |> observed_vs_theoretical()
```

## *plot\_inertia\_from\_tree*

*Plot inertia, absolute loss and relative loss from a classification tree*

## Description

Plot inertia, absolute loss and relative loss from a classification tree

## Usage

```
plot_inertia_from_tree(tree, k_max = 15)
```

```
get_inertia_from_tree(tree, k_max = 15)
```

## Arguments

<code>tree</code>	A dendrogram, i.e. an <code>stats::hclust</code> object, an <code>FactoMineR::HCPC</code> object or an object that can be converted to an <code>stats::hclust</code> object with <code>stats::as.hclust()</code> .
<code>k_max</code>	Maximum number of clusters to return / plot.

## Value

A ggplot2 plot or a tibble.

## Examples

```
hc <- hclust(dist(USArrests))
get_inertia_from_tree(hc)
plot_inertia_from_tree(hc)
```

---

proportion	<i>Compute proportions</i>
------------	----------------------------

---

## Description

`proportion()` lets you quickly count observations (like `dplyr::count()`) and compute relative proportions. Proportions are computed separately by group (see examples).

## Usage

```
proportion(data, ..., .by = NULL)

## S3 method for class 'data.frame'
proportion(
  data,
  ...,
  .by = NULL,
  .na.rm = FALSE,
  .weight = NULL,
  .scale = 100,
  .sort = FALSE,
  .drop = FALSE,
  .conf.int = FALSE,
  .conf.level = 0.95,
  .options = list(correct = TRUE)
)

## S3 method for class 'survey.design'
proportion(
  data,
  ...,
  .by = NULL,
  .na.rm = FALSE,
  .scale = 100,
  .sort = FALSE,
  .conf.int = FALSE,
  .conf.level = 0.95,
  .options = NULL
)
```

## Arguments

<code>.data</code>	A data frame, data frame extension (e.g. a tibble), or a survey design object.
<code>...</code>	< <a href="#">data-masking</a> > Variable(s) for those computing proportions.
<code>.by</code>	< <a href="#">tidy-select</a> > Optional additional variables to group by (in addition to those eventually previously declared using <code>dplyr::group_by()</code> ).
<code>.na.rm</code>	Should NA values be removed (from variables declared in <code>...</code> )?
<code>.weight</code>	< <a href="#">data-masking</a> > Frequency weights. Can be NULL or a variable.
<code>.scale</code>	A scaling factor applied to proportion. Use 1 for keeping proportions unchanged.
<code>.sort</code>	If TRUE, will show the highest proportions at the top.
<code>.drop</code>	If TRUE, will remove empty groups from the output.
<code>.conf.int</code>	If TRUE, will estimate confidence intervals.
<code>.conf.level</code>	Confidence level for the returned confidence intervals.
<code>.options</code>	Additional arguments passed to <code>stats::prop.test()</code> or <code>srvyr::survey_prop()</code> .

## Value

A tibble.

A tibble with one row per group.

## Examples

```
# univariable table
titanic |> proportion(Class)
titanic |> proportion(Class, .sort = TRUE)
titanic |> proportion(Class, .conf.int = TRUE)
titanic |> proportion(Class, .conf.int = TRUE, .scale = 1)

# bivariable table
titanic |> proportion(Class, Survived) # proportions of the total
titanic |> proportion(Survived, .by = Class) # row proportions
titanic |> # equivalent syntax
  dplyr::group_by(Class) |>
  proportion(Survived)

# combining 3 variables or more
titanic |> proportion(Class, Sex, Survived)
titanic |> proportion(Sex, Survived, .by = Class)
titanic |> proportion(Survived, .by = c(Class, Sex))

# missing values
dna <- titanic
dna$Survived[c(1:20, 500:530)] <- NA
dna |> proportion(Survived)
dna |> proportion(Survived, .na.rm = TRUE)

## SURVEY DATA -----
```

```

ds <- srvyr::as_survey(titanic)

# univariable table
ds |> proportion(Class)
ds |> proportion(Class, .sort = TRUE)
ds |> proportion(Class, .conf.int = TRUE)
ds |> proportion(Class, .conf.int = TRUE, .scale = 1)

# bivariable table
ds |> proportion(Class, Survived) # proportions of the total
ds |> proportion(Survived, .by = Class) # row proportions
ds |> dplyr::group_by(Class) |> proportion(Survived)

# combining 3 variables or more
ds |> proportion(Class, Sex, Survived)
ds |> proportion(Sex, Survived, .by = Class)
ds |> proportion(Survived, .by = c(Class, Sex))

# missing values
dsna <- srvyr::as_survey(dna)
dsna |> proportion(Survived)
dsna |> proportion(Survived, .na.rm = TRUE)

```

`round_preserve_sum`      *Round values while preserve their rounded sum in R*

## Description

Sometimes, the sum of rounded numbers (e.g., using `base::round()`) is not the same as their rounded sum.

## Usage

```
round_preserve_sum(x, digits = 0)
```

## Arguments

<code>x</code>	Numerical vector to sum.
<code>digits</code>	Number of decimals for rounding.

## Details

This solution applies the following algorithm

- Round down to the specified number of decimal places
- Order numbers by their remainder values
- Increment the specified decimal place of values with  $k$  largest remainders, where  $k$  is the number of values that must be incremented to preserve their rounded sum

**Value**

A numerical vector of same length as *x*.

**Source**

<https://biostatmatt.com/archives/2902>

**Examples**

```
sum(c(0.333, 0.333, 0.334))
round(c(0.333, 0.333, 0.334), 2)
sum(round(c(0.333, 0.333, 0.334), 2))
round_preserve_sum(c(0.333, 0.333, 0.334), 2)
sum(round_preserve_sum(c(0.333, 0.333, 0.334), 2))
```

*step\_with\_na*

*Apply step(), taking into account missing values*

**Description**

When your data contains missing values, concerned observations are removed from a model. However, then at a later stage, you try to apply a descending stepwise approach to reduce your model by minimization of AIC, you may encounter an error because the number of rows has changed.

**Usage**

```
step_with_na(model, ...)

## Default S3 method:
step_with_na(model, ..., full_data = eval(model$call$data))

## S3 method for class 'svyglm'
step_with_na(model, ..., design)
```

**Arguments**

- |                  |   |
|------------------|---|
| <i>model</i>     | A model object.   |
| <i>...</i>       | Additional parameters passed to <a href="#">stats::step()</a> .       |
| <i>full_data</i> | Full data frame used for the model, including missing data.           |
| <i>design</i>    | Survey design previously passed to <a href="#">survey::svyglm()</a> . |

## Details

`step_with_na()` applies the following strategy:

- recomputes the models using only complete cases;
- applies `stats::step()`;
- recomputes the reduced model using the full original dataset.

`step_with_na()` has been tested with `stats::lm()`, `stats::glm()`, `nnet::multinom()` and `survey::svyglm()`. It may be working with other types of models, but with no warranty.

In some cases, it may be necessary to provide the full dataset initially used to estimate the model.

`step_with_na()` may not work inside other functions. In that case, you may try to pass `full_data` to the function.

## Value

The stepwise-selected model.

## Examples

```
set.seed(42)
d <- titanic |>
  dplyr::mutate(
    Group = sample(
      c("a", "b", NA),
      dplyr::n(),
      replace = TRUE
    )
  )
mod <- glm(as.factor(Survived) ~ ., data = d, family = binomial())
# step(mod) should produce an error
mod2 <- step_with_na(mod)
mod2

## WITH SURVEY -----
library(survey)
ds <- d |>
  dplyr::mutate(Survived = as.factor(Survived)) |>
  svyr::as_survey()
mods <- survey::svyglm(
  Survived ~ Class + Group + Sex,
  design = ds,
  family = quasibinomial()
)
mod2s <- step_with_na(mods, design = ds)
mod2s
```

titanic	<i>Titanic data set in long format</i>
---------	--

## Description

This `titanic` dataset is equivalent to `datasets::Titanic |> dplyr::as_tibble() |> tidyverse::uncount(n)`.

## Usage

```
titanic
```

## Format

An object of class `tbl_df` (inherits from `tbl`, `data.frame`) with 2201 rows and 4 columns.

## See Also

[datasets::Titanic](#)

unrowwise	<i>Remove row-wise grouping</i>
-----------	---------------------------------

## Description

Remove row-wise grouping created with `dplyr::rowwise()` while preserving any other grouping declared with `dplyr::group_by()`.

## Usage

```
unrowwise(data)
```

## Arguments

`data` A data frame, data frame extension (e.g. a tibble), or a lazy data frame.

## Value

A tibble.

## Examples

```
titanic |> dplyr::rowwise()
titanic |> dplyr::rowwise() |> unrowwise()

titanic |> dplyr::group_by(Sex, Class) |> dplyr::rowwise()
titanic |> dplyr::group_by(Sex, Class) |> dplyr::rowwise() |> unrowwise()
```

# Index

```
* datasets
  titanic, 12
* logic
  is_different, 3
* manip
  long_to_periods, 4
  unrowwise, 12
* models
  observed_vs_theoretical, 5
  step_with_na, 10
* tree
  plot_inertia_from_tree, 6
* univar
  proportion, 7
  round_preserve_sum, 9
* utilities
  install_dependencies, 2
  leading_zeros, 4

base::formatC(), 4
base::round(), 9
base::sprintf(), 4

cumdifferent(is_different), 3

datasets::Titanic, 12
dplyr::count(), 7
dplyr::group_by(), 8, 12
dplyr::rowwise(), 12

FactoMineR::HCPC, 6

get_inertia_from_tree
  (plot_inertia_from_tree), 6

install_dependencies, 2
is_different, 3
is_equal(is_different), 3

leading_zeros, 4
long_to_periods, 4

nnet::multinom(), 11
num_cycle(is_different), 3
observed_vs_theoretical, 5
pak::pkg_install(), 2
plot_inertia_from_tree, 6
proportion, 7

renv::dependencies(), 2
round_preserve_sum, 9

srvyr::survey_prop(), 8
stats::as.hclust(), 6
stats::glm(), 6, 11
stats::hclust, 6
stats::lm(), 6, 11
stats::prop.test(), 8
stats::step(), 10, 11
step_with_na, 10
survey::svyglm(), 10, 11

titanic, 12

unrowwise, 12
```