

Package ‘vvauditor’

February 26, 2024

Title Creates Assertion Tests

Version 0.6.0

Description

Offers a comprehensive set of assertion tests to help users validate the integrity of their data. These tests can be used to check for specific conditions or properties within a dataset and help ensure that data is accurate and reliable. The package is designed to make it easy to add quality control checks to data analysis workflows and to aid in identifying and correcting any errors or inconsistencies in data.

License MIT + file LICENSE

Encoding UTF-8

RoxygenNote 7.2.3

Imports checkmate, cli, dplyr, findR, janitor, kit, lubridate, magrittr, purrr, stats, stringr, tibble

Suggests knitr, rmarkdown, testthat (>= 3.0.0)

Config/testthat/edition 3

VignetteBuilder knitr

NeedsCompilation no

Author Tomer Iwan [aut, cre, cph]

Maintainer Tomer Iwan <t.iwan@vu.nl>

Repository CRAN

Date/Publication 2024-02-26 13:30:05 UTC

R topics documented:

| | |
|--|---|
| assertion_message | 2 |
| assert_date_named | 3 |
| assert_logical_named | 4 |
| assert_no_duplicates_in_group | 4 |
| calculate_category_percentages | 5 |
| check_double_columns | 6 |
| check_duplicates | 6 |

| | |
|---|----|
| check_na_columns | 7 |
| check_non_zero_rows | 8 |
| check_no_duplicates_in_group | 8 |
| check_no_duplicate_rows | 9 |
| check_numeric_or_integer_type | 10 |
| check_posixct_type | 11 |
| check_rows | 11 |
| check_zero_columns | 12 |
| count_more_than_1 | 13 |
| create_dataset_summary_table | 13 |
| drop_na_column_names | 14 |
| duplicates_in_column | 14 |
| find_common_columns | 15 |
| find_maximum_value | 15 |
| find_minimum_value | 16 |
| find_pattern_r | 17 |
| get_distribution_statistics | 17 |
| get_first_element_class | 18 |
| get_values | 18 |
| identify_join_pairs | 19 |
| identify_outliers | 19 |
| is_unique_column | 20 |
| md_complete_cases | 21 |
| regex_content_parameter | 21 |
| regex_time | 22 |
| regex_year_date | 23 |
| remove_duplicates_and_na | 24 |
| retrieve_functions_and_packages | 24 |
| retrieve_function_calls | 25 |
| retrieve_package_usage | 25 |
| retrieve_sourced_scripts | 26 |
| retrieve_string_assignments | 26 |
| return_assertions_message | 27 |
| str_detect_in_file | 28 |
| test_all_equal | 28 |
| unique_id | 29 |

Index **30**

assertion_message *Assert Message Based on Type*

Description

This function asserts a message based on the type specified. It can either push the message to an AssertCollection, print a warning, or stop execution with an error message.

Usage

```
assertion_message(message, assertion_fail = "stop")
```

Arguments

`message` A character string representing the message to be asserted.

`assertion_fail` A character string indicating the action to take if the assertion fails. Can be an `AssertCollection`, "warning", or "stop" (default).

Value

None

| | |
|--------------------------------|------------------------------------|
| <code>assert_date_named</code> | <i>Assert Date Value in Column</i> |
|--------------------------------|------------------------------------|

Description

This function asserts that the values in a specified column of a data frame are of Date type. It uses the `checkmate::assert_date` function to perform the assertion.

Usage

```
assert_date_named(column, df, prefix_column = NULL, ...)
```

Arguments

`column` A character vector or string with the column name to be tested.

`df` The data frame that contains the column.

`prefix_column` A character string that will be prepended to the column name in the assertion message. Default is `NULL`.

`...` Additional parameters are passed to the `checkmate::assert_date` function.

Value

None

assert_logical_named *Assert Logical Value in Column*

Description

This function asserts that the values in a specified column of a data frame are logical. It uses the `checkmate::assert_logical` function to perform the assertion.

Usage

```
assert_logical_named(column, df, prefix_column = NULL, ...)
```

Arguments

| | |
|----------------------------|---|
| <code>column</code> | A character vector or string with the column name to be tested. |
| <code>df</code> | The data frame that contains the column. |
| <code>prefix_column</code> | A character string that will be prepended to the column name in the assertion message. Default is <code>NULL</code> . |
| <code>...</code> | Additional parameters are passed to the <code>checkmate::assert_logical</code> function. |

Value

None

Examples

```
# Create a data frame
df <- data.frame(a = c(TRUE, FALSE, TRUE, FALSE), b = c(1, 2, 3, 4))
# Assert that the values in column "a" are logical
assert_logical_named("a", df)
```

assert_no_duplicates_in_group
Assert No Duplicates in Group

Description

This function asserts that there are no duplicate rows in the specified columns of a data frame. It groups the data frame by the specified columns, counts the number of unique values for each group, and checks if there are any groups with more than one row. If there are, it prints an error message and stops the execution (unless `assertion_fail` is set to "warn").

Usage

```
assert_no_duplicates_in_group(df, group_vars, assertion_fail = "stop")
```

Arguments

`df` A data frame.

`group_vars` A character vector of column names.

`assertion_fail` A character string indicating the action to take if the assertion fails. Can be "stop" (default) or "warn".

Value

The input data frame.

calculate_category_percentages

Calculate the percentage of categories in a data vector

Description

This function calculates the percentage of each category in a given data vector and returns the top 10 categories along with their percentages. If the data vector is of Date class, it is converted to POSIXct. If the sum of the percentages is not 100%, an "Other" category is added to make up the difference, but only if the number of unique values exceeds 10. If the data vector is of POSIXct class and the smallest percentage is less than 1%, the function returns "Not enough occurrences."

Usage

```
calculate_category_percentages(data_vector)
```

Arguments

`data_vector` A vector of categorical data.

Value

A character string detailing the top 10 categories and their percentages, or a special message indicating not enough occurrences or unsupported data type.

Examples

```
# Example with a character vector
data_vector <- c("cat", "dog", "bird", "cat", "dog", "cat", "other")
calculate_category_percentages(data_vector)

# Example with a Date vector
data_vector <- as.Date(c("2020-01-01", "2020-01-02", "2020-01-03"))
calculate_category_percentages(data_vector)
```

check_double_columns *check double columns*

Description

Check whether two dataframes have intersecting column names.

Usage

```
check_double_columns(x, y, connector = NULL)
```

Arguments

| | |
|-----------|--|
| x | Data frame x. |
| y | Data frame y. |
| connector | The connector columns as strings. Also possible as vector. |

Value

Message informing about overlap in columns between the dataframes.

See Also

Other tests: [check_no_duplicates_in_group\(\)](#), [check_numeric_or_integer_type\(\)](#), [check_posixct_type\(\)](#), [duplicates_in_column\(\)](#), [test_all_equal\(\)](#)

Examples

```
check_double_columns(mtcars, iris)
```

check_duplicates *Check for Duplicate Rows in Selected Columns*

Description

This function checks if there are any duplicate rows in the specified columns of a data frame. It prints the unique rows and returns a boolean indicating whether the number of rows in the original data frame is the same as the number of rows in the data frame with duplicate rows removed.

Usage

```
check_duplicates(data, columns)
```

Arguments

data A data frame.
columns A character vector of column names.

Value

A logical value indicating whether the number of rows in the original data frame is the same as the number of rows in the data frame with duplicate rows removed.

Examples

```
# Create a data frame
df <- data.frame(a = c(1, 2, 3, 1), b = c(4, 5, 6, 4), c = c(7, 8, 9, 7))
# Check for duplicate rows in the first two columns
check_duplicates(df, c("a", "b"))
```

check_na_columns *Check for columns with only NA values*

Description

This function checks if there are any columns in the provided dataframe that contain only NA values. If such columns exist, their names are added to the provided collection.

Usage

```
check_na_columns(df, collection)
```

Arguments

df A dataframe.
collection A list to store the names of the columns with only NA values.

Value

The updated collection.

Examples

```
# Create a dataframe with some columns containing only NA values
df <- data.frame(a = c(1, NA, 3), b = c(NA, NA, NA), c = c(4, 5, 6))
collection <- checkmate::makeAssertCollection()
check_na_columns(df, collection)
```

check_non_zero_rows *Check for Non-Zero Rows*

Description

This function checks if there are more than 0 rows in the provided dataframe. If there are 0 rows, a message is added to the provided collection.

Usage

```
check_non_zero_rows(dataframe, collection)
```

Arguments

dataframe A dataframe.
collection A list to store the message if there are 0 rows.

Value

The updated collection.

Examples

```
# Create an empty dataframe  
dataframe <- data.frame()  
collection <- checkmate::makeAssertCollection()  
check_non_zero_rows(dataframe, collection)
```

check_no_duplicates_in_group
 Check for No Duplicates in Group

Description

This function checks if there is exactly one row per group in the provided dataframe. If there are multiple rows per group, the assertion fails.

Usage

```
check_no_duplicates_in_group(  
  dataframe,  
  group_variables = NULL,  
  assertion_fail = "stop"  
)
```


Arguments

`dataframe` The dataframe to be checked.

`group_variables` The group variables as a character vector. The default is NULL.

`assertion_fail` How the function reacts to a failure. This can be a "warning", where only a warning is given on the failure, or a "stop", where the function execution is stopped and the message is displayed, or an "AssertCollection", where the failure message is added to an assertion collection.

See Also

Other assertions: [check_numeric_or_integer_type\(\)](#), [check_posixct_type\(\)](#)

Other tests: [check_double_columns\(\)](#), [check_numeric_or_integer_type\(\)](#), [check_posixct_type\(\)](#), [duplicates_in_column\(\)](#), [test_all_equal\(\)](#)

Examples

```
# Create a dataframe with some groups having more than one row
dataframe <- data.frame(a = c(1, 1, 2), b = c(2, 2, 3), c = c("x", "x", "y"))
# Check the uniqueness of rows per group
check_no_duplicates_in_group(dataframe)
```

check_no_duplicate_rows

Check for No Duplicate Rows

Description

This function checks if there are any duplicate rows in the provided dataframe. If there are duplicate rows, a message is added to the provided collection.

Usage

```
check_no_duplicate_rows(dataframe, collection, unique_columns = NULL)
```

Arguments

`dataframe` A dataframe.

`collection` A list to store the message if there are duplicate rows.

`unique_columns` Default is NULL. If provided, these are the columns to check for uniqueness.

Value

The updated collection.

Examples

```
# Create a dataframe with some duplicate rows
dataframe <- data.frame(a = c(1, 1, 2), b = c(2, 2, 3))
collection <- checkmate::makeAssertCollection()
check_no_duplicate_rows(dataframe, collection, c("a", "b"))
```

```
check_numeric_or_integer_type
```

```
    Check for Numeric or Integer Type
```

Description

This function checks if the specified column in the provided dataframe has a numeric or integer type. It uses the `checkmate::assert_numeric` or `checkmate::assert_integer` function to perform the assertion, depending on the value of the `field_type` parameter.

Usage

```
check_numeric_or_integer_type(
  column_name,
  dataframe,
  column_prefix = NULL,
  field_type = "numeric",
  ...
)
```

Arguments

| | |
|----------------------------|---|
| <code>column_name</code> | A character vector or string with the column name to be tested. |
| <code>dataframe</code> | The dataframe that contains the column. |
| <code>column_prefix</code> | Default is <code>NULL</code> . If provided, this text is prepended to the variable name in the assertion message. |
| <code>field_type</code> | Default is <code>"numeric"</code> . Specify <code>"integer"</code> to check if the column has an integer type. This parameter must be either <code>"integer"</code> or <code>"numeric"</code> . |
| <code>...</code> | The remaining parameters are passed to the function <code>assert_numeric</code> or <code>assert_integer</code> . |

See Also

Other assertions: [check_no_duplicates_in_group\(\)](#), [check_posixct_type\(\)](#)

Other tests: [check_double_columns\(\)](#), [check_no_duplicates_in_group\(\)](#), [check_posixct_type\(\)](#), [duplicates_in_column\(\)](#), [test_all_equal\(\)](#)

Examples

```
# Create a dataframe with a numeric column
dataframe <- data.frame(a = c(1, 2, 3))
# Check the numeric type of the 'a' column
check_numeric_or_integer_type("a", dataframe)
```

| | |
|--------------------|-------------------------------|
| check_posixct_type | <i>Check for POSIXct Type</i> |
|--------------------|-------------------------------|

Description

This function checks if the specified column in the provided dataframe has a POSIXct type. It uses the `checkmate::assert_posixct` function to perform the assertion.

Usage

```
check_posixct_type(column_name, dataframe, column_prefix = NULL, ...)
```

Arguments

| | |
|----------------------------|---|
| <code>column_name</code> | A character vector or string with the column name to be tested. |
| <code>dataframe</code> | The dataframe that contains the column. |
| <code>column_prefix</code> | Default is NULL. If provided, this text is prepended to the variable name in the assertion message. |
| <code>...</code> | The remaining parameters are passed to the function <code>assert_posixct</code> . |

See Also

Other assertions: [check_no_duplicates_in_group\(\)](#), [check_numeric_or_integer_type\(\)](#)

Other tests: [check_double_columns\(\)](#), [check_no_duplicates_in_group\(\)](#), [check_numeric_or_integer_type\(\)](#), [duplicates_in_column\(\)](#), [test_all_equal\(\)](#)

Examples

```
# Create a dataframe with a POSIXct column
dataframe <- data.frame(date = as.POSIXct("2023-10-04"))
# Check the POSIXct type of the 'date' column
check_posixct_type("date", dataframe)
```

| | |
|------------|-------------------|
| check_rows | <i>Check rows</i> |
|------------|-------------------|

Description

This function prints the number of rows of a data frame. This function is used to check that rows are not deleted or doubled unless expected.

Usage

```
check_rows(df, name = NULL)
```

Arguments

df The data frame whose rows are to be counted
name The name of the data file (this will be printed)

Value

A message is printed to the console with the number of rows of the data

Examples

```
check_rows(mtcars)
```

check_zero_columns *Check for Columns with Only 0s*

Description

This function checks if there are any columns in the provided dataframe that contain only 0 values. If such columns exist, their names are added to the provided collection.

Usage

```
check_zero_columns(dataframe, collection)
```

Arguments

dataframe A dataframe.
collection A list to store the names of the columns with only 0 values.

Value

The updated collection.

Examples

```
# Create a dataframe with some columns containing only 0 values  
dataframe <- data.frame(a = c(0, 0, 0), b = c(1, 2, 3), c = c(0, 0, 0))  
collection <- checkmate::makeAssertCollection()  
check_zero_columns(dataframe, collection)
```

| | |
|-------------------|--------------------------|
| count_more_than_1 | <i>Count more than 1</i> |
|-------------------|--------------------------|

Description

Function to count the number of values greater than 1 in a vector This function is used in the function `Check_columns_for_double_rows` to count duplicate values.

Usage

```
count_more_than_1(x)
```

Arguments

| | |
|---|--------------------|
| x | The vector to test |
|---|--------------------|

Value

Number of values greater than 1.

Examples

```
count_more_than_1(c(1, 1, 4))
```

| | |
|------------------------------|--|
| create_dataset_summary_table | <i>Create dataset summary statistics table</i> |
|------------------------------|--|

Description

This function creates a summary statistics table for a dataframe, providing insights into the nature of the data contained within. It includes detailed statistics for each column, such as column types, missing value percentages, minimum and maximum values for numeric columns, patterns for character columns, uniqueness of identifiers, and distributions.

Usage

```
create_dataset_summary_table(df_input)
```

Arguments

| | |
|----------|---|
| df_input | A dataframe for which to create a summary statistics table. |
|----------|---|

Value

A tibble with comprehensive summary statistics for each column in the input dataframe.

drop_na_column_names *Drop NA column names*

Description

Deletes columns whose name is NA or whose name is empty

Usage

```
drop_na_column_names(x)
```

Arguments

x dataframe

Value

dataframe without columns that are NA

duplicates_in_column *Duplicates in column*

Description

Searches for duplicates in a data frame column.

Usage

```
duplicates_in_column(df, col)
```

Arguments

df Data frame.
col Column name.

Value

Rows containing duplicated values.

See Also

Other tests: [check_double_columns\(\)](#), [check_no_duplicates_in_group\(\)](#), [check_numeric_or_integer_type\(\)](#), [check_posixct_type\(\)](#), [test_all_equal\(\)](#)

Examples

```
duplicates_in_column(mtcars, "mpg")
```

find_common_columns *Find Common Columns Between Data Frames*

Description

This function identifies common column names between multiple data frames. It takes a variable number of data frames as input and returns a character vector containing the common column names.

Usage

```
find_common_columns(...)
```

Arguments

... A variable length list of data frames.

Value

A character vector of column names found in common between all data frames.

Examples

```
df1 <- data.frame(a = c(1, 2, 3), b = c(4, 5, 6))
df2 <- data.frame(a = c(7, 8, 9), b = c(10, 11, 12), c = c(13, 14, 15))
common_columns <- find_common_columns(df1, df2)
print(common_columns)
```

find_maximum_value *Find the maximum numeric value in a vector, ignoring non-numeric values*

Description

Find the maximum numeric value in a vector, ignoring non-numeric values

Usage

```
find_maximum_value(numeric_vector)
```

Arguments

numeric_vector A vector from which to find the maximum numeric value.

Value

The maximum numeric value in the input vector, or NA if none exist.

Examples

```
# Find the maximum of a numeric vector
find_maximum_value(c(3, 1, 4, 1, 5, 9)) # Returns 9

# Find the maximum of a mixed vector with non-numeric values
find_maximum_value(c(3, 1, 4, "two", 5, 9)) # Returns 9

# Attempt to find the maximum of a vector with only non-numeric values
find_maximum_value(c("one", "two", "three")) # Returns NA
```

| | |
|--------------------|--|
| find_minimum_value | <i>Find the minimum numeric value in a vector, ignoring non-numeric values</i> |
|--------------------|--|

Description

Find the minimum numeric value in a vector, ignoring non-numeric values

Usage

```
find_minimum_value(numeric_vector)
```

Arguments

`numeric_vector` A vector from which to find the minimum numeric value.

Value

The minimum numeric value in the input vector, or NA if none exist.

Examples

```
# Find the minimum of a numeric vector
find_minimum_value(c(3, 1, 4, 1, 5, 9)) # Returns 1

# Find the minimum of a mixed vector with non-numeric values
find_minimum_value(c(3, 1, 4, "two", 5, 9)) # Returns 1

# Attempt to find the minimum of a vector with only non-numeric values
find_minimum_value(c("one", "two", "three")) # Returns NA
```

| | |
|----------------|----------------------------------|
| find_pattern_r | <i>Find pattern in R scripts</i> |
|----------------|----------------------------------|

Description

Function to search for a pattern in R scripts.

Usage

```
find_pattern_r(pattern, path = ".", case.sensitive = TRUE, comments = FALSE)
```

Arguments

| | |
|----------------|--|
| pattern | Pattern to search |
| path | Directory to search in |
| case.sensitive | Whether pattern is case sensitive or not |
| comments | whether to search in commented lines |

Value

Dataframe containing R script paths

| | |
|-----------------------------|---|
| get_distribution_statistics | <i>Compute distribution statistics for a numeric vector</i> |
|-----------------------------|---|

Description

This function computes summary statistics such as quartiles, mean, and standard deviation for a numeric vector.

Usage

```
get_distribution_statistics(data_vector)
```

Arguments

| | |
|-------------|---|
| data_vector | A numeric vector for which to compute summary statistics. |
|-------------|---|

Value

A character string describing the summary statistics of the input vector.

Examples

```
# Compute summary statistics for a numeric vector
data_vector <- c(1, 2, 3, 4, 5, 6, 7, 8, 9, 10)
get_distribution_statistics(data_vector)
```

`get_first_element_class`*Retrieve the class of the first element of a vector*

Description

Retrieve the class of the first element of a vector

Usage

```
get_first_element_class(input_vector)
```

Arguments

`input_vector` A vector whose first element's class is to be retrieved.

Value

The class of the first element of the input vector.

Examples

```
# Get the class of the first element in a numeric vector
get_first_element_class(c(1, 2, 3)) # Returns "numeric"

# Get the class of the first element in a character vector
get_first_element_class(c("apple", "banana", "cherry")) # Returns "character"
```

`get_values`*Get values of column*

Description

A function to determine what kind of values are present in columns.

Usage

```
get_values(df, column)
```

Arguments

`df` The dataframe
`column` Column to get values from.

Value

The class of the column values

Examples

```
get_values(mtcars, "mpg")
```

identify_join_pairs *Identify Possible Join Pairs Between Data Frames*

Description

This function identifies potential join pairs between two data frames based on the overlap between the distinct values in their columns. It returns a data frame showing the possible join pairs.

Usage

```
identify_join_pairs(..., similarity_cutoff = 0.2)
```

Arguments

... A list of two data frames.
similarity_cutoff The minimal percentage of overlap between the distinct values in the columns.

Value

A data frame showing candidate join pairs.

Examples

```
identify_join_pairs(iris, iris3)
```

identify_outliers *Identify Outliers in a Data Frame Column*

Description

This function identifies outliers in a specified column of a data frame. It returns a tibble containing the unique values, tally, and whether it is an outlier or not.

Usage

```
identify_outliers(df, var)
```

Arguments

df The data frame.
var The column to check for outliers.

Value

A tibble containing the unique values, tally, and whether each value is an outlier or not.

Examples

```
df <- data.frame(a = c(1, 2, 3, 100, 101), b = c(4, 5, 6, 7, 8), c = c(7, 8, 9, 100, 101))
outliers <- identify_outliers(df, "a")
print(outliers)
```

| | |
|------------------|---|
| is_unique_column | <i>Check if a column in a dataframe has unique values</i> |
|------------------|---|

Description

Check if a column in a dataframe has unique values

Usage

```
is_unique_column(column_name, data_frame)
```

Arguments

| | |
|-------------|---|
| column_name | The name of the column to check for uniqueness. |
| data_frame | A dataframe containing the column to check. |

Value

TRUE if the column has unique values, FALSE otherwise.

Examples

```
# Create a dataframe with a unique ID column
data_frame <- tibble::tibble(
  id = c(1, 2, 3, 4, 5),
  value = c("a", "b", "c", "d", "e")
)
is_unique_column("id", data_frame) # Returns TRUE

# Create a dataframe with duplicate values in the ID column
data_frame <- tibble::tibble(
  id = c(1, 2, 3, 4, 5, 1),
  value = c("a", "b", "c", "d", "e", "a")
)
is_unique_column("id", data_frame) # Returns FALSE
```

| | |
|-------------------|--------------------------|
| md_complete_cases | <i>MD complete cases</i> |
|-------------------|--------------------------|

Description

Print the complete cases of the data.

Usage

```
md_complete_cases(data, digits = 1)
```

Arguments

| | |
|--------|--|
| data | The data frame. |
| digits | Default: 1. number of digits for rounding. |

Value

Message with the number of rows, number of rows with missing values and the percentage of complete rows.

Examples

```
# example code
md_complete_cases(iris)

iris$Sepal.Length[5] <- NA_character_
md_complete_cases(iris)
```

| |
|-------------------------|
| regex_content_parameter |
|-------------------------|

Construct Regex for Matching Function Parameter Content

Description

This function constructs a regex pattern for matching the content of a parameter in a function. It uses the `base::paste0` function to construct the regex pattern.

Usage

```
regex_content_parameter(parameter)
```

Arguments

| | |
|-----------|--|
| parameter | The parameter whose value is to be searched in a function. |
|-----------|--|

Value

A regex pattern as a character string.

Examples

```
# Create a parameter name
parameter <- "my_parameter"
# Construct a regex pattern for matching the content of the parameter
pattern <- regex_content_parameter(parameter)
```

regex_time

Generate regular expression of a time.

Description

This function generates a regular expression for time based on the input format.

Usage

```
regex_time(format = "hh:mm")
```

Arguments

| | |
|--------|--|
| format | The format of the time. Possible values are: <ul style="list-style-type: none">• "hh:mm": to generate "09:05".• "h:m": to generate "9:5".• "hh:mm:ss": to generate "09:05:00".• "h:m:s": to generate "9:5:0".• "hh:mm:ss AM/PM": to generate "09:05:00 AM".• "h:m:s AM/PM": to generate "9:5:0 AM". |
|--------|--|

Value

A regular expression.

Examples

```
regex_time("hh:mm")
regex_time("h:m")
regex_time("hh:mm:ss")
regex_time("h:m:s")
regex_time("hh:mm:ss AM/PM")
regex_time("h:m:s AM/PM")
```

| | |
|-----------------|--|
| regex_year_date | <i>Generate regular expression of a year date.</i> |
|-----------------|--|

Description

This function generates a regular expression for year date based on the input format.

Usage

```
regex_year_date(format = "yyyy")
```

Arguments

| | |
|--------|--|
| format | The format of the year date. Possible values are: <ul style="list-style-type: none">• "yyyy": to generate "2023".• "yyyy-MM-dd": to generate "2023-09-29".• "yyyy/MM/dd": to generate "2023/09/29".• "yyyy.MM.dd": to generate "2023.09.29".• "yyyy-M-d": to generate "2023-9-29".• "yyyy/M/d": to generate "2023/9/29".• "yyyy.M.d": to generate "2023.9.29".• "yyyy-MM-dd HH:mm:ss": to generate "2023-09-29 12:34:56".• "yyyy/MM/dd HH:mm:ss": to generate "2023/09/29 12:34:56".• "yyyy-MM-dd HH:mm": to generate "2023-09-29 12:34".• "yyyy/MM/dd HH:mm": to generate "2023/09/29 12:34". |
|--------|--|

Value

A regular expression.

Examples

```
regex_year_date("yyyy")
regex_year_date("yyyy-MM-dd")
regex_year_date("yyyy/MM/dd")
regex_year_date("yyyy.MM.dd")
regex_year_date("yyyy-M-d")
regex_year_date("yyyy/M/d")
regex_year_date("yyyy.M.d")
regex_year_date("yyyy-MM-dd HH:mm:ss")
regex_year_date("yyyy/MM/dd HH:mm:ss")
regex_year_date("yyyy-MM-dd HH:mm")
regex_year_date("yyyy/MM/dd HH:mm")
```

remove_duplicates_and_na

Remove Duplicates and NA Values from Input

Description

This function removes duplicate values and NA values from the input. It first removes NA values from the input using the `na.omit` function from the `stats` package. Then it removes duplicate values from the result using the `unique` function.

Usage

```
remove_duplicates_and_na(input)
```

Arguments

`input` A vector or data frame.

Value

A vector or data frame with duplicate values and NA values removed.

Examples

```
# Create a vector with duplicate values and NA values
input <- c(1, 2, NA, 2, NA, 3, 4, 4, NA, 5)
# Remove duplicate values and NA values
output <- remove_duplicates_and_na(input)
print(output)
```

retrieve_functions_and_packages

Retrieve functions and packages

Description

Retrieves functions and their corresponding packages used in a given script.

Usage

```
retrieve_functions_and_packages(path)
```

Arguments

`path` The complete path of the script.

Value

Used_functions

retrieve_function_calls
retrieve_function_calls

Description

retrieve_function_calls

Usage

retrieve_function_calls(script_name)

Arguments

script_name The script to search functions in

Value

dataframe

retrieve_package_usage
Retrieve packages that are loaded in a script

Description

Retrieve packages that are loaded in a script

Usage

retrieve_package_usage(script_name)

Arguments

script_name The path to the R script

Value

dataframe

```
retrieve_sourced_scripts  
    retrieve_sourced_scripts
```

Description

retrieve_sourced_scripts

Usage

```
retrieve_sourced_scripts(script_name)
```

Arguments

script_name The main script to search

Value

dataframe

```
retrieve_string_assignments  
    retrieve_string_assignments
```

Description

retrieve_string_assignments

Usage

```
retrieve_string_assignments(script_name)
```

Arguments

script_name The script to search objects in

Value

dataframe

`return_assertions_message`*Return Assertion Messages*

Description

This function returns a message indicating whether an assertion test has passed or failed. An "assertion collection" from the checkmate package must be provided. The message can be returned as an error or a warning. For some assertions, only warnings are allowed, as an error would stop the script from running. This is done for the following assertions: percentage missing values, duplicates, subset, and set_equal.

Usage

```
return_assertions_message(  
  collection,  
  collection_name,  
  fail = "stop",  
  silent = FALSE,  
  output_map = NULL  
)
```

Arguments

| | |
|------------------------------|---|
| <code>collection</code> | An object with the class "AssertCollection". |
| <code>collection_name</code> | The name of the collection. This name is mentioned in the messages. |
| <code>fail</code> | "stop" or "warning". If the assertions fail, an error is returned and the script output is stopped. If "warning", only a warning is returned. |
| <code>silent</code> | If FALSE (default), the success message is printed in the console. If TRUE, it is not shown. |
| <code>output_map</code> | A map, like 1. Read data, where the file is stored. |

Value

The message indicating whether the assertion test has passed or failed.

str_detect_in_file *Detect string in file*

Description

Detect string in file

Usage

```
str_detect_in_file(file, pattern, only_comments = FALSE, collapse = FALSE)
```

Arguments

| | |
|---------------|---|
| file | Path to file. |
| pattern | Pattern to match. |
| only_comments | default FALSE. Whether to only search in commented lines. |
| collapse | default: FALSE: search file line by line. If true, then pattern is search in the entire file at once after collapsing. (only_comments does not work when collapse is set to TRUE) |

Value

Boolean whether pattern exists in file.

test_all_equal *Test all equal*

Description

Test whether all values in a vector are equal.

Usage

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

Arguments

| | |
|-------|--|
| x | Vector to test. |
| na.rm | default: FALSE. exclude NAs from the test. |

Value

Boolean result of the test

See Also

Other tests: [check_double_columns\(\)](#), [check_no_duplicates_in_group\(\)](#), [check_numeric_or_integer_type\(\)](#), [check_posixct_type\(\)](#), [duplicates_in_column\(\)](#)

Examples

```
test_all_equal(c(5, 5, 5))
```

```
test_all_equal(c(5, 6, 3))
```

| | |
|-----------|------------------|
| unique_id | <i>unique id</i> |
|-----------|------------------|

Description

Check if parsed variable is a unique identifier. This function was adapted from: Source: <https://edwinth.github.io/blog/unique>.

Usage

```
unique_id(x, ...)
```

Arguments

x vector or dataframe.
... optional variables, e.g. name of column or a vector of names.

Value

Boolean whether variable is a unique identifier.

Examples

```
unique_id(iris, Species)
```

```
mtcars$name <- rownames(mtcars)
```

```
unique_id(mtcars, name)
```

Index

* assertions

- check_no_duplicates_in_group, 8
- check_numeric_or_integer_type, 10
- check_posixct_type, 11

* tests

- check_double_columns, 6
- check_no_duplicates_in_group, 8
- check_numeric_or_integer_type, 10
- check_posixct_type, 11
- duplicates_in_column, 14
- test_all_equal, 28

* vector calculations

- count_more_than_1, 13

- assert_date_named, 3
- assert_logical_named, 4
- assert_no_duplicates_in_group, 4
- assertion_message, 2

- calculate_category_percentages, 5
- check_double_columns, 6, 9–11, 14, 29
- check_duplicates, 6
- check_na_columns, 7
- check_no_duplicate_rows, 9
- check_no_duplicates_in_group, 6, 8, 10, 11, 14, 29
- check_non_zero_rows, 8
- check_numeric_or_integer_type, 6, 9, 10, 11, 14, 29
- check_posixct_type, 6, 9, 10, 11, 14, 29
- check_rows, 11
- check_zero_columns, 12
- count_more_than_1, 13
- create_dataset_summary_table, 13

- drop_na_column_names, 14
- duplicates_in_column, 6, 9–11, 14, 29

- find_common_columns, 15
- find_maximum_value, 15

- find_minimum_value, 16
- find_pattern_r, 17

- get_distribution_statistics, 17
- get_first_element_class, 18
- get_values, 18

- identify_join_pairs, 19
- identify_outliers, 19
- is_unique_column, 20

- md_complete_cases, 21

- regex_content_parameter, 21
- regex_time, 22
- regex_year_date, 23
- remove_duplicates_and_na, 24
- retrieve_function_calls, 25
- retrieve_functions_and_packages, 24
- retrieve_package_usage, 25
- retrieve_sourced_scripts, 26
- retrieve_string_assignments, 26
- return_assertions_message, 27

- str_detect_in_file, 28

- test_all_equal, 6, 9–11, 14, 28

- unique_id, 29