

# Package ‘sdLog’

October 8, 2021

**Title** Tools for Statistical Disclosure Control in Research Data Centers

**Version** 0.4.0

**Description** Tools for researchers to explicitly show that their results comply to rules for statistical disclosure control imposed by research data centers. These tools help in checking descriptive statistics and models and in calculating extreme values that are not individual data. Also included is a simple function to create log files. The methods used here are described in the “Guidelines for the checking of output based on microdata research” by Bond, Brandt, and de Wolf (2015) <[https://ec.europa.eu/eurostat/cros/system/files/dwb\\_standalone-document\\_output-checking-guidelines.pdf](https://ec.europa.eu/eurostat/cros/system/files/dwb_standalone-document_output-checking-guidelines.pdf)>.

**License** GPL-3

**URL** <https://github.com/matthiasgomolka/sdLog>

**BugReports** <https://github.com/matthiasgomolka/sdLog/issues>

**Depends** R (>= 3.5)

**Imports** broom (>= 0.5.5),  
checkmate (>= 2.0.0),  
crayon (>= 1.3.4),  
data.table (>= 1.12.8),  
mathjaxr,  
stats,  
utils

**Suggests** knitr,  
lfe,  
rmarkdown,  
skimr,  
spelling,  
testthat (>= 3.0.0),  
tibble

**VignetteBuilder** knitr

**RdMacros** mathjaxr

**Config/testthat/edition** 3

**Encoding** UTF-8**Language** en-US**LazyData** true**Roxygen** list(markdown = TRUE)**RoxygenNote** 7.1.2

## R topics documented:

common_arguments . . . . .	2
sdv_descriptives . . . . .	3
sdv_descriptives_DT . . . . .	4
sdv_log . . . . .	5
sdv_min_max . . . . .	6
sdv_min_max_DT . . . . .	7
sdv_model . . . . .	7
sdv_model_DT . . . . .	8

<b>Index</b>	<b>9</b>
--------------	----------

---

common_arguments	<i>arguments</i>
------------------	------------------

---

### Description

arguments

### Arguments

data	<a href="#">data.frame</a> from which the descriptive statistics are calculated.
id_var	<a href="#">character</a> The name of the id variable. Defaults to <code>getOption("sdv.id_var")</code> so that you can provide <code>options(sdv.id_var = "my_id_var")</code> at the top of your script.
val_var	<a href="#">character</a> vector of value variables on which descriptive statistics are computed.
by	<a href="#">character</a> vector of grouping variables.
zero_as_NA	<a href="#">logical</a> If TRUE, zeros in 'val_var' are treated as NA.
model	The estimated model object. Can be a model type like <a href="#">lm</a> , <a href="#">glm</a> and various others (anything which can be handled by <a href="#">broom::augment()</a> ).
min_obs	<a href="#">integer</a> The minimum number of observations used to calculate the minimum and maximum. Defaults to <code>getOption("sdv.n_ids", 5L)</code> . <i>This is not the number of distinct entities.</i>
max_obs	<a href="#">integer</a> The maximum number of observations used to calculate the minimum and maximum. Defaults to <code>nrow(data)</code> . <i>This is not the number of distinct entities.</i>

---

sdc\_descriptives      *Disclosure control for descriptive statistics*

---

## Description

Checks the number of distinct entities and the  $(n, k)$  dominance rule for your descriptive statistics.

That means that `sdc_descriptives()` checks if there are at least 5 distinct entities and if the largest 2 entities account for 85% or more of `val_var`. The parameters can be changed using options. For details see `vignette("options", package = "sdcLog")`.

## Usage

```
sdc_descriptives(
  data,
  id_var = getOption("sdc.id_var"),
  val_var = NULL,
  by = NULL,
  zero_as_NA = NULL
)
```

## Arguments

`data`            [data.frame](#) from which the descriptive statistics are calculated.

`id_var`           [character](#) The name of the id variable. Defaults to `getOption("sdc.id_var")` so that you can provide `options(sdc.id_var = "my_id_var")` at the top of your script.

`val_var`          [character](#) vector of value variables on which descriptive statistics are computed.

`by`                [character](#) vector of grouping variables.

`zero_as_NA`      [logical](#) If TRUE, zeros in 'val\_var' are treated as NA.

## Details

The general form of the  $(n, k)$  dominance rule can be formulated as:

$$\sum_{i=1}^n x_i > \frac{k}{100} \sum_{i=1}^N x_i$$

where  $x_1 \geq x_2 \geq \dots \geq x_N$ .  $n$  denotes the number of largest contributions to be considered,  $x_n$  the  $n$ -th largest contribution,  $k$  the maximal percentage these  $n$  contributions may account for, and  $N$  is the total number of observations.

If the statement above is true, the  $(n, k)$  dominance rule is violated.

## Value

A [list](#) of class `sdc_descriptives` with detailed information about options, settings, and compliance with the criteria distinct entities and dominance.

**Examples**

```
sdc_descriptives(  
  data = sdc_descriptives_DT,  
  id_var = "id",  
  val_var = "val_1"  
)
```

```
sdc_descriptives(  
  data = sdc_descriptives_DT,  
  id_var = "id",  
  val_var = "val_1",  
  by = "sector"  
)
```

```
sdc_descriptives(  
  data = sdc_descriptives_DT,  
  id_var = "id",  
  val_var = "val_1",  
  by = c("sector", "year")  
)
```

```
sdc_descriptives(  
  data = sdc_descriptives_DT,  
  id_var = "id",  
  val_var = "val_2",  
  by = c("sector", "year")  
)
```

```
sdc_descriptives(  
  data = sdc_descriptives_DT,  
  id_var = "id",  
  val_var = "val_2",  
  by = c("sector", "year"),  
  zero_as_NA = FALSE  
)
```

---

sdc\_descriptives\_DT    *Example data for sdc\_descriptives()*

---

**Description**

Utilized in the vignette.

**Usage**

```
data("sdc_descriptives_DT")
```

**Format**

A `data.table` with 20 rows and 5 columns.

**Details**

The `data.table` contains the following columns:

- `id` **factor** random identifier
- `sector` **factor** economic sector
- `year` **integer** time variable
- `val_1`, `val_2` **numeric** value variables

---

sdc\_log

*Create Stata-like log files from R Scripts*

---

**Description**

This function creates Stata-like log files from R Scripts. It can handle several files (in a **character** vector) at once.

**Usage**

```
sdc_log(r_script, destination, replace = FALSE, append = FALSE, local = FALSE)
```

**Arguments**

<code>r_script</code>	<b>character</b> Path of the R script to be run with logging.
<code>destination</code>	One of: <ul style="list-style-type: none"> <li>• <b>character</b> Path of the log file to be used.</li> <li>• <b>file</b> connection to which the log should be written. This is especially useful, when you have nested calls to <code>sdc_log()</code> and want to write everything into the same log file. Then, create a single <b>file</b> connection and provide this connection to all calls to <code>sdc_log()</code> (and close it afterwards).</li> </ul>
<code>replace</code>	<b>logical</b> Indicates whether to replace an existing log file.
<code>append</code>	<b>logical</b> Indicates whether to append an existing log file.
<code>local</code>	One of: <ul style="list-style-type: none"> <li>• <b>logical</b> Indicates whether to evaluate within the global environment (FALSE) or the calling environment (TRUE).</li> <li>• <b>environment</b> A specific evaluation environment. Determines the evaluation environment. Useful whenever <code>sdc_log()</code> is called from within a function, or for nested <code>sdc_log()</code> calls. By default (FALSE) evaluation occurs in the global environment. See also <b>source</b>.</li> </ul>

**Value**

**character** vector holding the path(s) of the written log file(s).

sdc\_min\_max

*Calculate RDC rule-compliant extreme values***Description**

Checks if calculation of extreme values comply to RDC rules. If so, function returns average min and max values according to RDC rules.

**Usage**

```
sdc_min_max(
  data,
  id_var = getOption("sdc.id_var"),
  val_var,
  by = NULL,
  max_obs = nrow(data)
)
```

**Arguments**

**data** [data.frame](#) from which the descriptive statistics are calculated.

**id\_var** [character](#) The name of the id variable. Defaults to `getOption("sdc.id_var")` so that you can provide `options(sdc.id_var = "my_id_var")` at the top of your script.

**val\_var** [character](#) vector of value variables on which descriptive statistics are computed.

**by** [character](#) vector of grouping variables.

**max\_obs** [integer](#) The maximum number of observations used to calculate the minimum and maximum. Defaults to `nrow(data)`. *This is not the number of distinct entities.*

**Value**

A list [list](#) of class `sdc_min_max` with detailed information about options, settings and the calculated extreme values (if possible).

**Examples**

```
sdc_min_max(sdc_min_max_DT, id_var = "id", val_var = "val_1")
sdc_min_max(sdc_min_max_DT, id_var = "id", val_var = "val_2")
sdc_min_max(sdc_min_max_DT, id_var = "id", val_var = "val_3", max_obs = 10)
sdc_min_max(sdc_min_max_DT, id_var = "id", val_var = "val_1", by = "year")
sdc_min_max(
  sdc_min_max_DT, id_var = "id", val_var = "val_1", by = c("sector", "year")
)
```

---

sdc_min_max_DT	<i>Example data for sdc_min_max()</i>
----------------	---------------------------------------

---

### Description

Utilized in the vignette

### Usage

```
data("sdc_min_max_DT")
```

### Format

A data.table with 20 rows and 6 columns.

### Details

The data.table contains the following columns:

- id [factor](#) random identifier
- sector [factor](#) economic sector
- year [integer](#) time variable
- val\_1 - val\_3 [numeric](#) value variables

---

sdc_model	<i>Disclosure control for models</i>
-----------	--------------------------------------

---

### Description

Checks if your model complies to RDC rules. Checks for overall number of entities and number of entities for each level of dummy variables.

### Usage

```
sdc_model(data, model, id_var = getOption("sdc.id_var"))
```

### Arguments

data	<a href="#">data.frame</a> which was used to build the model.
model	The estimated model object. Can be a model type like <a href="#">lm</a> , <a href="#">glm</a> and various others (anything which can be handled by <a href="#">broom::augment()</a> ).
id_var	<a href="#">character</a> The name of the id variable. Defaults to <code>getOption("sdc.id_var")</code> so that you can provide <code>options(sdc.id_var = "my_id_var")</code> at the top of your script.

**Value**

A [list](#) of class `sdc_model` with detailed information about options, settings, and compliance with the distinct entities criterion.

**Examples**

```
# Check simple models
model_1 <- lm(y ~ x_1 + x_2, data = sdc_model_DT)
sdc_model(data = sdc_model_DT, model = model_1, id_var = "id")

model_2 <- lm(y ~ x_1 + x_2 + x_3, data = sdc_model_DT)
sdc_model(data = sdc_model_DT, model = model_2, id_var = "id")

model_3 <- lm(y ~ x_1 + x_2 + dummy_3, data = sdc_model_DT)
sdc_model(data = sdc_model_DT, model = model_3, id_var = "id")
```

---

`sdc_model_DT`*Example data for sdc\_model()*

---

**Description**

Utilized in the vignette

**Usage**

```
data("sdc_model_DT")
```

**Format**

A `data.table` with 80 rows and 9 columns.

**Details**

The `data.table` contains the following columns:

- `id` [factor](#) random identifier
- `y - x_4` [numeric](#) value variables
- `dummy_1 - dummy_3` [factor](#) dummy variables

# Index

## \* datasets

- sdc\_descriptives\_DT, 4
- sdc\_min\_max\_DT, 7
- sdc\_model\_DT, 8

broom::augment(), 2, 7

character, 2, 3, 5–7  
common\_arguments, 2

data.frame, 2, 3, 6, 7

environment, 5

factor, 5, 7, 8  
file, 5

glm, 2, 7

integer, 2, 5–7

list, 3, 6, 8  
lm, 2, 7  
logical, 2, 3, 5

numeric, 5, 7, 8

sdc\_descriptives, 3  
sdc\_descriptives\_DT, 4  
sdc\_log, 5  
sdc\_min\_max, 6  
sdc\_min\_max\_DT, 7  
sdc\_model, 7  
sdc\_model\_DT, 8  
source, 5