

Statistical disclosure control and synthetic data generation using R

Alexander Kowarik,
Head of Statistical Methods and Survey Methodology

Johannes Gussenbauer
Statistical Methods and Survey Methodology

Vienna, 21 November 2023

www.statistik.at

Independent statistics for evidence-based decision making

Content

- SDC methods vs. synthetic data
- Micro data protection with sdcMicro
 - sdcApp(): Graphical User Interface
 - Target record swapping in R
- Anonymization of tabular data
- SDC for tabular data

Motivation

The image features a blue-tinted background of a modern building's interior, showing multiple levels with glass railings and potted plants. On the right side, there is a clear view through a window with a white frame, looking out at a modern building facade with a grid of windows and balconies.

Who is this person?

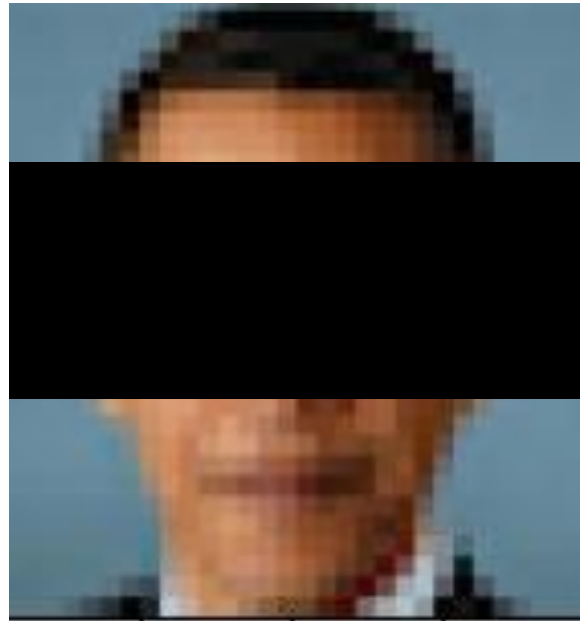
Synthetic data



Sources: <https://twitter.com/Chicken3gg/status/1274314622447820801>, <https://pbs.twimg.com/media/EbBrAKNX0A0NSO8?format=jpg&name=medium>

Who is this person?

SDC applied



Sources: <https://twitter.com/Chicken3gg/status/1274314622447820801> ; https://en.wikipedia.org/wiki/Social_policy_of_the_Barack_Obama_administration#/media/File:Official_portrait_of_Barack_Obama.jpg ; <https://pbs.twimg.com/media/EbBrAKNX0A0NSO8?format=jpg&name=medium>

SDC Methods vs. Synthetic Data

Which approach is better?

Original data

/

SDC applied

/

Synthetic data



Sources: <https://twitter.com/Chicken3gg/status/1274314622447820801> ; https://en.wikipedia.org/wiki/Social_policy_of_the_Barack_Obama_administration#/media/File:Official_portrait_of_Barack_Obama.jpg

SDC Methods vs. Synthetic Data

Which one is better?

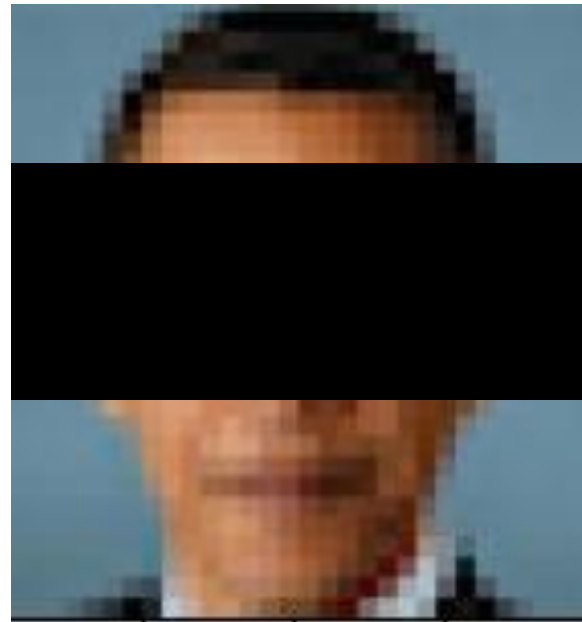
Original data

/

SDC applied

/

Synthetic data



Sources: <https://twitter.com/Chicken3gg/status/1274314622447820801> ; https://en.wikipedia.org/wiki/Social_policy_of_the_Barack_Obama_administration#/media/File:Official_portrait_of_Barack_Obama.jpg ; <https://pbs.twimg.com/media/EbBrAKNX0A0NSO8?format=jpg&name=medium>

SDC methods vs. Synthetic data

Both: Reducing risk / Maximizing utility

Both

/

SDC applied

/

Synthetic data

- More complicated for hierarchical data, e.g. households
- Re-identification risk $\neq 0$

- Loss of granularity
- Individual/manual approach

- High granularity
- Complex generation process
- Re-identification risk harder to assess (attribute disclosure, membership inference)

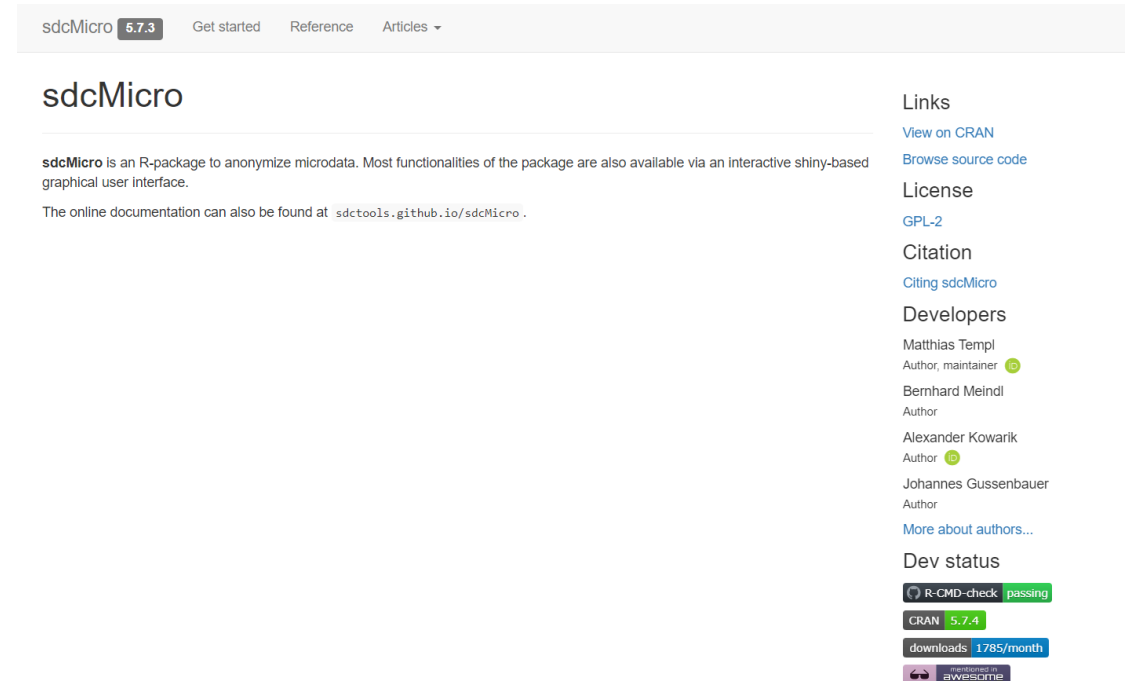
Micro data protection



Micro data protection in R

sdcMicro

- Released versions are on CRAN: <https://cran.r-project.org/web/packages/sdcMicro>
- Development version and issue tracking on Github: <https://github.com/sdcTools/sdcMicro>
- Collaborative development with colleagues from Statistics Austria and other institutions
- Documentation:
<http://sdctools.github.io/sdcMicro/>



The screenshot shows the CRAN page for the sdcMicro package. The page title is "sdcmicro 5.7.3". The navigation menu includes "Get started", "Reference", and "Articles". The main content area describes sdcMicro as an R-package for anonymizing microdata, with a link to the online documentation at sdctools.github.io/sdcMicro. The right sidebar contains sections for "Links" (View on CRAN, Browse source code), "License" (GPL-2), "Citation" (Citing sdcMicro), "Developers" (Matthias Templ, Bernhard Meindl, Alexander Kowarik, Johannes Gussenbauer), "Dev status" (R-CMD-check: passing, CRAN 5.7.4, downloads: 1785/month, featured on awesome), and "More about authors..."

sdcMicro

Categorical variables

- Risk estimation:
 - Concept of k-anonymity and
 - various methods to estimate the disclosure risk (individual risk, global risk, suda2)
- Deterministic methods
 - Non-perturbative protection methods
 - Top- and bottom coding
 - Recoding
 - Local suppression
 - Probabilistic methods
 - Perturbative protection methods
 - (Rank) swapping
 - Post-randomization (pram)

sdcMicro

Numerical variables

- Deterministic protection methods
 - Top- and bottom coding
 - Microaggregation
 - Rounding
- Perturbative protection methods based on randomness
 - Noise addition
 - (Rank)swapping
 - Shuffling

sdcMicroObj

- The S4-class sdcMicroObj Slots
 - original Data (@origData)
 - modified (key) variables (@manipKeyVars, @manipNumVars, ...)
 - categorical key variables (@keyVars)
 - numerical key variables (@numVars)
 - computed disclosure risk (@risk)
 - previous object (@prev)

sdcMicroObj methods

- The S4-class sdcMicroObj (some) Methods:
 - addNoise(): add noise to numerical key variables
 - Pram(): Post RAndomisation Method (PRAM)
 - rankSwap(): Rank Swapping
 - shuffle(): Shuffling and EGADP
 - report(): Create a report about the anonymization process
 - undolast(): undo last calculation
- . . .

sdcMicro

Tiny example

```
data(testdata)
sdc <- createSdcObj(testdata,
  keyVars=c('urbrur','roof','walls','water','electcon','relat'),
  numVars=c('expend','income','savings'),
  w='sampling_weight')
### Display Risk
sdc@risk$numeric [1] 1
### use addNoise without Parameters
sdc <- addNoise(sdc, variables = c("expend","income"))
### risk changed
sdc@risk$numeric [1] 0.07729258
```


sdCApp (trs not included)

- Online-Demo <https://sdctools.shinyapps.io/sdcapp/>

sdCApp GUI

About/Help Microdata Anonymize Risk/Utility Export Data Reproducibility Undo

What do you want to do?

- Display microdata
- Explore variables
- Reset variables
- Use subset of microdata
- Convert numeric to factor
- Convert variables to numeric
- Modify factor variable
- Create a stratification variable
- Set specific values to NA
- Hierarchical data

Reset inputdata

Loaded microdata

The loaded dataset is `testdata` and consists of 4580 observations and 15 variables. No variables were dropped because of all missing values.

Show 20 entries Search:

urbrur	roof	walls	water	electcon	relat	sex	age	hhcivil	expend	income	savings	ori_hid	sampling_weight	household_
2	4	3	3	1	1	1	46	2	90929693	57800000	116258.5	1	100	
2	4	3	3	1	2	2	41	2	27338058	25300000	279345	1	100	
2	4	3	3	1	3	1	9	1	26524717	69200000	5495381	1	100	
2	4	3	3	1	3	1	6	1	18073948	79600000	8695862	1	100	
2	4	2	3	1	1	1	52	2	6713247	90300000	203620.2	2	100	16.66666
2	4	2	3	1	2	2	47	2	49057636	32900000	1021268	2	100	16.66666
2	4	2	3	1	3	2	13	1	63386309	22700000	8119166	2	100	16.66666
2	4	2	3	1	3	2	19	1	1106874	89100000	9881406	2	100	16.66666
2	4	2	3	1	3	1	9	1	32659507	2087324	7043642	2	100	16.66666
2	4	2	3	1	3	2	16	1	34347609	44100000	4783134	2	100	16.66666
2	4	3	3	1	1	1	65	2	71883547	55500000	7942221	3	100	33.33333
2	4	3	3	1	2	2	60	2	55174345	41200000	4318171	3	100	33.33333
2	4	3	3	1	5	2	6	1	46002021	99600000	2680967	3	100	33.33333
2	4	3	3	1	1	1	34	2	33042094	98400000	3662611	4	100	33.33333

Showing 1 to 20 of 4,580 entries

Previous 1 2 3 4 5 ... 229

**View/Analyze existing
sdcProblem**[Show summary](#)

Explore variables

Add linked variables

Create new IDs

**Anonymize categorical
variables**

Recoding

k-Anonymity

PRAM (simple)

PRAM (expert)

Supress values with high risks

**Anonymize numerical
variables**

Top/bottom coding

Microaggregation

Adding noise

Rank swapping

Summary of dataset and variable selection

The loaded dataset consists of **4580** records and **15** variables.

Categorical key variable(s): **urbrur** **roof** **walls**

Numerical key variable(s): **expend**

Sampling weight: **sampling_weight**

Hierarchical identifier: **ori_hid**

Computation time

The current computation time was ~ **0.15 seconds** .

Information on categorical key variables

Reported is the number of levels, average frequency of each level and frequency of the smallest level for categorical key variables. In parentheses, the same information is reported for the original data. Note that NA (missing) is counted as a separate category.

Variable name	Number of levels	Average frequency	Frequency of smallest level
urbrur	3 (3)	2285.000 (2285.000)	643 (643)
roof	5 (5)	916.000 (916.000)	16 (16)
walls	3 (3)	1526.667 (1526.667)	50 (50)

Risk measures for categorical key variables

We expect **0.16** (**0.00%**) re-identifications in the population, as compared to **0.16** (**0.00%**) re-identifications in the original data.

0 observations have a higher risk than the risk in the main part of the data, as compared to **0** observations in the original data. **i**

Risk measures

Information of risk

[Suda2 risk measure](#)

I-Diversity risk measure

Visualizations

Barplot/Mosaicplot

Tabulations

Information loss

Obs. violating k-anon

Numerical risk measures

Compare summary statistics

Disclosure risk

Information loss

SUDA2 risk measure

The SUDA algorithm is used to search for Minimum Sample Uniques (MSU) in the data among the sample uniques to determine which sample uniques are also special uniques i.e., have subsets that are also unique. See the help files for more information on SUDA scores.

Reset to choose a different sampling fraction parameter

Suda scores (sampling fraction is 0.1)

The table below shows the frequencies of the records with a suda score in the specified intervals.

Interval	Number of records
== 0	4580
(0.0, 0.1]	0
(0.1, 0.2]	0
(0.2, 0.3]	0
(0.3, 0.4]	0
(0.4, 0.5]	0
(0.5, 0.6]	0
(0.6, 0.7]	0
> 0.7	0

Attribute contributions

The table below shows the contribution of each categorical key variable to the SUDA scores. The contribution of a variable is

Variable selection

Variable name	Type	Additional suppressions by local suppression algorithm
urbrur	cat. key variable	0
roof	cat. key variable	0
walls	cat. key variable	0
expend	num. key variable	
sampling_weight	sampling weight	

Additional parameters

Parameter	Value
-----------	-------

What do you want to do?

[View the current script](#)

Import a previously saved problem

Export/Save the current
sdcProblem

View the current generated script

Browse and download the script used to generate your results. These can be used later as a reminder of what you did or entered into R from command-line to reproduce results.

Save Script to File

```
require(sdcMicro)
obj$inputdata <- readMicrodata(path="testdata", type="rdf", convertCharToFac=FALSE, drop_all_missings=FALSE)
inputdataB <- inputdata

## Convert a numeric variable to factor (each distinct value becomes a factor level)
inputdata <- varToFactor(obj=inputdata, var=c("urbrur","water","electcon","relat"))
## Set up sdcMicro object
sdcObj <- createSdcObj(dat=inputdata,
  keyVars=c("urbrur","roof","walls"),
  numVars=c("expend"),
  weightVar=c("sampling_weight"),
  hhId=c("ori_hid"),
  strataVar=NULL,
  pramVars=NULL,
  excludeVars=NULL,
  seed=0,
  randomizeRecords=FALSE,
  alpha=c(1))

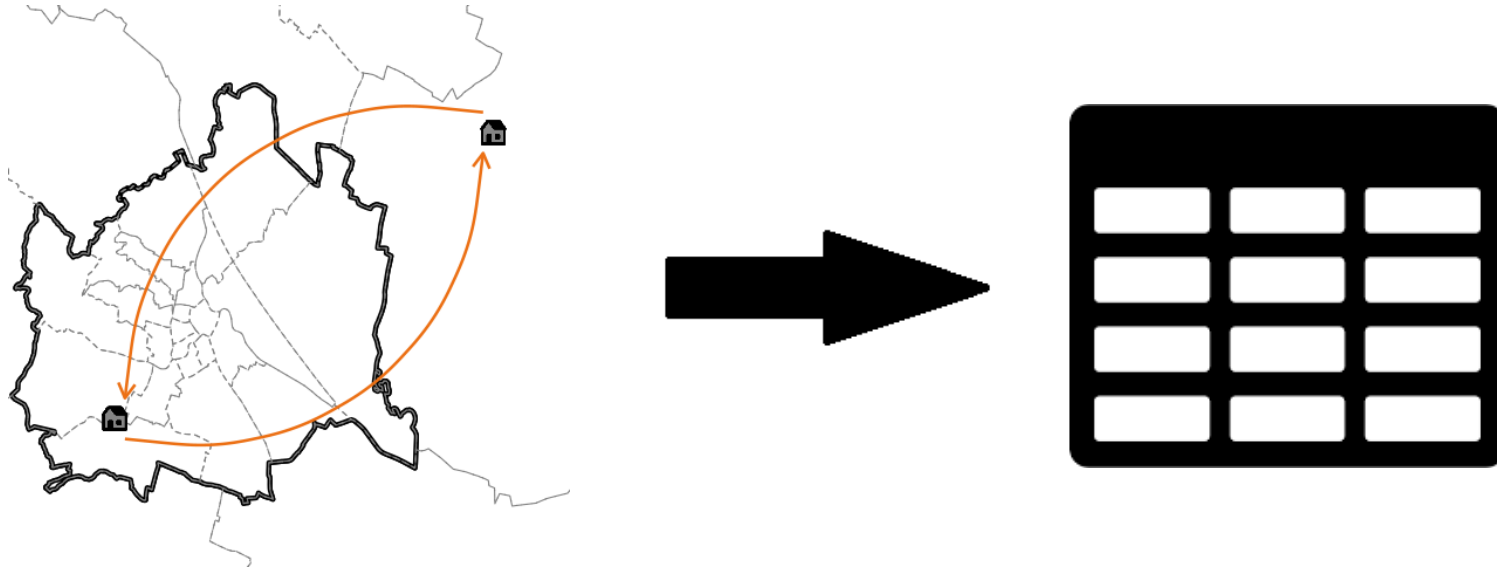
## Store name of uploaded file
opts <- get.sdcMicroObj(sdcObj, type="options")
opts$filename <- "testdata"
sdcObj <- set.sdcMicroObj(sdcObj, type="options", input=list(opts))

## calculating suda2 riskmeasure
sdcObj <- suda2(obj=sdcObj, DisEnaction=0.1, missing=NA)
```

Target record swapping

Recommended methodology for the census

- Data-swapping technique applied on micro data
- TRS: Swap households across administrative/geographic regions
- Swapped microdata used for all outputs



Target record swapping

- High risk → small frequency counts of individuals on a set of key variables (usually)
 - **Geographic hierarchy x risk variables**
 - k-anonymity in our algorithm
 - (own risk can be provided)
- Household at high risk ↔ individual at high risk
- “Similarity” variables
- Variables on which the swapped households must agree on (hh size, ...)
 - Preserves marginal distributions of those variables



Target record swapping with sdcMicro

- Data preparation

- Data needs to contain only integer/numeric columns
- no decimal places

- Convert column to integer using, e.g. `factor`

```
dat[, AGE.M := as.integer(factor(AGE.M)) ]
```

- Generate additional variables, e.g. for the similarity measure, e.g.
- A truncate household size

```
dat[, Size := pmin(5, Size) ]
```

```
dat[!duplicated(HID), .N, by=.(Size)] [order(Size) ]
```

Target record swapping parameters

- geographic hierarchy
 - *hierarchy* = `c("NUTS1", "NUTS2")`
- Column name of household id
 - *hid* = `"HID"`
- Variables for internal risk calculation
 - *risk_variables* = `c("COC.M", "POB.M")`
- Threshold for k-anonymity
 - *k_anonymity* = 3
- (Minimal) swap rate
 - *swaprate* = 0.05
- Similarity profile(s)
 - *similar* = `list(c("Size"), c("Size", "NationalityHead"))`

Target record swapping

Call and output

```
swapped <- recordSwap(data = dat, ... ,  
return_swapped_id = TRUE,  
seed = 123)
```

- `return_swapped_id = TRUE` get household ID of swapped household
- Seed to fix random seed
- Number of swapped households

```
dat_swapped[HID!=HID_swapped, uniqueN(HID) ]  
## [1] 520
```

Synthetic data generation

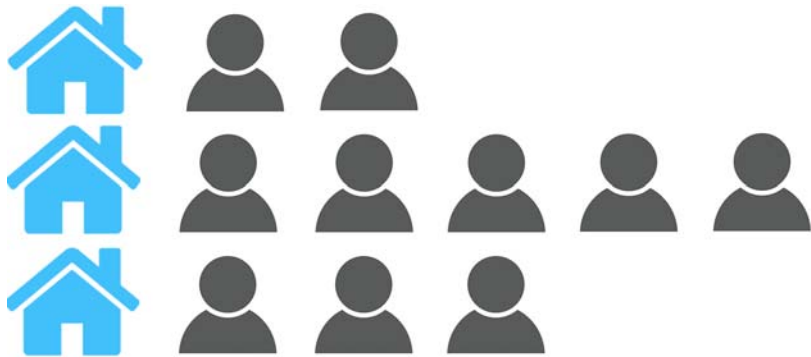
Johannes Gussenbauer

Overview

- Quick intro to R-package `simPop`
- Initialise and extend synthetic data
- Calibration of synth. population

simPop

- **simPop** R-Package to generate synthetic micro data ~ household data



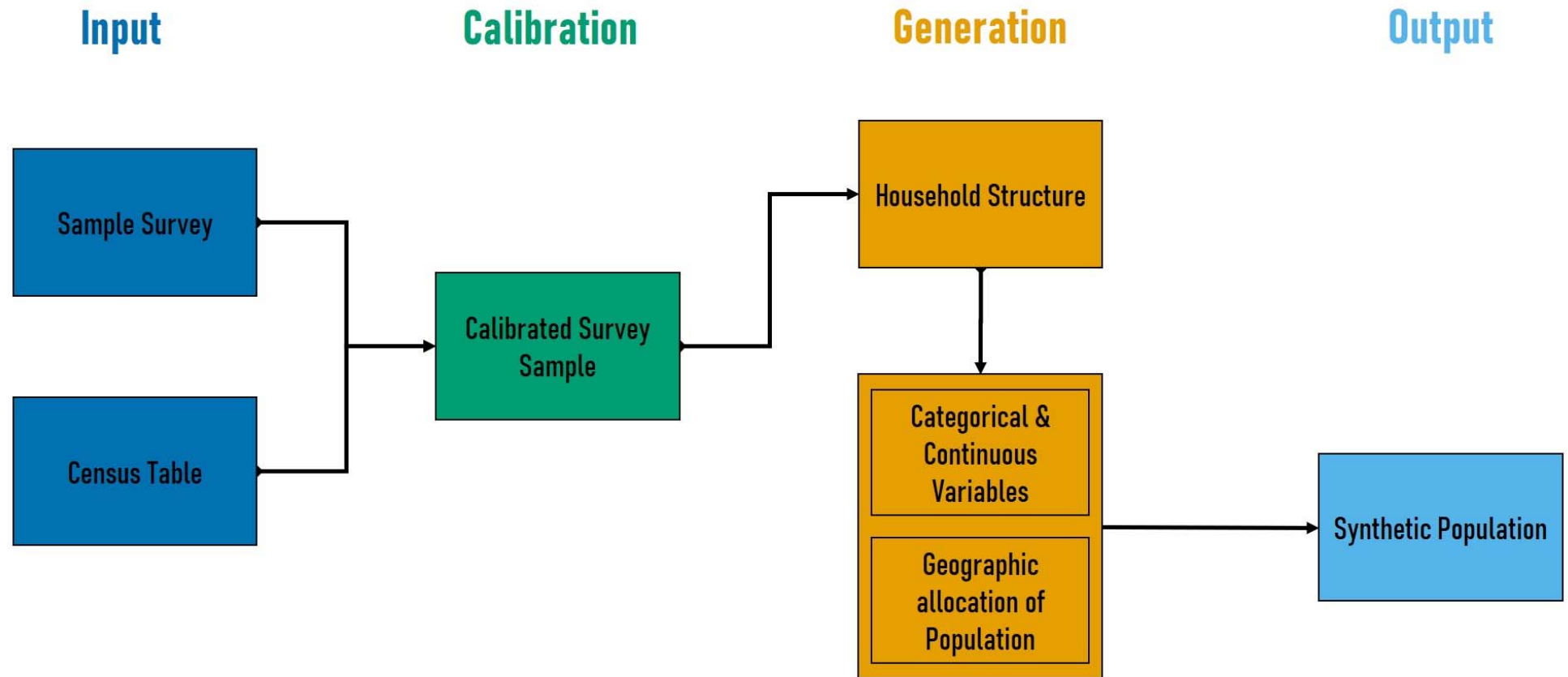
- Why generate synthetic micro data?
 - Microdata more and more needed by scientific community / the public
 - Often not possible to disseminate micro data (GDPR, national data protection laws)
 - Available micro data highly censored/takes long time to get access

simPop

- Synthetic data should
 - not reveal sensitive information
 - preserve correlation structure / Quasi-identical distribution

- Synthetic data can be used for
 - Prototype development
 - Micro-Simulation

simPop Workflow



Initialise population

- simPop → specifically designed for synthesising populations (persons living in households)
- Specify inputs

```
1 data(eusilcS)
2 # specify input
3 inp <- specifyInput(data=eusilcS, hhid="db030", hhsiz="hsize",
4                     weight="rb050", strata="db040")
```

- Initialise synthetic population by defining household structure

```
1 simPopObj <- simStructure(data=inp, method="direct",
2                           basicHHvars=c("age", "rb090"))
```

Simulate variables

- Sequentially add categorical or continuous variables

```
1 simPopObj <- simCategorical(simPopObj,  
2                             additional=c("p1030", "pb220a"),  
3                             method="multinom", nr_cpus=1)
```

- methods = c("multinom", "distribution", "ctree", "cforest", "ranger", "xgboost")

```
1 simPopObj <- simContinuous(simPopObj, additional="netIncome",  
2                             method = "lm",  
3                             regModel = ~rb090+hsize+p1030+pb220a,  
4                             nr_cpus=1)
```

- methods = c("multinom", "lm", "poisson", "xgboost")

Additional modelling functions

- Additional, more specific, modelling functions

1. Simulate categorical variables taking relationships between household members into account

```
1 ghanaP <- simRelation(simPopObj = ghanaP, relation = "relate",  
2                       head = "head",  
3                       additional = c("nation", "ethnic", "religion"),  
4                       nr_cpus = 1)
```

Additional modelling functions

2. Simulate components of continuous variables

```
1 simPopObj <- simComponents(simPopObj=simPopObj, total="netIncome",  
2   components=c("py010n", "py050n", "py090n", "py100n",  
3     "py110n", "py120n", "py130n", "py140n"),  
4   conditional = c("pl030"), replaceEmpty = "sequential", seed=1 )
```

3. Introduce smaller regions to already existing broader regions →
`simInitSpatial()`

4. Fix age heaping → `correctHeaps()`, `correctSingleHeap()`

simPop Workflow

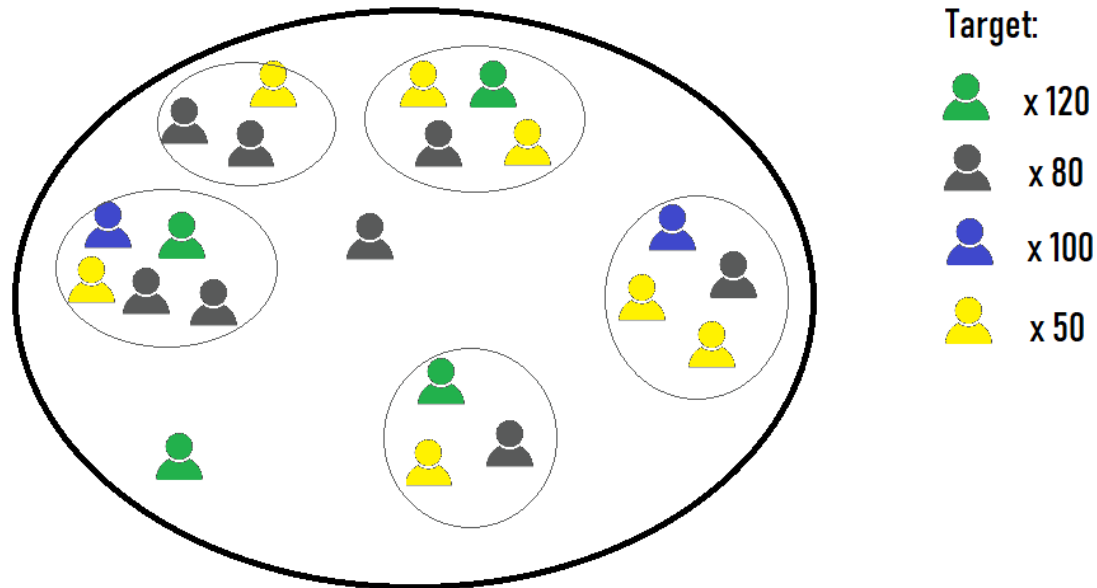
- After variables have been synthesised *calibrate* synthetic population to *fit* selected distributions
- Simulated annealing algorithm implemented in simPop

```
1 data("eusilcP")
2 # add margins
3 margins <- as.data.frame(
4   xtabs(rep(1, nrow(eusilcP)) ~ eusilcP$region +
5     eusilcP$gender + eusilcP$citizenship))
6 colnames(margins) <- c("db040", "rb090", "pb220a", "freq")
7 simPopObj <- addKnownMargins(simPopObj, margins)
```

```
1 # run calibration
2 simPop_adj2 <- calibPop(simPopObj, split="db040",
3   temp=1, epsP.factor=0.1,
4   epsMinN=10, nr_cpus = 1)
```

Calibrate synth. Population

- Calibrate synthetic population to “fit” selected distributions - how?



- Apply simulated annealing to try to find a local optimum

Simulated annealing in **simPop**

1. Multiply synthetic population
2. Make initial selection of households
3. Compare target distribution against distribution of selection using an objective function
4. Discard and redraw households with certain probabilities
5. Check objective function again and accept/reject new solution
6. Repeat 4. and 5. until distributions “differ” by ϵ

Recent improvements

- Multiple distributions allowed on household and personal level

```
1  calibPop(  
2    inp, # <- simPopObj  
3    ...  
4    hhTables = NULL,  
5    persTables = NULL  
6  )
```

- Objective / sampling probabilities / termination condition adjusted accordingly

Objective and sampling probability

- Objective, p target margins with k_1, \dots, k_p number of cells

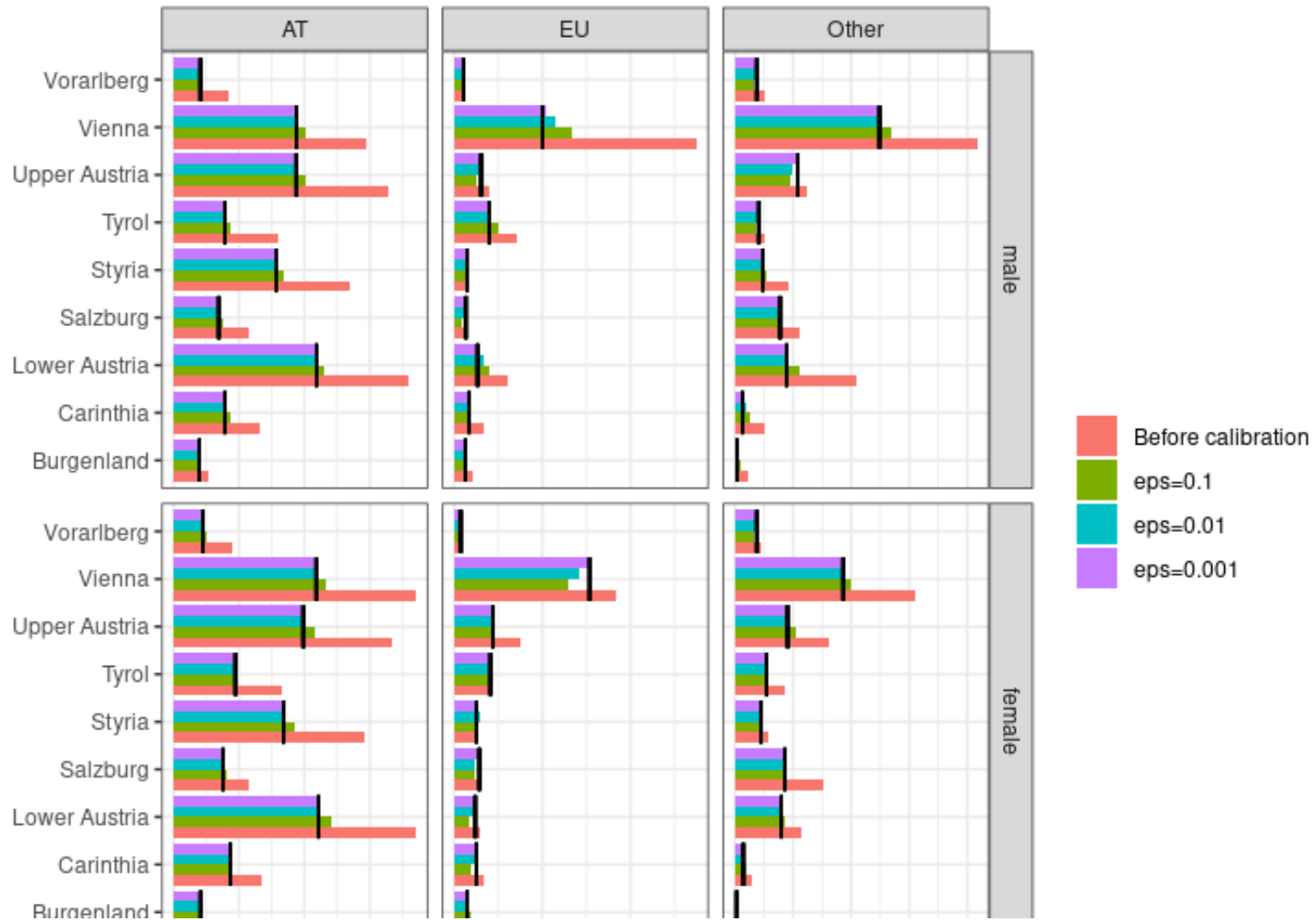
$$\sqrt{\frac{1}{p} \sum_{i=1}^p \left(\sum_{j=1}^{k_i} |n_{i,j} - \tilde{n}_{i,j}| \right)^2}$$

- Sampling prob. p_l for individual l

$$p_l = \begin{cases} \frac{f_l}{N_l} & \text{if } f_l > 0 \\ \exp\left(-\sum_{f_l \leq 0} \frac{f_l}{N_l}\right) & \end{cases}; \quad f_l = \frac{1}{p} \sum_i |e_{i,j(i,l)}| \cdot \text{sign}\left(\sum_i e_{i,j(i,l)}\right)$$

$$e_{i,j} = n_{i,j} - \tilde{n}_{i,j}$$

Improvements of calibration



References

- Münnich, Ralf and Josef Schürle (2003). “On the the simulation of complex universes in the case of applying the German Microcensus”
- Templ, Matthias, Alexander Kowarik, and Peter Filzmoser (2011). “Iterative stepwise regression imputation using standard & robust methods”. In: Computational Statistics & Data Analysis 55.10. DOI:10.1016/j.csda.2011.04.012, ISSN: 0167-9473, pp. 2793–2806. issn:0167-9473.
- Templ, Matthias, Bernhard Meindl, et al. (Aug. 2017). “Simulation of Synthetic Complex Data: The R Package simPop”. In: Journal of Statistical Software 79. doi:10.18637/jss.v079.i10.

Anonymization of tabular data

Prepared by Bernhard Meindl



Why?

Reasons for tabular data control

- **Disclosure risk** exists in aggregated (tabular) data
- **Goals:**
 - Protect statistical units that contribute to table
 - Take into account **trade-off** between risk and data-utility
- **Why?**
 - Publication requirements
 - Legal reasons
- **Statistical tables:**
 - frequency vs. magnitude tables
 - cells relate to each other
 - attackers can make use of those relations

What to protect against?

Group Attribute Disclosure

	female	male	Total
a	4	12	16
b	3	0	3
c	3	0	3
Total	10	12	22

- “**Definition**”: Attacker can learn an attribute about an individual
- **Example**:
 - only one **female** that is single → if additional tables exist (e.g **gender** x **marital_status** x **income**), attackers can derive the income-group for the this unit
- *Further Issue*: → **Linked tables** (identical cells can appear in multiple tables)

Packages

The image features a blue-tinted background of a modern office building's interior. On the left, there are multiple levels of balconies with glass railings, some containing potted plants. On the right, a window with white frames and horizontal blinds is partially open, showing a view of another building. The word "Packages" is written in white, bold, sans-serif font on the left side of the image.

Some relevant R-packages

We are going to introduce the following R packages:

- `sdcHierarchies`
- `sdcTable`
- `cellKey` / `pTable`
- `sdcSpatial`

sdcHierarchies

General Information

- **Goal / Idea:** Define (nested) hierarchies that are used to define statistical tables
- **Features:**
 - allows to programmatically construct (nested) hierarchies
 - Information about (individual) nodes can be extracted ' Hierarchies can be imported/exported from/to various formats
 - Shiny-App is provided to interactively create / modify / export hierarchies
- **Web:** github.com/bernhard-da/sdcHierarchies

sdcHierarchies

Create hierarchy

```
1 library(sdcHierarchies)
2 h <- hier_create(root = "Tot", nodes = LETTERS[1:3])
3 h <- hier_add(h, root = "A", nodes = paste0("a", 1:2))
4 h <- hier_add(h, root = "a1", nodes = "a11"); hier_display(h)
```

```
Tot
├── A
│   ├── a1
│   │   └── a11
│   └── a2
├── B
└── C
```

sdcHierarchies

Information about a node

```
1 hier_info(h, "A")
```

```
$name
```

```
[1] "A"
```

```
$is_rootnode
```

```
[1] FALSE
```

```
$level
```

```
[1] 2
```

```
$is_leaf
```

```
[1] FALSE
```

```
$siblings
```

```
[1] "B" "C"
```

sdCTable

General Information

- **Goal / Idea:** Protect statistical tables
- **Features:**
 - Setup of complex statistical tables
 - Identification primary sensitive table cells
 - Protect sensitive cells using different algorithms
 - Allow export of problem instances to τ -Argus
- **Web:** github.com/sdcTools/sdcTable

sdcTable

Setup Problem

```
1 library(sdcTable)
2 p <- sdc_testproblem(); print(str(p))
```

```
Formal class 'sdcProblem' [package "sdcTable"] with 8 slots
 ..@ dataObj          :Formal class 'dataObj' [package "sdcTable"] with 7
slots
 .. .. ..@ rawData   :Classes 'data.table' and 'data.frame':   100 obs.
of 5 variables:
 .. .. .. ..$ region : chr [1:100] "A" "A" "A" "A" ...
 .. .. .. ..$ gender : chr [1:100] "female" "female" "male"
"male" ...
 .. .. .. ..$ freq   : num [1:100] 1 1 1 1 1 1 1 1 1 1 ...
 .. .. .. ..$ tmpsamplingweights: num [1:100] 1 1 1 1 1 1 1 1 1 1 ...
 .. .. .. ..$ val    : num [1:100] 9 11 10 11 5 7 13 15 13 6 ...
 .. .. .. ..- attr(*, ".internal.selfref")=<externalptr>
 .. .. .. ..- attr(*, "sorted")= chr [1:2] "region" "gender"
 .. .. ..@ dimVarInd  : int [1:2] 1 2
 .. .. ..@ freqVarInd : int 3
```

sdcTable

Extract Table

```
1 sdcProb2df(p, addDups = FALSE, dimCodes = "original", addNumVars = TRUE)
```

```
# A data frame: 15 × 6
  strID   freq sdcStatus   val region gender
  <chr> <dbl> <chr>      <dbl> <chr> <chr>
1 0000     100 s          1284 total total
2 0001      55 s           802 total male
3 0002      45 s           482 total female
4 0100      20 s           198 A      total
5 0101      18 s           178 A      male
6 0102       2 s            20 A      female
7 0200      33 s           344 B      total
8 0201      14 s           140 B      male
9 0202      19 s           204 B      female
10 0300      22 s           224 C      total
11 0301      12 s           118 C      male
12 0302      10 s           106 C      female
```

sdcTable

Identify primary sensitive cells

```
1 p <- primarySuppression(p, type = "freq", maxN = 3)
```

- Table of different "cell-stati"
 - "s": "safe" for publication
 - "u": primary unsafe
 - "x": secondary suppression

```
1 table(getInfo(p, "sdcStatus"))
```

```
s  u
14 1
```

sdcTable

Protect table

```
1 p <- protectTable(p, method = "GAUSS")
2 getInfo(p, "finalData")
```

```
# A data frame: 15 × 5
  region gender  Freq  val sdcStatus
  <chr>  <chr> <dbl> <dbl> <chr>
1 total  total   100  1284 s
2 total  male    55   802 s
3 total  female  45   482 s
4 A      total   20   198 s
5 A      male    18   178 x
6 A      female   2    20 u
7 B      total   33   344 s
8 B      male    14   140 s
9 B      female  19   204 s
10 C     total   22   224 s
11 C     male    12   118 s
12 C     female  10   106 s
```

cellKey - ptable

General Information

- **Goal / Idea:** Persistent perturbation of statistical tables
- **Features:**
 - Allows generation of tables (similar to `sdcTable`)
 - Implements a perturbation algorithm that depends on record- and cell keys
 - Makes use of look-up tables (using `ptable` Package)
 - Different methods for frequency- and magnitude tables
 - Very useful for perturbation of large, linked tables
 - Drawback: Impact of perturbation visible (non-additive results)
- **Web:**
 - github.com/sdcTools/ptable | github.com/sdcTools/cellKey

cellKey + ptable

Setup

```
1 library(cellKey)
2 x <- ck_create_testdata()
3 tab <- ck_setup(
4   x = x,
5   rkey = 6, # digits
6   dims = list(
7     sex = hier_create(root = "Total", nodes = c("male", "female")),
8     age = hier_create(root = "Total", nodes = paste0("age_group", 1:6)),
9     w = "sampling_weight"
10  )
11 tab$print()
```

— Table Information

✓ 21 cells in 2 dimensions ('sex', 'age')

✓ weights: yes

— Tabulated / Perturbed countvars

'total'

cellKey + ptable

Define perturbation parameters

```
1 p_cnts <- ck_params_cnts(  
2   ptab = ptable::pt_ex_cnts()  
3 )  
4 print(head(p_cnts$params$ptable))
```

	i	j	p	v	lb	ub	type
1:	0	0	1.000000000	0	0.00000000	1.00000000	all
2:	1	0	0.508333333	-1	0.00000000	0.50833333	all
3:	1	2	0.475000000	1	0.50833333	0.98333333	all
4:	1	3	0.016666667	2	0.98333333	1.00000000	all
5:	2	0	0.16155827	-2	0.00000000	0.1615583	all
6:	2	2	0.55565037	0	0.1615583	0.7172086	all

- Assign parameters to specific variable

```
1 tab$params_cnts_set(v = "total", val = p_cnts)
```

cellKey + ptable

Perturb and evaluate

```
1 tab$perturb(v = "total")
2 tab$print()
```

— Table Information

✓ 21 cells in 2 dimensions ('sex', 'age')
✓ weights: yes
— Tabulated / Perturbed countvars

'total' (perturbed)

- Evaluate / extract results

```
1 tab$freqtab(v = "total")
```

A data frame: 21 × 7

	sex	age	vname	uwc	wc	puwc	pwc
	<chr>	<chr>	<chr>	<dbl>	<dbl>	<dbl>	<dbl>
1	Total	Total	total	<u>4580</u>	<u>275677</u>	<u>4581</u>	<u>275737.</u>
2	Total	age_group1	total	<u>1969</u>	<u>118905</u>	<u>1969</u>	<u>118905</u>
3	Total	age_group2	total	<u>1143</u>	<u>68788</u>	<u>1144</u>	<u>68848.</u>
4	Total	age_group3	total	864	<u>52136</u>	866	<u>52257.</u>
5	Total	age_group4	total	423	<u>25028</u>	425	<u>25146.</u>

sdcsSpatial

General Information

- **Goal / Idea:** Protect spatial data
- **Features:**
 - Based on functionality from `raster` package
 - Allows to identify sensitive cells
 - Multiple algorithms (removal, smoothing, aggregation) to protect raster-cells are implemented
- **Web:**
 - github.com/edwindj/sdcSpatial
 - Helpful vignette: [Intro to sdcSpatial](#)

Please address queries to
alexander.kowarik@statistik.gv.at
johannes.gussenbauer@statistik.gv.at
bernhard.meindl@statistik.gv.at

STATISTIK AUSTRIA
Guglgasse 13, 1110 Wien

Independent statistics for evidence-based decision making



What is the awesome list?

Curated list of software for
official statistics



awesome

www.awesomeofficialstatistics.org

Criteria

An item on this list is awesome because

1. it is free, open source, and available for download and
2. it is confirmed to be used in the production of official statistics by at least one institute or it provides access to official statistics publications.

We prefer packages that are easy to install and use, have at least one stable version, and are actively maintained. [Contributions](#) are welcome.

License



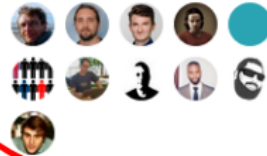
This work is licensed under a [Creative Commons Attribution 4.0 International License](#).

Open license

Social interactions

An awesome list of statistical
AI software for creating and
accessing official statistics

Contributors 15



+ 4 contributors

Working together

Contributions

Awesome contributions are welcome, here are ways to do it:

- The GitHub way: send us a [pull request](#) to add directly to this list.
- Add an item to the [issue tracker](#) issue tracker. (you need a GH account)
- Send an e-mail to [mark dot vanderloo at gmail dot com](#) or [olav dot tenberch at gmail dot com](#) or tweet [@markvdloo](#)

