

Chapter 4

Methods for Data Perturbation

Abstract Methods for perturbation of data differ for categorical and continuous variables. The risk for categorical key variables is dependent on the frequency counts of keys, whereas keys with only few observations are problematic. Categories of categorical key variables with low frequency counts are therefore often recoded and combined with other categories. However, a still too high disclosure risk may be present for some individuals. Local suppression is one method to further reduce the disclosure risk. In order to find a well-balanced, suitable solution, global recoding is usually applied in an explorative manner to observe with which (reasonable) recodings one achieves the best effect in terms of reducing the disclosure risk and providing high data utility. Especially with a large amount of key variables, swapping methods, such as PRAM, are good alternatives. Methods for continuous scaled variables are combining values (microaggregation) or adding noise to the values. Advanced methods such as shuffling allow to preserve certain statistics.

4.1 Kind of Methods

The SDC techniques may be categorized into three kind of methods:

- non-perturbative techniques, such as recoding and local suppression, which suppress or reduce the detail without altering the original data;
- perturbative techniques, such as adding noise, Post-Randomization Method (PRAM), micro-aggregation and shuffling, which distort the original micro-data set before release;
- techniques that generate a synthetic microdata file that preserves certain statistics or relationships of the original files.

This chapter focuses on the non-perturbative and perturbative techniques. As with disclosure risk measures, different SDC methods are applicable to categorical variables and continuous variables.

Note that generating synthetic data is a more complicated approach and the methods differ completely. This is the reason why these methods are not discussed in this chapter, and we refer to Chap. 6 that is dedicated to synthetic data simulation.

4.2 Methods for Categorical Key Variables

4.2.1 Recoding

Global recoding is a non-perturbative method that can be applied to both categorical and continuous key variables. For a categorical variable, the idea of recoding is to combine several categories into fewer categories with higher frequency counts and less detailed information. In other words, a global recoding achieves anonymity by mapping the values of the categorical key variables to generalized or altered categories. For example, one could combine multiple levels of schooling (e.g., secondary, tertiary, postgraduate) into one (e.g., secondary and above). For a continuous variable, recoding means to discretize the variable; for example, recoding a continuous income variable into a categorical variable of income levels. In both cases, the goal is to reduce the total number of possible values of a variable. Typically, recoding is applied to categorical variables to collapse categories with few observations into a single category with larger frequency counts. For example, if there are only two respondents with tertiary level of education, the tertiary can be combined with the secondary level into a single category of “secondary and above”.

Remember, we have the following variables in our `sdcMicroObj`

```
require("sdcMicro")
data(testdata, package="sdcMicro")
sdc <- createSdcObj(testdata,
  keyVars=c('urbrur', 'water', 'sex', 'age', 'relat'),
  numVars=c('expend', 'income', 'savings'),
  pramVars=c("walls"),
  w='sampling_weight',
  hhId='ori_hid')
```

The categorical key-variable `age` selected before has a lot of categories

```
table(testdata$age)
```

```
##
##      0      1      2      3      4      5      6      7      8      9     10     11     12     13     14
##  98    90   134   112   128   133   136   125   144   126   151   127   108   143   103
##  15    16    17    18    19    20    21    22    23    24    25    26    27    28    29
## 111   106   101    96    64    88    61    64    55    55    69    56    68    69    50
##   30   31   32   33   34   35   36   37   38   39   40   41   42   43   44
##   90   51   72   49   59   86   61   68   51   42   67   43   44   49   40
##   45   46   47   48   49   50   51   52   53   54   55   56   57   58   59
##   65   37   31   43   28   44   28   28    8   19   31   28   17   27   20
##   60   61   62   63   64   65   66   67   68   69   70   71   72   73   74
##   36   24   28   14   12   40    8   16   14    8   20    6   14    4    6
##   75   76   77   78   79   80   82   83   84   85   88   90   95
##    6    4    5    3    2    5    1    1    1    1    1    2    1
```

and some categories are sparse. Recoding `age` into `age` classes will cause that the number of observations in the keys increase and less observations are violating the 2- and 3-anonymity assumption.

```
labs <- c("1-9", "10-19", "20-29", "30-39",
          "40-49", "50-59", "60-69", "70-79", "80-130")
sdc <- globalRecode(sdc, column="age",
                   breaks=c(0,9,19,29,39,49,59,69,79,130),
                   labels=labs)
print(sdc)
```

```
## Infos on 2/3-Anonymity:
##
## Number of observations violating
## - 2-anonymity: 111 (2.424%) | in original data: 653 (14.258%)
## - 3-anonymity: 184 (4.017%) | in original data: 1087 (23.734%)
## - 5-anonymity: 345 (7.533%) | in original data: 1781 (38.886%)
##
##
-----
```

Let's have a look at the frequency counts of the categorized `age` variable. Note that we use function `extractManipData()` to extract the current data from the *sdcMicroObj* object.

```
table(extractManipData(sdc)$age)
```

```
##
##      1-9    10-19    20-29    30-39    40-49    50-59    60-69    70-79
##   1128    1110     635     629     447     250     200     70
##  80-130
##      13
```

Still some categories have only few cases, so we further modify the age variable by joining the last two groups. Here a general function could be applied, called `groupAndRename()`.

```
sdc <- groupAndRename(sdc, var = "age",
                      before = c("60-69", "70-130"),
                      after = "60-130")

print(sdc)

## Infos on 2/3-Anonymity:
##
## Number of observations violating
## - 2-anonymity: 85 (1.856%) | in original data: 653 (14.258%)
## - 3-anonymity: 162 (3.537%) | in original data: 1087 (23.734%)
## - 5-anonymity: 294 (6.419%) | in original data: 1781 (38.886%)
##
## -----

table(extractManipData(sdc)$age)

##
##      1-9    10-19    20-29    30-39    40-49    50-59    60-130
##      1128     1110      635      629      447      250      283
```

Alternatively we could also re-run the code using function `undolast()` and performing `globalRecode()` again with new breaks.

A special case of global recoding is top and bottom coding. This method can be applied to ordinal or continuous variables. The idea for this approach is that all values above (i.e., top coding) and/or below (i.e., bottom coding) a pre-specified threshold value are combined into a new category. Top coding sets an upper limit on all values of a variable and replaces any value greater than this limit by the upper limit; for example, top coding would replace the age value for any individual aged above 80 with 80. Similarly, bottom coding replaces any value below a pre-specified lower limit by the lower limit; for example, bottom coding would replace the age value for any individual aged under 5 with 5. Instead of 80 or 5, respectively, one can also choose any number, e.g. the arithmetic mean of the high (or low) scores so that the arithmetic mean is not changed after top and bottom coding. As already seen in previous code blocks, in **sdcMicro** function `globalRecode()` can be used to perform both global recoding and the top/bottom coding, but there is also a specialized function called `topBotCoding` available. A help file with more examples as seen above is accessible using `?globalRecode` and `?topBotCoding`.

To replace high incomes in our test data set with the mean of the high incomes we can use the following code. Again the *sdcMicroObj* `sdc` updates its slots (e.g. risk and utility) after applying an anonymization method.

```
highIncomes <- testdata$income[testdata$income > 100000000]
sdc <- topBotCoding(sdc,
  value = 100000000,
  replacement = mean(highIncomes),
  column = "income")
```

We note that **sdcMicroGUI** and the forthcoming browser-based version of the GUI offers a more user-friendly way of applying global recoding in general.

Exercises:

Question 4.1 Global recoding:

Take the `eusilc` data set from package **laeken** again. Assume the following disclosure scenario that defines `age`, `pb220a` (citizenship), `p1030` (education level), `rb090` (gender) and `hsize` (household size) as categorical key variables. Use the package **sdcMicro** to create an object of class `sdcMicroObj` considering the sampling weights (function argument `weightVar` in `createSdcObj`) and the household ID (function argument `hhId`). Reduce the disclosure risk through some reasonable recodings. Apply the recoding carefully, i.e. the analytical quality of the data should mostly remain the same while the disclosure risk should lower considerably.

4.2.2 Local Suppression

If unique combinations of categorical key variables remain after recoding, local suppression could be applied to the data with the aim to achieve k -anonymity. For the detailed explanation of k -anonymity, have a look in Sect. 3.3. Local suppression is a non-perturbative method typically applied to categorical variables. In this approach, missing values are created to replace certain values of individual variables to increase the number of key variables sharing the same pattern, thus reducing the record-level disclosure risks. There are two approaches to implementing local suppression. One approach sets the parameter k and tries to achieve k -anonymity (typically 3-anonymity) with minimum suppression of values. For example, in **sdcMicroGUI** (Kowarik et al. 2013), the user sets the value for k and orders key variables by the likelihood they will be suppressed. Then the application calls a heuristic algorithm to suppress a minimum number of values in the key variables to achieve k -anonymity. The second approach sets a record-level risk threshold. This method involves examining unsafe records with individual disclosure risks higher than the threshold and suppressing values of the selected key variable(s) for all the unsafe records.

Local suppression can be done univariate such as in Sect. 4.2.2.2, but in general it is a multivariate problem that is NP-hard, not solvable optimally in a reasonable time.

Say we have a problem P . It consists of achieving k -anonymity with the constraint to suppress as few values. In addition, a weighting determining the importance of key variables may just involve another additional parameter that can be incorporated in the constraints.

Some algorithms have been written which provide heuristic solutions.

Mondrian one solution is the algorithm Mondrian (LeFevre et al. 2006). This algorithm tries to combine categories to achieve k -anonymity, i.e. it is a sort of recoding based on the median counts of categories. However, this is a too oversimplistic approach of sorting and splitting, we do not obtained very promising results. More theoretical descriptions on the Mondrian algorithm can be found in Sect. 4.2.2.6.

all- M approach in μ -Argus and **sdcmicro** an algorithm is implemented that is often referred as the *all- M* approach. Using this algorithm, k -anonymity cannot be provided, but k -anonymity in all subsets of size M of the key variables only (in μ -Argus with the maximum of $M = 4$ variables). More precisely, the algorithm will provide k -anonymity for each combination of M key variables.

k -anonymity approach the default approach of **sdcmicro** in function `kAnon` ensures k -anonymity for the combination of all selected key variables. Naturally, this lead to k -anonymity but also to oversuppressions whenever the amount of key variables is too high, say larger than 7–10 key variables. Then often it is suitable to provide k -anonymity on a subset of variables (e.g. 7 variables) and apply a swapping method such as PRAM to the other key variables. Practical applications are shown in Sect. 4.2.2.3.

The real problem is the computation time and to find a good heuristic algorithm that solves our problem P to suppress as few values. Computation time increases a lot as soon as missing values are present in the data, since one has to treat them adequately. To decrease computation time, a lot of tricks must be used, based on (1) filtering the data in advance (2) ordering the data and (3) using C++ code and (4) thinking of good heuristics to find a good local optimal solution for this multivariate problem.

4.2.2.1 Illustrative Examples for Frequency Calculation with Missing Values

Here we show the whole problem of sample frequency calculation in case of missing values in the key variables, achieving k -anonymity and local suppression.

Remember, in Sect. 3.2.2 five possible methods how to calculate frequencies have been noted. For readability these five methods are defined once more:

1. (default method) Missing values increase frequencies in other categories.
2. (conservative method) Missing values do not increase frequencies in other categories but in those observations where a missing occurs.
3. (category size) Missing values do increase frequencies in other categories by a factor c .
4. (conservative method 2) Missing values do not increase frequencies in other categories.
5. (own category) Same as method 4, but missings are treated like an own category.

Table 4.1 Simple toy data set to explain different methods for sample frequency counts including missings/local suppressions in the following tables. Different solutions are displayed that show how to achieve 2-anonymity with local suppression based on different methods to calculate sample frequency counts

Original data					Method 1 (default)				
ID	Region	Status	Age group	f_k	ID	Region	Status	Age group	f_k
1	A	Single	30–49	2	1	A	Single	30–49	3
2	A	Married	30–49	2	2	A	Married	30–49	3
3	A	Married	30–49	2	3	A	Married	30–49	3
4	A	Single	30–49	2	4	A	Single	30–49	3
5	A	Widow	30–49	1	5	A	*	30–49	5
Method 2 (conservative)					Method 3 (category size)				
ID	Region	Status	Age group	f_k	ID	Region	Status	Age group	f_k
1	A	Single	30–49	2	1	A	Single	30–49	2.4
2	A	Married	30–49	2	2	A	Married	30–49	2.4
3	A	Married	30–49	2	3	A	Married	30–49	2.4
4	A	Single	30–49	2	4	A	Single	30–49	2.4
5	A	*	30–49	5	5	A	*	30–49	5 (or 3.2)
					2.4 = ratio single · 1 missing 2.4 = ratio married · 1 missing				
Method 4 (conservative 2)					Method 5 (own category)				
ID	Region	Status	Age group	f_k	ID	Region	Status	Age group	f_k
1	A	*	30–49	5	1	A	*	30–49	3
2	A	Married	30–49	2	2	A	Married	30–49	2
3	A	Married	30–49	2	3	A	Married	30–49	2
4	A	*	30–49	5	4	A	*	30–49	3
5	A	*	30–49	5	5	A	*	30–49	3

To motivate and explain these different concepts, a toy data example is used in the following, see Table 4.1 (upper left), with different values only in the second variable. It can be seen that 2-anonymity is not reached in the original data. For our scenarios, we assume that 2-anonymity should be achieved by local suppressions, i.e. the necessary suppressions under each frequency count method are made and sample frequencies are calculated. The results are displayed in Table 4.1. The default method, the conservative method and the category size method lead to the lowest numbers of necessary suppression to achieve 2-anonymity, i.e. with one suppression

Table 4.2 Simple toy data set to explain different methods for sample frequency counts including missings/local suppressions in the following tables. Different solutions are displayed to achieve 3-anonymity with local suppression based on different methods to calculate sample frequency counts

Original data					Method 1 (default)				
ID	Region	Status	Age group	f_k	ID	Region	Status	Age group	f_k
1	A	Single	30–49	2	1	A	Single	30–49	3
2	A	Married	30–49	2	2	A	Married	30–49	3
3	A	Married	30–49	2	3	A	Married	30–49	3
4	A	Single	30–49	2	4	A	Single	30–49	3
5	A	Widow	30–49	1	5	A	*	30–49	5
Method 2 (conservative)					Method 3 (category size)				
ID	Region	Status	Age group	f_k	ID	Region	Status	Age group	f_k
1	A	*	30–49	5	1	A	*	30–49	5
2	A	*	30–49	5	2	A	Married	30–49	3.2
3	A	*	30–49	5	3	A	Married	30–49	3.2
4	A	*	30–49	5	4	A	*	30–49	5
5	A	*	30–49	5	5	A	*	30–49	5
					3.2 = 2 + ratio of married · 3 missings				
Method 4 (conservative 2)					Method 5 (own category)				
ID	Region	Status	Age group	f_k	ID	Region	Status	Age group	f_k
1	A	*	30–49	5	1	A	*	30–49	5
2	A	*	30–49	5	2	A	*	30–49	5
3	A	*	30–49	5	3	A	*	30–49	5
4	A	*	30–49	5	4	A	*	30–49	5
5	A	*	30–49	5	5	A	*	30–49	5

2-anonymity is reached. However, the sample frequencies are different. The method using an own category for the missing values is most strict. It leads to the lowest frequency counts with three suppressions.

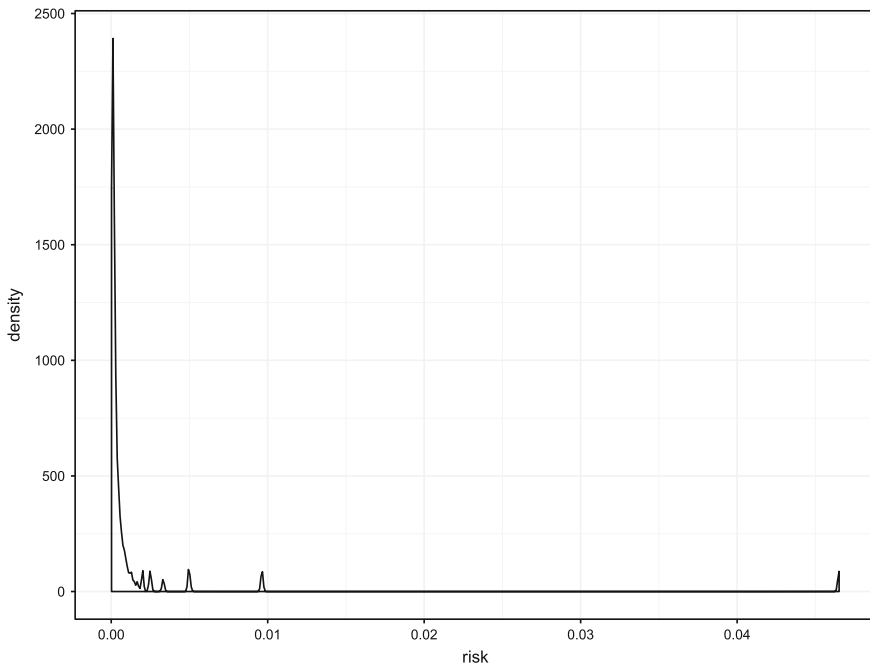
As from the previous table, from the upper left table in Table 4.2 it can be seen that 3-anonymity is not reached. For our scenarios, we assume that this time 3-anonymity should be achieved by local suppressions. The results are displayed in Table 4.2. The default method count a missing as a any possible category. Thus one suppression is already enough to achieve 3-anonymity. The category size methods counts the amount of values of category *married* plus the ratio of category *married* to the amount of categories times the number of missing values. All other approaches are more conservative and all values of *status* have to be suppressed to achieve 3-anonymity.

The choice of method for sample frequency calculation is best done with discussion with the law department or the management of the data holder. In the following, the so called default method is used. This method is used for years at Statistics Austria and every organisation who uses the package **sdcMicro**. The default method is amongst possible choices the less conservative one. It assumes that an intruder cannot be sure about the correct category of the suppressed value.

4.2.2.2 Univariate Local Suppression for Observations with High Risk

Using function `localSupp()` of **sdcMicro**, it is possible to suppress values of a key variable for all units with individual risks above a pre-defined threshold, given a disclosure scenario. This procedure requires user intervention by setting a threshold. This is illustrated in the following code listing. First, a density plot of the individual risks is plotted. This plot helps to determine the threshold. The majority of the data have very low risk, below 0.005, which already may serve as the value of the threshold. For those observations with higher individual risk than 0.005, the values of variable “*urbur*” are suppressed.

```
## extract vector of individual risks
risk <- get.sdcMicroObj(sdc, "risk")$individual[, "risk"]
## have a look on the risk
ggplot(data.frame(risk), aes(x=risk)) + geom_density() + theme_bw()
```



```
## suppress values in urbrur for risk higher than 0.05
sdc <- localSupp(sdc, keyVar = "urbrur", threshold = 0.05)
```

4.2.2.3 Ensuring k -Anonymity—The Optimal Approach

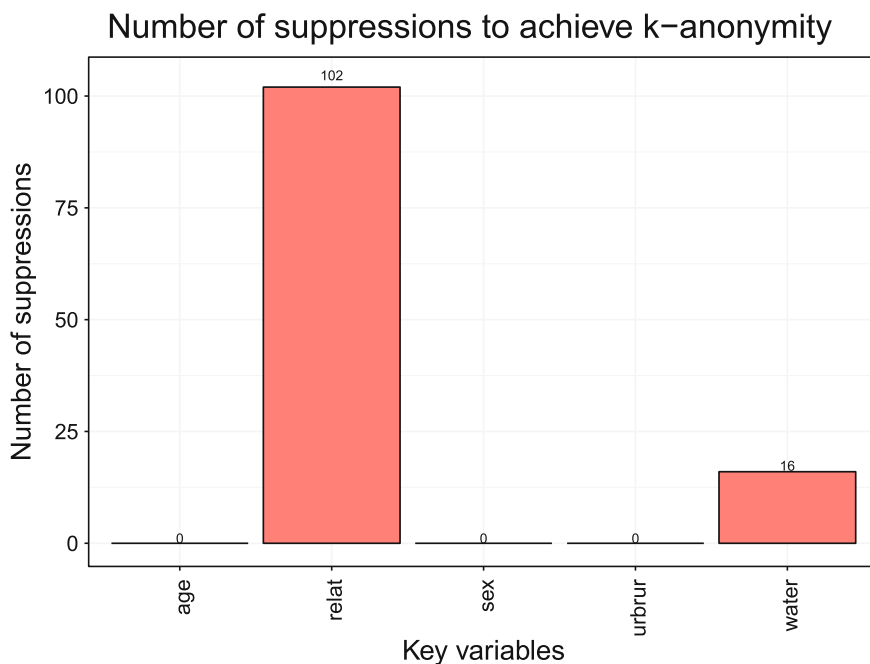
To automatically suppress a minimum amount of values in the key variables to achieve k -anonymity, one can use function `kAnon()`. In the code listing below, first the previous local suppression is reverted. Local suppression to achieve 3-anonymity is then done by setting parameter `k = 3`.

```
## undo last step (localSupp(...))
sdc <- undolast(sdc)

## local suppression to ensure 3-anonymity
## for given key variables
sdc <- kAnon(sdc, k=3)
```

The amount of suppression can be reported via a plot of the object `sdc`.

```
plot(sdc, "ls")
```



Also the print of this object gives information about the number and percentages of suppressions for each variable.

```
print(sdc, "ls")
```

```
## Local suppression:
## KeyVar | Suppressions (#) | Suppressions (%)
## urbrur | 0 | 0.000
## water | 16 | 0.349
## sex | 0 | 0.000
## age | 0 | 0.000
## relat | 102 | 2.227
## -----
```

In this implementation, a heuristic algorithm is called to suppress as few values as possible. We see that 3-anonymity is ensured. In brackets the number of observations violating 2- and 3-anonymity in the original data set are reported.

```
print(sdc, "kAnon")

## Infos on 2/3-Anonymity:
##
## Number of observations violating
##   - 2-anonymity: 0 (0.000%) | in original data: 653 (14.258%)
##   - 3-anonymity: 0 (0.000%) | in original data: 1087 (23.734%)
##   - 5-anonymity: 138 (3.013%) | in original data: 1781 (38.886%)
##
## -----
```

Importance of variables: It is possible to specify a desired ordering of key variables. If you have a look at the functions arguments in `kAnon()`, it can be seen that the function argument `importance` is responsible for it. The aim is that the higher the importance of a variable, the less suppressions are taken for this variable.

In other words, the algorithm allows the specification of a preference-vector that determines which key variables should be preferred when selecting values to be suppressed. By specifying this importance of variables as a parameter in `kAnon()`, for key variables with high importance, suppression will only take place if no other choices are possible. Still, it is possible to achieve k -anonymity for selected key variables.

```
## undo last local suppression
sdc <- undolast(sdc)
```

Remember, the key variables are

```
str(testdata[, sdc@keyVars])

## 'data.frame':   4580 obs. of  5 variables:
##  $ urbrur: int   2 2 2 2 2 2 2 2 2 ...
##  $ water  : int   3 3 3 3 3 3 3 3 3 ...
##  $ sex    : int   1 2 1 1 1 2 2 2 1 2 ...
##  $ age    : int  46 41 9 6 52 47 13 19 9 16 ...
##  $ relat  : int   1 2 3 3 1 2 3 3 3 3 ...
```

Another importance is assigned to those variables. In general, the lower the number, the less local suppressions. For any reason, we assume that variable *relat* is very important, giving the importance of 1. Variable *age* seems more important than *urbrur*, *water* and *sex*, thus giving age importance of 2. The rest of the variables are assigned by the numbers 5 (*urbrur*), 4 (*water*) and 3 (*sex*). The result is now different from before.

```
sdc <- kAnon(sdc, importance = c(5,4,3,2,1), k = 3)
print(sdc, "ls")
```

```
## Local suppression:
## KeyVar | Suppressions (#) | Suppressions (%)
## urbrur | 44 | 0.961
## water | 52 | 1.135
## sex | 13 | 0.284
## age | 9 | 0.197
## relat | 4 | 0.087
## -----
```

It can be seen that now other variables are mainly used for suppressions. In variable *relat*, for example, only 3 instead of 114 local suppression are made.

Stratification:

The methods can also be applied on each strata separately as long as the strata is specified in `createSdcObj`. Automatically then the algorithm ensures k -anonymity in all strata.

4.2.2.4 An Approach for Ensuring k -anonymity in Subsets

Due to computational issues another approach exists in μ -Argus, often referred as all- M approach. It is also implemented in **sdcMicro** (with no computational constraints), but mainly for comparison reasons. With this approach k -anonymity is usually not reached considering all key variables, but it can ensure k -anonymity on subsets of key variables.

It is also possible to provide k -anonymity for subsets of key-variables with varying parameter k . In the follow-up example we want to provide 10-anonymity for all combinations of 4 key variables, 20-anonymity for all combinations with 3 key variables and 30-anonymity for all combinations of 2 key variables. Note that strata are automatically considered.

```
sdc <- undolast(sdc) combs
<- 4:2 k <-
c(5,10,20)
sdc <- kAnon(sdc,
k=k, combs=combs)
print(sdc, "ls")

## Local suppression:
## KeyVar | Suppressions (#) | Suppressions (%)
## urbrur | 0 | 0.000 ## water |
96 | 2.096 ## sex | 0 |
0.000 ## age | 98 | 2.140 ## relat |
147 | 3.210 ##
-----
```

```
print(sdc, "kAnon")
```

```
## Infos on 2/3-Anonymity:
##
## Number of observations violating
##   - 2-anonymity: 0 (0.000%) | in original data: 653 (14.258%)
##   - 3-anonymity: 3 (0.066%) | in original data: 1087 (23.734%)
##   - 5-anonymity: 30 (0.655%) | in original data: 1781 (38.886%)
##
## -----
```

Since both, k -anonymity is not ensured by this approach and a lot of parameters have to be set (k 's and combinations), we strongly suggest to use the approach from the previous Sect. 4.2.2.3. In addition, more suppressions are made using this approach based on subsets.

However, if for any reason the SDC specialist decides to have a large amount of key variables this approach becomes important.

4.2.2.5 Choice of the Local Suppression Algorithm

The “Argus” subset approach and the optimal approach are equal if `combs = #` number of key variables. In any other case, the optimal approach ensures k -anonymity while the subset approach cannot ensure this. In addition, several parameters have to be determined when using the subset approach.

In general, it is advisable to provide k -anonymity for the most important indirect identifiers, i.e. to apply the optimal approach. The number of key variables is usually a number between 4 and 8. Scenarios with more than 8 key variables are seldom and rarely reported in literature. The optimal approach is preferable in this case since it ensures k -anonymity and does not lead to much over-suppression.

With a very high number of key variables, say 15, the optimal approach will reach k -anonymity but, because the high number of possible keys, many values might be suppressed.

With such a large number of key variables, one alternative is to apply PRAM (see Sect. 4.2.3) on some (e.g. 7 variables) and the optimal local suppression approach guarantees k -anonymity for the remaining variables (e.g. 8 variables). However, since it is not easy to determine the disclosure risk for prammed variables, alternatives are suggested in the following.

If PRAM should not be applied, k anonymity should be ensured in subsets of variables. Ideally, the number of variables in a subset should be high, e.g. 8. Then, 6435 combinations of 8 out of 15 variables exist. Thus, k -anonymity is ensured in each of the 6435 combinations by using the subset local suppression approach.

In the previous section, it can be seen that three values of k for three combinations are given. But using `combs <- 4:2` and `k <- c(5, 10, 20)` was somehow a subjective decision. On the one hand, one wants to guarantee low disclosure risk. With this aim, the values for `combs` should be high as well as the corresponding values of k . On the other hand, one wants to result in a low number of suppressions, thus low numbers of `comb` and k should be chosen. Currently, there is no algorithm available which would solve this optimization problem.

4.2.2.6 Mondrian

A very simple approach which is often cited is the Mondrian algorithm (LeFevre et al. 2006) which is based on multidimensional partitioning. For each split/partitioning, the key variable with the highest number of categories is chosen. The data are split according to the median in each partition whenever a split is possible. A split is possible when the splitted data contains enough observations. More precisely, this greedy (strict) median-partitioning algorithm results in a set of multidimensional regions, each containing between k and $2p(k - 1) + m$, with p the number of key variables and m the number of distinct categories (for more details, see LeFevre et al. 2006).

This algorithm ensures k -anonymity but due to its simple procedure it may lead to over-suppressions. The algorithm is not included in **sdcmicro** but available upon request.

4.2.2.7 Linked Variables in Local Suppression

As mentioned in Chap. 2 after applying local suppression (`kAnon` or `localSuppression`) ghost/linked variables should have the same suppression pattern as the variable that they are linked to.

We give an illustrative artificial example to show how this concept works. Another variable with the same content is produced, in the following this is variable `electcon2`, `electcon3` (linked to `electcon` and `water2` (linked to `water`)).

```
data("testdata", package = "sdcMicro")
testdata$electcon2 <- testdata$electcon
testdata$electcon3 <- testdata$electcon
testdata$water2 <- testdata$water
keyVars <- c("urbrur", "roof", "walls", "water", "electcon", "relat", "sex")
numVars <- c("expend", "income", "savings")
w <- "sampling_weight"
```

We want to make sure that some variables not used as key-variables will have the same suppression pattern as variables that have been selected as key variables. Thus, we are using *ghost*-variables. As indicated above, in our example we want variables `electcon2` and `electcon3` to be linked to key-variable `electcon` and we want variable `water2` to be linked to key-variable `water`.

```
ghostVars <- list()
ghostVars[[1]] <- list()
ghostVars[[1]][[1]] <- "electcon"
ghostVars[[1]][[2]] <- c("electcon2", "electcon3")
ghostVars[[2]] <- list()
ghostVars[[2]][[1]] <- "water"
ghostVars[[2]][[2]] <- "water2"
```

Next we create the *sdcMicroObj* object.

```
obj <- createSdcObj(testdata, keyVars=keyVars,
  numVars=numVars, w=w, ghostVars=ghostVars)
```

We apply 3-anonymity to selected key variables

```
obj <- kAnon(obj, k=3)
obj
```



```

obj <- kAnon(obj, k=3)
obj

## Data set with 4580 rows and 17 columns.
## --> Categorical key variables: urbrur, roof, walls, water, electcon,
##                               relat, sex
## --> Numerical key variables: expend, income, savings
## --> Weight variable: sampling_weight
## -----
##
## Information on categorical Key-Variables:
##
## Reported is the number, mean size and size of the smallest category
## for recoded variables.
## In parenthesis, the same statistics are shown for the unmodified data.
## Note: NA (missings) are counted as separate categories!
##
## Key Variable Number of categories      Mean size
##   urbrur                2 (2) 2290.000 (2290.000)
##   roof                  6 (5)  913.600 (916.000)
##   walls                 3 (3) 1526.667 (1526.667)
##   water                 9 (8)  569.250 (572.500)
##   electcon              4 (3) 1525.333 (1526.667)
##   relat                 9 (9)  554.000 (508.889)
##   sex                   2 (2) 2290.000 (2290.000)
## Size of smallest
##      646 (646)
##      15 (16)
##      50 (50)
##      25 (26)
##     103 (107)
##       3 (1)
##    2284 (2284)
## -----
##
## Infos on 2/3-Anonymity:
##
## Number of observations violating
## - 2-anonymity: 0 (original data: 157)
## - 3-anonymity: 0 (original data: 281)
##
## Percentage of observations violating
## - 2-anonymity: 0.000 % (original data: 3.428 %)
## - 3-anonymity: 0.000 % (original data: 6.135 %)
## -----
##
## Numerical key variables: expend, income, savings
##
## Disclosure risk (~100.00% in original data):
##   modified data: [0.00%; 100.00%]
##
## Current Information Loss in modified data (0.00% in original data):
##   IL1: 0.00
## Difference of Eigenvalues: 0.000%
## -----
##
## Local Suppression:
##   KeyVar | Suppressions (#) | Suppressions (%)
##   urbrur |          0 |          0.000
##   roof   |         12 |          0.262
##   walls  |          0 |          0.000
##   water  |         26 |          0.568
##   electcon |         4 |          0.087
##   relat  |        148 |          3.231
##   sex    |          0 |          0.000
## -----

```

In the following we can see that the suppression patterns of the key variable and its ghost variables are identical.

```
manipGhostVars <- get.sdcMicroObj(obj, "manipGhostVars")
manipKeyVars <- get.sdcMicroObj(obj, "manipKeyVars")
all(is.na(manipKeyVars$selectcon) == is.na(manipGhostVars$selectcon2))

## [1] TRUE

all(is.na(manipKeyVars$selectcon) == is.na(manipGhostVars$selectcon3))

## [1] TRUE

all(is.na(manipKeyVars$water) == is.na(manipGhostVars$water2))

## [1] TRUE
```

Exercises:

Question 4.2 Local suppression:

Take the `eusilc` data set from package `laeken` and the `sdcMicroObj` object produced in the previous example including all already done recodings. The task is now to ensure k -anonymity.

4.2.3 Post-randomization Method (PRAM)

If there are a larger number of categorical key variables (e.g., more than 5), recoding might not sufficiently reduce disclosure risks, or local suppression might lead to great information loss. In this case, the PRAM (Gouweleeuw et al. 1998) may be a more efficient alternative.

PRAM (Gouweleeuw et al. 1998) is a probabilistic, perturbative method for protecting categorical variables. The method swaps the categories for selected variables based on a pre-defined transition matrix, which specifies the probabilities for each category to be swapped with other categories.

To illustrate, consider the variable `location`, with three categories: `location = 1` “east”, `location = 2` “middle”, `location = 3` “west”. We define a 3-by-3 transition matrix, where p_{ij} is the probability of changing category i to j . For example, in the following matrix,

$$\mathbf{P} = \begin{pmatrix} 0.1 & 0.9 & 0 \\ 0.2 & 0.1 & 0.7 \\ 0.9 & 0 & 0.1 \end{pmatrix}$$

the probability that the value of the variable will stay the same after perturbation is 0.1, since we set $p_{11} = p_{22} = p_{33}$. The probability of east being changed into middle is p_{12} , while east will not be changed into west because p_{13} is set to be 0. PRAM protects the records by perturbing the original data file, while at the same time, since the probability mechanism used is known, the characteristics of the original data can be estimated from the perturbed data file. PRAM can be applied to each record independently, allowing the flexibility to specify the transition matrix as a function parameter according to desired effects. For example, it is possible to prohibit changes from one category to another by setting the corresponding probability in the transition matrix to 0, as shown in the example above. It is also possible to apply PRAM to subsets of the microdata independently.

```
set.seed(1234)
A <- as.factor(rep(c("A1", "A2", "A3"), each=5)); A

## [1] A1 A1 A1 A1 A1 A2 A2 A2 A2 A2 A3 A3 A3 A3 A3
## Levels: A1 A2 A3
```

We apply `pram()` on vector `A` and print the result:

```
Apramed <- pram(A); Apramed

## Number of changed observations:
## - - - - -
## x != x_pram : 1 (6.67%)
```

The summary provides more detail. It shows a table of original frequencies and the corresponding table after applying PRAM. All transitions that took place are also listed:

```
summary(Apramed)

## Variable:  x
## -----
## Frequencies in original and perturbed data:
##              x A1 A2 A3 NA
## 1:          Original Frequencies  5  5  5  0
## 2: Frequencies after Perturbation  6  5  4  0
##
## Transitions:
## transition Frequency
## 1:  1 --> 1          5
## 2:  2 --> 2          5
## 3:  3 --> 1          1
## 4:  3 --> 3          4
```

PRAM is applied to each observation independently and randomly. This means that different solutions are obtained for every run of PRAM if no seed is specified for the random number generator. A main advantage of the PRAM procedure is the flexibility of the method. Since the transition matrix can be specified freely as a function

parameter, all desired effects can be modeled. For example, it is possible to prohibit changes from one category to another by setting the corresponding probability in the transition matrix to 0.

In **sdcMicro**, `pram()` allows PRAM to be performed. The corresponding help file can be accessed by typing `?pram` into an R console. When using the function it is possible to apply PRAM to sub-groups of the micro-data set independently. In this case, the user needs to select the stratification variable defining the sub-groups. If the specification of this variable is omitted, the PRAM procedure is applied to all observations of the data set.

We again create an object of class *sdcMicroObj*, also specifying a strata variable to apply PRAM within these strata independently. Note that when applying PRAM to variables, these variables should be saved as a factor.

```
require("sdcMicro")
data(testdata, package="sdcMicro")
# categorical variables should be saved as a factor
vars <- c('urbrur', 'water', 'sex', 'age', 'relat', 'walls', 'roof')
testdata[, vars] <- lapply(testdata[, vars], as.factor)

# setting up the SDC problem
sdc <- createSdcObj(testdata,
  keyVars = c('urbrur', 'water', 'sex', 'age', 'relat'),
  numVars = c('expend', 'income', 'savings'),
  pramVars = c("walls"),
  w = 'sampling_weight',
  hhId = 'ori_hid',
  strataVar = 'hhcivil'
)
```

```
sdc <- pram(sdc)
print(sdc, "pram")

## Post-Randomization (PRAM):
## Variable: walls
## --> final Transition-Matrix:
##           2           3           9
## 2 0.84447887 0.1499160 0.005605102
## 3 0.05420769 0.9410131 0.004779201
## 9 0.13485876 0.3180081 0.547133184
##
## Changed observations:
##   variable nrChanges percChanges
## 1   walls       398         8.69
## -----
```

Sometimes it is useful to apply `pram` not on all categories. To give an illustrative example, children with age below 16 should not be prammed in variable education. Or governmental organisations should not be prammed in variable ownership.

We show an example for our `testdata` set. We want to apply `pram` to variable `urbrur` for each group of variable `urbrur`. However, no value should be changed where `roof`

equals 4. Thus we are creating a new value for these observations. First the previous application of PRAM is undone, then a new category is made for those observations that should not change.

```
sdc <- undolast(sdc)
sv <- testdata$urbrur
levels(sv) <- c("1", "2", "3")
sv[testdata$roof == 4] <- 3
```

Next PRAM is applied. For a later check, the original data as well as the prammed variable *roof* is extracted.

```
sdc <- pram(sdc, variables = "roof", strata_variables = sv)
orig <- get.sdcMicroObj(sdc, "origData")$roof
pramed <- get.sdcMicroObj(sdc, "manipPramVars")$roof
```

Now it can be validated if any of the category 4 in the prammed variable differs from the original category.

```
all(pramed[orig == 4] == 4)

## [1] FALSE
```

It can be seen that nothing is changed for this category.

4.3 Methods for Continuous Key Variables

4.3.1 Microaggregation

Micro-aggregation (Defays and Anwar 1998) is a perturbing method typically applied to continuous variables. It is also a natural approach to achieving *k*-anonymity. The method first partitions records into groups, then assigns an aggregate value (typically the arithmetic mean, but other robust methods are also possible) to each variable in the group. As an example, in Table 4, records are first partitioned into groups of two, and then the values are replaced by the group means. Note that in the example, by setting group size of two, micro-aggregation automatically achieves 2-anonymity with respect to the three key variables.

Individual values of the records for each variable are replaced by the group aggregation value, which is often the mean; as an example, see Table 4.3, where two values that are most similar are replaced by their column-wise means.

To preserve the multivariate structure of the data, the most challenging part of micro-aggregation is grouping records by how “similar” they are. The simplest method is to sort data based on a single variable in ascending or descending order. Another option is to cluster data first, and sort by the most influential variable in each cluster. These methods, however, might not be optimal for multivariate data (Templ and Meindl 2008).

Table 4.3 Example of micro-aggregation. Columns 1–3 contain the original variables, columns 4–6 the micro-aggregated values (rounded on two digits)

	Num1	Num2	Num3	Mic1	Mic2	Mic3
1	0.30	0.400	4	0.65	0.85	8.5
2	0.12	0.220	22	0.15	0.51	15.0
3	0.18	0.800	8	0.15	0.51	15.0
4	1.90	9.000	91	1.45	5.20	52.5
5	1.00	1.300	13	0.65	0.85	8.5
6	1.00	1.400	14	1.45	5.20	52.5
7	0.10	0.010	1	0.12	0.26	3.0
8	0.15	0.500	5	0.12	0.26	3.0

Individual ranking method

The individual ranking methods (Defays and Anwar 1998) is often applied because of its simplicity. The method replaces values by its aggregates column by column independently. First, the first column is sorted and the index of sorting is memorized to be able to sort the values back in the original order. Then the first k values are replaced by their aggregate (usually the arithmetic mean), the next k values are replaced by their aggregate, and so on, until all values are aggregated from the first variable. The variable is then back-sorted. This procedure is then applied on the other variables independently.

Microaggregation based on PCA

The Principle Component Analysis method sorts data on the first principal components (see, e.g., Templ and Meindl 2008). A robust version of this method can be applied to clustered data for small- or medium-sized datasets (see, e.g., Templ and Meindl 2008). This approach is fast and performs well whenever the first principal component explains a high percentage of the variance for the variables considered for micro-aggregation. If this is not the case, other algorithms are preferable, such as the MDAV algorithm explained next.

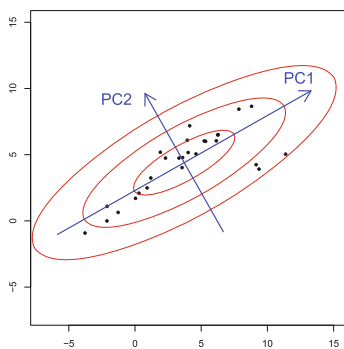
Figure 4.1 shows the procedure. First, the first principal component must be estimated and along the first principal component the values are then aggregated. As for almost any microaggregation procedure, the remaining $(2k - 1)$ are microaggregated whenever $n \bmod k$ is unequal zero.

Microaggregation by MDAV

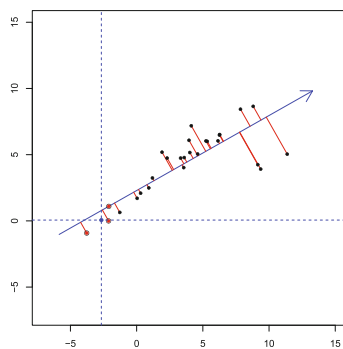
The Maximum Distance to Average Vector (MDAV) method is a standard method that groups records based on classical Euclidean distances in a multivariate space (Domingo-Ferrer and Mateo-Sanz 2002).

Figure 4.2 shows the MDAV procedure, which works as follows.

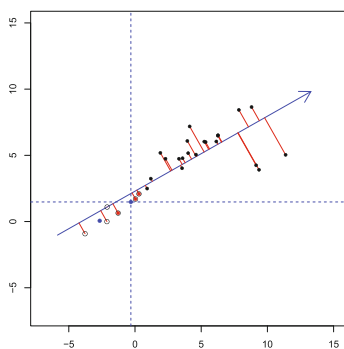
1. The center of the data is estimated using column-wise arithmetic means.
2. The farthest observation (Euclidean distance), say \mathbf{x}_r from the center is then chosen.



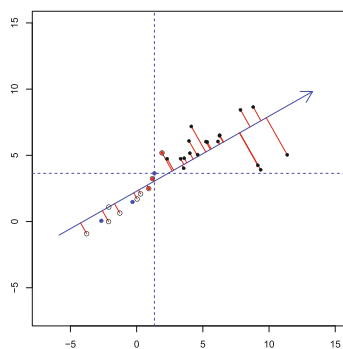
(a) First two principal components.



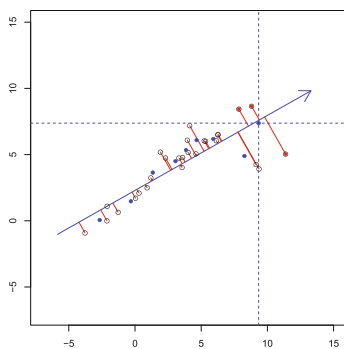
(b) Projection/sorting on first component. Aggregation along the sorted values (first three observations).



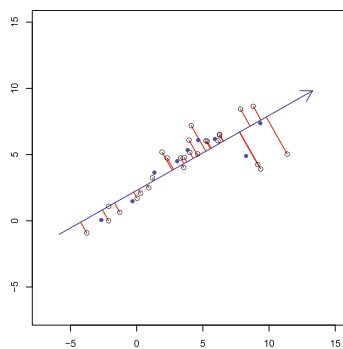
(c) Aggregation along the sorted values (next three observations).



(d) Aggregation along the sorted values (next three observations).



(e) Aggregation along the sorted values (remaining last observations).



(f) Microaggregated values in blue.

Fig. 4.1 Schematic workflow of microaggregation using principal component analysis on two-dimensional toy data. Aggregation is done along the first principal component

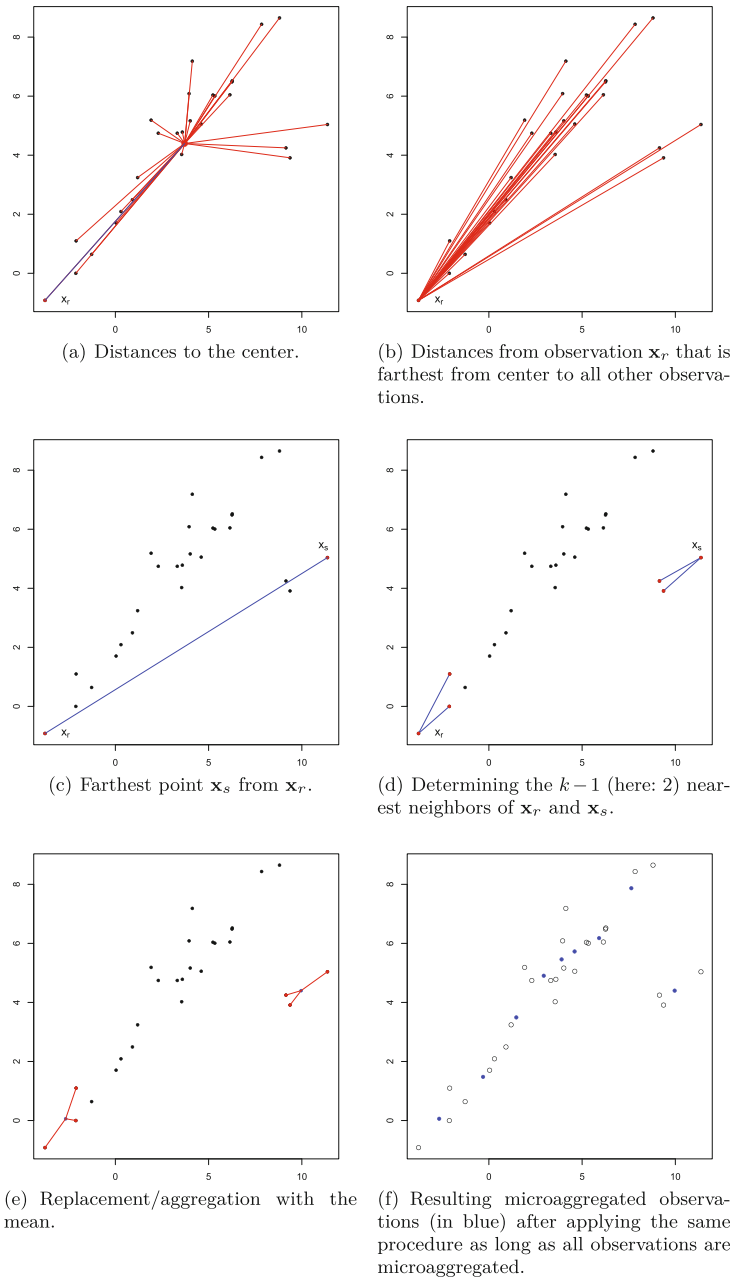


Fig. 4.2 Schematic workflow of microaggregation using MDAV on two-dimensional toy data. Aggregation is stepwise applied after all observations are microaggregated. The first step is shown in (b)–(e)

3. The farthest observation (Euclidean distance) from \mathbf{x}_s is also selected.
4. $k - 1$ nearest neighbors (using Euclidean distances) are chosen for \mathbf{x}_r , and also for \mathbf{x}_s and the arithmetic column-wise means are estimated for both groups of observations.
5. The corresponding observations are replaced by their arithmetic column-wise means.
6. Start at (1) for the remaining observations. A variant of the procedure is to start not at step (1) but at step (2), i.e. not to estimate the centers in each step but held them fixed. This iteration is done until less than or equal $2k - 1$ observations are left.
7. Aggregate the remaining observations.

Microaggregation by robust Mahalanobis distances

The MDAV method was further improved by replacing Euclidean distances with robust multivariate (Mahalanobis) distance measures (Templ and Meindl 2008). All of these methods are implemented in **sdcMicro** (Templ et al. 2015) and **sdcMicroGUI** (Kowarik et al. 2013).

Microaggregation for mixed scale data

To deal not only with continuous key variables but also with categorical (and ordered categorical or semi-continuous) in the same time, two central issues have to be changed. First, a more general definition of distance is required and secondly, for non-continuous variables a strategy for aggregation has to be developed.

For the first task, this microaggregation method is using distances computed similar to the Gower distance (Gower 1971). The distance function makes distinction between the variable types categorical, ordered categorical, continuous and semi-continuous (variables with a fixed probability mass at a constant value, e.g. 0). The distance calculation is explained in Sect. 5.2.1.

After calculating the distances and to microaggregate, it is necessary to sample a category in each group whenever a variable is non-(semi-)continuous. This probabilities corresponding to the occurrence of the level in the groups. The level with the most occurrences is then chosen, or if the maximum is not unique it is selected randomly.

Applying microaggregation in **sdcMicro** is straightforward, we again apply the corresponding function to our *sdcMicroObj* object, and the utility and risk will be automatically updated as well as other slots such as *manipData*.

```
sdc <- microaggregation(sdc)
print(sdc, "numrisk")

## Numerical key variables: expend, income, savings
##
## Disclosure risk is currently between [0.00%; 36.09%]
##
## Current Information Loss:
##   - IL1: 0.11
##   - Difference of Eigenvalues: 0.060%
## -----
```

Depending on the chosen method in function `microaggregation()`, additional parameters can be specified. For example, it is possible to specify the number of observations that should be aggregated (parameter `aggr`) as well as the statistics used to calculate the aggregation (parameter `measure`). This statistics defaults is the arithmetic mean. It is also possible to perform micro-aggregation independently to pre-defined strata (the parameter is called `strata_variables`) or to use cluster methods to achieve the grouping (depending on parameter `method`). Let us undo the last action and perform microaggregation using different parameters. As mentioned before, the disclosure risk updates automatically.

```
sdc <- undolast(sdc)
sdc <- microaggregation(sdc,
                        aggr = 4,
                        strata_variables="age",
                        method="mdav")
print(sdc, "numrisk")

## Numerical key variables: expend, income, savings
##
## Disclosure risk is currently between [0.00%; 25.35%]
##
## Current Information Loss:
##   - IL1: 0.13
##   - Difference of Eigenvalues: 0.100%
## -----
```

We can observe that the disclosure risk decreased considerably – \sim only 2.88% of the original values do not fall within intervals calculated around the perturbed values, compare Sect. 3.11 on distance-based risk estimation. We can see that the information loss criteria increased slightly. All of the previous settings (and many more) can be applied in **sdcMicro**. The corresponding help file can be viewed with command `?microaggregation`.

Not very commonly used is the microaggregation of categorical and continuous variables. This can be achieved with function `microaggrGower` that uses the Gower distance (Gower 1971) to measure the similarity between observations. For details have a look at the help function `?microaggrGower` in R. its application is also straightforward, just apply the function to an object of class *sdcMicroObj*. Note that in such an object, every important information is already stored and the variables to be microaggregated are chosen automatically.

```
sdc <- microaggrGower(sdc)
```

4.3.2 Noise Addition

Adding noise is a perturbative method typically applied to continuous variables. The idea is to add or multiply a stochastic or randomized number to the original values to protect data from exact matching with external files. While this approach sounds simple in principle, many different algorithms can be used. In this section, we introduce the uncorrelated and correlated additive noise methods.

Adding noise should be used with caution, as the results depend greatly on the amount of noise chosen.

4.3.2.1 Uncorrelated Additive Noise

Uncorrelated additive noise (see, e.g., Hundepool et al. 2007) can be expressed as the following:

$$\mathbf{z}_j = \mathbf{x}_j + \epsilon_j \quad , \quad (4.1)$$

where vector \mathbf{x}_j represents the original values of variable j , \mathbf{z}_j represents the perturbed values of variable j and ϵ_j (uncorrelated noise, or white noise) denotes normally distributed errors with $Cov(\epsilon_l, \epsilon_k) = 0$ for all $k \neq l$. In matrix notation this looks like

$$\mathbf{Z} = \mathbf{X} + \epsilon \quad , \quad (4.2)$$

with $X \sim (\mu, \Sigma)$, $\epsilon \sim N(0, \Sigma_\epsilon)$ and $\Sigma_\epsilon = c \cdot \text{diag}(\sigma_1^2, \sigma_2^2, \dots, \sigma_p^2)$, for a constant $c > 0$.

In the following code, synthetic toy data are simulated and noise is added. The original but also the perturbed data are displayed in Fig. 4.3. Noise can be added to variables using function `addNoise()` and its parameter `noise`. The corresponding help file can be accessed with `?addNoise`. In addition, a 97.5% tolerance ellipse showing the covariance structure is drawn for both the original and the perturbed data. More noise than probably needed is added to the original variables in order to show the flaws of this method. It is clearly visible that the covariance of the data set is not respected in the noise addition. The perturbed data have a different covariance structure than the original data.

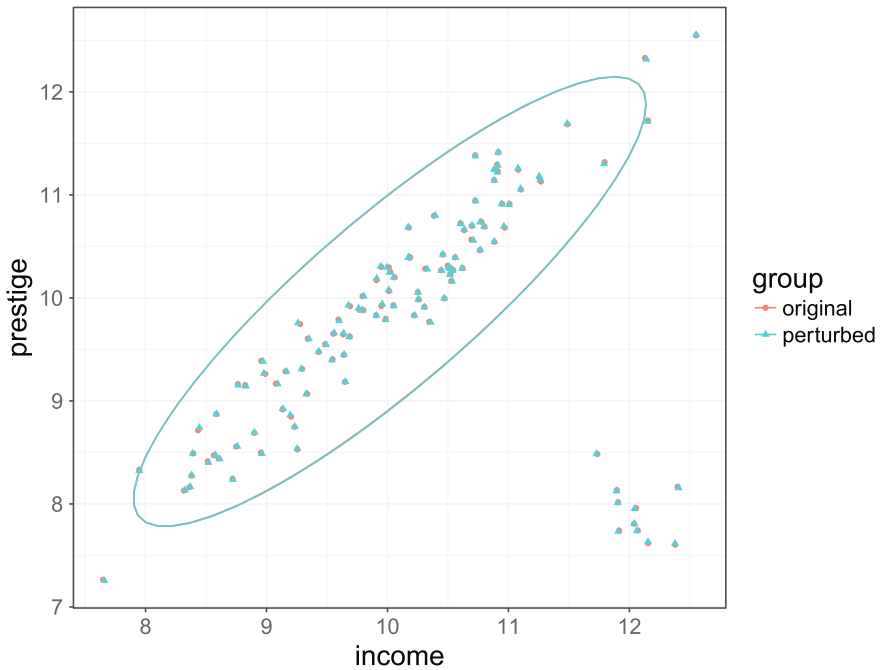


Fig. 4.3 Adding additive noise. More noise than probably needed is added to see the flaws of the method

```
## we generate some synthetic bivariate toy data to illustrate noise addition
X <- MASS::mvrnorm(100,
  mu = c(10,10),
  Sigma = matrix(c(1,0.95,0.95,1),
    byrow = TRUE, ncol=2))
X2 <- MASS::mvrnorm(10,
  mu = c(12,8),
  Sigma = matrix(c(0.1,0,0,0.1),
    byrow = TRUE, ncol=2))
X <- data.frame(rbind(X, X2))
colnames(X) <- c("income", "prestige")
head(X)

##      income  prestige
## 1 10.887193 11.249357
## 2 11.101999 11.053242
## 3 12.128663 12.329934
## 4  8.585157  8.873169
## 5  8.899352  8.689755
## 6  8.564348  8.472442

## now add noise:
set.seed(123)
Y <- addNoise(X, method="additive", noise=0.7)$xm
```

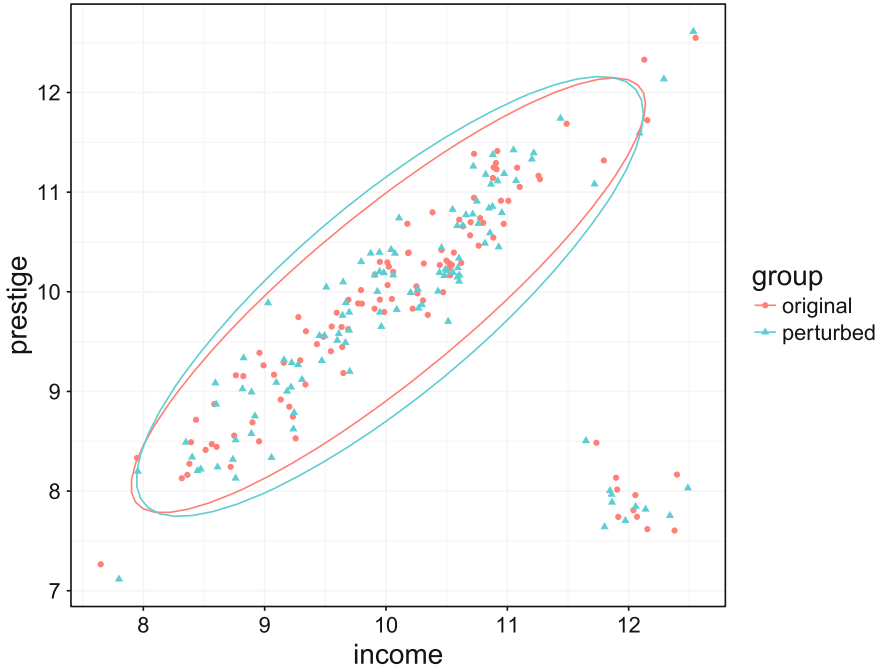


Fig. 4.4 Adding correlated noise

4.3.2.2 Correlated Additive Noise

While adding uncorrelated additive noise preserves the means and variance of the original data, co-variances and correlation coefficients are not preserved. It is preferable to apply correlated noise because the co-variance matrix of the errors is proportional to the co-variance matrix of the original data (Brand 2002).

The difference to the uncorrelated noise method is that the covariance matrix of the errors is now designed to be proportional to the covariance of the original data, i.e. $\epsilon \sim N(0, \Sigma_\epsilon = c\Sigma)$. The following holds

$$\Sigma_Z = \Sigma + c\Sigma = (1 + c)\Sigma \quad . \quad (4.3)$$

Whenever the constant c is known, the covariance of the original data can be estimated from the perturbed data set.

Figure 4.4 presents the original data as well as the perturbed data set with the additive correlated noise method. Again, more noise as probably necessary is added to better see the effect of the method. Nevertheless, at least the covariance structure of this toy data set is well preserved.

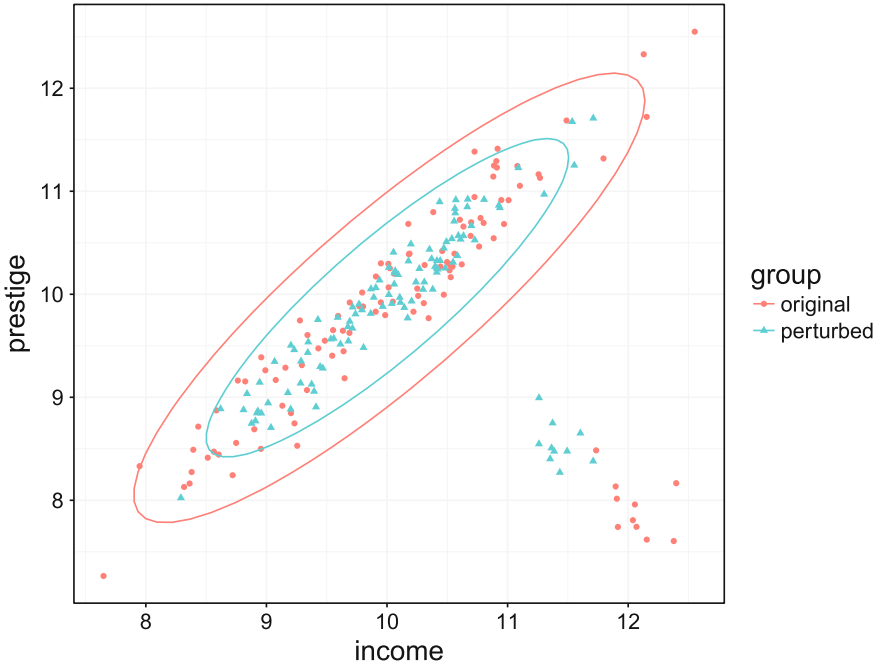


Fig. 4.5 Adding correlated noise based on transformations

4.3.2.3 Adding Correlated Noise Based on Transformations

Adding correlated noise based on transformations (Kim 1986) is another method that is also implemented in **sdcMicro**. Here we calculate $d = (1 - c^2)\epsilon$ and then $x_j d + cz_j$ whereas y_j are random numbers from $N(\frac{(1-d)\bar{x}_j}{c}, s_j)$, with \bar{x}_j and s_j the mean standard deviation of x_j .

Another similar method which takes the sample size into account is described, for example, in Brand (2002). It is a method which is often denoted as the restricted correlated noise method and which is implemented in **sdcMicro** as well (method `restr.`).

Figure 4.5 presents the original data as well as the perturbed data set with the transformation based additive correlated noise method. Again, more noise as probably necessary is added to better see the effect of the method. The tails of the distribution are wider for the perturbed data and also a bias is introduced for the outlying group in the bottom-right of the graphic.

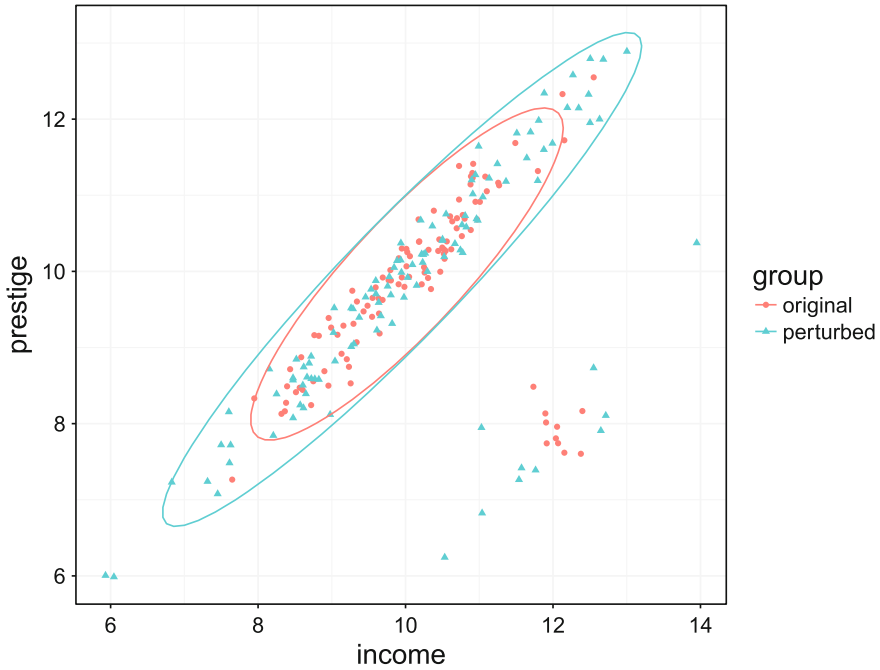


Fig. 4.6 Adding noise by Random Orthogonal Matrix Masking

4.3.2.4 Random Orthogonal Matrix Masking

Furthermore, for method ROMM (Random Orthogonal Matrix Masking) (Ting et al. 2005) perturbed data are obtained by $\mathbf{Z} = \mathbf{A}\mathbf{X}$, whereas \mathbf{A} is randomly generated and fulfils $\mathbf{A}^{-1} = \mathbf{A}^T$ (orthogonality condition). To obtain an orthogonal matrix as described in Ting et al. (2005) the Gram-Schmidt procedure (see, e.g., Golub and Van Loan 1996) is used.

Figure 4.6 presents the original data as well as the perturbed data set using ROMM. Again, more noise as probably necessary is added (parameter p) to better see the potential flaws of the method. The tails of the distribution are more narrow for the perturbed data.

4.3.2.5 Robustness Issues

The distribution of the original variables \mathbf{x}_j may not follow a normal distribution and, especially the correlated noise methods may be influenced by this fact.

In this case, a robust version of the correlated noise method is described in detail by Templ and Meindl (2008). Hereby, the covariance matrices are estimated with robust covariance estimators such as the MCD estimator (Rousseeuw and Van Driessen 1999).

In **sdcMicro**, many algorithms are implemented that can be used to add noise to continuous variables. For example, it is possible to add noise only to outlying observations. In this case it is assumed that such observations possess higher risks than non-outlying observations.

4.3.2.6 Applying Adding Noise Methods

Of course, any mentioned noise methods can be applied to objects from class *sdcMicroObj* using function `addNoise`, e.g.

We now want to apply a method for adding correlated noise based on non-perturbed data after we undo microaggregation on our previously generated `sd` object:

```
sd <- undolast(sdc)
sd <- addNoise(sdc, method="correlated2")
print(sdc, "numrisk")

## Numerical key variables: expend, income, savings
##
## Disclosure risk is currently between [0.00
##
## Current Information Loss:
##   - IL1: 0.11
##   - Difference of Eigenvalues: 0.060
## -----
```

We see that the data utility measure is comparable with the microaggregation on strata in the previous code chunk but the risk is higher than before using microaggregation.

Other methods ensure that the amount of noise added takes into account the underlying sample size and sampling weights.

4.3.3 Shuffling

Shuffling (Muralidhar and Sarathy 2006) generates new values for selected sensitive variables based on the conditional density of sensitive variables given non-sensitive variables. As a rough illustration, assume we have two sensitive variables, income and savings, which contain confidential information. We first use age, occupation, race and education variables as predictors in a regression model to simulate a new set of values for income and savings. We then apply reverse mapping (i.e., shuffling) to replace ranked new values with the ranked original values for income and savings. This way, the shuffled data consists of the original values of the sensitive variables. Moreover, Muralidhar and Sarathy (2006) showed that since we only need the rank of

the perturbed value in this approach, instead of generating a new value, shuffling can be implemented using only the rank order correlation matrix (measuring the strength of the association between the ranked sensitive variables and ranked non-sensitive variables) and the ranks of values of non-sensitive variables.

In the following code we do not use default values because we want to show how to specify the form of the model. We first restore the previous results and remove the effect of adding noise using `undolast()`.

```
sdc <- undolast(sdc)
form <- formula(expend + income + savings ~ age + urbrur + water +
  electcon + age + sex, data=testdata)
sdc <- shuffle(sdc, form)
print(sdc, "numrisk")

## Numerical key variables: expend, income, savings
##
## Disclosure risk is currently between [0.00%; 0.22%]
##
## Current Information Loss:
##   - IL1: 1.63
##   - Difference of Eigenvalues: 2.090%
## -----
```

To find a good model is essential to provide good results. If we would for example add variable walls to the model, the results are no longer of high quality.

Exercises:

Question 4.3 Adding noise:

Take the data set EIA. Assume that variables RESSALES, COMSALES, INDSALES, and OTHRSALES are continuous key variables, and UTILNAME and STATE the categorical key variables. Define your *sdcMicroObj* object. Compare additive noise and the method called `correlated2` for adding noise in terms of disclosure risk. Remark: In the next chapter, we will also compare the data utility.

Question 4.4 Microaggregation

Apply microaggregation on each state independently.

Question 4.5 Microaggregation

Assume that in general your management defines the disclosure risk on the basis of 3-anonymity for your categorical key variables. Is your data set *safe* if you apply microaggregation on your continuous key variables and ensure 3-anonymity for your categorical key variables?

References

- Kowarik, A., Templ, M., Meindl, B., & Fonteneau, F. (2013). *sdMicroGUI: Graphical user interface for package sdMicro*. R package version 1.1.1. <http://CRAN.R-project.org/package=sdMicroGUI>.
- LeFevre, K., DeWitt, D. J., & Ramakrishnan, R. (2006). Mondrian multidimensional k-anonymity. In *ICDE '06: Proceedings of the 22nd international conference on data engineering* (p. 25).
- Gouweleew, J., Kooiman, P., Willenborg, L., & De Wolf, P.-P. (1998). Post randomisation for statistical disclosure control: Theory and implementation. *Journal of Official Statistics* 14(4), 463–478.
- Defays, D., & Anwar, M. N. (1998). Masking microdata using micro-aggregation. *Journal of Official Statistics*, 14(4), 449–461.
- Templ, M., & Meindl, B. (2008b). Robustification of microdata masking methods and the comparison with existing methods. In *Privacy in statistical databases*. Lecture Notes in Computer Science (Vol. 5262, pp. 113–126). Springer.
- Domingo-Ferrer, J., & Mateo-Sanz, J. M. (2002). Practical data-oriented microaggregation for statistical disclosure control. *IEEE Transactions on Knowledge and Data Engineering*, 14(1), 189–201.
- Templ, M., Meindl, B., & Kowarik, A. (2015). Statistical disclosure control for micro-data using the R package sdMicro. *Journal of Statistical Software*, 67(1), 1–37.
- Gower, J. C. (1971a). A general coefficient of similarity and some of its properties. *Biometrics*, 27(4), 857–871.
- Hundepool, A., et al. (2007). *Handbook on statistical disclosure control*.
- Brand, R. (2002). *Microdata protection through noise addition*. Lecture Notes in Computer Science London: Springer.
- Kim, J. J. (1986). A method for limiting disclosure in microdata based on random noise and transformation. In *Proceedings of the section on survey research methods* (pp. 303–308). American Statistical Association.
- Ting, D., Fienberg, S., & Trottini, M. (2005). Romm methodology for microdata release. In *Monographs of official statistics: Work session on statistical data confidentiality*. Luxembourg: Eurostat.
- Golub, G. H., & Van Loan, C. F. (1996). *Matrix computations* (3rd ed.). Baltimore: Johns Hopkins University Press.
- Rousseeuw, P. J., & Van Driessen, K. (1999). A fast algorithm for the minimum covariance determinant estimator. *Technometrics*, 41, 212–223.
- Muralidhar, K., & Sarathy, R. (2006). Data shuffling—A new masking approach for numerical data. *Management Science*, 52(2), 658–670.

Chapter 5

Data Utility and Information Loss

Abstract Once SDC methods have been applied to modify the original data set and to lower the disclosure risk, it is critical to measure the resulting information loss and data utility. Basically, two different kinds of complementary approaches exist to assess information loss: (i) direct measuring of distances/frequencies between the original data and perturbed data, and (ii) comparing statistics computed on the original and perturbed data. The first concept is common but often of limited use. The latter concept is closer to the users and data sets since its aim is to measure the differences for the most important indicators/estimates.

5.1 Element-Wise Comparisons

5.1.1 Comparing Missing Values

Missing values (NA) might be accounted for using a simple utility measure that counts the missing values in the original data \mathbf{X} and the anonymized data \mathbf{Y} . Let $\mathbf{R}^{(X)}$ and $\mathbf{R}^{(Y)}$ be (indicator) matrices of the same size as \mathbf{X} and \mathbf{Y} . A cell/element of $\mathbf{R}^{(X)}$ is 1 when \mathbf{X} has a missing value on that position, and zero if the corresponding value is not missing. The same is with $\mathbf{R}^{(Y)}$ regarding \mathbf{Y} . Thus the matrices $\mathbf{R}^{(X)}$ and $\mathbf{R}^{(Y)}$ consist of zeros and ones depending on the position of missings in \mathbf{X} and \mathbf{Y} . Let \mathbf{R} a matrix of the same size with elements

$$r_{ij} = \begin{cases} 0 & \text{if } r_{ik}^{(X)} = r_{ik}^{(Y)} = 0, \\ 1 & \text{if } r_{ik}^{(X)} = 1 \wedge r_{ik}^{(Y)} = 1, \\ 0 & \text{if } r_{ik}^{(X)} = 0 \wedge r_{ik}^{(Y)} = 1, \\ 0 & \text{if } r_{ik}^{(X)} = 1 \wedge r_{ik}^{(Y)} = 0. \end{cases} \quad (5.1)$$

Hereby, it is assumed that an original data set X always has more or equal information as the anonymized data set Y . Thus, a possible imputed value in Y will not effect the following measure. We can now count the number of additional missings per variable caused by anonymizing the data using the indicator matrix \mathbf{R} with n

observations and p variables,

$$m_j = \sum_i^n r_{ij} \quad , \quad j \in \{1, \dots, p\} \quad . \quad (5.2)$$

This can also be seen relatively by dividing by the number of observations, or as percentages of new missing values in each variable.

$$mp_j = 100 \cdot \frac{m_j}{n} \quad . \quad (5.3)$$

The higher m_j (or mp_j) the higher the information loss.

Let us illustrate this again on the EU-SILC data. For demonstration purposes we only define four key variables and apply local suppression on the categorical key variables.

```
library("laeken")
data("eusilc")
sdc <- createSdcObj(eusilc,
  keyVars = c("db040", "hsize", "pb220a",
    "rb090"),
  weightVar = "rb050", hhId = "db030")
sdc <- kAnon(sdc)
print(sdc, "ls")
```

```
## Local suppression:
## KeyVar | Suppressions (#) | Suppressions (%)
## db040 | 0 | 0.000
## hsize | 9 | 0.061
## pb220a | 0 | 0.000
## rb090 | 0 | 0.000
## -----
```

From this result it can be seen the number of missing values is the same for `hsize`, `pb220a` and `rb090`, but for key variable `db040` there are 4 additional missing values in the anonymized data set.

5.1.2 Comparing Aggregated Information

Instead of the direct comparison of the values of categorical variables, an alternative is to compare a contingency table $\mathbf{T}^{(X)}$ calculated from categorical variables of the original data \mathbf{X} and the contingency table $\mathbf{T}^{(Y)}$ from the anonymized data \mathbf{Y} . More precisely, the normed sum of the absolute distances between the cells of the tables with n_1 rows and n_2 columns (see also the description from Domingo-Ferrer 2009),

$$UT = \frac{1}{n_1 n_2} \sum_{i=1}^{n_1} \sum_{j=1}^{n_2} \left| T_{ij}^{(X)} - T_{ij}^{(Y)} \right| . \quad (5.4)$$

The higher UT the lower the data quality.

However, another normalization approach could be to consider only the relative change in each cell (in percentages).

$$UT2 = 100 \cdot \frac{1}{n_1 n_2} \sum_{i=1}^{n_1} \sum_{j=1}^{n_2} \left| \frac{T_{ij}^{(X)} - T_{ij}^{(Y)}}{T_{ij}^{(X)}} \right| . \quad (5.5)$$

In the following code, a contingency table of $rb090 \times db040$ (gender \times federal state) is computed for the anonymized data and original data. First, an object of class *sdcMicroObj* is created, and then we apply PRAM on federal state.

Let's start from the beginning. First we create the **sdcMicro** object, then we apply PRAM and compare the original and table considering the prammed variable `db040` (region).

```
library("laeken")
data("eusilc")
X <- Y <- eusilc
sdc <- createSdcObj(X,
  keyVars = c("db040", "hsize", "pb220a",
    "rb090", "pl030", "age"),
  pramVars = "db040",
  weightVar = "rb050", hhId = "db030")
sdc <- pram(sdc)
Y <- extractManipData(sdc)
```

We now compare the tables according to Eq. (5.5).

```
ct <- c("rb090", "db040")
Tx <- table(X[, ct])
Ty <- table(Y[, ct])
Tx

##          db040
## rb090  Burgenland Carinthia Lower Austria Salzburg Styria
## male      261      517      1417      440    1128
## female    288      561      1387      484    1167
##
##          db040
## rb090  Tyrol Upper Austria Vienna Vorarlberg
## male      650      1363    1132      359
## female    667      1442    1190      374

Ty

##          db040
## rb090  Burgenland Carinthia Lower Austria Salzburg Styria
```

```
##   male      266      514      1425      426      1116
##   female    290      559      1397      474      1177
##           db040
##   rb090    Tyrol Upper Austria Vienna Vorarlberg
##   male      649      1368      1129      374
##   female    654      1434      1198      377

n1 <- nrow(Ty)
n2 <- ncol(Ty)
## UT
sum(abs(Tx - Ty)) / (n1 * n2)

## [1] 7.333333

## UT2
sum(abs(Tx - Ty)/Tx) / (n1 * n2) * 100

## [1] 1.163519
```

We see that the mean difference in the cell values of the tables is approximately 1%.

However, contingency tables are of compositional nature (Egozcue et al. 2015) since no negative cell values are possible and not any cell can be larger than the sum of cells or the corresponding marginals in the table. To consider this, one can also compare the cells using Aitchison distances, i.e. since compositional data are represented only in the simplex sample space, we have to use a different distance measure, like the Aitchison distance. It is defined for two compositions $\mathbf{x}_i = (x_1, \dots, x_{n_2})$ and $\mathbf{y}_i = (y_1, \dots, y_{n_2})$ as

$$d_a(\mathbf{x}, \mathbf{y}) = \sqrt{\sum_{i=1}^{n_2-1} \sum_{j=i+1}^{n_2} \left(\frac{x_i}{x_j} - \frac{y_i}{y_j} \right)^2} . \quad (5.6)$$

Thus, the Aitchison distance takes care of the property that compositional data include their information only in the ratios between the parts.

We define the relative Aitchison distance as a measure of data utility

$$UTA = \frac{1}{n_1 n_2} \sum_{i=1}^{n_1} d_A \left(\mathbf{T}_i^{(X)}, \mathbf{T}_i^{(Y)} \right) , \quad (5.7)$$

where \mathbf{T}_i again denotes the i -th row of the contingency table obtained from the original and the anonymized data.

This is easy to be computed using the R package **robCompositions** (Templ et al. 2011).

```
library("robCompositions")
## UTA
```

```
aDist(Tx, Ty)
```

```
## [1] 0.09068296
```

UTA is theoretically sound, but the disadvantage is that the interpretation is more difficult than for *UT2*, which is computed without considering the simplex space of a contingency table.

Remark: For recoding variables, it is quite usual to join categories in categorical key variables. Doing so, it is no longer possible to compare the whole contingency tables since they are of different size as soon as less categories are present in the anonymized file.

Let's start from the beginning. First we create the **sdcMicro** object, then we do a recoding of *hsize* (household size).

```
library("laeken")
data("eusilc")
X <- eusilc
X$hsize <- factor(X$hsize)
sdc <- createSdcObj(X,
  keyVars = c("db040", "hsize", "pb220a",
    "rb090", "pl030", "age"),
  numVars = c("age", "eqIncome"),
  weightVar = "rb050", hhId = "db030")
before <- paste(6:9)
after <- "6-9"
sdc <- groupAndRename(sdc, var="hsize", before=before,
  after=after)
Y <- extractManipData(sdc)
```

In the following the data utility measures (UT, UT2 and UTA) are calculated.

```
ct <- c("rb090", "hsize")
Tx <- table(X[, ct])
Ty <- table(Y[, ct])
n1 <- nrow(Ty)
n2 <- ncol(Ty)
## UT
sum(abs(Tx[1:n1, 1:n2] - Ty)) / (n1 * n2)

## [1] 29.83333

## UT2
sum(abs(Tx[1:n1, 1:n2] - Ty)/Tx[1:n1, 1:n2]) / (n1 * n2) * 100

## [1] 9.478475

## UTA
aDist(Tx[1:n1, 1:n2], Ty)

## [1] 0.8215698
```

The values in the cells differ about 10% and the Aitchison-based differences are much higher than in the previous example.

Contingency tables can also be visualized using mosaic plots and the multivariate dependencies of categorical data can be shown. Mosaic plots, introduced by Hartigan and Kleiner (1981), are graphical representations of multi-way contingency tables. The frequencies of the different cells of categorical variables are visualized by area-proportional rectangles (tiles). For constructing a mosaic plot, a rectangle is first split vertically at positions corresponding to the relative frequencies of the categories of a corresponding variable. Then the resulting smaller rectangles are again subdivided according to the conditional probabilities of a second variable. This can be continued for further variables accordingly. Hofmann (2003) provides an excellent description of the construction of mosaic plots and the underlying mathematical theory, and **vcd** package in R with its **strucplot** framework (Meyer et al. 2006) gives a well-performing implementation of mosaic plots into the hand of the users.

Figure 5.1 shows the relative frequencies for all combinations of **rb090** (sex), **pb220a** (citizenship) and **hsize** (household size). Typically, fewer larger households exist than small households as well as most respondents are from Austria for the Austrian EU-SILC data. The mosaic plots in Fig. 5.1 are produced as follows.

```
require(vcd)
ct <- c("rb090", "pb220a", "hsize")
Tx <- table(X[, ct])
Ty <- table(Y[, ct])
par(mfrow=c(1,2))
## mosaic for original data
mosaic(Tx)
## mosaic for perturbed data
mosaic(Ty)
```

However, in case of complex survey data sampled from finite populations, the population frequencies should be calculated instead of sample frequencies, i.e. Horwitz-Thompson estimates of the counts are needed. Here the function **spTable** and **spMosaic** from package **simPop** (Templ et al. 2017) could be used. The result for population frequency counts for household size \times region is shown in Fig. 5.2. They are produced with the following code.

```
ct <- c("hsize", "db040")
Tx_pop <- tableWt(X[, ct], weights = X$rb050)
Ty_pop <- tableWt(Y[, ct], weights = X$rb050)
par(mfrow=c(1,2))
## mosaic for original data
mosaic(Tx_pop)
## mosaic for perturbed data
mosaic(Ty_pop)
```

In both Figures the recoding of the large households is clearly visible.

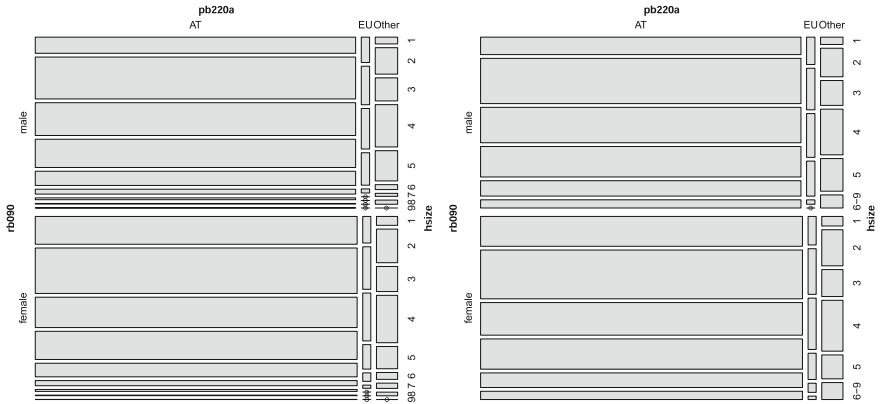


Fig. 5.1 Mosaic plot of gender (rb090) \times citizenship (pb220a) \times household size (hsize) showing the original sample frequencies (*left plot*) and the sample frequencies from the perturbed data (*right plot*)

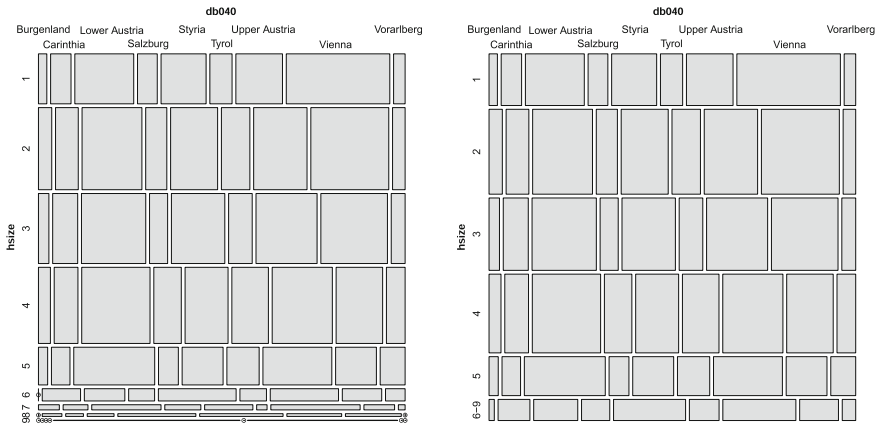


Fig. 5.2 Mosaic plot of gender (rb090) \times citizenship (pb220a) \times household size (hsize) showing the original sample frequencies (*left plot*) and the sample frequencies from the perturbed data (*right plot*)

5.2 Element-Wise Measures for Continuous Variables

In literature, element-wise measures are not often mentioned. For example, Winkler (1998) defines a microdata set to be analytically valid if (1) means and covariances on a small set of subdomains are approximately preserved, (2) marginal values for a few tabulations of the data are approximately the same and (3) at least one distributional characteristic is preserved. In addition, he proposed that a microdata file is analytically interesting if six variables on important subdomains can be validly analyzed. Also Hundepool et al. (2012) suggest to compare means, covariances, correlations,

loadings from a principal component analysis, etc. In **sdcMicro**, a measure called **eigen** is estimated, calculating relative absolute differences between eigenvalues of the co-variances from standardized continuous key variables of the original and perturbed variables. Eigenvalues can be estimated from a robust or classical version of the covariance matrix. However, all these suggestions and proposed measures are quite general, and may be not suitable for every data set. For example, covariance-based measures are only suitable in the multivariate context without any missing values and zeros in the data.

After discussing element-wise measures for continuous variables, this book focuses on specific utility measures in Sect. 5.5 that are focusing on quality of indicators that are data-dependent and subject matter dependent. General methods as like covariance-based methods are not mentioned because of limited use.

In any case, element-wise comparisons of values from the original and the anonymized data requires the definition of a distance. Element-wise measures of information loss and data utility can be, for example, be based on the classical or robust distances between original and perturbed values. Classical distances are covered in the definition of the generalized Gower distance in Eq. 5.11. Multivariate (robust) distances can also be considered by calculating (robust) Mahalanobis distances for continuous scaled variables of the original and the anonymized data, defined for an observation \mathbf{x}_i as

$$MD(\mathbf{x}_i) = [(\mathbf{x}_i - U)'C^{-1}(\mathbf{x}_i - U)]^{1/2}, i = 1, \dots, n, \quad (5.8)$$

where U and C are estimators of location and covariance. Clearly, for reliable consideration if the main bulk of the data is well preserved in the sense of robust statistics, both T and C have to be estimated in a robust way, and not in the traditional way by arithmetic mean vector and sample covariance matrix. Robust estimates of location and covariance can be obtained for instance from the MCD (Minimum Covariance Determinant) estimator (Rousseeuw and Driessen 1998).

A utility measure be based on (robust) covariance matrices can be defined as follows

```
X <- eusilc[, c("age", "eqIncome")]
## add noise to age and eqIncome
sdc <- addNoise(sdc, method = "correlated2")
Y <- extractManipData(sdc[, c("age", "eqIncome")])
Y$age <- round(as.numeric(as.character(Y$age)))
require(robustbase)
covX <- covMcd(X)
covY <- covMcd(Y)
mdX <- sqrt(mahalanobis(X,
                        center = covX$center,
                        cov = covX$cov))
mdY <- sqrt(mahalanobis(Y,
                        center = covY$center,
                        cov = covY$cov))
cat("\n relative difference in percentages:",
    round(1 / length(mdX) * sum(abs(mdX - mdY) / mdX, na.rm = TRUE) *
        100, 4), "% \n")

##
## relative difference in percentages: 11.3035 %
```

In average the Mahalanobis distances differs around 11% after recoding of age.

Following are three common proposals to measure the information loss and data utility:

- ILIs, proposed by Yancey et al. (2002), can be interpreted as the scaled distances between original and perturbed values. Again let $\mathbf{X} = \{x_{ij}\}$ be the original data set, $\mathbf{Y} = \{y_{ij}\}$ is a perturbed version of \mathbf{X} , and x_{ij} is the j -th variable in the i -th original observation. Both data sets consist of n observations and p variables each. The measure of information loss is defined by

$$ILI = \frac{1}{pn} \sum_{j=1}^p \sum_{i=1}^n \frac{|x_{ij} - y_{ij}|}{\sqrt{2}S_j}, \quad (5.9)$$

where S_j is the standard deviation of the j -th variable in the original data set.

- **prediction quality** measures the differences between estimates obtained from fitting a pre-specified regression model on the original data and the perturbed data:

$$|(\tilde{y}_w^o - \tilde{y}_w^m)/\tilde{y}_w^o|, \quad (5.10)$$

with \tilde{y}_w being fitted values from a pre-specified model obtained from the original (index o) and the modified data (index m). Index w indicates that the survey weights should be considered when fitting the model.

Note that two of these measures are automatically estimated in **sdcMicro** when an object of class `sdcMicroObj` is generated or whenever continuous key variables are modified in such an object. Thus, no user input is needed. The data utility measures are shown when printing the risk (see previous chunks) but they can also be extracted. We define our *sdcMicroObj* first and then already apply a method to anonymize the continuous key variables (Fig. 5.3).

```
df <- data.frame(mdX = mdX, mdY = mdY)
require(ggplot2)
gg <- ggplot(df, aes(x = mdX, y = mdY)) + geom_point()
gg <- gg + xlab("Mahalanobis distances - original data") +
  ylab("Mahalanobis distances - perturbed data") print(gg)
```

```
library("laeken")
data("eusilc")
sdc <- createSdcObj(eusilc,
  keyVars = c("db040", "hsize", "pb220a",
    "rb090", "pl030", "age"),
  numVars = "eqIncome",
  weightVar = "rb050",
  hhId = "db030")
sdc <- microaggregation(sdc)
```

```
get.sdcMicroObj(sdc, "utility")
```

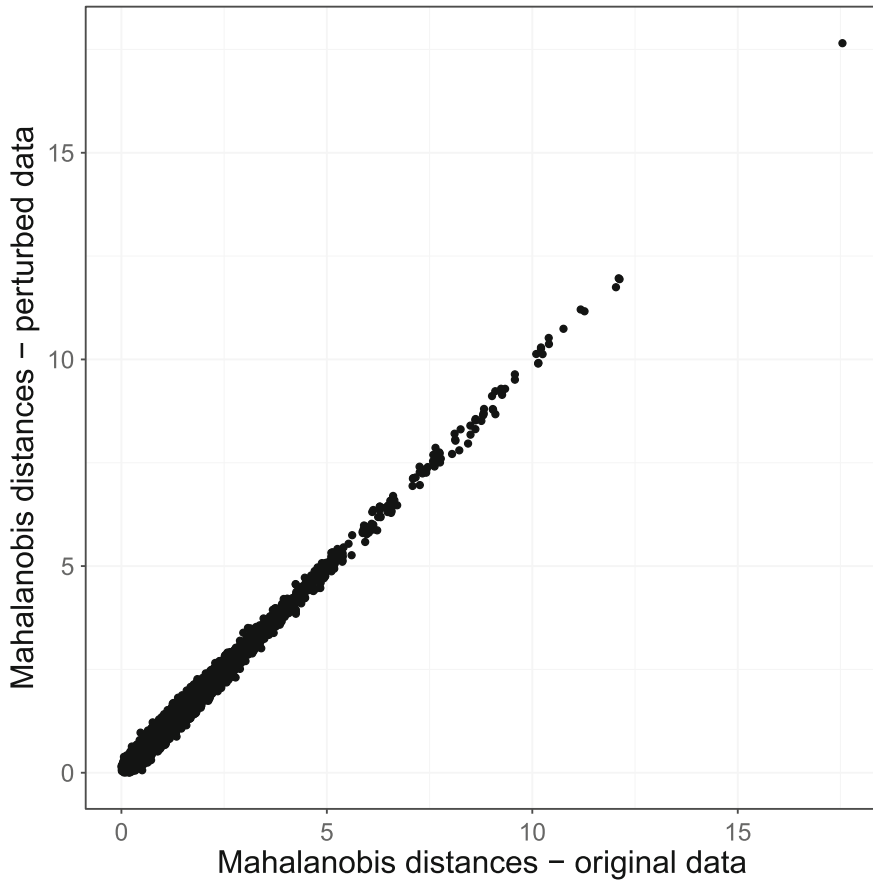


Fig. 5.3 Mahalanobis distances from original data (EU-SILC) versus the perturbed data after adding correlated noise to age and income

```
## $i11
## [1] 0.001442564
##
## $eigen
## [1] 0
```

Let's have a look at another method.

```
sdc <- undolast(sdc)
sdc <- addNoise(sdc, method = "correlated2")
get.sdcMicroObj(sdc, "utility")

## $i11
## [1] 0.1138542
##
```

```
## $eigen
## [1] 2.220446e-16
```

Since the higher the values of these measures, the lower is the data utility, we see that the data utility is lower for adding correlated noise than for microaggregation (method *mdav*). In the following code listing we see that additive noise gives much worse results.

```
sdc <- undolast(sdc)
sdc <- addNoise(sdc, method = "additive")
get.sdcMicroObj(sdc, "utility")

## $i11
## [1] 1.687381
##
## $eigen
## [1] 2.220446e-16
```

5.2.1 Element-Wise Comparisons of Mixed Scaled Variables

A distance measure that takes different scales of variables into account is based on an extension of the Gower distance (Gower 1971), which can handle distance variables of the type binary, categorical, ordered, continuous and semi-continuous. The distance between two observations is the weighted mean of the contributions of each variable, where the weight should represent the importance of the variable.

The observations of a data matrix \mathbf{X} (the original data) and the observations of the anonymized data \mathbf{Y} are collected in the rows $\mathbf{x}_i = (x_{i1}, \dots, x_{ip})^t$ and $\mathbf{y}_i = (y_{i1}, \dots, y_{ip})^t$ respectively (for $i = 1, \dots, n$ and p the number of variables).

The distance between the i -th observation \mathbf{x}_i and \mathbf{y}_i can be defined as

$$d(\mathbf{x}_i, \mathbf{y}_i) = \frac{\sum_{k=1}^p w_k \delta_{i,k}}{\sum_{k=1}^p w_k}, \quad (5.11)$$

where w_k is the weight and $\delta_{i,k}$ is the contribution of the k -th variable.

For continuous variables the absolute distance divided by the total range is used

$$\delta_{i,k} = |x_{i,k} - y_{i,k}| / r_k, \quad (5.12)$$

where $x_{i,k}$ and $y_{i,k}$ are the values of the k -th variable in the original and anonymized data of the i -th observation. r_k is the range of the k -variable in the original data set. Ordinal variables are converted to integer variables and then the absolute distance divided by the range is computed. The categories are therefore treated as if they were equidistant.

For nominal and binary variables a simple 0/1 distance is used

$$\delta_{i,k} = \begin{cases} 0 & \text{if } x_{i,k} = y_{i,k} , \\ 1 & \text{if } x_{i,k} \neq y_{i,k} . \end{cases} \quad (5.13)$$

Another special type of variables are semi-continuous variables, consisting of continuously distributed part and probability mass at one point. An example for such a variable might be an income component, which is 0 for some observations and continuously distributed in the remaining observations. The contributions for semi-continuous variables is computed as mixture of the contribution for nominal and continuous variables

$$\delta_{i,k} = \begin{cases} 0 & \text{if } x_{i,k} = s_k \wedge y_{i,k} = s_k \\ 1 & \text{if } x_{i,k} \neq s_k \wedge y_{i,k} = s_k \\ 1 & \text{if } x_{i,k} = s_k \wedge y_{i,k} \neq s_k \\ |x_{i,k} - y_{i,k}|/r_k^{(x_k)} & \text{if } x_{i,k} \neq s_k \wedge y_{i,k} \neq s_k \end{cases} , \quad (5.14)$$

where s_k is the special value for the k -th variable, e.g., the 0 in the income variable example.

As described above, all $\delta_{i,k}$ are in $[0, 1]$, as a consequence the computed distances $d(\mathbf{x}_i, \mathbf{y}_i)$ between two observations are also inside this interval.

In general, the higher the distance in Eq. (5.11), the lower the data utility.

```
library("laeken")
library("VIM")
data("eusilc")
X<-Y <- eusilc
sdc<- createSdcObj(X,
  keyVars = c("db040", "hsize", "pb220a",
    "rb090", "pl030", "age"),
  numVars = c("eqIncome"),
  pramVars = "db040",
  weightVar = "rb050", hhId = "db030")
sdc <- pram(sdc)
sdc <- microaggregation(sdc, strata="db040")
Y <- extractManipData(sdc)
gd <- VIM::gowerD(X[, c("eqIncome", "db040")], Y[, c("eqIncome", "db040")],
  numerical="eqIncome", factors="db040",
  orders = NULL, mixed=NULL, levOrders = NULL,
  weights=rep(1,2), mixed.constant = 0)
sum(abs(gd)) / nrow(eusilc) * 100
```

5.3 Entropy

The alternative option is to use an entropy function. Given c_1, c_2, \dots, c_k categories of a variable \mathbf{X}_j , the entropy E_{c_j} is defined as

$$E_{c_j} = -\frac{1}{n} \sum_{c_j \in \mathbf{X}_j} f_{c_j} \log \left(\frac{f_{c_j}}{n} \right) , \quad (5.15)$$

where f_{c_j} is the frequency of category c_j of variable \mathbf{X}_j and n the total number of observations.

This could also be used to suppress variables, e.g. the variable with the lowest value of the entropy function is chosen for further global recoding (*age* in our case, see below).

```
## entropy of key variables on original data X
entropy <- function(fk, n){
  1 / n * sum(fk * log(fk / n))
}
## for hsize
n <- nrow(eusilc)
fk <- as.numeric(table(eusilc$hsize))
entropy(fk, n)

## [1] -1.765339

## for age
entropy(as.numeric(table(eusilc$age)), n)

## [1] -4.440551

## for pb220a
entropy(as.numeric(table(eusilc$pb220a)), n)

## [1] -0.4446661
```

5.4 Propensity Score Methods

Woo et al. (2009) propose the use of a method that is based on the idea of propensity scores. Propensity scores are usually used in medical studies where the propensity score is the probability of being assigned to treatment given covariate variables. From this theory, it can be said that treatment assignment and covariates are conditionally independent given the propensity scores. The groups should have similar distributions of covariates, if they have the same distributions of propensity scores.

The idea is to merge/join the original and the perturbed data sets and then create a new index variable with ones for the original data and zeros for observations from anonymized data. A logistic regression model is then fitted using the new index variable as the response variable. Predictions from this model are then compared with the proportion of observations of the perturbed data to the original data (usually 1/2).

Woo also describes two other measures, one based on cluster analysis (evaluating the cluster sizes) and another which compares the empirical cumulative distribution function. They concentrate only on data utility measures and do not account for disclosure risk. Karr et al. (2006) proposes measures based on differences between inferences on original and perturbed data that are tailored to normally distributed data and they used also the propensity score method in Oganian and Karr (2006).

Let us have a closer look at (Woo et al. 2009) with a practical example. Again, the `eusilc` data set used is

```
eusilc$hsize <- factor(eusilc$hsize)
## create sdcMicroObj
sdc <- createSdcObj(eusilc,
  keyVars = c("db040", "hsize", "pb220a",
              "rb090", "pl030", "age"),
  pramVars = "db040",
  numVars = "eqIncome",
  weightVar = "rb050", hhId = "db030")
## apply first anonymization
sdc <- pram(sdc)
Y <- extractManipData(sdc)
## apply second and third anonymization
before <- paste(6:9)
after <- c("6-9")
sdc <- groupAndRename(sdc, var="hsize", before=before,
  after=after)
sdc <- microaggregation(sdc)
Y2 <- extractManipData(sdc)
```

The original and the perturbed data are joined together (row binding).

```
## to benchmark, original and original
Z <- rbind(eusilc, eusilc)
## original and perturbed version 1
ZA <- rbind(eusilc, Y)
## original and perturbed version 2
ZA2 <- rbind(eusilc, Y2)
```

Next the index vector indicating is created, determining if an observation belongs to the original data or the perturbed data.

```
Z$index <- ZA$index <- ZA2$index <- rep(0:1, each=nrow(eusilc))
```

We then fit the following model. For simplicity we choose a very simply model. However, we recommend to improve the model by carefully searching for the best model.

```
form <- as.formula("index ~ db040 + hsize + pb220a +
  rb090 + pl030 + eqIncome")
res <- glm(form, data=Z, family = binomial())
resA <- glm(form, data=ZA, family = binomial())
resA2 <- glm(form, data=ZA2, family = binomial())
```


First we can look on how many one's and zero's are predicted. If the model fits perfect, the ratio between one's and zero's should be one.

```
ps <- function(mod){
  p <- predict(mod, type="response")
  t1 <- as.numeric(table(p < 0.5))
  return(t1[1] / t1[2])
}
p <- ps(res)
p

## [1] 0.8302343

pa <- ps(resA)
pa

## [1] 1.19668

pa2 <- ps(resA2)
pa2

## [1] 1.185971
```

Ideally `ps(res)` should be 1 or very close to one. In this example we see that the more variables are modified, the higher the ratio between (0/1) in the predictions. Note that the results may get better when selecting the best possible model by possible also considering also interaction terms.

Woo et al. (2009) compares also the distribution of the predicted values by

$$UP = \frac{1}{n_m} \sum_{i=1}^{n_m} (p_i - c)^2, \quad (5.16)$$

where n_m is the number of observations in the merged data set, p_i is the estimated probabilities being in group 1 (original data) or group 2 (perturbed data). c is usually determined as 0.5 whenever the perturbed data have the same amount of observations as the original data. If UP is close to zero, the data utility is high. The worse case is if $UP \sim 1/4$, where the two data sets are completely distinguishable (see also Woo et al. 2009).

```
1 / nrow(Z) * sum((predict(res, type="response") - 0.5)^2)

## [1] 9.055325e-30

1 / nrow(Z) * sum((predict(resA, type="response") - 0.5)^2)
```

```
## [1] 1.097417e-05

1 / nrow(Z) * sum((predict(resA2, type="response") - 0.5)^2)

## [1] 0.01095378
```

We see that the grouping of the large households and the microaggregation of the equivalized income has a large effect on *UP*, while only to apply PRAM on variable *db040* (federal states) has no big effect and the data utility is reported to be high.

Exercises:

Question 5.1 As done already before, produce an object of class *sdcMicroObj* using the data set *eusilc*. Apply some recodings and local suppression for categorical key variables, and apply a method to anonymize the continuous key variable *eqIncome*. Cross tabulate *hsize*, *rb090* (gender) and *pb220a* (citizenship) and compare this table with the result from the raw survey data.

Question 5.2 Plot the equivalized income from the raw survey data against the modified equivalized income. Split the plot by conditioning on *p1030*.

5.5 Quality Indicators

Although, in practice, it is not possible to create a file with the exact same structure as the original file after applying SDC methods, an important goal of SDC should be to minimize the difference in the statistical properties of the perturbed data and the original data. It is a fact that not every estimate can be preserved. Therefore the aim is to preserve the most important estimates, i.e. instead of calculating generally defined utility measures as in the previous sections, evaluation on context/data dependent indicators is proposed. Such an approach to measuring data utility is based on benchmarking indicators (Ichim and Franconi 2010; Templ 2011a, 2015).

5.5.1 General Procedure

The first step in quality assessment is to evaluate what users of the underlying data are analyzing and then try to determine the most important estimates, or *benchmarking indicators* (see, e.g., Templ 2011a, b). Special emphasis should be put on benchmarking indicators that take into account the most important variables of the micro data set. Indicators that refer to the most sensitive variables within the microdata should also be calculated. The general procedure is quite simple and can be described in the following steps:

- selection of a set of (benchmarking) indicators;
- choice of a set of criteria as to how to compare the indicators;
- calculation of all benchmarking indicators of the original micro data;
- calculation of the benchmarking indicators on the protected micro data set;
- comparison of statistical properties such as point estimates, variances or overlaps in confidence intervals for each benchmarking indicator;
- assessment as to whether the data utility of the protected micro data set is good enough to be used by researchers.

If the main indicators calculated from the protected data differ significantly from those estimated from the original data set, the SDC procedure should be restarted. It is possible to either change some parameters of the applied methods or start from scratch and completely change the choice of SDC methods. The benchmarking indicator approach is usually applied to assess the impact of SDC methods on continuous variables. But it is also applicable to categorical variables. In addition, the benchmarking indicators approach can be applied to subsets of the data. In this case, benchmarking indicators are evaluated for each of the subsets and the results are evaluated by reviewing differences between indicators for original and modified data within each subset.

In addition, it is interesting to evaluate the set of benchmarking indicators not only for the entire data set but also independently for subsets of the data. In this case, the micro-data are partitioned into a set of h groups. The evaluation of benchmarking indicators is then performed for each of the groups and the results are evaluated by reviewing differences between indicators for original and modified data in each group.

5.5.2 Differences in Point Estimates

For an object of class *sdcMicroObj*, the slot `additionalResults` can be used to store additional results such as self-defined indicators in an object of class *sdcMicroObj*. This is sometimes useful because all results corresponding to an anonymization process are stored in one single object.

For the EU-SILC data, one of the most important indicators is the Gini coefficient. Therefore, this indicator is used to evaluate the quality of the anonymized data set. Given a vector x_1, \dots, x_n with sample weights w_1, \dots, w_n the Gini coefficient can be estimated by

$$\widehat{Gini} = \frac{2 \sum_{i=1}^n \left(w_i x_i \sum_{j=1}^i w_j \right) - \sum_{i=1}^n w_i^2 x_i}{\left(\sum_{i=1}^n w_i \right) \sum_{i=1}^n w_i x_i} - 1 \quad .$$

The Gini coefficient takes on values between 0 and 1. A value of 0 stand for perfect equality, meaning that every data point in the sample has the same value or every

individual has the same income or volume of a certain good. With a value of 1 the Gini coefficient would indicate perfect inequality, meaning that all but one data point are equal to zero or that one individual has all the income or volume of certain good.

The following line of code estimates the Gini coefficient and stores the result within the `sdcMicroObj` `sdc`.

```
sdc@additionalResults$gini <- gini(inc = "eqIncome",
  weights = "rb050",
  breakdown = "db040",
  data = extractManipData(sdc)$valueByStratum$value
```

This result can then be simple compared for each category of `db040` (federal state) with the result of the original data to get an indicator of data utility expressed as relative absolute errors (in percentages):

```
res <- gini(inc = "eqIncome",
  weights = "rb050",
  breakdown = "db040",
  data = eusilc)$valueByStratum$value
100*abs((res - sdc@additionalResults$gini)/res)

## [1] 2.173689 2.386133 0.270414 1.289537 1.055491 1.080695
## [7] 1.925351 1.468171 4.996778
```

It can be seen that in few countries the change in the Gini coefficient is about 3% of its absolute value, in other regions, the differences are smaller.

5.5.3 Differences in Variances and MSE

However, it is also recommended to not look only on point estimates but also estimate the variances and compare the variances obtained from the original and the perturbed data.

The estimation of variance is not always an easy task when working with complex survey collected based on complex sample designs. A calibrated bootstrap (Alfons and Templ 2013) to estimate the variances is applied in the following.

```
res <- gini(inc = "eqIncome",
  weights = "rb050",
  breakdown = "db040",
  data = eusilc)
resVar <- variance("eqIncome", weights = "rb050",
  design = "db040", breakdown = "db040",
  data = eusilc, indicator = res, R = 50,
  X = calibVars(eusilc$db040), seed = 123)
res <- resVar$valueByStratum$value
resVar <- resVar$varByStratum$var
eusilcA <- extractManipData(sdc)
resA <- gini(inc = "eqIncome",
```

```

weights = "rb050",
breakdown = "db040",
data = eusilcA)
resVarA <- variance("eqIncome", weights = "rb050",
  design = "db040", breakdown = "db040",
  data = eusilcA, indicator = resA, R = 50,
  X = calibVars(eusilc$db040), seed = 123)
resA <- resVarA$valueByStratum$value
resVarA <- resVarA$varByStratum$var
100*abs((res - resA) / res)

## [1] 2.173689 2.386133 0.270414 1.289537 1.055491 1.080695
## [7] 1.925351 1.468171 4.996778

100*abs((resVar - resVarA) / resVar)

## [1] 35.123843 13.152475 1.601867 26.820150 45.971562
## [6] 11.283797 28.450694 18.381783 31.100778

```

A mean squared error (MSE) kind of measure would then be the differences in the point estimates (bias) to the square plus the differences in variances.

```

(res - resA)^2 + (resVar - resVarA)

## [1] 1.164910938 0.315197915 0.007902761 0.226787385
## [5] -0.013603540 0.131555209 0.3100s45368 0.218618271
## [9] 1.867190441

```

5.5.4 Overlap in Confidence Intervals

However, also the overlap of confidence intervals is of interest.

Let $[l_{c_j}^{(X)}, u_{c_j}^{(X)}]$ the 95%-confidence interval for a given parameter in the original data \mathbf{X} in strata c_j and let $[l_{c_j}^{(Y)}, u_{c_j}^{(Y)}]$ the corresponding interval in the anonymized data. The intersection of the two intervals is denoted by

$$[li_{c_j}, ui_{c_j}] = [l_{c_j}^{(X)}, u_{c_j}^{(X)}] \cap [l_{c_j}^{(Y)}, u_{c_j}^{(Y)}] \quad (5.17)$$

The utility measure is then given by

$$UC = \frac{1}{2K} \sum_{j=1}^K \left(\frac{ui_{c_j} - li_{c_j}}{u_{c_j}^{(X)} - l_{c_j}^{(X)}} + \frac{ui_{c_j} - li_{c_j}}{u_{c_j}^{(Y)} - l_{c_j}^{(Y)}} \right), \quad (5.18)$$

Table 5.1 Lower (l) and upper (u) limits of the confidence intervals for the GDP for each category of *education* for the original data, for recoded and local suppressed data, for recoded, local suppressed and noise addition to data, for prammed and microaggregated data and for recoded, local suppressed and shuffled data

Data	ISCED 0-1	ISCED 2	ISCED 3-4	ISCED 5A	ISCED 5B
original (l)	0.15938	0.12102	0.22572	0.29568	0.21744
original (u)	0.26525	0.15023	0.23944	0.35010	0.25835
rec+ls+ma (l)	0.16123	0.12144	0.22624	0.28891	0.21290
rec+ls+ma (u)	0.27062	0.15211	0.23970	0.34381	0.25904
rec+ls+noise (l)	0.17012	0.12106	0.22399	0.29135	0.21152
rec+ls+noise (u)	0.27011	0.15172	0.23776	0.34551	0.25805
pram+ma (l)	0.17682	0.12200	0.22554	0.29064	0.21946
pram+ma (u)	0.27230	0.15065	0.24197	0.33822	0.26172
rec+ls+shuffle (l)	-0.01865	0.09365	0.18510	0.19294	0.19859
rec+ls+shuffle (u)	0.24584	0.12496	0.20950	0.25071	0.26183

Table 5.2 Coverage rates for confidence intervals of the gender pay gap in each educational sector between the original and perturbed data

Data	ISCED 0 and 1	ISCED 2	ISCED 3 and 4	ISCED 5A	ISCED 5B
rec+ls+ma	98.25	98.55	96.21	88.45	88.65
rec+ls+noise	89.85	99.86	87.81	91.58	99.26
pram+ma	83.52	96.63	83.45	78.18	95.08
rec+ls+shuffle	81.67	13.51	0.00	0.00	64.67

with K the number of strata in which the confidence intervals are being estimated. When the intervals are identical in both \mathbf{X} and \mathbf{X} , $UC = 1$. In the other extreme case, when the intervals do not overlap at all, $UC = 0$.

As an example, the upper and lower confidence intervals for the GDP in domain *education* are given in Table 5.1. It is easy to see that the length of the confidence intervals are shorter for category ISCED 3-4 and largest for ISCED 0-1.

Again, the shuffling method does not seem to be able to give approximately the same confidence intervals.

A clearer picture is supported by Table 5.2 where the overlap of the confidence intervals for the GDP—estimated from the perturbed and the original data—is reported.

The coverage rates are relatively high for all methods except recoding+local suppression+shuffling. Differences in some categories are visible when comparing the other methods whereas no clear ranking of them in terms of quality can be made.

The coverage rates for the gender pay gap in domain *age* (Table 5.3) are similar. Mostly the recoding+local suppression+microaggregation methods performs slightly better than recoding+local suppression+adding noise and pram+microaggregation.

Table 5.3 Coverage rates for confidence intervals of the GDP in each age class between the original and perturbed data

Data	(0,19]	(19,29]	(29,39]	(39,49]	(49,59]	(59,120]
rec+ls+ma	98.81	76.40	99.28	82.41	95.82	91.45
rec+ls+noise	94.90	80.27	94.31	89.60	89.70	96.76
pram+ma	84.26	88.92	95.02	88.55	92.58	86.94
rec+ls+shuffle	0.00	32.75	0.00	0.00	0.00	0.00

Table 5.4 Coverage rates for confidence intervals of the Gini indices in each age \times gender domain between the original and perturbed data

Data	(0,19]:f	(0,19]:m	(19,29]:f	(19,29]:m	(29,39]:f	(29,39]:m
rec+ls+ma	93.64	81.66	96.83	94.93	89.24	95.63
rec+ls+noise	52.71	0.00	22.12	37.18	63.09	87.92
pram+ma	88.29	82.49	88.39	93.05	85.36	94.50
rec+ls+shuffle	82.61	0.00	0.00	0.00	20.38	0.00
	(39,49]:f	(39,49]:m	(49,59]:f	(49,59]:m	(59,120]:f	(59,120]:m
rec+ls+ma	84.69	75.33	99.21	94.59	95.40	92.22
rec+ls+noise	88.49	83.07	80.52	89.70	93.94	96.03
pram+ma	97.89	85.00	96.78	82.25	88.31	94.94
rec+ls+shuffle	12.55	55.93	0.00	0.00	0.00	0.00

However, a completely different picture is seen for the absolute relative bias of the Gini index in Table 5.4. Recoding+local suppression+microaggregation outperforms all other methods. PRAM+microaggregation also gives acceptable results but recoding+local suppression+adding noise gives low coverage rates for age classes below 29 years. Shuffling results in the highest biased estimates.

5.5.5 Differences in Model Estimates

As already mentioned, to compare the regression coefficients of original and anonymized data sets, the same categories in the explanatory variables of the model must be present. Thus the recoded 12 categories of *economic activity* are used also for the original data set, keeping in mind that this means a certain kind of information loss.

In Table 5.5 the regression coefficients for the original and the anonymized data sets are shown.

Recoding+local suppression+microaggregation again performs best and the confidence intervals obtained from the anonymized data cover the confidence intervals obtained from the original data almost always completely. Almost as good is the quality of data anonymized by recoding+local suppression+adding correlated noise.

Table 5.5 Regression coefficients

	original	rec+ls+ma	rec+ls+noise	pram+ma	rec+ls+shuffle
(Intercept)	1.50454	1.52627	1.51374	1.40474	1.63726
Sexmale	0.20478	0.20484	0.20433	0.20970	0.19733
age(19,29]	0.57210	0.57190	0.58560	0.57659	0.76536
age(29,39]	0.73750	0.73745	0.75186	0.74388	0.91469
age(39,49]	0.81758	0.81746	0.83260	0.82634	0.96199
age(49,59]	0.85660	0.85597	0.87072	0.86754	0.89338
age(59,120]	0.81553	0.81067	0.82604	0.82169	0.49264
educationISCED 2	0.03692	0.02006	0.01011	0.03834	-0.25102
educationISCED 3 and 4	0.28314	0.26646	0.25737	0.28667	-0.16874
educationISCED 5A	0.73406	0.71508	0.70647	0.74198	0.09813
educationISCED 5B	0.44484	0.42802	0.41959	0.45337	-0.01251
LocationAT2	-0.07516	-0.07523	-0.07528	-0.06368	-0.00673
LocationAT3	-0.01230	-0.01207	-0.01132	-0.00900	-0.00098
NACE1D- Manufacturing	-0.05542	-0.06029	-0.05441	0.01740	-0.01600
NACE1E- Electricity	0.09709	0.09018	0.09264	0.12244	-0.02588
NACE1F- Construction	-0.12280	-0.12891	-0.12260	-0.03775	-0.01806
NACE1G-Trade	-0.18916	-0.19422	-0.18848	-0.09872	-0.02576
NACE1H-Hotels	-0.37478	-0.37962	-0.37589	-0.24398	-0.02269
NACE1I- Transport	-0.17130	-0.17632	-0.17061	-0.07943	-0.00939
NACE1J- FinancInt	0.14921	0.14532	0.15055	0.19273	-0.01993
NACE1K- RealEstate	-0.13433	-0.13901	-0.13517	-0.05156	-0.02072
NACE1M- Education	-0.16289	-0.16650	-0.16300	-0.07845	-0.02505
NACE1N-Health	-0.11299	-0.11734	-0.11360	-0.02939	-0.01838
NACE1O-Other	-0.19113	-0.19585	-0.19353	-0.10283	-0.01054

The results from invariant pram+microaggregation are good for all coefficients except those related to *economic activity*. This is not surprising since this variable was one of the variables which was changed using PRAM. Some few coefficients are well preserved from the recoding+local suppression+shuffling anonymized data, but others are not. The reason is that even if the distribution of the continuous shuffled

variables are well preserved, the relation to other variables that are not included in the shuffling model might be not preserved. A better model would probably lead to better results.

Excercises:

Question 5.3 Undo-the last step of the anonymization procedure above, apply a different disclosure limitation technique and re-calculate the Gini coefficients. Compare it with previous results.

References

- Alfons, A., & Templ, M. (2013). Estimation of social exclusion indicators from complex surveys: The R package Laeken. *Journal of Statistical Software*, 54(15), 1–25.
- Domingo-Ferrer, J. (2009). Information loss measures. In L. Liu & M. Tamerzsu (Eds.), *Encyclopedia of database systems* (pp. 1499–1501). Springer US. ISBN 978-0-387-35544-3.
- Egozcue, J. J., Pawlovsky, V., Templ, M., & Hron, K. (2015). Independence in contingency tables using simplicial geometry. *Communications in Statistics*, 44(18), 3978–3996.
- Gower, J. C. (1971). A general coefficient of similarity and some of its properties. *Biometrics*, 27, 857–871.
- Hartigan, J. A., & Kleiner, B. (1981). Mosaics for contingency tables. In W. F. Eddy (Ed.), *Computer Science and Statistics: Proceedings of the 13th Symposium on the Interface* (pp. 268–273). New York: Springer.
- Hofmann, H. (2003). Constructing and reading mosaicplots. *Computational Statistics & Data Analysis*, 43(4), 565–580.
- Hundepool, A., Domingo-Ferrer, J., Franconi, L., Giessing, S., Nordholt, E. S., Spicer, K., et al. (2012). *Statistical disclosure control*. Wiley Series in Survey Methodology. Wiley. ISBN 9781118348222. <https://books.google.at/books?id=BGa3zKkFm9oC>.
- Ichim, D., & Franconi, L. (2010). Strategies to achieve SDC harmonisation at european level: Multiple countries, multiple files, multiple surveys. In *Privacy in statistical databases* (pp. 284–296).
- Karr, A. F., Kohnen, C. N., Oganian, A., Reiter, J. P., & Sanil, A. P. (2006). A framework for evaluating the utility of data altered to protect confidentiality. *The American Statistician*, 60(3), 224–232.
- Meyer, D., Zeileis, A., & Hornik, K. (2006). The strucplot framework: Visualizing mult-way contingency tables with vcd. *Journal of Statistical Software*, 17(3), 1–48.
- Oganian, A., & Karr, A. F. (2006). Combinations of SDC methods for microdata protection. In J. Domingo-Ferrer & L. Franconi (Eds.), *Privacy in statistical databases* (Vol. 4302, pp. 102–113), Lecture Notes in Computer Science Heidelberg: Springer.
- Rousseeuw, P. J., & Van Driessen, K. (1998). A fast algorithm for the minimum covariance determinant estimator. *Technometrics*, 41, 212–223.
- Templ, M. (2011a). Estimators and model predictions from the structural earnings survey for benchmarking statistical disclosure methods. Research Report CS-2011-4, Department of Statistics and Probability Theory, Vienna University of Technology.
- Templ, M. (2011b). Comparison of perturbation methods based on pre-defined quality indicators. In *Joint UNECE/Eurostat work session on statistical data confidentiality*. Tarragona, Spain. Invited paper.
- Templ, M. (2015, accepted for publication). Quality indicators for statistical disclosure methods: A case study on the structural earnings survey. *Journal of Official Statistics*.

- Templ, M., Hron, K., Filzmoser, P. (2011) *robCompositions: An R-package for robust statistical analysis of compositional data* (pp. 341–355). Wiley. ISBN 9781119976462. <http://dx.doi.org/10.1002/9781119976462.ch25>.
- Templ, M., Meindl, B., Kowarik, A., & Dupriez, O. (2017, accepted for publication in December 2015). Simulation of synthetic complex data: The R-package simPop. *Journal of Statistical Software*, 1–38.
- Winkler, W. E. (1998). Re-identification methods for evaluating the confidentiality of analytically valid microdata. *Research in Official Statistics*, 1, 50–69.
- Woo, M., Reiter, J. P., Oganian, A., & Karr, A. F. (2009). Global measures of data utility for microdata masked for disclosure limitation. *Journal of Privacy and Confidentiality*, 1(1), 111–124.
- Yancey, W. E., Winkler, W. E., & Creecy, R. H. (2002). Disclosure risk assessment in perturbative microdata protection. In *Inference control in statistical databases*. Lecture Notes in Computer Science (pp. 49–60). Springer.