



Statistics Netherlands

Division of Methodology and Quality
Department of Methodology, The Hague

Methods Series:
Statistical disclosure control

Anco Hundepool and Peter-Paul de Wolf

18 July 2011

Version management

Version history				
Version	Date	Description	Authors	Reviewers
0.5	02-04-2007	First Dutch version	Anco Hundepool Peter-Paul de Wolf	John Schalen Eric Schulte Nordholt
1.1	23-01-2008	Minor modifications to layout	Anco Hundepool Peter-Paul de Wolf	
1.2	15-02-2010	Chapters on frequency tables and analysis results added	Anco Hundepool Peter-Paul de Wolf	John Schalen Eric Schulte Nordholt
1.2E	19-07-2011	First English version	Anco Hundepool Peter-Paul de Wolf	

Table of Contents

1.	Statistical Disclosure Control	4
1.1	General description.....	4
1.2	Scope and relationship with other themes	5
1.3	Place in the statistical process	5
1.4	Definitions	5
2.	Statistical Disclosure Control of Microdata.....	7
2.1	General description and reading guide	7
2.2	Scope and relationship with other themes and subthemes.....	8
2.3	Global recoding	8
2.4	Local suppression	10
2.5	Top-coding	12
2.6	Adding noise to weights	13
2.7	PRAM.....	15
2.8	Conclusion.....	19
3.	Statistical Disclosure Control of Quantitative Tables.....	20
3.1	General description and reading guide	20
3.2	Scope and relationship with other themes and subthemes.....	21
3.3	<i>P</i> % rule	22
3.4	Table restructuring.....	24
3.5	Cell suppression	25
3.6	Additive rounding.....	31
3.7	Conclusion.....	34
4.	Statistical Disclosure Control of Frequency Tables.....	35
4.1	General description and reading guide	35
4.2	Scope and relationship with other themes and subthemes.....	36
4.3	Temporary standardisation of a frequency table	37
4.4	Table restructuring.....	39
4.5	Suppression	40
4.6	Additive rounding.....	43
5.	Statistical Disclosure Control of Analysis Results	46
5.1	General description and reading guide	46
5.2	Scope and relationship with other themes and subthemes.....	47
5.3	Disclosure control of analysis results	47
6.	References.....	49

1. Statistical Disclosure Control

1.1 General description

When publishing statistical information, Statistics Netherlands must achieve a balance between the interests of its data suppliers and the interests of its users. On the one hand, the Statistics Netherlands users want as much information as possible, and as detailed as possible. On the other, the data suppliers (people and companies, as well as the registration holders and the Dutch Data Protection Authority) require that their privacy is guaranteed. *Private lives and public policies: confidentiality and accessibility of government statistics* (Duncan et al., 1993) is the very relevant title of an American book about this problem.

What Statistics Netherlands may and may not publish follows from its statistical disclosure control policy, as set down in the Statistical Disclosure Control Handbook (Hundepool et al., 2006). Here, statistical disclosure control means preventing that content-related conclusions about recognisable units are made based on published or otherwise available Statistics Netherlands data.

It must not be possible to make such conclusions based on the statistical publications from Statistics Netherlands (StatLine tables, web articles, press releases, scientific articles). However, also if Statistics Netherlands makes microdata available for scientific analysis, this basic rule of statistics must remain in force.

The Statistical Disclosure Control Handbook describes the policy and other rules that individual publications must comply with. However, not all publications satisfy these rules in and of themselves. On the contrary, frequently a publication will have to be “protected”. Different methods are available to protect microdata, table data and analysis results. The theme of Statistical Disclosure Control in the Methods Series can thus be broken down into a number of subthemes:

- Statistical disclosure control methods for microdata,
- Statistical disclosure control methods for quantitative tables,
- Statistical disclosure control methods for frequency tables,
- Statistical disclosure control methods for analysis results.

The conflicting interests of privacy protection and information retention play an ongoing role in statistical disclosure control. When using the different methods for statistical disclosure control, these two aspects must be taken into account. The statistical disclosure control policy of Statistics Netherlands sets down a minimum level of protection. The real skill of the person protecting the data is to use different disclosure control methods in such a way that the minimum required level of protection is achieved and that the information loss is as small as possible. This will be different in every situation, as the concept of “information loss” can have different meanings for the different users of the Statistics Netherlands data.

The methods mentioned in this theme will each be explained separately. However, in practice, for each situation, multiple methods will often be used at the same time to create “safe” publications. The interaction between the different methods when used simultaneously will not be described in this Methods Series.

1.2 Scope and relationship with other themes

The statistical disclosure control policy will not be described here. That policy is laid down in the Statistical Disclosure Control Handbook referred to above. However, various available methods will be described that can be used to apply that policy to Statistics Netherlands publications. Some of the methods described are actively used at Statistics Netherlands, while other methods are, at present, only used at statistical bureaus abroad.

When applying statistical disclosure control methods, both the level of protection and the information loss of the publications must be examined. Since the concept of “information” is subjective and therefore can be defined differently by each user (even in a single publication), it is not possible to prescribe a specific method for each specific situation. The methods will therefore be described along with their advantages and disadvantages, along with their effects on the level of protection and information loss. A Statistics Netherlands staff member who is in charge of the statistical disclosure control of a publication (in whatever form), will subsequently have to choose the most suitable method for the publication in question.

1.3 Place in the statistical process

Statistical disclosure control traditionally takes place at the end of the statistical process: statistical disclosure control is applied immediately before publication (in whatever form). Ideally, account should be taken during the entire statistical process of the fact that, ultimately, the publication will have to satisfy the statistical disclosure control policy. However, measures can also be taken at the start of the statistical process, such as formulating the cover letter for participation in a survey (“informed consent”).

The concept of statistical disclosure control therefore plays a role during the entire statistical process. However, the specific methods as described in this document are only used at the end of the statistical process, immediately before publication.

1.4 Definitions

Term	Definition
μ -ARGUS	Software for the statistical disclosure control of microdata files
τ -ARGUS	Software for the statistical disclosure control of tables
Disclosure	The obtaining of information from statistical data about a recognisable specific person, household, company or institution
Identifying variable	Variable of which the value can contribute to the identification of a specific person, household, company or institution

Primary risky cell	Cell in a table that does not satisfy the disclosure control rules
Secondary risky cell	Cell in a table that does satisfy the disclosure control rules, but which must be suppressed to protect the primary risky cells
Structural zero cell	A cell for which it is generally known that, logically, this cell <i>cannot</i> have a contribution

An extensive glossary for statistical disclosure control can be found at:

<http://neon.vb.cbs.nl/casc/Glossary.htm>.

2. Statistical Disclosure Control of Microdata

2.1 General description and reading guide

2.1.1 General description

The statistical disclosure control of microdata refers to the creation of microdata that complies with the disclosure control policy of Statistics Netherlands and which may be released as such by Statistics Netherlands. This therefore expressly does not include microdata files that remain at Statistics Netherlands, including files for onsite and remote access. The disclosure control policy for microdata is set down in chapter 3 of the Statistical Disclosure Control Handbook (Hundepool et al., 2006).

The methods described in this subtheme are used to create protected microdata files. The extent to which the methods are used (or how strictly they are applied) depends partly on the type of file that is going to be released. This is described in detail in the Statistical Disclosure Control Handbook, where a distinction is made between Public Use Files and Microdata Files under Contract.

The methods described in this subtheme are easy to apply with the μ -ARGUS software package. This package was developed by DMV in a European context.

2.1.2 Reading guide

As the first step in the statistical disclosure control of microdata, it will have to be determined whether disclosure is possible: is there any information about individual respondents in the microdata that may not be disclosed? This “sensitive” information usually concerns respondents that can be recognised as unique or rare cases in the microdata. Such respondents must be protected.

Several of the disclosure control methods that we will discuss here can be applied to categorical variables: global recoding (section 2.3) and PRAM (section 2.7). Top (and bottom) coding is mainly intended for continuous variables; see section 2.5. Local suppression (section 2.4) can be used for both categorical and continuous variables. There is also the possibility of adding noise to raising weights; see section 2.6.

Which method or combination of methods will ultimately be used in a specific situation cannot be determined in advance. The department responsible for the construction of the microdata file is also responsible for adequate statistical protection. When selecting the method or methods to be used, two competing aspects must be taken into account:

- disclosure risk
- information loss.

In general, it can be said that reducing the disclosure risk will lead to increased information loss. The converse is also true: the smaller the information loss, the larger the disclosure risk. In certain cases, an assessment will have to be made, in which, at a minimum, the rules in the Statistical Disclosure Control Handbook (Hundepool et al., 2006) will always have to be satisfied.

2.2 Scope and relationship with other themes and subthemes

The methods that will be described in this subtheme are applied directly on the microdata itself. As a result, different levels of protection can arise, which correspond to those of the Public Use Files or the Microdata Files under Contract.

Users of unprotected files (including files for onsite and remote access) and Microdata Files under Contract can generate output that does not necessarily satisfy the disclosure control policy of Statistics Netherlands. In these situations, other methods will have to be used to protect the output. For such methods, see the subtheme “Statistical Disclosure Control of Analysis Results”.

2.3 Global recoding

2.3.1 Short description

In the statistical disclosure control of microdata files that are released by Statistics Netherlands, we mainly examine the variables that can potentially be used to identify a respondent. These types of variables are called *identifying* variables. Identifying variables are generally categorical variables. Combinations of categories of identifying variables tend to lead to unique or rare people. Consider, for example, “Mayor in Amsterdam” (unique) or “Female neurosurgeon older than 55 years of age from Staphorst” (rare). The rules for Microdata Files under Contract (see chapter 3 from the Statistical Disclosure Control Handbook (Hundepool et al., 2006)) state that such combinations must occur sufficiently often in the target population.

By combining categories of identifying variables, rare combinations can be made less rare.

2.3.2 Applicability

In the disclosure control of microdata files that are released by Statistics Netherlands, certain combinations of identifying variables must occur sufficiently often in the population. In particular, if an identifying variable is present in a very detailed form in the file, global recoding can in many cases be used to sufficiently protect the file, while the information loss remains limited.

For some researchers, however, global recoding will remove too much detail, as a result of which they will no longer be able to perform their analyses. It is therefore the task of the Statistics Netherlands staff member who is charged with the statistical disclosure control of the file to assess whether global recoding is a suitable protection method for the case in question.

However global recoding does not have to be limited to identifying variables. Non-identifying variables can also be globally recoded, as long as they are categorical variables. In such an application, with respect to the possible identification of a respondent, only less detailed (and therefore probably generally known) information would be disclosed.

2.3.3 Detailed description

Global recoding involves the adaptation of the code list of an identifying variable. If the variable is hierarchical (for example, region), an obvious aspect of the recoding is to delete some detail levels. For example, in a recoding of the City/Town variable, all cities could be replaced by the associated province.

After the code list of a variable is adapted, for *each* record, the score on that variable is adapted to the new code list. This is therefore done not only for the risky records, but also for the safe records.

2.3.4 Example

Figure 1 shows several records from a fictitious microdata file. The records are numbered for easy reference.

	Occupation	City/Town	Gender	Education	...
1	Mayor	Amsterdam	Man	High	...
2	Fisherman	Urk	Man	Low	...
3	Teacher	Amsterdam	Woman	High	...
4	Plumber	Papendrecht	Man	Medium	...

Figure 1: Several records from a fictitious microdata file

The mayor from Amsterdam is, obviously, unique. The variable “City/Town” is now globally recoded by replacing the city names by the associated province. This generates the records as shown in Figure 2.

	Occupation	City/Town	Gender	Education	...
1	Mayor	Noord-Holland	Man	High	...
2	Fisherman	Flevoland	Man	Low	...
3	Teacher	Noord-Holland	Woman	High	...
4	Plumber	Zuid-Holland	Man	Medium	...

Figure 2: Records from Figure 1 after global recoding of “City”

Now, in record 1, the mayor is no longer unique. Because the recoding is applied *globally*, the city/town variable in the safe records 2 to 4 is also adapted.

2.4 Local suppression

2.4.1 Short description

In the statistical disclosure control of microdata files that are released by Statistics Netherlands, we mainly examine the variables that can potentially be used to identify a respondent. These types of variables are called *identifying* variables. Identifying variables are generally categorical variables. Combinations of categories of identifying variables tend to lead to unique or rare people. Consider, for example, “Mayor in Amsterdam” (unique) or “Female neurosurgeon older than 55 years of age from Staphorst” (rare). The rules for Microdata Files under Contract (see chapter 3 from the Statistical Disclosure Control Handbook (Hundepool et al., 2006)) state that such combinations must occur sufficiently often in the population.

In local suppression, the score on at least one of the variables in a combination that occurs insufficiently often in the target population is suppressed (or it is assigned the score of “Unknown”). As a result, the combination of the remaining variables describes a potentially larger group in the target population.

2.4.2 Applicability

In the disclosure control of microdata files that are released by Statistics Netherlands, certain combinations of identifying variables must occur sufficiently often in the target population. In particular, if an identifying variable occurs in a very detailed state in the file, oftentimes local suppression can be used to sufficiently protect the file, while the information loss remains limited.

Local suppression is often used as the final disclosure control method. At this point, most of the protection has already been provided by other methods, and local suppression is used to protect the last risky records.

Local suppression leads to missing values in the file. The way in which these missing values are selected, however, is certainly not random: the goal is to protect records that belong to small, identifiable groups. The effect of these missing values on the analyses to be conducted is different from the effect of missing values as a result of non-response.

For that matter, local suppression does not have to be limited to identifying variables. Non-identifying variables can also be locally suppressed. In the event that a respondent is identified, this ensures that no sensitive information would be disclosed.

2.4.3 Detailed description

In local suppression, the value of an identifying or other variable is set to “Unknown”. According to the disclosure control rules from the Statistical Disclosure Control Handbook (Hundepool et al., 2006), combinations of identifying variables must occur sufficiently often in the target population. By suppressing the score for at least one variable from such a combination, a lower dimensional combination is, in

fact, created. The result of this is that the combination will potentially describe a larger group of respondents in the target population.

Local suppression is only applied to risky records. It is possible that multiple risky combinations of identifying variables occur in a single record. By suppressing the right variable or variables in an intelligent manner, multiple risky combinations can sometimes be protected simultaneously.

If a microdata file has multiple records of people from the same household, this must be taken into account in local suppression. Such records may contain so-called household variables. These are variables for which each member of the household has the same score, for example, household income, household size and city/town. If a risky combination with a household variable occurs for at least one person from a household and this household variable is locally suppressed, then this variable must be suppressed for all the people in that household. In that case, values may therefore also be suppressed in safe records.

When selecting the variable that is going to be suppressed from a rare combination of scores on identifying variables, this choice is, in principle, free. However, two options are possible in μ -ARGUS.

First, the user can indicate, by assigning weights to variables, the extent to which the suppression of the score on that variable is desired (or not). μ -ARGUS then chooses to suppress those variables for which the sum of the weights is as small as possible. Consequently, it is possible, for example, to refrain (to a certain extent) from locally suppressing those variables that have already been adapted through other protection methods.

In the second option, μ -ARGUS uses a type of entropy argument to select the variable or variables to be suppressed. Each variable is then assigned the following weight:

$$w_X = -\sum_{i=1}^{K_X} \frac{f_X(i)}{n} \log \frac{f_X(i)}{n}, \quad (2.4.1)$$

where K_X is the number of categories of variable X , n the number of records in the microdata file and $f_X(i)$ the number of records with score i on variable X . As a result, variables with larger numbers of categories are suppressed less frequently than variables with only a few categories.

2.4.4 Example

Figure 3 shows some records from a fictitious microdata file. The records are numbered for easy reference.

	Occupation	City/Town	Gender	Education	...
1	Mayor	Amsterdam	Man	High	...
2	Fisherman	Urk	Man	Low	...
3	Teacher	Amsterdam	Woman	High	...
4	Plumber	Papendrecht	Man	Medium	...

Figure 3: Some records from a fictitious microdata file

The mayor from Amsterdam is, obviously, unique. The variable “City/Town” is now locally suppressed by replacing the city names by the score “Unknown” in the risky records. This generates the records as shown in Figure 4.

	Occupation	City/Town	Gender	Education	...
1	Mayor	Unknown	Man	High	...
2	Fisherman	Urk	Man	Low	...
3	Teacher	Amsterdam	Woman	High	...
4	Plumber	Papendrecht	Man	Medium	...

Figure 4: Records from Figure 3 after local suppression of City

Now, in record 1, the mayor is no longer unique. Since the suppression is applied *locally*, the city/town variable in the safe records 2 to 4 is not suppressed.

2.5 Top-coding

2.5.1 Short description

When protecting microdata, most attention is paid to the treatment of the identifying variables. They play an important role in the disclosure control. The numerical variables are often the variables that are of interest to the data user (and also a possible discloser), such as income, etc. The actual income of an average Dutch person does not identify this person to a significant extent, but that is not the case for people with an extremely high income. Suddenly, the variable of income has become an identifying variable, and therefore the need for extra protection must be assessed.

Top-coding is a suitable method in this situation. It is a simple method, in which values above a certain threshold are replaced by the same standard value. This can be an indication such as (‘many’) or (‘> threshold’). However, the mean of all records with a value above that threshold can also be used. The advantage of this last choice is that the mean for the top-coded variable remains the same for all records.

In addition, bottom-coding can be used in an equivalent manner.

It is clear that top-coding is only useful for numerical variables. For qualitative variables, global recoding (see section 2.3) can be used to obtain a sort of top-coding.

2.5.2 Applicability

This method can be used as additional protection in those situations where some extremes of numerical variables must be considered as identifying.

2.5.3 Detailed description

An implementation of this method is available in μ -ARGUS. Using the menu option Modify|ModifyNumericalVariables, the user can easily indicate in a dialogue box (see Figure 5) the variable on which top or bottom-coding should be applied and which replacement value should be included in the file.

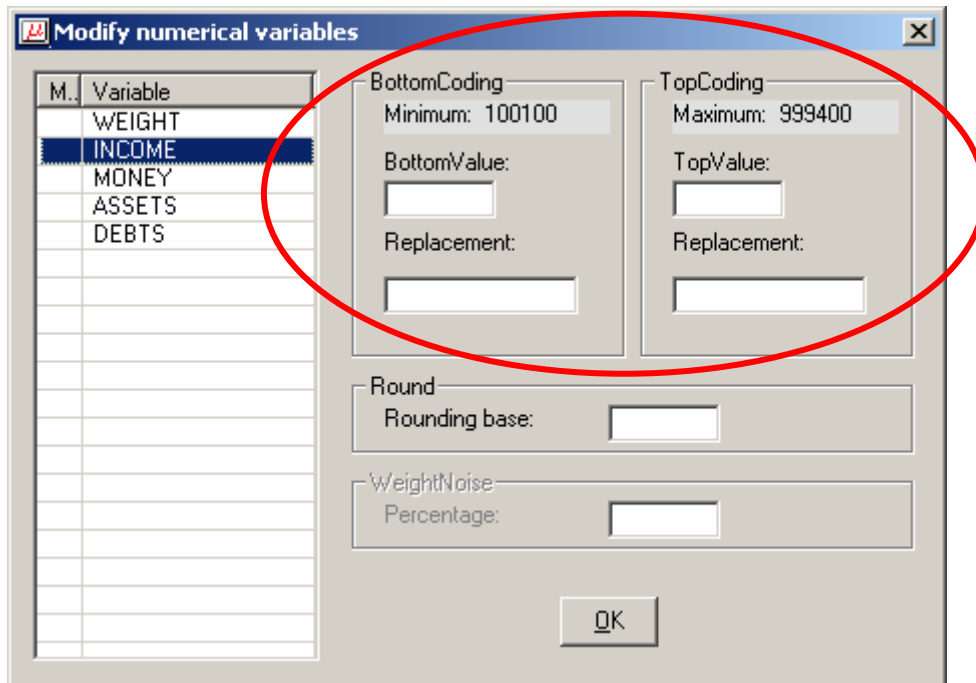


Figure 5: Dialogue box from μ -ARGUS for top/bottom-coding

Once the information has been entered, μ -ARGUS will only save the specification in this phase. Only when a protected file will actually be saved the top or bottom-coding is actually performed.

2.6 Adding noise to weights

2.6.1 Short description

If the file contains raising weights (to correct for the sample and/or non-response), the data protector must consider whether, using the information about the sample design, certain information could be retrieved from those raising weights that could lead to disclosure. A well-known example is that the region is often used as a stratification variable. If, in the protection with global recoding (see section 2.3), the region information is limited or possibly completely removed, a consideration should be made as to whether information can still be derived about the region from the value of the raising variable. If, for example, the city is replaced by the province, it is still possible that the raising weight could reveal that it concerns a large city. And therefore it is clear which (suppressed) city information this relates to.

This type of disclosure can be avoided by adding sufficient noise to the raising weight. By adding random noise, the raising weight will generally still be useful in analyses.

2.6.2 Applicability

In such cases as indicated above, knowledge about the sample design could divulge information which could contribute to the disclosure of data. This method can help to prevent information about individual respondents from being disclosed from raising weights.

2.6.3 Detailed description

An implementation of this method is available in μ -ARGUS. Using the menu option Modify|ModifyNumericalVariables, the user can easily indicate in a dialogue box (see Figure 6) how much noise should be added to the raising weight.

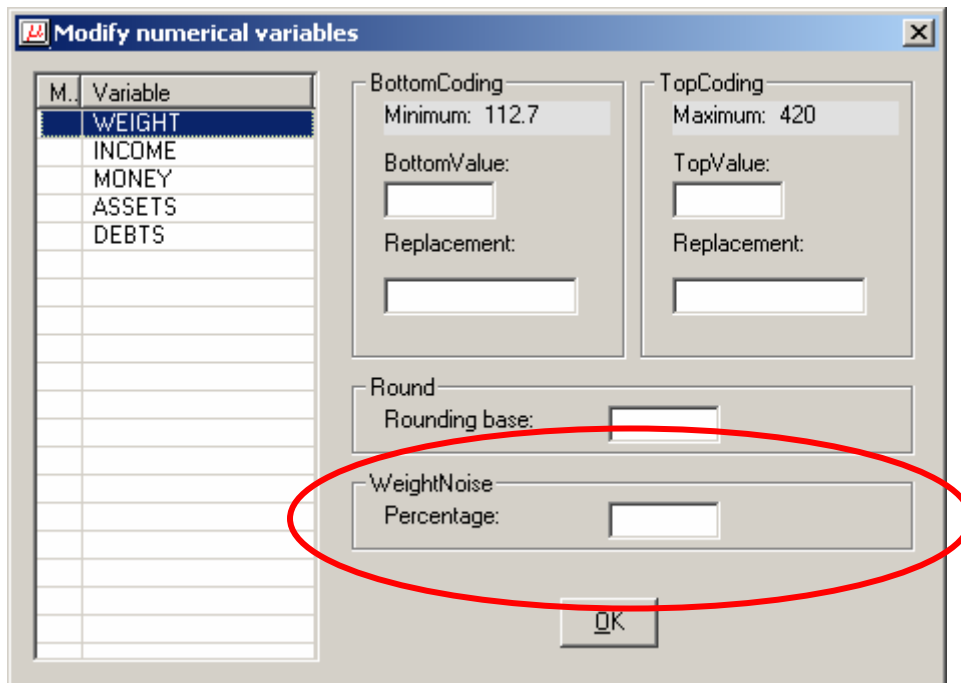


Figure 6: Dialogue box from μ -ARGUS for adding noise to weights

A percentage p can be indicated, so that μ -ARGUS will replace the weight w_i by a random value from the interval

$$\left[\frac{(100 - p)}{100} w_i, \frac{(100 + p)}{100} w_i \right]. \quad (2.6.1)$$

Once the information has been entered, μ -ARGUS will only save the specification in this phase. Only when a protected file will actually be saved the noise will be added to the raising variable.

2.7 PRAM

2.7.1 *Short description*

The Post Randomisation Method (PRAM) is a method for the statistical disclosure control of categorical variables. PRAM can be considered as an intentional misclassification, for which the misclassification probabilities are recorded by the data protector. PRAM is also related to the Randomised Response (RR) technique. However, RR is performed when the questions are being asked, while PRAM is only applied after the answer has been provided.

When PRAM is used, for each record in a microdata file, the score on one or more categorical variables is changed – or not – based on a certain probability. This is done independently on all the records. The probability mechanism that determines the transition of the scores is recorded in advance in a so-called Markov matrix.

Because PRAM is a stochastic method, the disclosure risk is directly affected: if a discloser believes that he or she recognises a record, there is a certain probability that this record does *not* correspond to the person that the discloser is thinking of. After all, several scores on identifying variables are changed with a certain probability.

The fact that the probability mechanism used is known when PRAM is applied means that it is possible, using the protected microdata and the Markov matrix, to construct unbiased estimators for certain statistical attributes of the original data. In addition, techniques from the misclassification and the Randomised Response can also be used.

For a detailed description of PRAM, we refer to Gouweleeuw et al. (1998a and 1998b).

2.7.2 *Applicability*

The Statistics Netherlands policy for the statistical disclosure control of microdata under contract states that the identification of individual people must be prevented (or, in any case, it must be made more difficult). To identify an individual person, a discloser will have to use identifying variables, such as gender, marital status, age and educational level. Naturally, this only works if the discloser is certain that the variables in the file provided are actually the true scores. By applying PRAM to identifying variables, this certainty is eliminated: there is now a positive probability that the score is no longer the original score.

In the statistical disclosure control of a microdata file under contract, it is generally not possible to include very detailed regional variables. This is particularly the case if other detailed identifying variables are present in the file. In this situation, the traditional statistical disclosure control methods, such as recoding, top-coding and local suppression, would produce a file that is virtually unusable for analyses in which the regional detail is important. PRAM would then be a possible alternative:

the detail level is maintained, but the actual score on an identifying variable can no longer be seen with certainty.

A user of a file that is protected using PRAM, however, must have sufficient statistical knowledge to be able to correct his or her desired analysis method for the changes made to the records. How these methods must be adapted is known for several analysis methods. See, for example, Gouweleeuw et al. (1998a and 1998b), Van den Hout (1999), Van den Hout and van der Heijden (2002) and Ronning et al. (2004).

Files that are protected using PRAM are thus mainly intended for theoretically or otherwise experienced statisticians. In addition, microdata files on which PRAM was applied can also be used as “test files”; for example, to test scripts or to determine research trends. The ultimate definitive analysis would then have to be performed on the original (unprotected) file by means of remote execution or an onsite session.

2.7.3 Detailed description

For a detailed theoretical description of the method, please refer to Gouweleeuw et al. (1998a and 1998b).

The Markov matrix with transition probabilities plays an important role in the application of PRAM. The transition probabilities determine the level of protection and affect the information loss. It is therefore important to properly select these probabilities. Each user will experience information loss in a different way. It is therefore preferable to keep the users’ wishes in mind when determining the transition probabilities. De Wolf (2006) provides different measures for information loss.

Because PRAM is a stochastic disclosure control method, the standard rules as described in the Statistical Disclosure Control Handbook (Hundepool et al., 2006) are not directly applicable. However, alternative rules are provided in, for example, De Wolf (2006), and these are related to the standard rules from the Statistical Disclosure Control Handbook.

It should be clear that the selection of the transition probabilities is not an easy task. There is no universal way to take the right decision in every situation. The following questions play a role in determining the transition probabilities:

- On which variables will PRAM be applied?
- Will PRAM be applied independently on these variables or on a subset thereof?
- Are there impossible combinations that must be prevented by setting the associated transition probabilities to zero?
- What effect does it have on the information loss?
- What effect does it have on the disclosure risk?

For each case, the answers to these questions will determine the selection of the specific transition probabilities. There is therefore no universal method available to determine the ideal transition probabilities.

For an empirical study into the consequences of different possibilities for the transition probabilities on both the disclosure risk and the information loss, please refer to De Wolf (2006).

When selecting a matrix of transition probabilities, a number of typical structures are possible. For example, a band matrix with bandwidth b can be useful for ordinal variables such as Age. In that case, an age can be replaced with a certain probability by an age within plus or minus b years. Completely filled matrices are mainly useful for nominal variables with a limited number of categories, such as the variable Marital status. See Figure 7 for a few examples.

$$\begin{matrix} \begin{pmatrix} 0.70 & 0.10 & 0.10 & 0.10 \\ 0.02 & 0.94 & 0.02 & 0.02 \\ 0.09 & 0.09 & 0.73 & 0.09 \\ 0.11 & 0.11 & 0.11 & 0.67 \end{pmatrix} & \begin{pmatrix} 0.80 & 0.20 & 0 & 0 \\ 0.10 & 0.80 & 0.10 & 0 \\ 0 & 0.20 & 0.60 & 0.20 \\ 0 & 0 & 0.10 & 0.90 \end{pmatrix} \\ \text{(a)} & \text{(b)} \end{matrix}$$

Figure 7: Examples of matrices with transition probabilities: (a) Completely filled matrix, (b) Band matrix with bandwidth 1

For other variables, a block matrix is a more obvious solution. For example, for a variable such as Region (at city/town level), we can consider a block matrix in which the blocks correspond to the Provinces. In this way, cities can only be replaced by other cities from the same province. See Figure 8 for an example of a block matrix with transition probabilities.

$$\begin{pmatrix} 0.90 & 0.10 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0.20 & 0.80 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0.70 & 0.20 & 0.10 & 0 & 0 & 0 \\ 0 & 0 & 0.10 & 0.80 & 0.10 & 0 & 0 & 0 \\ 0 & 0 & 0.15 & 0.15 & 0.70 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0.75 & 0.15 & 0.10 \\ 0 & 0 & 0 & 0 & 0 & 0.09 & 0.82 & 0.09 \\ 0 & 0 & 0 & 0 & 0 & 0.05 & 0.05 & 0.90 \end{pmatrix}$$

Figure 8: Example of a block matrix with three blocks with transition probabilities.

2.7.4 Example

At present, μ -ARGUS only offers a limited facility to perform statistical disclosure control using PRAM. In that package, it is possible to apply PRAM per variable, for

which the Markov matrix may either be a band matrix or a completely filled matrix. The bandwidth of a band matrix is adjustable, the same as the diagonal probabilities (the probabilities that certain categories do *not* change).

Because PRAM is a stochastic disclosure control method (only the transition *probabilities* are recorded), a protected file may look different after every application of PRAM: such a protected file is, in any case, the outcome of a probability experiment. Analyses can therefore only be corrected *in expectation* for the fact that they are used on a file that is protected using PRAM. This means that, for example, the expectation for corrected estimated parameters will be the same as the parameter estimations based on the original file.

To obtain an impression of possible adaptations of analyses, consider the simple case of PRAM applied to the variable Gender (two categories), in which we want to estimate the frequency table for the number of men and the number of women. We notate the variable Gender by ξ , where $\xi = 1 = \text{Man}$ and $\xi = 2 = \text{Woman}$. We notate the associated frequency table by \mathbf{T}_ξ . Suppose that the original file contains 100 men and 100 women, so $\mathbf{T}_\xi = (100, 100)^t$. PRAM is applied to the variable Gender using the following matrix with transition probabilities:

$$\mathbf{P} = \begin{pmatrix} 0.90 & 0.10 \\ 0.20 & 0.80 \end{pmatrix}. \quad (2.7.1)$$

Written out: the probability that the gender Man will be changed to Woman is 10%, the probability that the gender Woman will be changed to Man is 20%. The variable ξ is notated as X after the application of PRAM. The frequency table of Gender based on the protected file is then written \mathbf{T}_X . It can easily be derived that

$$E(\mathbf{T}_X | \xi) = \mathbf{P}' \mathbf{T}_\xi, \quad (2.7.2)$$

where the expectation is conditional on the original file. In the example, this means that it is expected that 110 men and 90 women will occur in the protected file. An unbiased estimator for the original frequency table follows directly from equation (2.7.2), i.e.

$$\hat{\mathbf{T}}_\xi = (\mathbf{P}')^{-1} \mathbf{T}_X. \quad (2.7.3)$$

The original frequency table will only be reproduced *in expectation* by this corrected estimator. In other words,

$$E(\hat{\mathbf{T}}_\xi | \xi) = \mathbf{T}_\xi. \quad (2.7.4)$$

Suppose that the protected file contains 112 men and 88 women (Please note: this is an example, because this can differ for each realisation of the probability experiment), then the unbiased estimation (rounded to whole numbers) would be represented by

$$\hat{\mathbf{T}}_\xi = (\mathbf{P}')^{-1} \mathbf{T}_X = \begin{pmatrix} 0.90 & 0.20 \\ 0.10 & 0.80 \end{pmatrix}^{-1} \begin{pmatrix} 112 \\ 88 \end{pmatrix} = \begin{pmatrix} 103 \\ 97 \end{pmatrix}. \quad (2.7.5)$$

Note that this corrected estimation of the frequency table is much closer to the original frequency than the uncorrected estimation (the direct count from the protected file), but that the exact original values were not obtained.

2.8 Conclusion

The package μ -ARGUS is available at Statistics Netherlands to protect microdata files. For a detailed description of the package, we refer to the associated manual (Hundepool et al., 2007).

When μ -ARGUS is used, a report is made after each session in which one or more files were protected. This report includes the methods and parameters used.

With μ -ARGUS, it is easy to see the effects of different statistical disclosure control methods. Methods can be applied and then also undone (within a single session). The automatically generated report can be used to easily derive for a subsequent version of the same file which method or combinations of methods were ultimately used.

Because μ -ARGUS was and is being developed in a European context, it also includes several methods that are not described in this Methods Series. These are methods that are used by several other EU countries, but not by Statistics Netherlands.

3. Statistical Disclosure Control of Quantitative Tables

3.1 General description and reading guide

3.1.1 General description

The statistical disclosure control of quantitative tables encompasses the production of quantitative tables that satisfy Statistics Netherlands policy for statistical disclosure control and, as such, can be published. The disclosure control policy for quantitative tables is set down in chapter 4 of the Statistical Disclosure Control Handbook (Hundepool et al., 2006). Quantitative tables are tables in which the cell values are composed by summation of a continuous variable over all the contributors to a cell. This is in contrast to frequency tables in which only the *number* of contributors per cell is given. Other rules apply to frequency tables, and other protection methods may be more suitable than those for quantitative tables. Disclosure control methods for frequency tables are discussed in the subtheme “Statistical Disclosure Control of Frequency Tables”.

If exactly one or two contributors produce a cell total, it is clear that this cell cannot be published. In the case of a single contributor, individual information is released directly, and in the case of two contributors, one contributor can exactly calculate the other contribution by subtracting his or her own contribution from the cell total.

However, undesirable situations can arise also if there are more than two contributors in a cell. In principle, in the statistical disclosure control of quantitative tables, we must prevent (or at least make it more difficult) that a any contribution can be estimated too accurately. This may occur, for example, also in the case that a very large contributor is present in a single cell along with several relatively small contributors. In this case, the second-largest contributor can calculate that the largest contribution does not contribute more than the cell total minus the second-largest contribution to the cell. A relatively good estimation of the contribution of the largest contributor can be obtained as a result, in conflict with the disclosure control rules of Statistics Netherlands.

The presence of empty cells also requires extra attention. In some cases, an empty cell will be a so-called *structural zero cell*. This means that it is generally known that, logically, it is *impossible* for this cell to have a contribution. Such cells can therefore also not be used in the disclosure control: whatever you do, everyone knows that they must be empty cells.

At the same time, reliable information can sometimes be disclosed using *non-structural zero cells*. If there are contributors in such a cell, there is actually a sort of group disclosure: it is immediately clear that all the contributors have provided a contribution of zero (assuming that the contributions are non-negative). If there are

no contributors in the cell, but it is not automatically impossible for a contributor to be in this cell, this in itself also reveals direct information.

The methods described in this subtheme can be easily applied using the software package τ -ARGUS. This package was developed by DMV in a European context.

3.1.2 Reading guide

As the initial step in determining the correct statistical disclosure control for a quantitative table, it will first have to be determined whether disclosure is possible. In the first instance, the basis for this is “common sense”: is there information present in the table that may not be disclosed about individual respondents? For quantitative tables, such information is generally a respondent’s individual contribution to the total of a specific cell in the table.

In addition, an objective method is needed to determine which cells in the table contain respondents that potentially run a risk of their individual contribution being disclosed. The p% rule (see section 3.3) is intended to identify such primary risky cells. This method can only be used for quantitative tables and not for frequency tables.

After the risky cells have been identified, the table will generally have to be further protected. There are three general methods available for this purpose: restructuring the table (see section 3.4), suppressing cells (see section 3.5) and rounding (see section 3.6).

Which method or combination thereof is ultimately used in a specific situation cannot be determined in advance. This depends to a significant extent on the intended users. For example, in Eurostat regulations, it is not always possible to restructure the table, and cell suppression will often have to be chosen. The department responsible for the quantitative table concerned is also responsible for the adequate statistical disclosure control of the table. When selecting the method or methods to be used, two competing aspects must be taken into account:

- disclosure risk;
- information loss.

In general, it can be said that reducing the disclosure will lead to increased information loss. The converse is also true: the smaller the information loss, the larger the disclosure risk. In certain cases, an assessment will have to be made, in which, at a minimum, the rules in the Statistical Disclosure Control Handbook (Hundepool et al., 2006) will always have to be satisfied.

3.2 Scope and relationship with other themes and subthemes

This subtheme discusses methods that can be used for the statistical disclosure control of quantitative tables. This chapter does *not* discuss any methods that can only be used for frequency tables. For such methods, see the subtheme “Statistical Disclosure Control of Frequency Tables”.

A few of the methods described in this chapter can, in principle, be used for quantitative tables as well as for frequency tables. Such methods will be repeated in the subtheme “Statistical Disclosure Control of Frequency Tables”.

The methods in this chapter can be divided into two variants: methods for determining the primary (or other) risky cells of a quantitative table and methods for making tables with risky cells suitable for publication.

3.3 *P* % rule

3.3.1 Short description

The goal of statistical disclosure control is to prevent the disclosure of information about individual contributors to a table, or at least to make this more difficult. To achieve this, the cells where there is a risk of possible disclosure will first have to be identified. An objective measure is needed for this, one that indicates how well an individual contribution to a cell can be estimated based on the published table. The p % rule provides for this. This is also the method to use to indicate to what extent the disclosure control rule has been violated and how large the measures to be taken must be.

3.3.2 Applicability

Before a quantitative table can be protected statistically, it must first be indicated where potential problems occur in that table. The p % rule indicates how well a contributor in a cell would be able to estimate another contributor in that same cell. This serves to determine the primary risky cells, and it also gives an indication how much protection must be provided to satisfy the Statistics Netherlands policy for the publication of quantitative tables.

With this method, account can also be taken of possible authorisations/waivers: contributors who have indicated that they do not object to publications from which their contribution can be derived. Such contributors are then simply excluded in the application of the p % rule.

The p % rule may only be used:

- in the case of quantitative tables;
- with non-negative contributors;
- for which the largest contributors are identifiable for the discloser;
- on non-empty cells with a positive cell total.

3.3.3 Detailed description

Let T_A be the cell value of cell A in the table in question. Denote the largest contributor without a waiver by X_s and the largest of the remaining contributors by X_r . Then cell T_A is risky if:

$$\frac{(T_A - X_r) - X_s}{X_s} < \frac{p}{100}. \quad (3.3.1)$$

That is, in the situation of no waivers, a cell is risky if the second-largest contributor can estimate the largest contributor with an accuracy exceeding p %.

It is simple to see that this is the worst scenario: if the second-largest contributor *cannot* estimate the largest contributor more accurately than p %, then no other contributor can estimate an arbitrary other contributor more accurately than that p %, and therefore the cell is safe. In other words: the most accurate estimation can be made by the second-largest contributor, when this party estimates the largest contribution.

The value of the difference between the left side and the right side of the inequality in formula (3.3.1) also indicates how much protection a risky cell needs. For more detail, please refer to Loeve (2001).

With the standard software at Statistics Netherlands for the protection of tables, τ -ARGUS, it is easy to apply the p % rule. This method is one of the standard built-in rules that can be used to identify the primary risky cells. Moreover, τ -ARGUS automatically calculates how much protection a risky cell needs and uses that in the further protection of the table concerned. To make it possible for τ -ARGUS to identify the primary risky cells using the p % rule, however, it is necessary that the input for τ -ARGUS consists of the microdata from which the table concerned is composed. To apply the p % rule, information is needed, in any case, about the individual contributors. For more information about the use of τ -ARGUS, please refer to the associated manual (Hundepool et al., 2003).

The value selected for p is determined by Statistics Netherlands policy. The Statistical Disclosure Control Handbook from Statistics Netherlands (Hundepool et al., 2006) gives an interval within which p should be selected ($5 \leq p \leq 15$). The exact value for p is determined by the statistical division and may never be revealed to external parties, because this could help them in the calculation of the suppressed cells.

A large value for p results in strict disclosure control, because, when estimating an arbitrary contribution in that case, not even a relatively “large” error may be made. A small value for p results in less strict disclosure control, because a cell is only risky in this situation if a contribution can be estimated very accurately.

3.3.4 Example

In this fictitious example, we look at a cell in a table with turnover according to SBI (the Dutch Standard Industrial Classification) and Region. Suppose that the cell with SBI = 32 and Region = Noord Brabant consists of four contributors with the values 324, 4, 2 and 10. Suppose that we want to use the p % rule where $p = 5$, then we must first sort the contributors: $X_1 = 324$, $X_2 = 10$, $X_3 = 4$ and $X_4 = 2$. The cell total T_A

is then 340. If we calculate the quotient from formula (3.3.1), we obtain the value 0.0185. This is clearly smaller than 5 %, and therefore the cell is risky.

3.4 Table restructuring

3.4.1 Short description

Section 3.3 describes when a cell in a table must be considered risky. In general, cells with a limited number of contributors or a cell with one or two large contributors are the obvious candidates to be characterised as risky. All risky cells must be protected. Before performing suppression on a large scale, restructuring the table can also be considered. By combining rows and/or columns, cells are pooled and the content per cell is increased. The result of this is that fewer cells are identified as risky by the p % rule, as described in section 3.3.

3.4.2 Applicability

This method will generally lead to fewer risky cells in the table. Combining cells creates cells that are safer than the individual cells that were combined.

There are no methodological conditions for using this method. However, externally imposed obligations sometimes specify what detail level of a table must be published. This may be a Eurostat obligation, but Statistics Netherlands policy can also mean that a certain detail level of a table must be published. In these cases, the method can be applied from a technical perspective, but its use is prevented by external policy decisions.

Furthermore, an assessment must be made between the information loss resulting from the larger number of crosses (suppressed cells) that are needed to protect the table, and the information loss resulting from combining columns/rows, for which fewer crosses are needed.

3.4.3 Detailed description

The software package τ -ARGUS has provisions for recoding rows and/or columns in tables. In this regard, a distinction is made between two situations:

- In the case of a hierarchical spanning variable, the recoding implies that certain splits are omitted at the lowest level.
- In the case of an unstructured spanning variable, users are free to combine the columns or rows of a table as they choose.

3.4.4 Example

Figure 9 presents a fictitious table of the turnover according to Region (hierarchical) and SizeClass. Figure 10 provides two possible restructuring possibilities for this table. The variable SizeClass is recoded such that the categories 2 to 6 are combined into the category MediumSmall, and that the categories 7, 8 and 9 are combined into

the category Large. Note that, in this way, all the primary risky cells are combined to create safe cells. In the recoding of the variable Region the smallest detail level has been removed. This restructuring does not resolve all the problems: the primary risky cells at region level (for North and East) are still present in the table.

	tot	2	4	5	6	7	8	9	99
tot	16,847,646.84	20.00	25.00	2,711,808.00	2,320,534.00	2,505,042.58	2,799,074.26	6,510,758.00	385.00
- North	4,373,664.00	X	X	719,049.00	659,680.00	688,962.00	756,529.00	1,549,049.00	385.00
.. 1	1,986,129.00	X	X	398,062.00	348,039.00	354,711.00	418,778.00	466,529.00	-
.. 2	1,809,246.00	0.00	-	223,990.00	221,332.00	241,913.00	258,233.00	863,393.00	385.00
.. 3	578,289.00	-	-	96,997.00	90,309.00	92,338.00	79,518.00	219,127.00	-
- East	3,703,896.00	15.00	X	642,238.00	515,003.00	534,147.00	620,392.00	1,392,096.00	-
.. 4	124,336.00	X	-	36,311.00	32,132.00	25,770.00	18,150.00	-	X
.. 5	526,279.00	-	-	93,589.00	94,957.00	110,930.00	81,799.00	145,004.00	-
.. 6	2,234,995.00	X	X	345,803.00	251,358.00	251,188.00	303,377.00	1,083,254.00	-
.. 7	818,286.00	-	-	166,535.00	136,556.00	146,259.00	217,066.00	151,870.00	-
- West	4,576,115.84	-	-	648,972.00	543,570.00	663,896.58	775,132.26	1,944,545.00	-
.. 8	485,326.00	-	-	63,767.00	75,442.00	87,305.00	59,953.00	198,859.00	-
.. 9	3,664,559.84	-	-	537,911.00	430,851.00	515,019.58	643,762.26	1,537,016.00	-
.. 10	426,230.00	-	-	47,294.00	37,277.00	61,572.00	71,417.00	208,670.00	-
- South	4,193,971.00	-	15.00	701,549.00	602,281.00	618,037.00	647,021.00	1,625,068.00	-
.. 11	2,752,743.00	-	15.00	488,613.00	392,395.00	363,490.00	402,925.00	1,105,305.00	-
.. 12	1,441,228.00	-	-	212,936.00	209,886.00	254,547.00	244,096.00	519,763.00	-
.99	-	-	-	-	-	-	-	-	-

Figure 9: Quantitative table for turnover according to region and size class

	tot	Large	SmallMedium	99
tot	16,847,646.84	11,814,874.84	5,032,387.00	385.00
- North	4,373,664.00	2,994,540.00	1,378,739.00	385.00
.. 1	1,986,129.00	1,240,018.00	746,111.00	-
.. 2	1,809,246.00	1,363,539.00	445,322.00	385.00
.. 3	578,289.00	390,983.00	187,306.00	-
- East	3,703,896.00	2,546,635.00	1,157,261.00	-
.. 4	124,336.00	55,888.00	68,448.00	-
.. 5	526,279.00	337,733.00	188,546.00	-
.. 6	2,234,995.00	1,637,819.00	597,176.00	-
.. 7	818,286.00	515,195.00	303,091.00	-
- West	4,576,115.84	3,383,573.84	1,192,542.00	-
.. 8	485,326.00	346,117.00	139,209.00	-
.. 9	3,664,559.84	2,695,797.84	968,762.00	-
.. 10	426,230.00	341,659.00	84,571.00	-
- South	4,193,971.00	2,890,126.00	1,303,845.00	-
.. 11	2,752,743.00	1,871,720.00	881,023.00	-
.. 12	1,441,228.00	1,018,406.00	422,822.00	-
.99	-	-	-	-

(a) Recoding of SizeClass (all primary risky cells have been protected)

	tot	2	4	5	6	7	8	9	99
tot	16,847,646.84	20.00	25.00	2,711,808.00	2,320,534.00	2,505,042.58	2,799,074.26	6,510,758.00	385.00
North	4,373,664.00	X	X	719,049.00	659,680.00	688,962.00	756,529.00	1,549,049.00	385.00
East	3,703,896.00	15.00	X	642,238.00	515,003.00	534,147.00	620,392.00	1,392,096.00	-
West	4,576,115.84	-	-	648,972.00	543,570.00	663,896.58	775,132.26	1,944,545.00	-
South	4,193,971.00	-	15.00	701,549.00	602,281.00	618,037.00	647,021.00	1,625,068.00	-
.99	-	-	-	-	-	-	-	-	-

(b) Recoding of Region (not all primary risky cells have been protected)

Figure 10: Two possible restructuring possibilities used on the table from Figure 9

3.5 Cell suppression

3.5.1 Short description

A frequently used method to protect primary risky cells is to suppress (not publish) certain cells. The cell value is then simply replaced by an X (×).

In a quantitative table when the marginals are also provided, however, it is often not sufficient to suppress only the primary risky cells. If a suppressed cell is the only suppressed cell in a row, the suppressed value can, after all, simply be calculated by subtracting the other cell values in that row from the corresponding marginal.

To sufficiently protect primary risky cells, it is therefore also necessary to suppress other cells which, in themselves, are safe. This is called *secondary suppression*. It is not easy to perform this in such a way such that the primary risky cells are protected sufficiently, while also ensuring that not too much information is removed from the table. Furthermore, account must also be taken of the fact that structural zero cells cannot be used as secondary suppressions: everyone knows that, by definition, these cells are empty.

To prevent a situation where suppressed, primary risky cells can be calculated exactly, secondary suppressions are therefore necessary. However, also a “too accurate” estimation for a suppressed cell is not desirable. Indeed, what is the difference between the following statements: “This suppressed cell actually has a value of 10000” and “This suppressed cell actually has a value of between 9998 and 10002”. Given a suppression pattern, it is always¹ possible to calculate an interval in which a suppressed cell must lie. The method of “Cell Suppression” must then also produce a suppression pattern, for which the intervals that can be calculated are sufficiently large. The size of these intervals is determined by the rule that is used to determine the primary risky cells.

Fischetti and Salazar (2000) have developed a method to solve the above problem in an optimal manner. Their method is, in theory, applicable to arbitrary, additive tables with non-negative contributors. In practice, however, their solution involves too much computing time if the tables become too large, either in size or complexity. This is why a number of suboptimal methods have been developed to find suitable suppression patterns for larger and/or more complex tables.

For example, the “modular approach” (HiTaS) splits a hierarchical table into a large number of non-hierarchical subtables and applies the optimal method to each individual subtable. By correctly combining the results, a suboptimal solution can be obtained for the entire table, with a significantly shorter computing time.

The “hypercube approach” can also protect large tables by protecting the subtables in a certain iterative way. The protection of each subtable also takes place suboptimally. Consequently, the approach is relatively fast, but, in general, more cells are suppressed than strictly necessary to obtain a protected table.

3.5.2 Applicability

This method can be used to adequately protect quantitative tables with cells that do not satisfy the requirements of the Statistics Netherlands statistical disclosure control

¹ In the case that the table is composed of non-negative contributors and the marginals are also given.

policy. In particular, if the table cannot be restructured further or at all, the cell suppression method can be used effectively.

The contributions to the table to be protected must not be negative and the table must be additive, and the marginals must also be provided.

In the modular approach, the table must be three-dimensional at a maximum. Each dimension may be hierarchical. Linked tables can be protected by copying the suppressions from one table to the other, and then protecting the tables. This should then possibly be performed in an iterative manner. Recent developments in τ -ARGUS make it possible to solve the linked tables problem automatically.

In the hypercube approach as implemented in τ -ARGUS, the table may be seven-dimensional at a maximum. The table may be hierarchical in every dimension. Linked tables are also possible in principle.

It should be mentioned that for both approaches, from a performance perspective, the recommendation is to avoid using long, unstructured (non-hierarchical) code lists.

3.5.3 Detailed description

The software package τ -ARGUS has a provision to apply cell suppression to quantitative tables. If the original microdata is used as input, τ -ARGUS will determine the primary risky cells with the associated safety intervals (see also section 3.3).

After this, τ -ARGUS will have to determine a suppression pattern that guarantees the necessary safety intervals. There are various options for this. We will discuss the two approaches that are the most interesting for Statistics Netherlands.

3.5.3.1 Modular approach

For a detailed description and an elaborated example of the modular approach, see De Wolf (2002).

Generally, the modular approach can be described as follows:

1. Split the hierarchical table into all logical non-hierarchical subtables.
2. Group the subtables in classes in such a way that all tables in a single class can be protected independently of each other. For a suitable classification, see De Wolf (2002).
3. Protect all tables in class K .
4. If no secondary suppressions are placed in the marginals of the subtables of class K , continue with class $K + 1$, including any secondary suppressions in the inside of a table as primary suppressions for class $K + 1$.
5. If secondary suppressions do have to be placed in a marginal of at least one subtable, go back to class $K - 1$, including only the secondary suppressions in the marginals as primary suppressions.

6. Repeat steps 4 and/or 5 until all subtables have been protected at the lowest (most detailed) hierarchical level.

All non-hierarchical subtables will be protected using the mixed integer approach from Fischetti and Salazar (2000). In this approach, the required safety intervals are guaranteed, while a certain cost function is minimised. This cost function can be selected in different ways, as a result of which various forms of information loss can be minimised. This minimisation takes place *locally*, so that the ultimate solution for the entire (hierarchical) table does not necessarily also have to be optimal.

In selecting the cost function in τ -ARGUS, several options can be selected, including:

- A variable from the dataset (such as the quantitative value on which tabulation takes place);
- A constant (so that the number of suppressions is minimised);
- The number of contributors per cell (so that the total number of suppressed contributions is minimised).

In the disclosure control of a subtable, also the so-called singletons problem must be taken into account: cells with only one contribution. If such cells are in a suppression pattern, the contributors involved can reverse part or all of the suppression pattern. After all, they know what their own contribution is and can therefore fill in that suppressed value, as a result of which it may also be possible to calculate other suppressed cells. In the current implementation of the mixed integer approach in τ -ARGUS, it is not possible to keep each conceivable combination of a singleton with another suppressed cell under control while searching for a suppression pattern. However, it is possible to take account of the combinations within a single row, column or layer² in the table. The combinations which must be taken into account consist of exactly two primary risky cells in a single row, column or layer, of which at least one cell is a singleton. By giving the larger of these two primary risky cells a safety interval that is just large enough that it cannot be satisfied by the other primary risky cell, at least one extra secondary suppression will always be made in the row, column or layer concerned.

In a similar way, it is ensured that, within a single row, column or layer, all the suppressed cells together contain more than the minimum required number of contributors for a safe cell.

3.5.3.2 Hypercube approach

For a more detailed description of the hypercube approach, see Giessing and Repsilber (2002).

² A row consists of the cells with coordinates (r, k, l) where k and l are fixed. A column consists of the cells with the coordinates (r, k, l) where r and l are fixed. A layer consists of the cells with coordinates (r, k, l) where r and k are fixed.

In this approach too, a hierarchical table is split into non-hierarchical subtables. The non-hierarchical subtables are then protected in a certain order, where the subtables at the highest level are dealt with first.

For each subtable, all possible hypercubes are constructed for each primary risky cell in which that primary risky cell is one of the corner points. For each hypercube, the interval is calculated around the primary risky cell if all other corner points of the hypercube are also suppressed. If that interval is large enough (depending on the protection rule used), the associated hypercube is designated as “feasible”. The information loss is then calculated for each feasible hypercube. Finally, the admissible hypercube with the smallest information loss is selected to protect the primary risky cell concerned.

No linear programming problem needs to be solved in order to calculate the safety intervals resulting from a hypercube. This significantly accelerates the procedure. The hypercube approach is therefore, in general, faster than the modular approach, for which a mixed integer programming problem needs to be solved.

After all subtables are protected in this way, the entire procedure is repeated. Secondary suppressed cells from a certain subtable that also occur in other subtables are considered as primary risky cells in those other subtables, and dealt with as such. This process is repeated until no more changes take place.

Note that the use of hypercubes to protect primary risky cells is a sufficient but not necessary condition for a safe suppression pattern. In other words, in some cases, the combination of the different hypercubes will not lead to an optimal suppression pattern, but it will always produce a safe suppression pattern. Consequently, this approach tends to suppress more cells than necessary for a safe suppression pattern.

This approach also takes account of the so-called singletons. A cell with only one contributor would indeed allow all suppressed corner points of a hypercube to be calculated. Therefore the extra requirement in the case of singletons is that this type of cell must be a corner point of at least two different hypercubes.

3.5.4 Example

Using τ -ARGUS, it is easy to apply cell suppression to a quantitative table. Both the modular approach and the hypercube approach are implemented in τ -ARGUS. It is also possible to select multiple information loss measures for the cost function that must be minimised. For the use of τ -ARGUS, please refer to the associated manual (Hundepool et al., 2003).

Figure 11 presents an example of a table in which only the primary risky cells are suppressed.

	tot	2	4	5	6	7	8	9	99
tot	16,847,646.84	20.00	25.00	2,711,808.00	2,320,534.00	2,505,042.58	2,799,074.26	6,510,758.00	385.00
- North	4,373,664.00	X	X	719,049.00	659,680.00	688,962.00	756,529.00	1,549,049.00	385.00
.. 1	1,986,129.00	X	X	398,062.00	348,039.00	354,711.00	418,778.00	466,529.00	-
.. 2	1,809,246.00	0.00	-	223,990.00	221,332.00	241,913.00	258,233.00	863,393.00	385.00
.. 3	578,289.00	-	-	96,997.00	90,309.00	92,338.00	79,518.00	219,127.00	-
- East	3,703,896.00	15.00	X	642,238.00	515,003.00	534,147.00	620,392.00	1,392,096.00	-
.. 4	124,336.00	X	-	36,311.00	32,132.00	25,770.00	18,150.00	-	X
.. 5	526,279.00	-	-	93,589.00	94,957.00	110,930.00	81,799.00	145,004.00	-
.. 6	2,234,995.00	X	X	345,803.00	251,358.00	251,188.00	303,377.00	1,083,254.00	-
.. 7	818,286.00	-	-	166,535.00	136,556.00	146,259.00	217,066.00	151,870.00	-
- West	4,576,115.84	-	-	648,972.00	543,570.00	663,896.58	775,132.26	1,944,545.00	-
.. 8	485,326.00	-	-	63,767.00	75,442.00	87,305.00	59,953.00	198,859.00	-
.. 9	3,664,559.84	-	-	537,911.00	430,851.00	515,019.58	643,762.26	1,537,016.00	-
..10	426,230.00	-	-	47,294.00	37,277.00	61,572.00	71,417.00	208,670.00	-
- South	4,193,971.00	-	15.00	701,549.00	602,281.00	618,037.00	647,021.00	1,625,068.00	-
..11	2,752,743.00	-	15.00	488,613.00	392,395.00	363,490.00	402,925.00	1,105,305.00	-
..12	1,441,228.00	-	-	212,936.00	209,886.00	254,547.00	244,096.00	519,763.00	-
.99	-	-	-	-	-	-	-	-	-

Figure 11: Quantitative table for turnover according to region and size class

It is clear that this is not sufficient: both the cell (East, 4) and the cell (4, 9) can be directly calculated: (East, 4) = 3 703 896 – 15 – 642 238 – 515 003 – 534 147 – 620 392 – 1 392 096 = 5 and (4, 9) = 1 392 096 – 145 004 – 1 083 254 – 151 870 = 11 968.

Figure 12 shows the suppression pattern that was determined with τ -ARGUS using the hypercube approach. Figure 13 shows the same based on the modular approach. Of course, in a publication, it should be impossible to make a distinction between primary and secondary suppressions.

	tot	2	4	5	6	7	8	9	99
tot	16,847,646.84	20.00	25.00	2,711,808.00	2,320,534.00	2,505,042.58	2,799,074.26	6,510,758.00	385.00
- North	4,373,664.00	X	X	719,049.00	-	X	688,962.00	756,529.00	1,549,049.00
.. 1	1,986,129.00	X	X	398,062.00	-	X	354,711.00	418,778.00	466,529.00
.. 2	1,809,246.00	0.00	-	223,990.00	221,332.00	241,913.00	258,233.00	863,393.00	385.00
.. 3	578,289.00	-	-	96,997.00	90,309.00	92,338.00	79,518.00	219,127.00	-
- East	3,703,896.00	X	X	642,238.00	-	X	534,147.00	620,392.00	1,392,096.00
.. 4	124,336.00	X	-	36,311.00	-	X	25,770.00	18,150.00	X
.. 5	526,279.00	-	-	93,589.00	94,957.00	110,930.00	81,799.00	145,004.00	-
.. 6	2,234,995.00	X	X	345,803.00	-	X	251,188.00	303,377.00	X
.. 7	818,286.00	-	-	166,535.00	136,556.00	146,259.00	217,066.00	151,870.00	-
- West	4,576,115.84	-	-	648,972.00	543,570.00	663,896.58	775,132.26	1,944,545.00	-
.. 8	485,326.00	-	-	63,767.00	75,442.00	87,305.00	59,953.00	198,859.00	-
.. 9	3,664,559.84	-	-	537,911.00	430,851.00	515,019.58	643,762.26	1,537,016.00	-
..10	426,230.00	-	-	47,294.00	37,277.00	61,572.00	71,417.00	208,670.00	-
- South	4,193,971.00	-	X	X	-	X	618,037.00	647,021.00	1,625,068.00
..11	2,752,743.00	-	X	488,613.00	-	X	363,490.00	402,925.00	1,105,305.00
..12	1,441,228.00	-	-	X	-	X	254,547.00	244,096.00	519,763.00
.99	-	-	-	-	-	-	-	-	-

Figure 12: Suppression pattern for the table from Figure 11, using the hypercube approach

	tot	2	4	5	6	7	8	9	99
tot	16,847,646.84	20.00	25.00	2,711,808.00	2,320,534.00	2,505,042.58	2,799,074.26	6,510,758.00	385.00
- ..North	4,373,664.00	5.00	5.00	719,049.00	659,680.00	688,962.00	756,529.00	1,549,049.00	385.00
.. 1	1,986,129.00	5.00	5.00	398,062.00	348,039.00	354,711.00	418,778.00	466,529.00	-
.. 2	1,809,246.00	0.00	-	223,990.00	221,332.00	241,913.00	258,233.00	863,393.00	385.00
.. 3	578,289.00	-	-	96,997.00	90,309.00	92,338.00	79,518.00	219,127.00	-
- ..East	3,703,896.00	15.00	5.00	642,238.00	515,003.00	534,147.00	620,392.00	1,392,096.00	-
.. 4	124,336.00	5.00	-	36,311.00	32,132.00	25,770.00	18,150.00	11,968.00	-
.. 5	526,279.00	-	-	93,589.00	94,957.00	110,930.00	81,799.00	145,004.00	-
.. 6	2,234,995.00	10.00	5.00	345,803.00	251,358.00	251,188.00	303,377.00	1,083,254.00	-
.. 7	818,286.00	-	-	166,535.00	136,556.00	146,259.00	217,066.00	151,870.00	-
- ..West	4,576,115.84	-	-	648,972.00	543,570.00	663,896.58	775,132.26	1,944,545.00	-
.. 8	485,326.00	-	-	63,767.00	75,442.00	87,305.00	59,953.00	198,859.00	-
.. 9	3,664,559.84	-	-	537,911.00	430,851.00	515,019.58	643,762.26	1,537,016.00	-
..10	426,230.00	-	-	47,294.00	37,277.00	61,572.00	71,417.00	208,670.00	-
- ..South	4,193,971.00	-	15.00	701,549.00	602,281.00	618,037.00	647,021.00	1,625,068.00	-
..11	2,752,743.00	-	15.00	488,613.00	392,395.00	363,490.00	402,925.00	1,105,305.00	-
..12	1,441,228.00	-	-	212,936.00	209,886.00	254,547.00	244,096.00	519,763.00	-
..99	-	-	-	-	-	-	-	-	-

Figure 13: Suppression pattern for the table from Figure 11, using the modular approach

3.5.5 Quality indicators

If a table is protected by means of cell suppression, it is possible to calculate the realised safety interval for each suppressed cell. Given the suppression pattern and the structure of the table, two LP-problems must be solved for each suppressed cell (minimising and maximising the value for the suppressed cell).

If τ -ARGUS is used for the protection of a quantitative table, at the end of the session, a report is generated that contains the steps taken and the associated results. It is also possible during the session to obtain information about the protected or unprotected table (for example: the number of primary risky cells, the number of secondary suppressions, information loss).

3.6 Additive rounding

3.6.1 Short description

When rounding cell values in a quantitative table, the exact cell values are only known within a certain interval. A table with primary risky cells can also be protected in this way. The extent to which rounding is performed will, of course, have an impact on the size of the intervals. If each cell is rounded independently, the additivity of the table will not necessarily be maintained.

Of course, there is a simple way to guarantee the additivity: by rounding the cells in the interior of the table independently of each other and then recalculating the marginals. As a result, however, the marginals can deviate significantly from the rounded original values.

In additive rounding, the table is rounded such that the additivity is maintained and that the rounded table deviates from the original as little as possible. Furthermore, it is possible to perform additive rounding in such a way that safety intervals specified in advance can also be guaranteed. Whether this can be achieved, however, depends on the size of the selected rounding base in relation to the safety intervals.

3.6.2 Applicability

Additive rounding can be used for the statistical disclosure control of both quantitative tables and frequency tables. Often, a presentation argument will also play a role: a large number of significant figures suggests a high degree of precision that is not always justified because of sampling errors and measurement errors. Rounding the table values reduces this false precision to a certain extent.

3.6.3 Detailed description

In additive rounding, the cell values in a table are rounded to multiples of a rounding base b , keeping the totals and subtotals in the table equal to the sum of the corresponding parts.

Oftentimes, additive rounding is performed in a “zero restricted” manner. In other words, cell values that are already a multiple of the rounding base are not changed, while the other cell values are rounded to one of the adjacent multiples of that rounding base. The rounded values are selected such that the sum of the absolute deviations of the cell values in the rounded table with respect to the cell values in the original table is minimised, under the restriction that the rounded table remains additive. As a result, it is possible that cell values are not rounded to the closest multiple of the rounding base.

In certain conditions, it is not possible to construct a rounded table under the scenario described above. In that case, the restriction that rounding is performed to one of the adjacent multiples of the rounding base is weakened by allowing a cell value to also be rounded to non-adjacent multiples of the rounding base. This weakening can be limited slightly by determining a maximum for the number of steps that may exist between the rounded value and the original value.

In the case of “zero restricted” additive rounding using rounding base $b > 0$ of the non-negative number $z = ub + r$, where $0 \leq r < b$, rounding is performed on the number a , such that

$$a \in \{ub, (u + 1_{(0,b)}(r))b\} \quad (3.6.1)$$

where $1_{(0,b)}(r)$ is equal to 1 if $r \in (0, b)$ and equal to 0 if $r = 0$.

This means that, in the case that $r = 0$, a is always rounded to ub and, in the case that $r \in (0, b)$, a is always rounded to ub or to $(u + 1)b$.

If, however, the restriction is weakened by a maximum of $K > 0$ steps further than the adjacent multiples of the rounding base, then rounding is performed on the number a , such that

$$a \in \{(0 \vee (u + j))b \mid j = -K, \dots, (K + 1_{(0,b)}(r))\} \quad (3.6.2)$$

where $x \vee y = \max(x, y)$.

Multiple additive rounded versions may exist for a given table. These are all *feasible* tables. The table closest to the original table can subsequently be selected from the feasible tables. In τ -ARGUS, the distance that is minimised is represented by

$$\sum_{i=1}^N |z_i - a_i| \quad (3.6.3)$$

where N is the number of cells in the table (including all totals and subtotals), z_i the cell values in the original table and a_i the corresponding rounded cell values.

Finding the optimal solution is a problem that requires intensive computation (NP-complete). For large tables, this can result in unacceptably long calculation times. Partitioning is built into τ -ARGUS for this reason: a large table can be split into a number of subtables that are rounded individually. After these subtables are rounded, they are combined, calculating (if necessary) the totals and subtotals in question from the rounded parts.

3.6.4 Example

τ -ARGUS can be used to easily perform additive rounding on quantitative tables, while the desired protection margins are guaranteed.

Figure 14 presents an example of a table that contains a number of primary risky cells. Figure 15 contains the associated additively rounded table, with a rounding base of 2000. Of course, in a publication, the primary risky cells are not allowed to be recognisable.

	tot	2	4	5	6	7	8	9	99
tot	16,847,647	20	25	2,711,808	2,320,534	2,505,043	2,799,074	6,510,758	385
- North	4,373,664	5	5	719,049	659,680	688,962	756,529	1,549,049	385
.. 1	1,986,129	5	5	398,062	348,039	354,711	418,778	466,529	-
.. 2	1,809,246	0	-	223,990	221,332	241,913	258,233	863,393	385
.. 3	578,289	-	-	96,997	90,309	92,338	79,518	219,127	-
- East	3,703,896	15	5	642,238	515,003	534,147	620,392	1,392,096	-
.. 4	124,336	5	-	36,311	32,132	25,770	18,150	11,968	-
.. 5	526,279	-	-	93,589	94,957	110,930	81,799	145,004	-
.. 6	2,234,995	10	5	345,803	251,358	251,188	303,377	1,083,254	-
.. 7	818,286	-	-	166,535	136,556	146,259	217,066	151,870	-
- West	4,576,116	-	-	648,972	543,570	663,897	775,132	1,944,545	-
.. 8	485,326	-	-	63,767	75,442	87,305	59,953	198,859	-
.. 9	3,664,560	-	-	537,911	430,851	515,020	643,762	1,537,016	-
..10	426,230	-	-	47,294	37,277	61,572	71,417	208,670	-
- South	4,193,971	-	15	701,549	602,281	618,037	647,021	1,625,068	-
..11	2,752,743	-	15	488,613	392,395	363,490	402,925	1,105,305	-
..12	1,441,228	-	-	212,936	209,886	254,547	244,096	519,763	-
.99	-	-	-	-	-	-	-	-	-

Figure 14: Quantitative table for turnover according to Region and Size class

	tot	2	4	5	6	7	8	9	99
tot	16,848,000	0	0	2,712,000	2,320,000	2,506,000	2,800,000	6,510,000	0
- .North	4,374,000	0	0	720,000	660,000	690,000	756,000	1,548,000	0
.. 1	1,986,000	0	0	398,000	348,000	356,000	418,000	466,000	-
.. 2	1,810,000	0	-	224,000	222,000	242,000	258,000	864,000	0
.. 3	578,000	-	-	98,000	90,000	92,000	80,000	218,000	-
- .East	3,704,000	0	0	642,000	514,000	534,000	622,000	1,392,000	-
.. 4	124,000	0	-	36,000	32,000	26,000	18,000	12,000	-
.. 5	526,000	-	-	94,000	94,000	110,000	82,000	146,000	-
.. 6	2,236,000	0	0	346,000	252,000	252,000	304,000	1,082,000	-
.. 7	818,000	-	-	166,000	136,000	146,000	218,000	152,000	-
- .West	4,576,000	-	-	648,000	544,000	664,000	776,000	1,944,000	-
.. 8	486,000	-	-	64,000	76,000	88,000	60,000	198,000	-
.. 9	3,664,000	-	-	538,000	430,000	514,000	644,000	1,538,000	-
..10	426,000	-	-	46,000	38,000	62,000	72,000	208,000	-
- .South	4,194,000	-	0	702,000	602,000	618,000	646,000	1,626,000	-
..11	2,752,000	-	0	488,000	392,000	364,000	402,000	1,106,000	-
..12	1,442,000	-	-	214,000	210,000	254,000	244,000	520,000	-
.99	-	-	-	-	-	-	-	-	-

Figure 15: Table from Figure 14, protectively additively rounded with a rounding base of 2000

3.7 Conclusion

The package τ -ARGUS is available at Statistics Netherlands for the protection of quantitative tables. For a detailed description of that package, please refer to its associated manual (Hundepool et al., 2007).

When τ -ARGUS is used, a report is generated after each session in which one or more tables are protected. That report includes the methods and parameters used.

Using τ -ARGUS, the effects of various statistical disclosure control methods on the tables can easily be made visible. The different methods can be applied, but they can also be ‘undone’ during the same session.

4. Statistical Disclosure Control of Frequency Tables

4.1 General description and reading guide

4.1.1 General description

The statistical disclosure control of frequency tables encompasses the production of frequency tables that satisfy the Statistics Netherlands policy on statistical disclosure control and can be published as such. The disclosure control policy for frequency tables is laid down in chapter 5 of the Statistical Disclosure Control Handbook (Hundepool et al., 2006). Frequency tables are tables in which the number of contributors per cell is given. This is in contrast to quantitative tables in which the cell values are created by summation of a continuous variable over all the contributors to a cell. Other rules apply to quantitative tables, and other protection methods may be more suitable than those for frequency tables. Disclosure control methods for quantitative tables are discussed in the subtheme “Statistical Disclosure Control of Quantitative Tables”.

Article 37 of the Statistics Netherlands Act (2004) requires the protection of recognisable data about statistical units. A violation of the statistical confidentiality (“disclosure”) boils down to the combination of two facts: the recognition of a unit and the disclosure of further details about that unit.

For frequency tables, this can be formulated as follows. The user must first recognise a contributor or group of contributors in the table. This is followed by a statement about these contributor(s) due to the frequency distribution over the cells. The statement that the table makes possible about this group must provide more information about the members of the group than just the group size. In this sense, knowledge that is needed to recognise the members of the group is not considered information about the members of the group.

The statutory requirement is satisfied if the table does not provide any information about an individual statistical unit as such. The statistical professional standards and Statistics Netherlands’ own interest in the continuity of reporting to Statistics Netherlands, however, leads in certain cases to the requirement that the table does not provide information about groups of statistical units (people or households, etc.). In particular, that is the case if the table contains variables that could provide harmful or potentially damaging information about these groups. Such data will be referred to hereinafter as “sensitive data”.

The methods described in this subtheme can be easily applied using the software package τ -ARGUS. This package was developed by DMV (and its predecessors) in a European context.

4.1.2 Reading guide

As the first step in determining the correct statistical disclosure control for a frequency table, it will have to be determined whether disclosure is possible. In the first instance, the basis for this is “common sense”: is there information present in the table that may not be disclosed about individual respondents? For frequency tables, such information can be hidden in the spanning variables. Part of the spanning variables can be considered as identifying variables and the rest as sensitive. With respect to sensitivity, an additional distinction is made in regard to the degree of sensitivity. See Hundepool et al. (2006) for more information.

Once risky situations have been identified, the table will generally have to be further protected. There are three general methods available for this purpose: restructuring the table (see section 4.4), suppression (see section 4.5) and rounding (see section 4.6).

Which method or combination thereof is ultimately used in a specific situation cannot be determined in advance. This depends significantly on the intended users. For example, in Eurostat regulations, it is not always possible to restructure the table, and cell suppression or rounding will often have to be selected. The department responsible for the frequency table concerned is also responsible for the adequate statistical disclosure control thereof. When selecting the method or methods to be used, two competing aspects must be taken into account:

- disclosure risk;
- information loss.

In general, it can be said that reducing the disclosure risk will lead to increased information loss. The converse is also true: the smaller the information loss, the larger the disclosure risk. In certain cases, an assessment will have to be made, in which, at a minimum, the rules in the Statistical Disclosure Control Handbook (Hundepool et al., 2006) will always have to be satisfied.

4.2 Scope and relationship with other themes and subthemes

This subtheme discusses methods that are used for the statistical disclosure control of frequency tables. This chapter does *not* discuss any methods that are only used for quantitative tables. For these methods, see the subtheme “Statistical Disclosure Control of Quantitative Tables”.

A few of the methods described in this chapter can, in principle, be used for both quantitative tables and frequency tables. Such methods will be repeated in the subtheme “Statistical Disclosure Control of Quantitative Tables”.

The methods in this chapter can be divided into two variants: methods for determining the primary (or other) risky cells of a frequency table and methods for making tables with risky cells suitable for publication.

4.3 Temporarily standardisation of a frequency table

4.3.1 Short description

The goal of statistical disclosure control is to prevent the disclosure of information about individual contributors to a table, or at least to make this more difficult. To achieve this, the cells will first have to be identified where there is risk of a possible disclosure. In frequency tables, at least two aspects play a role in this: recognisable groups and sensitive variables. Stated briefly, a risky situation occurs in a frequency table either if a cell that corresponds to a recognisable group of respondents contains too few respondents, or if the distribution of the respondents from a recognisable group is too concentrated in one or two categories. The Statistical Disclosure Control Handbook (Hundepool et al., 2006) sets out what Statistics Netherlands policy means by “too small” and “too concentrated”.

To do this, it is useful to look at the frequency table in a standard format. In some cases, the frequency table will already be available in this format. In others, it will be necessary to temporarily convert it. Once the frequency table has been protected, it can be restructured in its original format.

4.3.2 Applicability

Before a frequency table can be statistically protected, it will first have to be indicated where possible problems arise in that table. For this purpose, it is convenient to temporarily look at the frequency table in a standard way, so that a clear distinction is visible between identifying and sensitive variables.

4.3.3 Detailed description

To detect risky situations in frequency tables, it is necessary to divide the spanning variables into identifying variables and sensitive variables. The qualification for each variable can, in principle, be determined by the department concerned. To promote coordination between the different frequency tables to be published, it is a good idea to keep track of this centrally.

Next, the frequency table can be restructured temporarily so that the identifying variables are included in the left-hand column and the categories of the sensitive variables are present in the other columns. This must also involve “hidden” variables that define the population or subpopulation that is the subject of the frequency table.

The rules as referred to in Hundepool et al. (2006) can then be applied to the frequency table standardised in this way.

4.3.4 Example

Suppose that Table 1 is a frequency table of the number of people who, in a certain year, have died from a non-natural death, as stated in a publication³.

³ The figures in the table are fictitious.

Type of non-natural death	Gender	Age						
		Total	<15	15-<20	20-<40	40-<60	60-<80	>=80
Suicide	Total	1530	8	43	418	674	298	89
	Man	1027	8	34	297	453	181	54
	Woman	503	-	9	121	221	117	35
Murder and manslaughter	Total	141	11	13	74	34	8	1
	Man	96	9	5	47	32	2	1
	Woman	45	2	8	27	2	6	-
Traffic accident	Total	880	47	120	380	67	179	87
	Man	636	23	87	315	52	98	61
	Woman	244	24	33	65	15	81	26
Workplace accident	Total	81	-	3	30	42	6	-
	Man	79	-	3	28	42	6	-
	Woman	2	-	-	2	-	-	-
Personal accident	Total	2013	64	6	120	60	481	1282
	Man	834	32	2	100	56	223	421
	Woman	1179	32	4	20	4	258	861
Other/unknown	Total	110	2	6	24	8	37	33
	Man	63	1	4	18	7	20	13
	Woman	47	1	2	6	1	17	20
Total	Total	4755	132	191	1046	885	1009	1492
	Man	2735	73	135	805	642	530	550
	Woman	2020	59	56	241	243	479	942

Table 1: Number of people who died from a non-natural death in year J

The standardised form of this frequency table for statistical disclosure control is created by placing the identifying variables in the left-hand column and the sensitive variables in the other columns. In Table 1, the identifying variables are “Gender” and “Age”. The sensitive variable is the variable “Type of non-natural death”. The standardised version of Table 1 is presented in Table 2.

Gender	Age	Type of non-natural death						Total
		Suicide	Murder and manslaughter	Traffic accident	Workplace accident	Personal accident	Other / Unknown	
Man	<15	8	9	23	-	32	1	73
	15-<20	34	5	87	3	2	4	135
	20-<40	297	47	315	28	100	18	805
	40-<60	453	32	52	42	56	7	642
	60-<80	181	2	98	6	223	20	530
	>=80	54	1	61	-	421	13	550
	Total	1027	96	636	79	834	63	2735
Woman	<15	-	2	24	-	32	1	59
	15-<20	9	8	33	-	4	2	56
	20-<40	121	27	65	2	20	6	241
	40-<60	221	2	15	-	4	1	243
	60-<80	117	6	81	-	258	17	479
	>=80	35	-	26	-	861	20	942
	Total	503	43	220	2	1147	46	1961
Total	<15	8	11	47	-	64	2	132
	15-<20	43	13	120	3	6	6	191
	20-<40	418	74	380	30	120	24	1046
	40-<60	674	34	67	42	60	8	885
	60-<80	298	8	179	6	481	37	1009
	>=80	89	1	87	-	1282	33	1492
	Total	1530	139	856	81	1981	109	4696

Table 2: Standardised version of Table 1

4.4 Table restructuring

4.4.1 Short description

Section 4.3 describes how risky situations can be discovered in frequency tables by temporarily looking at the table in standardised form. If risky cells are subsequently found, the table will have to be protected before it can be published. An initial option to make a frequency table with risky cells suitable for publication is to restructure the table. By combining categories, the content per cell is increased. This affects the distribution of the sensitive spanning variable(s) among the various categories. This method can also be used to increase the content per recognisable group.

4.4.2 Applicability

This method will generally lead to fewer risky cells occurring in the table. By combining rows and/or columns, cells are combined and the content per cell is increased. The distribution of the sensitive spanning variable(s) among the various categories is also affected as a result.

There are no methodological conditions for using this method. However, externally imposed delivery obligations sometimes specify what detail level of a table must be published. This may be a Eurostat obligation, but Statistics Netherlands policy can also imply that a certain detail level of a table must be published. In these cases, the method can be used from a technical perspective, but this is prevented by external or other policy decisions.

4.4.3 Detailed description

The standardised version of the frequency table is used to determine whether a risky situation is present. The restructuring can take place in two ways:

- a. Restructuring the original table
- b. Restructuring the standardised version of the table.

If option a is selected, the table will have to be examined again in the standardised form after restructuring to determine whether the disclosure control rules have been satisfied. The standardised form will also have to be used to determine which risky cells must be dealt with. In the case of option b, it is immediately clear which cells must be dealt with, but the restructuring will still have to be converted to the original table.

4.4.4 Example

Table 2 shows that the distribution of the respondents among the various categories of the sensitive variable does not satisfy the disclosure control rules in two rows. These rules require that a cell may not contain a concentration of nearly all the

respondents. The cell (Woman, 80+, Personal accident) contains 91% of the total group of women aged 80+ who died a non-natural death, and the cell (Woman, 40-60, Suicide) also contains 91% of the total group of women between 40 and 60 years of age who died a non-natural death.

Based on the original table, the choice could be made to not split the causes of death “Suicide” and “Personal accident” in terms of gender.

Based on the standardised form, the choice could be made to condense the ages categories to “<15”, “15-<20”, “20-<60” and “>=60”. This creates a table where there is no longer a strongly concentrated distribution of recognisable groups in a single cell. See Table 3 for the associated table in its original form.

Type of non-natural death	Gender	Age				
		Total	<15	15-<20	20-<60	>=60
Suicide	Total	1530	8	43	1092	387
	Man	1027	8	34	750	235
	Woman	503	-	9	342	152
Murder and manslaughter	Total	141	11	13	108	9
	Man	96	9	5	79	3
	Woman	45	2	8	29	6
Traffic accident	Total	880	47	120	447	266
	Man	636	23	87	367	159
	Woman	244	24	33	80	107
Workplace accident	Total	81	-	3	72	6
	Man	79	-	3	70	6
	Woman	2	-	-	2	-
Personal accident	Total	2013	64	6	180	1763
	Man	834	32	2	156	644
	Woman	1179	32	4	24	1119
Other/unknown	Total	110	2	6	32	70
	Man	63	1	4	25	33
	Woman	47	1	2	7	37
Total	Total	4755	132	191	1931	2501
	Man	2735	73	135	1447	1080
	Woman	2020	59	56	484	1421

Table 3: Protected version of Table 1

4.5 Suppression

4.5.1 Short description

A frequently used method to protect primary risky cells is to suppress (not publish) certain cells. The cell value is then simply replaced by an X.

In a frequency table where the totals and subtotals are also provided, however, it is often not sufficient to suppress only the primary risky cells. After all, if a suppressed cell is the only suppressed cell in a row, the suppressed value is simple to calculate by subtracting the other cell values in that row from the associated marginal.

To sufficiently protect primary risky cells, it is therefore also necessary to suppress other cells which, in themselves, are safe. This is called *secondary suppression*. It is not easy to perform this in such a way such that the primary risky cells are protected sufficiently, while also ensuring that not too much information is removed from the

table. Furthermore, account must also be taken of the fact that structural zero cells cannot be used as secondary suppressions: everyone knows that, by definition, these cells are empty.

To prevent a situation where suppressed, primary risky cells can be calculated exactly, secondary suppressions are therefore necessary. However, what again plays a role here is that a “too accurate” estimation for a suppressed cell is not desirable. There is little difference between the statements “This suppressed cell actually has a value of 10000” and “This suppressed cell actually has a value of between 9998 and 10002”. Given a suppression pattern, it is always⁴ possible to calculate an interval within which a suppressed cell must occur. The method of “Cell Suppression” must therefore produce a suppression pattern for which the intervals to be calculated are sufficiently large.

Fischetti and Salazar (2000) have developed a method to solve the above problem in an optimal manner. Their method is, in theory, applicable to arbitrary, additive tables with non-negative contributors. In practice, however, their solution involves too much computing time if the tables become too large, either in size or complexity. This is why a number of suboptimal methods have been developed to find suitable suppression patterns for larger and/or more complex tables.

For example, the “modular approach” (HiTaS) splits a hierarchical table into a large number of non-hierarchical subtables and applies the optimal method to each individual subtable. By correctly combining the results, a suboptimal solution can be obtained for the entire table, with a significantly shorter computing time.

The “hypercube approach” can also protect large tables by protecting the subtables in a certain iterative way. The protection of each subtable also takes place suboptimally. Consequently, the approach is relatively fast, but, in general, more cells are suppressed than strictly necessary to obtain a protected table.

4.5.2 Applicability

Risky situations in frequency tables can be divided into two cases:

- a. The recognisable group is too small;
- b. The distribution of the recognisable group among the sensitive variable(s) is too concentrated in a single sensitive cell.

To determine a suitable suppression pattern, it is necessary to know how one can comply with the disclosure control rules imposed. In many algorithms, so-called safety intervals are used for this purpose. These are the minimum intervals for primary suppressed cells that should arise from the suppression pattern. At present, contrary to the case of quantitative tables, no method is available to calculate the minimum intervals for primary risky cells in frequency tables. The method as described in Fischetti and Salazar (2000) is therefore not directly applicable as yet.

⁴ In the case that the table is composed of contributors within a certain range (for example, non-negative) and the marginals are also provided.

4.5.3 Detailed description

If a row total in the standardised form of the table is too small (the recognisable group is too small), this cell will have to be suppressed. Of course, multiple cells will have to be suppressed to prevent the row total from being calculated. In general, this will mean that the total row will have to be suppressed, including a second possibly “safe” row.

A second situation that may arise is a sufficiently large row total which, however, is too concentrated in a single sensitive category of the variable. In that case, the row total is, in principle, suitable for publishing. The cell associated with the category of the sensitive variable in which the respondents are concentrated can then be viewed as the primary cell to be suppressed. In a table with totals and subtotals, one must also look for secondary cells to be suppressed. In many algorithms, safety intervals are used for this purpose. These are the minimum intervals for primary suppressed cells that should follow from the suppression pattern. At present, contrary to the case for quantitative tables, no method is available to calculate minimum intervals for primary risky cells in frequency tables. The method as described in Fischetti and Salazar (2000) is therefore also not directly applicable as yet.

An additional problem is formed by what is called “meaningful aggregates” in the disclosure control rules. If multiple cells in a row are suppressed, the total of these suppressed cells is actually published. If the suppressed cells form a meaningful aggregate, then the respondents may also not be too concentrated in that combined cell. Account should therefore be taken of this when determining secondary suppressions. It is not yet clear if the Fischetti and Salazar model (2000) is general enough to take this into account.

4.5.4 Example

Table 4 shows a suppression pattern in which it is assumed that the aggregate “Suicide” + “Personal accident” is not a “meaningful aggregate”. Both problematic cells are suppressed by placing Xs in the cells.

Type of non-natural death	Gender	Age						
		Total	<15	15-<20	20-<40	40-<60	60-<80	>=80
Suicide	Total	1530	8	43	418	674	298	89
	Man	1027	8	34	297	×	181	×
	Woman	503	-	9	121	×	117	×
Murder and manslaughter	Total	141	11	13	74	34	8	1
	Man	96	9	5	47	32	2	1
	Woman	45	2	8	27	2	6	-
Traffic accident	Total	880	47	120	380	67	179	87
	Man	636	23	87	315	52	98	61
	Woman	244	24	33	65	15	81	26
Workplace accident	Total	81	-	3	30	42	6	-
	Man	79	-	3	28	42	6	-
	Woman	2	-	-	2	-	-	-
Personal accident	Total	2013	64	6	120	60	481	1282
	Man	834	32	2	100	×	223	×
	Woman	1179	32	4	20	×	258	×
Other/unknown	Total	110	2	6	24	8	37	33
	Man	63	1	4	18	7	20	13
	Woman	47	1	2	6	1	17	20
Total	Total	4755	132	191	1046	885	1009	1492
	Man	2735	73	135	805	642	530	550
	Woman	2020	59	56	241	243	479	942

Table 4: Suppression pattern for the protection of Table 1

4.6 Additive rounding

4.6.1 Short description

In frequency tables, rounding is a rather natural method of disclosure control. First of all, the exact cell values are only known in a certain interval when rounding is used. The extent to which rounding is performed will, of course, have an impact on the size of the intervals. Second, an unrounded frequency table creates the impression of great precision: the counting has been performed down to the individual units. In the case of estimated frequencies, this is false precision. Rounding also cuts down on this false precision.

If each cell is rounded independently, the additivity of the table will not necessarily be maintained. Of course, there is a simple way to guarantee the additivity: by rounding the cells in the interior of the table independently of one other and then recalculating the marginals. As a result, however, the marginals can deviate significantly from the rounded or unrounded original values.

In additive rounding, the table is rounded such that the additivity is maintained and that the rounded table deviates from the original as little as possible. The size of the rounding base determines the extent to which the frequency table is protected: the larger the rounding base, the greater the protection will generally be. At present, no method is available to determine the correct rounding base.

4.6.2 Applicability

Additive rounding can be used for the statistical disclosure control of both quantitative tables and frequency tables. Often, a presentation argument will also

play a role: a large number of significant figures suggests a high degree of precision that is not always justified because of sampling errors and measurement errors. Rounding the table values reduces this false precision to a certain extent.

4.6.3 Detailed description

In additive rounding, the cell values in a table are rounded to multiples of a rounding base b , keeping the totals and subtotals in the table equal to the sum of the corresponding parts.

Oftentimes, additive rounding is performed in a “zero restricted” manner. In other words, cell values that are already a multiple of the rounding base are not changed, while the other cell values are rounded to one of the adjacent multiples of that rounding base. The rounded values are selected such that the sum of the absolute deviations of the cell values in the rounded table with respect to the cell values in the original table is minimised, under the restriction that the rounded table remains additive. As a result, it is possible that cell values are not rounded to the closest multiple of the rounding base.

In certain conditions, it is not possible to construct a rounded table under the scenario described above. In that case, the restriction that rounding is performed to one of the adjacent multiples of the rounding base is weakened by allowing a cell value to also be rounded to non-adjacent multiples of the rounding base. This weakening can be limited slightly by determining a maximum for the number of steps that may exist between the rounded value and the original value.

In the case of “zero restricted” additive rounding using rounding base $b > 0$ for the non-negative number $z = ub + r$, where $0 \leq r < b$, rounding is performed on the number a , such that

$$a \in \{ub, (u + 1_{(0,b)}(r))b\} \quad (4.6.1)$$

where $1_{(0,b)}(r)$ is equal to 1 if $r \in (0, b)$ and equal to 0 if $r = 0$.

Therefore, in the case that $r = 0$, a is always rounded to ub and, in the case that $r \in (0, b)$, a is always rounded to ub or to $(u + 1)b$.

If, however, the restriction is weakened by a maximum of $K > 0$ steps further than the adjacent multiples of the rounding base, then rounding is performed on the number a , such that

$$a \in \{(0 \vee (u + j))b \mid j = -K, \dots, (K + 1_{(0,b)}(r))\} \quad (4.6.2)$$

where $x \vee y = \max(x, y)$.

Multiple additive rounded versions may exist for a given table. These are all *feasible* tables. The table closest to the original table can subsequently be selected from the feasible tables. In τ -ARGUS (see Hundepool et al. 2007), the distance that is minimised is represented by

$$\sum_{i=1}^N |z_i - a_i| \quad (4.6.3)$$

where N is the number of cells in the table (including all totals and subtotals), z_i the cell values in the original table and a_i the corresponding rounded cell values.

Finding the optimal solution is a problem that requires intensive calculation (NP-complete). For large tables, this can result in unacceptably long calculation times. Partitioning is built into τ -ARGUS for this reason: a large table can be split into a number of subtables that are rounded individually. After these subtables are rounded, they are combined, calculating (if necessary) the totals and subtotals in question from the rounded parts.

4.6.4 Example

Table 5 shows a rounded version of Table 1, where additive rounding is performed using rounding base 50.

Type	Gender	Age						
		Total	<15	15-<20	20-<40	40-<60	60-<80	>=80
Suicide	Total	1550	0	50	400	700	300	100
	Man	1050	0	50	300	450	200	50
	Woman	500	-	0	100	250	100	50
Murder and manslaughter	Total	150	0	0	100	50	0	0
	Man	100	0	0	50	50	0	0
	Woman	50	0	0	50	0	0	-
Traffic accident	Total	850	50	150	350	50	200	50
	Man	600	0	100	300	50	100	50
	Woman	250	50	50	50	0	100	0
Workplace accident	Total	100	-	0	50	50	0	-
	Man	100	-	0	50	50	0	-
	Woman	0	-	-	0	-	-	-
Personal accident	Total	2000	50	0	150	50	450	1300
	Man	850	50	0	100	50	200	450
	Woman	1150	0	0	50	0	250	850
Other/unknown	Total	100	0	0	0	0	50	50
	Man	50	0	0	0	0	50	0
	Woman	50	0	0	0	0	0	50
Total	Total	4750	100	200	1050	900	1000	1500
	Man	2750	50	150	800	650	550	550
	Woman	2000	50	50	250	250	450	950

Table 5: Rounded version of Table 1, with rounding base 50. An "0" is a rounded 0, a "-" is an empty cell

5. Statistical Disclosure Control of Analysis Results

5.1 General description and reading guide

5.1.1 General description

In addition to problems and methods for protecting microdata, quantitative tables and frequency tables described in the previous paragraphs, there is still a very broad, diverse group of statistical output. This concerns the results of various types of statistical analyses and model estimations. In principle, these results also run a risk of disclosing the data for individual respondents and must therefore be treated with care. The risk of disclosure is present particularly in the case of outliers. When determining whether these results are sufficiently safe, the underlying frequency tables are often examined. There is often a strong correlation between the analysis model and an underlying frequency table. In the Statistical Disclosure Control Handbook (Hundepool et al., 2006), a start has been made for the protection of analysis results.

5.1.2 Reading guide

The problem of determining whether the results of statistical analyses are sufficiently safe arises mainly when checking the output of OnSite working and Remote Access. This is where many statistical analyses are performed on unprotected data, while the users would very much like to use and publish the results of their research outside of Statistics Netherlands. Checking the output is a necessary part of this valued service from Statistics Netherlands and statistical offices in general. With respect to checking the output, it does not matter at all whether the output is obtained through OnSite or Remote Access. In both cases, the same analyses are performed on the same data files using the same tools (SPSS, SAS etc.).

Because this problem does not only occur at Statistics Netherlands, but actually at every statistical bureau in Europe, it was decided to make this a subject of the Statistical Disclosure Control ESSnet project (2008-2009). The ESSnet is subsidised by Eurostat. Statistics Netherlands provided the project manager for this project. One of the tasks in the ESSnet project was to draw up guidelines for checking output. For this subject in the Methods Series, use is also made of these “Guidelines for Output Checking”, which can be found on the ESSnet website (http://neon.vb.cbs.nl/casc/ESSnet/guidelines_on_outputchecking.pdf).

The subject is still under development. As such, the project group does not regard the current version as the final one, but instead as a very useable first version. The project group hopes that it will be able to continue its work in this area in a subsequent project.

Due to the diversity of the problem, both in terms of the number of possible analysis methods and the number of different statistical packages, each with their own forms of output, it is not possible to develop ready-made software for this purpose.

5.2 Scope and relationship with other themes and subthemes

This subtheme discusses methods that can be used to determine whether the results of statistical analyses are sufficiently safe. To a large extent, use is made of the results of a European project group that drew up these guidelines. These guidelines also discuss tables. However, because these subjects have already been covered in the previous chapters, these subjects from the guidelines are less relevant here.

5.3 Disclosure control of analysis results

The methods for the protection of analysis results tie in with these European guidelines.

A number of considerations played a role when drawing up the guidelines for output checking. Of course, it is not possible to fully discuss all possible forms of output. The number of different methods available in SAS and SPSS is so large that it is impossible to assess all of these methods with respect to their possible disclosure risks. Just consider the size of the SPSS and SAS documentation.

Another aspect that plays an important role in the guidelines is their feasibility in practice. In assessing output, we must take account of two possible errors: first, incorrectly approving risky results and, second, incorrectly holding back safe results.

In the guidelines, two methods are provided for each subject. A “Rule of Thumb”, which primarily minimises the first error, and a “principles-based” rule that tries to minimise both errors.

The idea behind this distinction is that much of the research output can easily be handled by the simple rule. If the output is not allowed due to the “Rule of Thumb”, and the researcher wants it approved anyway, extra work must be performed (also by the researcher) to demonstrate that the results are indeed safe.

Here is a list of the types of output that are currently discussed in the guidelines:

Descriptive statistics	Frequency tables
	Magnitude tables
	Maxima, minima and percentiles (incl. median)
	Mode
	Means, indices, ratios, indicators
	Concentration ratios
	Higher moments of distributions (incl. variance, covariance, kurtosis, skewness)
	Graphs: pictorial representations of actual data
Correlation and Regression Analysis	Linear regression coefficients
	Non-linear regression coefficients
	Estimation residuals
	Summary and test statistics from estimates (R^2 , χ^2 etc.)
	Correlation coefficients

For the rest, please refer to the European Guidelines.

6. References

- CBS (2004), *Wet op het Centraal Bureau voor de Statistiek*, Staatsblad 2004, 695.
See also:
<http://www.cbs.nl/nl-NL/menu/organisatie/corporate-informatie/default.htm>
- De Wolf, P.P. (2002), *HiTaS: A Heuristic Approach to Cell Suppression in Hierarchical Tables*, In: 'Inference Control in Statistical Databases' Domingo-Ferrer (Ed.), LNCS 2316, Springer-Verlag Berlin Heidelberg, pp. 74-82.
- De Wolf, P.P. (2006), *Risk, Utility and PRAM*, in 'Privacy in Statistical Databases 2006', Domingo-Ferrer, J. and Franconi, L. (Eds.), LNCS 4302, Springer-Verlag, Berlin Heidelberg, pp. 189-204.
- Duncan, G.T., Jabine, T.B. and V.A. de Wolf (Eds.) (1993), *Private lives and public policies: confidentiality and accessibility of government statistics*. The National Academies Press, ISBN 0309086515.
- Fischetti, M. and Salazar Gonzales, J.J. (2000), *Models and Algorithms for Optimizing Cell Suppression in Tabular Data with Linear Constraints*, Journal of the American Statistical Association, vol. 95, pp. 916-928.
- Giessing, S. and Repsilber, D. (2002), *Tools and Strategies to Protect Multiple Tables with the GHQUAR Cell Suppression Engine*, In: 'Inference Control in Statistical Databases' Domingo-Ferrer (Ed.), LNCS 2316, Springer-Verlag Berlin Heidelberg, pp. 181-192.
- Gouweleeuw, J.M., Kooiman, P., Willenborg, L.C.R.J. and De Wolf, P.P. (1998a), *Post Randomisation for Statistical Disclosure Control: Theory and Implementation*, Journal of Official Statistics, vol. 14, 4, pp. 463-478.
- Gouweleeuw, J.M., Kooiman, P., Willenborg, L.C.R.J. and De Wolf, P.P. (1998b), *The post randomisation method for protecting microdata*, Qüestió, Quaderns d'Estadística i Investigació Operativa, vol. 22, 1, pp. 145-156.
- Hundepool, A., Jonker, J., Nobel, J., Schulte Nordholt E. and De Wolf, P.P. (2006), *Handboek Statistische Beveiliging*, Statistics Netherlands, Voorburg.
- Hundepool, A., Van de Wetering, A., Ramaswamy, R., Franconi, L., Poletini, S., Capobianchi, A., de Wolf, P.P., Domingo, J., Torra, V., Brand, R. and Giessing, S. (2007), *μ -ARGUS user manual 4.1*, Statistics Netherlands, Voorburg.
- Hundepool, A., Van de Wetering, A., Ramaswamy, R., De Wolf, P.P., Giessing, S., Fischetti, M., Salazar, J.J., Castro, J. and Lowthian, P. (2007), *τ -ARGUS user manual 3.2*, Statistics Netherlands, Voorburg.
- Loeve, A. (2001), *Notes on sensitivity measures and protection levels*, Internal report, Statistics Netherlands, Voorburg.

- Ritchie, F., Welpton, R., Franconi, L., Lucarelli, M., Seri, G., Brandt, M., Guerke, C., Hundepool, A.J. and Mol, J. (2010), *Guidelines for the checking of output based on microdata research*, ESSnet-SDC-project.
http://neon.vb.cbs.nl/casc/ESSnet/guidelines_on_outputchecking.pdf
- Ronning, G., Rosemann, M. and Strotmann, H. (2004), *Estimation of the probit model using anonymized micro data*, Paper prepared for the 'European Conference on Quality and Methodology in Official Statistics (Q2004)', Mainz, 24–26 May 2004.
- Van den Hout, A. (1999), *The analysis of data perturbed by pram*, Delft University Press, Delft.
- Van den Hout, A. and Van der Heijden, P.G.M. (2002), *Randomized response, statistical disclosure control and misclassification: a review*, *International Statistical Review* 70(2), pp. 269-288.

Version history

Version	Date	Description	Authors	Reviewers
Dutch version: Statistische beveiliging				
0.5	02-04-2007	First Dutch version	Anco Hundepool Peter-Paul de Wolf	John Schalen Eric Schulte Nordholt
1.1	23-01-2008	Minor modifications to layout	Anco Hundepool Peter-Paul de Wolf	
1.2	15-02-2010	Chapters on frequency tables and analysis results added	Anco Hundepool Peter-Paul de Wolf	John Schalen Eric Schulte Nordholt
English version: Statistical Disclosure Control				
1.2E	19-07-2011	First English version	Anco Hundepool Peter-Paul de Wolf	