

# $\tau$

---

# ARGUS

Version 3.5



## User's Manual

Project: ESSnet-project  
Date: October 2011  
Revised version  
BPA no: 769-02-TMO

Statistics Netherlands  
P.O. Box 24500  
2490 HA The Hague  
The Netherlands  
email: [aj.hundepool@cbs.nl](mailto:aj.hundepool@cbs.nl)

Contributors: Anco Hundepool, Aad van de Wetering and Ramya Ramaswamy.  
Peter-Paul de Wolf (Modular), Sarah Giessing (GHMiter, audit)  
Matteo Fischetti, and Juan-José Salazar (Optimisation)  
Jordi Castro (Network solutions), Philip Lowthian (manual)



## Contents

1	Introduction.....	3
1.1	Preface.....	3
1.2	About the name ARGUS.....	3
1.3	Contact.....	4
1.4	Acknowledgments.....	4
1.5	The CASC-project.....	5
1.6	CASC-partners.....	5
1.7	The CENEX project.....	6
1.8	The ESSNet project.....	7
1.9	Latest improvements.....	7
1.10	The structure of this manual.....	7
2	Producing Safe tables.....	8
2.1	Introduction.....	8
2.2	Sensitive cells in magnitude table.....	8
2.3	Sensitive cells in frequency count tables.....	10
2.4	Table redesign.....	10
2.5	Secondary cell suppression.....	11
2.6	Information loss in terms of cell costs.....	11
2.7	Series of tables.....	11
2.8	The Hypercube/GHMITER method.....	12
2.8.1	The hypercube method.....	12
2.8.2	The ARGUS implementation of GHMITER.....	13
2.8.3	References on GHMiter.....	15
2.9	Optimisation models for secondary cell suppression.....	16
2.10	The Modular approach.....	17
2.11	The modular approach for linked tables.....	21
2.12	Network solution for large 2 dimensional tables with one hierarchy.....	22
2.13	Controlled rounding.....	24
2.13.1	Restricted and non-restricted controlled rounding.....	24
2.13.2	Optimal, first feasible and RAPID solutions.....	25
2.13.3	Protection provided by controlled rounding.....	25
2.13.4	Choosing the parameters for Controlled Rounding.....	26
2.14	Audit.....	26
2.15	Functional design of $\tau$ -ARGUS.....	29
3	A tour of $\tau$ -ARGUS.....	30
3.1	Preparation.....	31
3.1.1	First steps.....	31
3.1.2	Open a microdata file.....	32
3.1.3	Specify metafile.....	32
3.1.4	Specify tables.....	35
3.2	The Process of disclosure control.....	38
3.2.1	View table.....	39
3.2.2	Save the safe table.....	47
4	Reference Section - Description of the Menu Items.....	49
4.1	Main Window.....	49
4.2	The File menu.....	50
4.2.1	File   Open Microdata.....	50
4.2.2	File   Open Table.....	52
4.2.3	File   Open Table Set.....	54
4.2.4	File   Open Batch Process.....	54
4.2.5	File   Exit.....	55
4.3	The Specify menu.....	55

4.3.1	Specify   Metafile [for microdata].....	55
4.3.2	Specify   Metafile [SPSS System files].....	58
4.3.3	Specify   Metafile [for tabular data].....	59
4.3.4	Specify   Specify Tables [for microdata] .....	61
4.3.5	Specify   Specify tables [for tabular data].....	69
4.4	The Modify menu.....	70
4.4.1	Modify   Select Table.....	70
4.4.2	Modify   View Table.....	71
4.4.2.1	The View table screen.....	71
4.4.2.2	Global recoding.....	75
4.4.2.3	Secondary suppression.....	79
4.4.2.4	Controlled rounding.....	83
4.4.2.5	The audit procedure.....	87
4.4.2.6	The Options at the Bottom of the table .....	88
4.4.3	Modify   Linked Tables.....	91
4.5	The Output menu.....	94
4.5.1	Output   Save Table.....	94
4.5.2	Output   View Report.....	96
4.5.3	Output   GenerateApriory.....	97
4.5.4	Output   Write Batch File.....	99
4.6	The Help menu.....	100
4.6.1	Help   Contents.....	100
4.6.2	Help   News.....	100
4.6.3	Help   Options .....	100
4.6.4	Help   About.....	101
5	Further decriptions.....	103
5.1	Meta data files .....	103
5.1.1	Meta data for fixed format micro data .....	103
5.1.2	Meta data for free format micro data .....	106
5.1.3	Meta data for SPSS system files .....	107
5.1.4	Meta data for tabular data files .....	107
5.2	Hierarchy file.....	109
5.3	Codelist file .....	110
5.4	Global recode file .....	110
5.5	The apriori file.....	111
5.6	The Batch command file .....	113
5.7	Log file .....	117
6	Index .....	118

# 1 INTRODUCTION

---

## 1.1 Preface

---

This is the user manual for  $\tau$ -ARGUS version 3.4.  $\tau$ -ARGUS is a software tool designed to assist a data protector in producing safe tables. This manual describes the version of  $\tau$ -ARGUS at the end of the ESSNet-SDC project. With respect to the previous release of  $\tau$ -ARGUS, we have made many steps forward, and  $\tau$ -ARGUS now has facilities to protect hierarchical and some linked tables via improved apriority files. Also controlled rounding is available.

The purpose of  $\tau$ -ARGUS is to protect tables against the risk of disclosure, i.e. the accidental or deliberate disclosure of information related to individuals from a statistical table. This is achieved by modifying the table so that it contains less detailed information.  $\tau$ -ARGUS allows for several modifications of a table: a table can be redesigned, meaning that rows and columns can be combined; sensitive cells can be suppressed and additional cells to protect these can be found in some optimum way (secondary cell suppression).

$\tau$ -ARGUS is one of a twin set of disclosure control packages. Within the CASC-project a tool for the protection of microdata - called  $\mu$ -ARGUS - has also been developed, which is the twin brother of  $\tau$ -ARGUS.<sup>1</sup> This is manifest not only when one looks at the user interfaces of both packages, but also when one looks at the source code: the bodies of the twins are so much combined that they in fact are like Siamese twins!

## 1.2 About the name ARGUS

---

Somewhat jokingly the name ARGUS can be interpreted as the acronym of ‘Anti-Re-identification General Utility System’<sup>2</sup>. As a matter of fact, the name ARGUS was inspired by a myth of the ancient Greeks. In this myth Zeus has a girl friend named Io. Hera, Zeus’ wife, did not approve of this relationship and turned Io into a cow. She let the monster ARGUS guard Io. ARGUS seemed to be particularly well qualified for this job, because it had a hundred eyes that could watch over Io. If it would fall asleep only two of its eyes were closed. That would leave plenty of eyes to watch Io. Zeus was eager to find a way to get Io back. He hired Hermes who could make ARGUS fall asleep by the enchanting music on his flute. When Hermes played his flute to ARGUS this indeed happened: all its eyes closed, one by one. When Hermes had succeeded in making ARGUS fall asleep, ARGUS was decapitated. ARGUS’ eyes were planted onto a bird’s tail - a type of bird that we now know under the name of peacock. That explains why a peacock has these eye-shaped marks on its tail. This also explains the picture on the cover of this manual.

---

<sup>1</sup> See Anco Hundepool et al., 2010,  $\mu$ -ARGUS version 4.3 user’s manual, Statistics Netherlands, Voorburg, The Netherlands.

<sup>2</sup> This interpretation is due to Peter Kooiman, former head of the methodology department at Statistics Netherlands

It is a copperplate engraving of Gerard de Lairesse (1641-1711) depicting the process where the eyes of ARGUS are being removed and placed on the peacock's tail.<sup>3</sup>

Like the mythological ARGUS, the software is supposed to guard something, in this case data. This is where the similarity between the myth and the package is supposed to end, as we believe that the package is a winner and not a loser as the mythological ARGUS is.

### **1.3 Contact**

---

Feedback from users will help improve future versions of  $\tau$ -ARGUS and is therefore greatly appreciated. The authors of this manual can be contacted directly for suggestions that may lead to improved versions of  $\tau$ -ARGUS in writing or otherwise; e-mail messages can also be sent to [aj.hundepool@cbs.nl](mailto:aj.hundepool@cbs.nl).

### **1.4 Acknowledgments**

---

$\tau$ -ARGUS was started as part of the fourth framework SDC-project and in its current form the first version has been developed as part of the CASC project that was partly sponsored by the EU under contract number IST-2000-25069. This support is highly appreciated. The CASC (Computational Aspects of Statistical Confidentiality) project is part of the Fifth Framework of the European Union. The main part of  $\tau$ -ARGUS has been developed at Statistics Netherlands by Aad van de Wetering and Ramya Ramaswamy (who wrote the kernel) and Anco Hundepool (who wrote the interface). However this software would not have been possible without the contributions of several others, both partners in the CASC-project and outsiders. Recent extensions of  $\tau$ -ARGUS have been made possible during the European CENEX-SDC-project (grant agreement 25200.2005.001-2005.619) and the ESSNet-SDC project (grant agreement 25200.2005.003-2007.670.)

The German partners Statistisches Bundesamt (Sarah Giessing and Dietz Repsilber) have contributed the GHMITER software, which offers a solution for secondary cell suppression based on hypercubes. Peter-Paul de Wolf has built a search algorithm based on non-hierarchical optimal solutions. This algorithm breaks down a large hierarchical table into small non-hierarchical subtables, which are then individually protected. A team led by JJ Salazar of the University La Laguna Tenerife, Spain, has developed the optimisation routines. Additionally Jordi Castro has developed a solution based on networks.

The controlled rounding procedure has been developed by JJ Salazar in a project sponsored by ONS. In order to enhance the usability  $\tau$ -ARGUS now also can handle SPSS-system files. For using  $\tau$ -ARGUS in combination with SAS, several reports have been produced during the ESSnet project and are available from the CASC/ESSNet website.

---

<sup>3</sup> The original copy of this engraving is in the collection of 'Het Leidsch Prentenkabinet' in Leiden, The Netherlands.

For solving these optimisation problems,  $\tau$ -ARGUS uses commercial LP-solvers. Traditionally we use Xpress as an LP-solver. This package is kindly made available for users of  $\tau$ -ARGUS in a special agreement between the  $\tau$ -ARGUS-team and FICO, the developers of Xpress. Alternatively  $\tau$ -ARGUS can also use the Cplex-package. Users can choose either solver to link to  $\tau$ -ARGUS (provided, of course, they purchase a license for the solver chosen). However users already having a licence for one of these packages for other applications can use their current licence for  $\tau$ -ARGUS as well.

## 1.5 *The CASC-project*

---

The CASC project is the initiative in the 5<sup>th</sup> framework to explore new possibilities of Statistical Disclosure Control and to extend further the existing methods and tools. A key issue in this project is an emphasis more on practical tools, and the research needed to develop them. For this purpose a new consortium has been brought together. It has taken over the results and products emerging from the SDC-project. One of the main tasks of this new consortium was to further develop the ARGUS-software. The main software developments in CASC are  $\mu$ -ARGUS, the software package for the disclosure control of microdata, while  $\tau$ -ARGUS handles tabular data.

The CASC-project has involved both research and software development. As far as research is concerned, the project has concentrated on those areas that were expected to result in practical solutions, which can then be built into the software. Therefore the CASC-project has been designed round this software twin ARGUS. This will make the outcome of the research readily available for application in the daily practice of statistical institutes.

## 1.6 *CASC-partners*

---

At first sight the CASC-project team had become rather large. However there is a clear structure in the project, defining which partners are working together for which tasks. Sometimes groups working closely together have been split into independent partners only for administrative reasons.

<b>Institute</b>	<b>Short</b>	<b>Country</b>
1. Statistics Netherlands	CBS	NL
2. Istituto Nazionale di Statistica	ISTAT	I
3. University of Plymouth	UoP	UK
4. Office for National Statistics	ONS	UK
5. University of Southampton	SOTON	UK
6. The Victoria University of Manchester	UNIMAN	UK
7. Statistisches Bundesamt	StBA	D
8. University La Laguna	ULL	ES
9. Institut d'Estadística de Catalunya	IDESCAT	ES
10. Institut National de Estadística	INE	ES
11. TU Ilmenau	TUilm	D

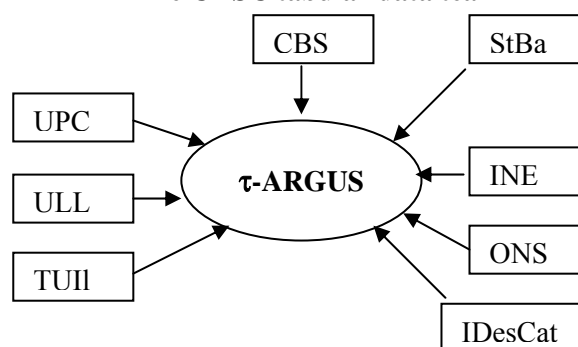
12. Institut d'Investigació en Intel·ligència Artificial-CSIC	CIS	ES
13. Universitat Rovira i Virgili	URV	ES
14. Universitat Politècnica de Catalunya	UPC	ES

Although Statistics Netherlands is the main contractor, the management of this project is a joint responsibility of the steering committee. This steering committee constitutes of 5 partners, representing the 5 countries involved and also bearing a responsibility for a specific part of the CASC-project:

#### CASC Steering Committee

Institute	Country	Responsibility
Statistics Netherlands	Netherlands	Overall manager Software development
Istituto Nazionale di Statistica	Italy	Testing
Office for National Statistics	UK	
Statistisches Bundesamt	Germany	Tabular data
Universitat Rovira i Virgili	Spain	Microdata

#### The CASC tabular data team



## 1.7 The CENEX project

In 2006 thanks to the CENEX project on SDC further extensions on the ARGUS software have been made. The CENEX (Centres of Excellence in Statistics) was an initiative by Eurostat to further foster the cooperation between the NSI's in Europe. The emphasis was more on using and promoting the current state of affairs by organising courses, meetings. Also further testing of ARGUS has led to several improvements. The batch-version is becoming more and more important, opening options to link  $\tau$ -ARGUS to packages like SuperCross, but also to facilitate the building blocks project of Eurostat. Many improvements (including the error checking) of the batch-version have made. A new, special output format for transmitting data to Eurostat (the SBS-format) has been added. Recently also the controlled rounding has been included in  $\tau$ -ARGUS.

The CENEX project was led by Statistics Netherlands. The NSI's of Italy, Germany, UK, Sweden, Slovenia, Estonia and Austria participated. The University Rovira i Virgili, Tarragona, also participated.

## 1.8 *The ESSNet project*

---

The ESSNet project, that can be seen as a continuation/follow-up of the CENEX project brought the development of  $\tau$ -ARGUS another step forward. In this first release within the ESSNet we have corrected several errors and inconsistencies and also introduced the option of reading data from a SPSS system file. Also the audit option has been updated by including a very much improved version of the audit/routine originally developed by the Illmenau-team during the CASC-project. The original version proved to be slow, but we achieved a big improvement here. This version now works with both CPLEX and XPress.

The testing of the members of the ESSNet team has contributed to several smaller improvements.

## 1.9 *Latest improvements*

---

The latest extensions in version 3.5 of  $\tau$ -ARGUS are :

- the functionality to protect a set of linked tables, both with the hypercube and with modular.
- The method for additional protection of singletons has been improved. The old method might lead to some overprotection.
- Tables with negative values can be protected now with the modular approach.

## 1.10 *The structure of this manual*

---

The remaining part of this manual consists of four chapters and an index.

In chapter 2 we will give a short introduction to the theory and methodology. However for a more fundamental description we refer to the handbook on Statistical Disclosure Control<sup>4</sup>. This handbook is the result of the CENEX and the ESSNet projects. Where appropriate we will refer to this handbook. The CENEX-project made a start with writing a handbook on Statistical Disclosure Control. In the ESSNet project the work on the handbook has been continued.

In Chapter 3 a short tour of ARGUS will be given as a first impression of the program.

Chapter 4 is the reference manual of  $\tau$ -ARGUS. It will describe in detail the program. This chapter is organized by the menu items of  $\tau$ -ARGUS.

Chapter 5 gives details of files used by  $\tau$ -ARGUS

The manual is concluded with an index.

---

<sup>4</sup> Hundepool, Anco, Josep Domingo-Ferrer, Luisa Franconi, Sarah Giessing, Rainer Lenz, Jane Longhurst, Eric Schulte Nordholt, Giovanni Seri, Peter-Paul De Wolf (2009), ESSNet handbook on Statistical Disclosure Control, <http://neon.vb.cbs.nl/casc/handbook.htm>

## 2 PRODUCING SAFE TABLES

---

### 2.1 Introduction

---

The growing demands from researchers, policy makers and others for more and more detailed statistical information lead to a conflict. Statistical offices collect large amounts of data for statistical purposes. The respondents are only willing to provide the statistical offices with the required information if they can be certain that these statistical offices will treat their data with the utmost care. This implies that respondents' confidentiality must be guaranteed. This imposes limitations on the amount of detail in the publications. Practice and research have generated insights into how to protect tables, but the problem is not yet definitively solved.

Before we go into more details, the basic ideas on which  $\tau$ -ARGUS is based, we give a sketch of the general ideas. At first sight one might find it difficult to understand that information presented in tabular form presents a disclosure risk. After all, one might say that the information is presented only in aggregate form.

Safe tables are produced from unsafe ones by applying certain SDC measures to the tables. These SDC measures - as far as they are implemented in  $\tau$ -ARGUS - are discussed in the present section. Some key concepts such as sensitive cells, information loss and the like are discussed as well.

### 2.2 Sensitive cells in magnitude table<sup>5</sup>

---

The well-known dominance rule is often used to find the sensitive cells in tables, i.e. the cells that can not be as they might reveal information on individual records. More particularly, this rule states that a cell of a table is unsafe for publication if a few ( $n$ ) major contributors to a cell are responsible for a certain percentage ( $k$ ) of the total of that cell. The idea behind this rule is that in that case at least the major contributors themselves can determine with sufficient precision the contributions of the other contributors to that cell. The choice  $n=3$  and  $k=70\%$  is not uncommon, but  $\tau$ -ARGUS will allow the users to specify their own values of  $n$  and  $k$ .

As an alternative the prior-posterior rule has been proposed. The basic idea is that a contributor to a cell has a better chance to estimate competitors in a cell than an outsider, and also that these kind of intrusions can occur rather often. The precision with which a competitor can estimate is a measure of the sensitivity of a cell. The worst case is that the second largest contributor will be able to estimate the largest contributor. If this precision is more than  $p\%$ , the cell is considered unsafe. An extension is that also the global knowledge about each cell is taken into account. In that case we assume that each intruder has a basic knowledge of the value of each contributor of  $q\%$ . Note, that it is actually the ratio  $p/q$  that determines which cells are considered safe, or unsafe. In this version of ARGUS, the  $q$ -parameter is fixed to 100. Literature refers to this rule as (minimum protection of)  $p\%$ -rule. If the

---

<sup>5</sup> See section 4.2.1 Sensitive Cells in Magnitude Tables of the Handbook.

intention is to state a prior-posterior rule with parameters  $p_0$  and  $q_0$ , where  $q_0 < 100$ , choose the parameter  $p$  of the  $p$  %-rule as  $p = p_0/q_0 * 100$ . See Loeve (2001)<sup>6</sup>

With these rules as a starting point it is easy to identify the sensitive cells, provided that the tabulation package has the facility not only to calculate the cell totals, but also to calculate the number of contributors and the  $n$  individual contributions of the major contributors. Tabulation packages like ABACUS (from Statistics Netherlands) and the package ‘SuperCross’ developed in Australia by Space-Time Research have that capacity. In fact  $\tau$ -ARGUS not only stores the sum of the  $n$  major contributions for each cell, but the individual major contributions themselves. The reason for this is that this is very handy in case rows and columns etc. in a table are combined. By merging and sorting the sets of individual contributions of the cells to be combined, one can quickly determine the major contributions of the new cell, without going back to the original file. This implies that one can quickly apply the dominance rule to the combined cells. Combining rows and columns (table redesign) is one of the major tools for reducing the number of unsafe cells.

This too is the reason why  $\tau$ -ARGUS reads microdata files. However due to continuous demands from users we have now also provide the option to read ready-made tables, but with the restriction that the options for table redesign will not then be available.

A problem, however, arises when also the marginals of the table are published. It is no longer enough to just suppress the sensitive cells, as they can be easily recalculated using the marginals. Even if it is not possible to exactly recalculate the suppressed cell, it is possible to calculate an interval that contains the suppressed cell. This is possible if some constraints are known to hold for the cell values in a table. A commonly found constraint is that the cell values are all nonnegative.

If the size of such an interval is rather small, then the suppressed cell can be estimated rather precisely. This is not acceptable either. Therefore it is necessary to suppress additional information to achieve sufficiently large intervals.

Several solutions are available to protect the information of the sensitive cells:

- Combining categories of the spanning variables (table redesign). Larger cells tend to protect the information about the individual contributors better.
- Suppression of additional (secondary) cells to prevent the recalculation of the sensitive (primary) cells.

The calculation of the optimal set (with respect to the loss of information) of secondary cells is a complex OR-problem.  $\tau$ -ARGUS has been built around this solution, and takes care of the whole process. A typical  $\tau$ -ARGUS session will be one in which the users will first be presented with the table containing only the primary unsafe cells. The user can then choose how to protect these cells. This can involve the combining of categories, equivalent to the global recoding of  $\mu$ -ARGUS. The result will be an update of the table with fewer unsafe cells (certainly not more) if the recoding has worked. At a certain stage the user requests the system to solve the remaining unsafe cells by finding secondary cells to protect the primary cells.

---

<sup>6</sup> Loeve, Anneke, 2001, Notes on sensitivity measures and protection levels, Research paper, Statistics Netherlands. Available at <http://neon.vb.cbs.nl/casc/related/marges.pdf>

At this stage the user can choose between several options to protect the primary sensitive cells. Either they choose the hypercube method or the optimal solution. In this case they also has to select the solver to be used, Xpress or Cplex. After this, the table can be stored for further processing if necessary, and eventual publication.

## 2.3 *Sensitive cells in frequency count tables*<sup>7</sup>

---

In the simplest way of using  $\tau$ -ARGUS, sensitive cells in frequency count tables are defined as those cells that contain a frequency that is below a certain threshold value. This threshold value is to be provided by the data protector. This way of identifying unsafe cells in a table is the one that is implemented in the current version of  $\tau$ -ARGUS. It should be remarked, however, that this is not always an adequate way to protect a frequency count table.<sup>8</sup> Yet it is applied a lot. Applying a dominance rule or a p% rule is useless in this context. One should think about possible disclosure risks that a frequency count table poses and possible disclosure scenarios in order to simulate the behaviour of an intruder. Such an analysis would probably come up with different insights than using a simple thresholding rule, *e.g.* like the one sketched in the reference just mentioned. We just mention here the risks of group-disclosure; when a (small) group of respondents have all the same score on a certain category. This risk is often also referred to as the problem of 100%-cells. Further research on this topic is being carried out at a.o. Statistics Netherlands.

## 2.4 *Table redesign*

---

If a large number of sensitive cells are present in a table, it might be an indication that the spanning variables are too detailed. In that case one could consider combining certain rows and columns in the table. (This might not always be possible because of publication policy.) Otherwise the number of secondary cell suppressions might just be too enormous. The situation is comparable to the case of microdata containing many unsafe combinations. Rather than eliminating them with local suppressions one can remove them by using global recodings. For tabular data we use the phrase “table redesign” to denote an operation analogous to global recoding in microdata sets. The idea of table redesign is to combine rows, columns etc., by adding the cell contents of corresponding cells from the different rows, columns etc. It is a property of the sensitivity rules that a joint cell is safer than any of the individual cells. So as a result of this operation the number of unsafe cells is reduced. One can try to eliminate all unsafe combinations in this way, but that might lead to an unacceptably high information loss. Instead, one could stop at some point, and eliminate the remaining unsafe combinations by using other techniques such as cell suppression.

---

<sup>7</sup> See section 5.2 Unsafe cells of the Handbook.

<sup>8</sup> See for instance Leon Willenborg and Ton de Waal, 1996, Statistical disclosure control in practice, Springer-Verlag, New York, Section 6.3.

## 2.5 *Secondary cell suppression*

---

Once the sensitive cells in a table have been identified, possibly following table redesign it might be a good idea to suppress these values. In case no constraints on the possible values in the cells of a table exist this is easy: one simply removes the cell values concerned and the problem is solved. In practice, however, this situation hardly ever occurs. Instead one has constraints on the values in the cells due to the presence of marginals and lower bounds for the cell values (typically 0). The problem then is to find additional cells that should be suppressed in order to protect the sensitive cells. The additional cells should be chosen in such a way that the interval of possible values for each sensitive cell value is sufficiently large. What is “sufficiently large” can be specified by the data protector in  $\tau$ -ARGUS by specifying the protection intervals.

In general the secondary cell suppression problem turns out to be a hard problem, provided the aim is to retain as much information in the table as possible, which, of course, is a quite natural requirement. The optimisation problems that will then result are quite difficult to solve and require expert knowledge in the area of combinatorial optimisation.

## 2.6 *Information loss in terms of cell costs<sup>9</sup>*

---

In case of secondary cell suppression it is possible that a data protector might want to differentiate between the candidate cells for secondary suppression. It is possible that they would strongly prefer to preserve the content of certain cells, and are willing to sacrifice the values of other cells instead. A mechanism that can be used to make such a distinction between cells in a table is that of cell costs. In  $\tau$ -ARGUS it is possible to associate different costs with the cells in a table. The higher the cost the more important the corresponding cell value is considered and the less likely it will be suppressed. We shall interpret this by saying that the cells with the higher associated costs have a higher information content. The aim of secondary cell suppression can be summarised by saying that a safe table should be produced from an unsafe one, by minimising the information loss, expressed as the sum of the costs associated with the cells that have secondarily been suppressed.

$\tau$ -ARGUS offers several ways to compute these costs. The first option is to compute the costs as the sum of the contributions to a cell. Alternatively another variable in the data file can be used as the cost function. Secondly this cost can be the frequency of the contributors to a cell, and finally each cell can have cost = 1, minimising the number of suppressed cells.

## 2.7 *Series of tables*

---

In  $\tau$ -ARGUS it is possible to specify a series of tables that will be protected one by one, and independently of each other. It is more efficient to choose this option since  $\tau$ -ARGUS requires only a single run through the microdata in order to

---

<sup>9</sup> See section 5.5 Information Loss of the Handbook.

produce the tables. But also for the user it is often more attractive to specify a series of tables and let  $\tau$ -ARGUS protect them in a single session, rather than have several independent sessions.

## 2.8 *The Hypercube/GHMITER method*<sup>10</sup>

---

In order to ensure tractability also of big applications,  $\tau$ -ARGUS interfaces with the GHMITER hypercube method of R. D. Repsilber of the Landesamt für Datenverarbeitung und Statistik in Nordrhein-Westfalen/Germany, offering a quick heuristic solution. The method has been described in depth in Repsilber (1994), Repsilber (1999) and Repsilber (2002), for a briefer description see Giessing and Repsilber (2002).

### 2.8.1 The hypercube method

The approach builds on the fact that a suppressed cell in a simple n-dimensional table without substructure cannot be disclosed exactly if that cell is contained in a pattern of suppressed, nonzero cells, forming the corner points of a hypercube.

The algorithm subdivides n-dimensional tables with hierarchical structure into a set of n-dimensional sub-tables without substructure. These sub-tables are then protected successively in an iterative procedure that starts from the highest level. Successively, for each primary suppression in the current sub-table, all possible hypercubes with this cell as one of the corner points are constructed.

If protection against inferential disclosure is requested, for each hypercube, a lower bound for the width of the suppression interval for the primary suppression that would result from the suppression of all corner points of the particular hypercube will be estimated. To estimate that bound, it is not necessary to implement the time consuming solution to the corresponding Linear Programming problem. Only if it turns out that the bound is sufficiently large, the hypercube becomes a feasible solution.

If no protection against inferential disclosure is requested, any hypercube will be considered feasible. This may of course lead to some cases of underprotection.

For any of the feasible hypercubes, the loss of information associated with the suppression of its corner points is computed. The particular hypercube that leads to minimum information loss is selected, and all its corner points are suppressed.

Note that the information loss concept of the hypercube method is slightly different from the one of the other, linear programming based methods for secondary cell suppression offered by  $\tau$ -ARGUS it operates rather like a two-stage concept. In the first way, the algorithm will look at the number of additional suppressions (additional to those that are already suppressed because they are primary unsafe, or because they were selected as secondary suppression in another subtable) that would be caused by the selection of a particular candidate hypercube. If there is more than one hypercube that would result in the same, smallest number of additional secondary suppressions, at second priority the method will select the one

---

<sup>10</sup> The section on GHMITER has been contributed by Sarah GIESSING, *Federal Statistical Office of Germany* 65180 Wiesbaden E-mail: [sarah.giessing@destatis.de](mailto:sarah.giessing@destatis.de). See section 4.4.4 of the Handbook.

with the smallest sum of costs associated to the suppression of the corresponding additional secondary suppressions. Cell costs associated to a cell are indeed a logarithmic transformation of the cell value plus eventually a large constant, if the cell is a marginal cell of the current sub-table.

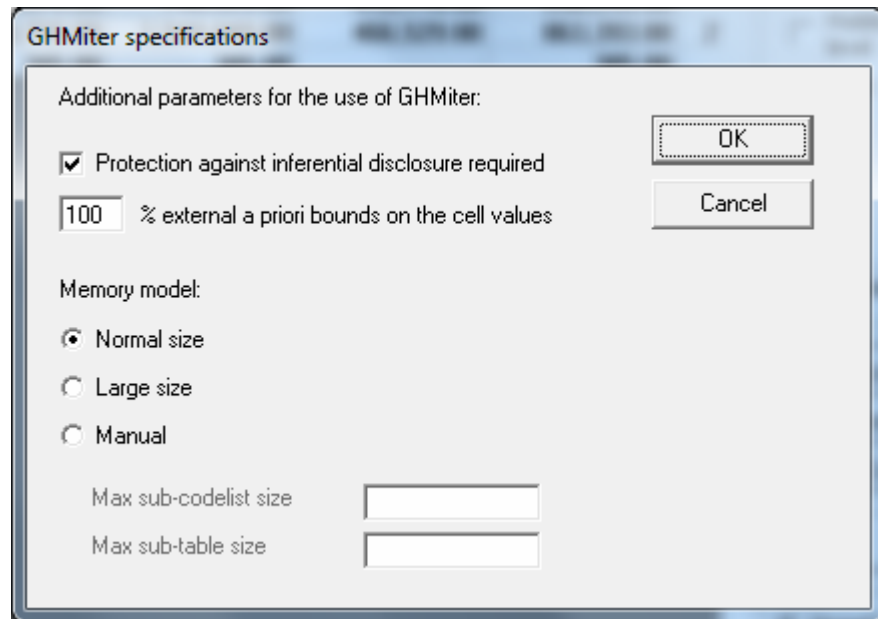
After all sub-tables have been protected once, the procedure is repeated in an iterative fashion. Within this procedure, when cells belonging to more than one sub-table are chosen as secondary suppressions in one of these sub-tables, in further processing they will be treated like sensitive cells in the other sub-tables they belong to. The same iterative approach is used for sets of linked tables.

It should be mentioned here that the ‘hypercube criterion’ is a sufficient but not a necessary criterion for a ‘safe’ suppression pattern. Thus, for particular subtables the ‘best’ suppression pattern may not be a set of hypercubes – in which case, of course, the hypercube method will miss the best solution and lead to some overprotection. Other simplifications of the heuristic approach that add to this tendency for over-suppression are the following: when assessing the feasibility of a hypercube to protect specific target suppressions against interval disclosure, the method

- is not able to consider protection maybe already provided by other cell suppressions (suppressed cells that are not corner points of this hypercube) within the same sub-table,
- does not consider the sensitivity of multi-contributor primary suppressions properly, that is, it does not consider the protection already provided in advance of cell suppression through aggregation of these contributions,
- attempts to provide the same *relative* ambiguity to (eventually large) secondary suppressions that have been selected to protect cells in a linked sub-table, as if they were single-respondent primary suppressions, while actually it would be enough to provide the same *absolute* ambiguity as required by the corresponding primary suppressions.

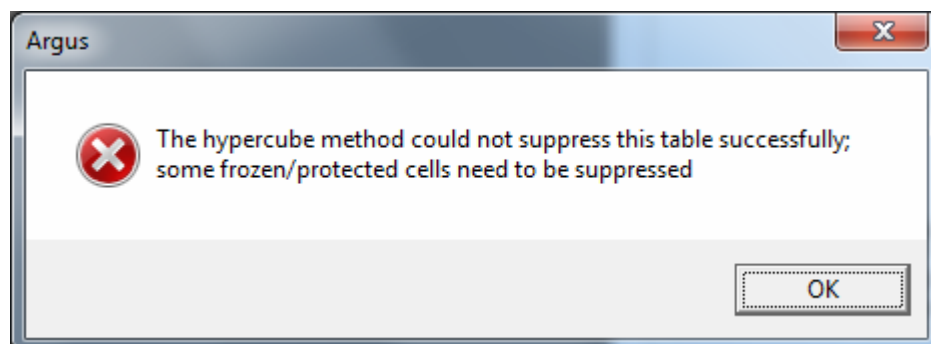
## 2.8.2 The ARGUS implementation of GHMITER

- In the implementation offered by ARGUS, GHMITER makes sure that a single respondent cell will never appear to be corner point of one hypercube only, but of two hypercubes at least. Otherwise it could happen that a single respondent, who often can be reasonably assumed to know that he is the only respondent, could use his knowledge on the amount of his own contribution to recalculate the value of any other suppressed corner point of this hypercube.
- As explained above, GHMITER uses an elaborate internal cost assignment mechanism which is essential to achieve an optimal performance (given the natural restrictions of the simple heuristic approach, of course). This mechanism should not be cast out of balance. Therefore, the user’s choice of the cell costs (c.f. 3.1.3, 4.3.3) does not have any impact, when using the hypercube method.
- For tables presenting magnitude data, if protection against inferential disclosure is requested (see the upper part of the pop-up window below)  $\tau$ -ARGUS will ensure that GHMITER selects secondary suppressions that protect the sensitive cells properly. Only cells will be considered feasible as secondary suppressions that are large enough to give enough protection to the target sensitive cell as explained in Giessing (2003).

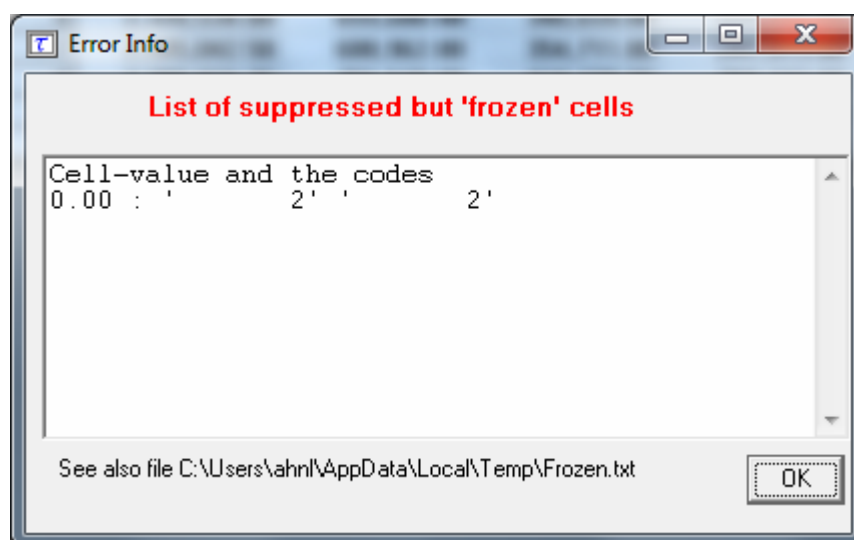


- In order to achieve this,  $\tau$ -ARGUS computes a suitable *sliding protection ratio* (for explanation see Giessing (2003),  $\tau$ -ARGUS will display the value of this ratio in the report file) to be used by GHMITER. If in the screen above the option “*Protection against inferential disclosure required*” is inactivated, GHMITER will not check whether secondary suppressions are sufficiently large.
- As mentioned above, GHMITER is unable to 'add' the protection given by multiple hypercubes. In certain situations, it is not possible to provide sufficient protection to a particular sensitive cell (or secondary suppression) by suppression of one single hypercube. In such a case, GHMITER is unable to confirm that this cell has been protected properly, according to the specified *sliding protection ratio*. It will then reduce the *sliding protection ratio* automatically, and individually, step by step for those cells, the protection of which the program cannot confirm otherwise. In steps 1 to 9 we divide the original ratio by  $k$ , values of  $k$  from 2 to 10, and if this still does not help, in step 10 we divide by an extremely large value, and finally, if even that does not solve the problem, step 11 will set the ratio to zero). The  $\tau$ -ARGUS report file will display the number of cases where the sliding protection range was reduced by finally confirmed sliding protection ratio.
- Note, that the number of cases with range reduction reported by this statistic in the report file is very likely to exceed the actual number of cells concerned, because cells belonging to multiple (sub-) tables are counted multiple times. In our experience this concerns particularly the cases, where the protection level was reduced to an, ‘infinitely’ small (positive) value (in step 10, see above). Step 10 is usually required to confirm protection of large, high level secondary suppressions, which are likely to appear in multiple tables, especially in processing of linked tables. By the way, terms “reduction of the *sliding protection ratio*” and “reduction of the *protection level*” are used synonymously in the report file.
- Note that step 11 will make cells eligible for secondary suppression that  $\tau$ -ARGUS considers as ‘protected’ (so called ‘frozen’ cells, for discussion of this option see for instance Giessing (2003)).

As this is inconsistent with the current view on protected cells in  $\tau$ -ARGUS this will lead to the following error message:



Codes and cell values of those suppressed *frozen* cells are then displayed by ARGUS:



When the status of these cells is changed into ‘unprotected’ before re-running the hypercube method, the solution will be a feasible solution for  $\tau$ -ARGUS.

### Negative values

The hypercube method has no problems when certain cells are negative.

## 2.8.3 References on GHMiter

- Repsilber, R. D. (1994), ‘Preservation of Confidentiality in Aggregated data’, paper presented at the Second International Seminar on Statistical Confidentiality, Luxembourg, 1994
- Repsilber, D. (1999), ‘Das Quaderverfahren’ - in Forum der Bundesstatistik, Band 31/1999: Methoden zur Sicherung der Statistischen Geheimhaltung, (in German)
- Repsilber, D. (2002), ‘Sicherung persönlicher Angaben in Tabellendaten’ - in Statistische Analysen und Studien Nordrhein-Westfalen, Landesamt für Datenverarbeitung und Statistik NRW, Ausgabe 1/2002 (in German)
- 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’

## 2.9 *Optimisation models for secondary cell suppression*<sup>11</sup>

---

$\tau$ -ARGUS applies different approaches to find optimal and near-optimal solutions. One of these approaches is based on a Mathematical Programming technique which consists of solving Integer Linear Programming programs modelling the combinatorial problems under different methodologies (Cell Suppression and Controlled Rounding). The main characteristic of these models is that they share the same structure, thus based only on a 0-1 variable for each cell. In the Cell Suppression methodology, the variable is 1 if and only if the cell value must be suppressed. In the Controlled Rounding methodology, the variable is 1 if and only if the cell value must be rounded up. No other variables are necessary, so the number of variables in the model is exactly the number of cells in the table to be protected. In addition, the model also imposes the protection level requirements (upper, lower and sliding) in the same way for the different methodologies (Cell Suppression and Controlled Rounding). These requirements ask for a guarantee that an attacker will not get too narrow an interval of potential values for a sensitive cell, which he/she will compute by solving two linear programming programs (called attacker problems). Even if a first model containing this two-attacker problem would lead to a bi-level programming model, complex to be solved in practice, a Benders' decomposition approach allows us to convert the attacker problems into a set of linear inequalities. This conversion provides a second model for each methodology that can be efficiently solved by a modern cutting-plane approach. Since the variables are 0-1, a branching phase can be necessary, and the whole approach is named "branch-and-cut algorithm".

Branch-and-cut algorithms are modern techniques in Operations Research that provide excellent results when solving larger and complicated combinatorial problems arising in many applied fields (like routing, scheduling, planning, telecommunications, etc.). Shortly, the idea is to solve a compact 0-1 model containing a large number of linear inequalities (as the ones above mentioned for the Cell Suppression and for the Controlled Rounding) through an iterative procedure that does not consider all the inequalities at the same time, but generates the important ones when needed. This dynamic procedure of dealing with large models allows the program to replace the resolution of a huge large model by a short sequence of small models, which is termed a "decomposition approach". The on-line generation of the linear inequalities (rows) was also extended in this work to the variables (columns), thus the algorithm can also works on tables with a large number of cells, and the overall algorithm is named "branch-and-cut-and-price" in the Operations Research literature.

To obtain good performance, the implementation has also considered many other ingredients, standard in branch-and-cut-and-price approaches. For example, it is

---

<sup>11</sup> The optimisation models have been built by a team of researchers headed by Juan-José Salazar-Gonzalez of the University La Laguna, Tenerife, Spain. Other members of the team were: G. Andreatta, M. Fischetti, R. Betancort Villalva, M.D. Montesdeoca Sanchez and M. Schoch

fundamentally the implementation of a pre-processing approach where redundant equations defining the table are eliminated, where variables associated to non-relevant cells are removed, and where dominated protection levels are detected. The pre-processing is fundamental to make the problem as small as possible before starting the optimization phase. Another fundamental ingredient is the heuristic routine, which allows the algorithm to start with an upper bound of the optimal loss of information. This heuristic routine ensures the production of a protected pattern if the algorithm is interrupted by the user before the end. In other words, thanks to the heuristic routine, the implemented algorithm provide a near-optimal solution if the execution is cancelled before having a proof of optimality. During the implicit enumeration approach (i.e., the branch-and-cut-and-price) the heuristic routine is called several times, thus providing different protected patterns, and the best one will be the optimal solution if its loss of information is equal to the lower bound. This lower bound is computed by solving a relaxed model, which consists of removing the integrability condition on the integer model. Since the relaxed model is a linear program, a linear programming solver must be called.

We have not implemented our own linear programming solver, but used a commercial solver which is already tested by other programmers for many years. A robust linear programming solver is a guarantee that no numerical trouble will appear during the computation.

That is the reason to requires either CPLEX (from ILOG) or XPRESS (from FICO). Because the model to be solved can be applied to all type of table structures (2-dim, 3-dim, 4-dim, etc), including hierarchical and linked tables, we cannot use special simplex algorithm implementations, like the min-cost flow computation which would required to work with tables that can be modelled as a network (e.g., 2-dimensional tables or collections of 2-dim tables linked by one link). On this special table, ad-hoc approaches (solving network flows or short path problems) could be implemented to avoid using general linear programming solvers.

In any case, future works will try to replace the commercial solvers by freely available linear-programming solvers.

---

## 2.10 *The Modular approach*<sup>12</sup>

---

The modular (HiTaS) solution is a heuristic approach to cell suppression in hierarchical tables. Hierarchical tables are specially linked tables: at least one of the spanning variables exhibits a hierarchical structure, *i.e.* contains (many) sub-totals.

In Fischetti and Salazar (1998) a theoretical framework is presented that should be able to deal with hierarchical and generally linked tables. In what follows, this will be called the mixed integer approach. In this framework, additional constraints to a linear programming problem are generated. The number of added constraints however, grows rapidly when dealing with hierarchical tables, since many dependencies exist between all possible (sub-)tables containing many (sub-)totals. The implemented heuristic approach (HiTaS) deals with a large set of (sub-)tables in a particular order. A non hierarchical table can be considered to be a hierarchical table with just one level. In that case, the approach reduces to the original mixed integer approach and hence provides the optimal solution. In case of a hierarchical

---

<sup>12</sup> See section 4.4.2 Modular of the Handbook.

table, the approach will provide a sub-optimal solution that minimises the information loss per sub-table, but not necessarily the global information loss of the complete set of hierarchically linked tables.

In the following section, a short description of the approach is given. For a more detailed description of the method, including some examples, see *e.g.*, De Wolf (2002).

HiTaS deals with cell suppression in hierarchical tables using a top-down approach. The first step is to determine the primary unsafe cells in the base-table consisting of all the cells that appear when crossing the hierarchical spanning variables. This way all cells, whether representing a (sub-)total or not, are checked for primary suppression. Knowing all primary unsafe cells, the secondary cell suppressions have to be found in such a way that each (sub-)table of the base-table is protected and that the different tables cannot be combined to undo the protection of any of the other (sub-)tables. The basic idea behind the top-down approach is to start with the highest levels of the variables and calculate the secondary suppressions for the resulting table. The suppressions in the interior of the protected table is then transported to the corresponding marginal cells of the tables that appear when crossing lower levels of the two variables. All marginal cells, both suppressed and not suppressed, are then ‘fixed’ in the calculation of the secondary suppressions of that lower level table, i.e., they are not allowed to be (secondarily) suppressed. This procedure is then repeated until the tables that are constructed by crossing the lowest levels of the spanning variables are dealt with.

A suppression pattern at a higher level only introduces restrictions on the marginal cells of lower level tables. Calculating secondary suppressions in the interior while keeping the marginal cells fixed, is then independent between the tables on that lower level, i.e., all these (sub-)tables can be dealt with independently of each other. Moreover, added primary suppressions in the interior of a lower level table are dealt with at that same level: secondary suppressions can only occur in the same interior, since the marginal cells are kept fixed.

However, when several empty cells are apparent in a low level table, it might be the case that no solution can be found if one is restricted to suppress interior cells only. Unfortunately, backtracking is then needed.

Obviously, all possible (sub)tables should be dealt with in a particular order, such that the marginal cells of the table under consideration have been protected as the interior of a previously considered table. To that end, certain groups of tables are formed in a specific way (see De Wolf (2002)). All tables within such a group are dealt separately, using the mixed integer approach.

The number of tables within a group is determined by the number of parent-categories the variables have one level up in the hierarchy. A parent-category is defined as a category that has one or more sub-categories. Note that the total number of (sub)-tables that have to be considered thus grows rapidly.

### Singletons

Singleton cells should be treated with extra care. The single respondent in this cell could easily undo the protection if no extra measures were taken. The most dangerous situation is that there are only two singletons in a row, or one and one other primary unsafe cell. These singletons could easily disclose the other cell.

We have added options for extra singleton protection in the following situations.

1. If on a row or column of a subtable there are only two singletons and no other primary suppressions.

2. If there is only one singleton and one multiple primary unsafe cell.
3. If a frequency rule is used, it could happen that two cells on a row/column are primary unsafe, but the sum of the two cells could still be unsafe. In that case it should be prevented that these two cells protect each other.

Cells within a table sometimes consist of exactly one contributor. Such a cell is called a singleton. Linear sensitivity rules will usually label this cell as (primary) unsafe. When cell suppression is used to protect a table with unsafe cells, these singletons need to be taken care of in a special way.

Within a suppression pattern, contributors in singletons may be able to recalculate other suppressed cells. Obviously, a contributor could always insert its own contribution and thereby recalculate its own suppressed cell. This could in turn lead to the possibility of recalculating other suppressed cells in the same suppression pattern. Whenever such a recalculated cell is (primary) unsafe, this means disclosure.

Within the current models used to determine suppression patterns, it is not possible to take all possible situations into account when singletons are part of a suppression pattern. However, an important group of instances of disclosure by singletons, is when a singleton is part of a row with exactly one additional (also primary) suppression.

4. If on a row or column of a subtable there are only two singletons and no other primary suppressions.
5. If there is only one singleton and one multiple primary unsafe cell.
6. If a frequency rule is used, it could happen that two cells on a row/column are primary unsafe, but the sum of the two cells could still be unsafe. In that case it should be prevented that these two cells protect each other.

Note that the last situation is not really a singleton problem, but this problem is handled in the same way.

To prevent this kind of disclosure, it would be sufficient to force an additional (third) suppression in the same row. In prior versions of  $\tau$ -ARGUS this was accomplished by increasing the protection levels of one of the (primary) unsafe cells in the row. In short, the protection level of one of the primary suppressed cells was raised in such a way that the other primary suppression would not be able to give sufficient protection. The largest primary unsafe cell in the row got the *cell value* of the other unsafe cell in the row, plus a small value, as protection level. Indeed, this forces a third suppression in the row.

However, since the *cell value* of one of the suppressed cells was involved, this meant that the increased protection level of this cell could become quite large, which would have an effect on the suppression pattern in one of the other dimensions. In certain situations this led to oversuppression.

To circumvent this problem, the newly implemented approach adds a virtual cell to the table. That virtual cell is assigned a value equal to the sum of the two primary suppressed cells in the row, and is given the status '(primary) unsafe'. That virtual cell then only has to be protected against exact disclosure, i.e., it suffices to impose a small protection interval.

Table~\ref{tab:SingletonExample} shows an example table, displaying the singleton problem. In Table~\ref{tab:SingletonExample}(a), the values of the cells are given, with in bold italic the (primary) unsafe cells. Table~\ref{tab:SingletonExample}(b) shows the names of the cells, where

$c_{ij}$  stands for the cell with coordinates  $(i, j)$ .

	<b>Total</b>	<b>X1</b>	<b>X2</b>	<b>X3</b>	<b>X4</b>
<b>Total</b>	227	73	33	93	25
<b>A</b>	146	52	<b>15</b>	62	<b>17</b>
<b>B</b>	81	24	18	31	8

	<b>Total</b>	<b>X1</b>	<b>X2</b>	<b>X3</b>	<b>X4</b>
<b>Total</b>	$c_{00}$	$c_{01}$	$c_{02}$	$c_{03}$	$c_{04}$
<b>A</b>	$c_{10}$	$c_{11}$	<b><math>c_{12}</math></b>	$c_{13}$	<b><math>c_{14}</math></b>
<b>B</b>	$c_{20}$	$c_{21}$	$c_{22}$	$c_{23}$	$c_{24}$

Table 1 Example table to explain Singleton Problem.

Bold and red means (primary) unsafe.

Now assume that cell  $c_{12} = (A, X2)$  is a singleton and cell  $c_{14} = (A, X4)$  is unsafe according to a p%-rule with  $p=10$ . Hence, cell  $c_{14}$  is the only other (primary) unsafe cell in that row. To protect cell  $c_{14}$  against disclosure by the contributor of singleton  $c_{12}$ , a `virtual cell  $c_v$  is defined with value 32. Moreover, that virtual cell is given a small protection interval, (32,33) say. The relations that define the table structure, including the virtual cell, are given below:

$$\begin{aligned}
 c_{00} &= c_{01} + c_{02} + c_{03} + c_{04} \\
 c_{10} &= c_{11} + c_{12} + c_{13} + c_{14} \\
 c_{20} &= c_{21} + c_{22} + c_{23} + c_{24} \\
 c_{00} &= c_{10} + c_{20} \\
 c_{01} &= c_{11} + c_{21} \\
 c_{02} &= c_{12} + c_{22} \\
 c_{03} &= c_{13} + c_{23} \\
 c_{04} &= c_{14} + c_{24} \\
 c_v &= c_{12} + c_{14}
 \end{aligned}$$

Table 2 Relations defining table structure of Table 1

Within  $\tau$ -ARGUS, this procedure is implemented in both the optimal approach as well as in the modular approach. For the modular approach, this procedure is applied to each subtable separately, whenever a subtable is dealt with within the modular approach.

This special attention to singletons is only given when the other suppressed cell in the same row is a `true' primary suppression. This is natural, since it has to be done prior to the search for secondary suppressions. In the modular approach, a hierarchical table is divided into many, non-hierarchical, subtables. Secondary suppressions in one table sometimes temporarily become primary suppressions in other tables during the process. I.e., those suppression are not `true' primary suppressions. It is therefore also natural not to construct virtual cells in case a singleton is in the same row with exactly one other primary suppression that was originally a secondary suppression. This is indeed the way it is implemented in the modular approach.

In previous versions of  $\tau$ -ARGUS a similar procedure was available. But then the additional protection was achieved by increasing the protection level of the singleton cell. This would lead however also in additional protection in other dimensions and would create over-protection

### **Negative values**

The implementation by Fischetti and Salazar does not allow for negative values. However it is not uncommon, that some cells in a table have negative values. Therefore additional measures have been taken. If in a subtable during the process negative values are found, all cell values are increased such that the lowest value becomes positive. Of course the margins have to be recalculated, but a safe protection pattern will be found.

### **References on the modular method**

Fischetti, M. and J.J. Salazar-González (1998). Models and Algorithms for Optimizing Cell Suppression in Tabular Data with Linear Constraints. Technical Paper, University of La Laguna, Tenerife.

P.P. de Wolf (2002). HiTaS: a heuristic approach to cell suppression in hierarchical tables. Proceedings of the AMRADS meeting in Luxembourg (2002).

Additional reading on the optimisation models can be found at the CASC-website (<http://neon.vb.cbs.nl/casc/Related/99wol-heu-r.pdf>)

## ***2.11 The modular approach for linked tables***

---

When tables are linked through simple linear constraints, cell suppressions must obviously be coordinated between tables. The most typical case is when tables share common cells (usually marginals), i.e., when they are linked through constraints saying literally that cell  $X$  of table  $A$  is identical to cell  $Y$  of table  $B$ .

Suppose a set of  $N$  tables,  $\{T_1, \dots, T_N\}$ , need to be protected. These tables are assumed to be linked. Each table has a hierarchical structure that may differ from the hierarchical structures of the other tables. However, it is assumed that tables using the same spanning variables have hierarchies that can be covered. Loosely speaking this means that a single hierarchy can be constructed such that all hierarchies of the same variable in the  $N$  tables are a sub hierarchy of the cover hierarchy. See De Wolf and Giessing (2009) for more details. In the context of pre-planned table production processes which are typically in place in statistical agencies for the production of certain sets of pre-specified standard tabulations, it is normally no problem to satisfy these conditions. Literally speaking, the assumption is that tables in a set of linked tables may present the data in a breakdown by the same spanning variable at various amounts of detail. But only under the condition that, if in one of the tables some categories of a spanning variable are grouped into a certain intermediate sum category, during SDC processing this intermediate sum category is considered in any other table presenting the data in a breakdown of the same spanning variable and at that much detail.

The idea is then as follows. Suppose that the  $N$  tables  $\{T_1, \dots, T_N\}$  that need to be protected simultaneously, contain  $M$  different spanning variables. Since the hierarchies are supposed to be coverable, an  $M$ -dimensional table exists having all

the specified tables as subtables. The spanning variables will be numbered 1 up to  $M$ .

Each spanning variable can have several hierarchies in the specified tables. Denote those hierarchies for spanning variable  $i$  by  $\mathcal{H}_1^i, \dots, \mathcal{H}_{I_i}^i$  where  $I_i$  is the number of different hierarchies of variable  $i$ .

Define the  $M$ -dimensional table by the table with spanning variables according to hierarchies  $\mathcal{G}_1, \dots, \mathcal{G}_M$  such that, for each  $i = 1, \dots, M$  hierarchy  $\mathcal{G}_i$  covers the set of hierarchies  $\{\mathcal{H}_j^i\}$  with  $j = 1, \dots, I_i$ . This  $M$ -dimensional table will be called the cover table. See De Wolf and Giessing (2009) for more details.

Then use the Modular approach (see section 2.10) on the cover table  $T_C$ , but only consider those subtables that are also subtables of at least one of the specified tables  $T_1, \dots, T_N$  and disregard the other subtables.

I.e., the procedure of the Modular approach is followed, but during that process any simple subtable that is not a subtable of any of the tables in the set  $\{T_1, \dots, T_N\}$  is skipped. I.e., the order the simple subtables will be protected, is the same as in the ‘complete’ Modular approach, only some subtables will be skipped.

See De Wolf and Hundepool (2010) for a practical application of the Adjusted Modular Approach.

### References on the modular approach for linked tables

- De Wolf, P.P. and S. Giessing (2009), Adjusting the  $\tau$ -ARGUS modular approach to deal with linked tables, *Data & Knowledge Engineering*, Volume 68, Issue 11, pp. 1160-1174.
- De Wolf, P.P. and A. Hundepool (2010), Three ways to deal with a set of linked SBS tables using  $\tau$ -ARGUS, *Privacy in Statistical Databases*, J. Domingo-Ferrer and E. Magkos (Eds.), Springer 2010, LNCS 6344 pp. 66-74.

## ***2.12 Network solution for large 2 dimensional tables with one hierarchy***

---

$\tau$ -ARGUS also contains a solution for the secondary cell suppression based on network flows. This contribution is by Jordi Casto of the Universitat Politècnica de Catalunya in Barcelona. The network flows solution for cell suppression implements a fast heuristic for the protection of statistical data in two-dimensional tables with one hierarchical dimension (1H2D tables). This new heuristic sensibly combines and improves ideas of previous approaches for the secondary cell suppression problem in two-dimensional general, see Castro(1994) and positive tables, see Kelly(1992) and Castro(2003) tables. Details about the heuristic can be found in Castro(1996) and Cox(1995). Unfortunately this approach is only possible for two-dimensional tables with only one hierarchy, due to the limitations of the network flows.

The heuristic is based on the solution of a sequence of shortest-path subproblems that guarantee a feasible pattern of suppressions (i.e., one that satisfies the protection levels of sensitive cells). Hopefully, this feasible pattern will be close to the optimal one.

The current package is linked with three solvers: CPLEX7.5/8.0 see ILOG(2000) PPRN see Castro(1996), and an efficient implementation of the bidirectional Dijkstra's algorithm for shortest-paths (that will be denoted as "Dijkstra") see Ahuja(1993). Later releases of CPLEX will also work if the interface routines are the same than for version 8.0. The heuristic can use any of the three solvers for the solution of the shortest path subproblems, although Dijkstra is recommended (and the default one) for efficiency reasons. CPLEX is needed if a lower bound of the optimal solution want to be computed. The auditing phase can be performed with either CPLEX or PPRN.

PPRN and Dijkstra were implemented at the Dept. of Statistics and Operations Research of the Universitat Politècnica de Catalunya, and are included in NF CSP. PPRN was originally developed during 1992–1995, but it had to be significantly improved within the CASC project to work with NF CSP. Dijkstra was completely developed in the scope of CASC. The third solver, CPLEX, is a commercial tool, and requires purchasing a license. However, PPRN is a fairly good replacement—although not so robust—for the network flows routines of CPLEX. Therefore, in principle, there is no need for an external commercial solver, unless lower bounds want to be computed.

Even though two of the three solvers are included in the distribution of NF CSP, this document only describes the features of the heuristic, and from the user's point of view. A detailed description of PPRN and Dijkstra's solvers can be found in Castro(1996) and Ahuja(1993), respectively.

The current implementation in  $\tau$ -ARGUS however only uses the Dijkstra and the PPRN solvers. We have restricted ourselves from commercial solvers here as the network flows give already a very fast solution.

### References on the network solution

- Ahuja, R.K, Magnanti, T.L., Orlin, J.B., Network Flows, Prentice Hall (1993).
- Castro, J., PPRN 1.0, User's Guide, Technical report DR 94/06 Dept. of Statistics and Operations Research, Universitat Politècnica de Catalunya, Barcelona, Spain, 1994.
- Castro, J., Network flows heuristics for complementary cell suppression: an empirical evaluation and extensions, in LNCS 2316, Inference Control in Statistical Databases, J. Domingo-Ferrer (Ed), (2002) 59–73.
- Castro, J., Nabona, N. An implementation of linear and nonlinear multicommodity network flows. European Journal of Operational Research 92, (1996) 37–53.
- Cox, L.H., Network models for complementary cell suppression. J. Am. Stat. Assoc. 90, (1995) 1453–1462.
- ILOG CPLEX, ILOG CPLEX 7.5 Reference Manual Library, ILOG, (2000).
- Kelly, J.P., Golden, B.L, Assad, A.A., Cell Suppression: disclosure protection for sensitive tabular data, Networks 22, (1992) 28–55.
- Castro, J. User's and programmer's manual of the network flows heuristics package for cell suppression in 2D tables Technical Report DR 2003-07, Dept. of Statistics and Operations Research, Universitat Politècnica de Catalunya, Barcelona, Spain, 2003;
- See [http://neon.vb.cbs.nl/casc/deliv/41D6\\_NF1H2D-Tau-ARGUS.pdf](http://neon.vb.cbs.nl/casc/deliv/41D6_NF1H2D-Tau-ARGUS.pdf)

## 2.13 Controlled rounding<sup>13</sup>

---

Controlled rounding is a rounding procedure that, differently from other rounding methods, yields additive rounded tables. That is to say that the rounded values add up to the rounded totals and sub-totals shown in the table. This property not only permits the release of realistic tables but also makes it impossible to reduce the protection by “unpicking” the original values by exploiting the differences in the sums of the rounded values. The CRP implemented in  $\tau$ -ARGUS also allows the specification hierarchical links.

Controlled rounding is a SDC method that is most effective for frequency tables. In fact, this method gives adequate protection to small frequencies by creating uncertainty also with respect to zero values (*i.e.* empty cells). The same cannot be said for suppression in the way it is implemented now in  $\tau$ -ARGUS.

### 2.13.1 Restricted and non-restricted controlled rounding

In Zero-restricted Controlled Rounding the rounded values are chosen leaving unaltered the original values that are already multiples of the rounding base, while rounding the others to one of the adjacent multiples of this base. The modified values are chosen so that the sum of the absolute differences between the original values and the rounded ones is minimized under the additivity constraint. Therefore, some values will be rounded up or down to the most distant multiple of the base in order to satisfy the constraints. In most cases such a solution can be found but in some cases it cannot. The zero-restriction constraint in CRP can be relaxed allowing the values to be rounded to a nonadjacent multiple of the base. This relaxation is controlled by allowing a maximum number of *steps*. For example, consider rounding the value 7 when the base equals 5. In zero-restricted rounding, the solution can be either 5 or 10. If 1 step is allowed, the solution can be 0, 5, 10 or 15. In general, let  $z$  be the integer to be rounded in base  $b$ , then this number can be written as

$$z = ub + r,$$

where  $ub$  is the lower adjacent multiple of  $b$  (hence  $u$  is the floor value of  $z/b$ ) and  $r$  is the remainder. In the zero-restricted solution the rounded value,  $a$ , can take values:

$$\begin{cases} a = ub & \text{if } r = 0; \\ a = \begin{cases} ub \\ (u+1)b \end{cases} & \text{if } r \neq 0. \end{cases}$$

If  $K$  steps are allowed, then  $a$ , can take values:

$$\begin{cases} a = \max\{0, (u+j)b\}, j = -K, \dots, K, & \text{if } r = 0; \\ a = \max\{0, (u+j)b\}, j = -K, \dots, (K+1), & \text{if } r \neq 0. \end{cases}$$

---

<sup>13</sup> See section 5.4 Rounding of the Handbook.

### 2.13.2 Optimal, first feasible and RAPID solutions<sup>14</sup>

For a given table there could exist more than one controlled rounded solutions; any of these solutions is a *feasible* solution. The Controlled Rounding Program embedded in  $\tau$ -ARGUS determines the *optimal* solution by minimising the sum of the absolute distances of the rounded values from the original ones. Denoting the cell values, including the totals and sub-totals, with  $z_i$  and the corresponding rounded values with  $a_i$ , the function that is minimised is

$$\sum_{i=1}^N |z_i - a_i|,$$

where  $N$  is the number of cells in a table (including the marginal ones). The optimisation procedure for controlled rounding is a rather complex one ( $NP$ -complete program), so finding the optimal solution may take a long time for large tables. In fact, the algorithm iteratively builds different rounded tables until it finds the optimal solution. In order to limit the time required to obtain a solution, the algorithm can be stopped when the first feasible solution is found. In many cases, this solution is quite close to the optimal one and it can be found in significantly less time.

The RAPID solution is produced by CRP as an approximated solution when not even a feasible one can be found. This solution is obtained by rounding the internal cells to the closest multiple of the base and then computing the marginal cells by addition. This means that the computed marginal values can be many jumps away from the original value. However, a RAPID solution is produced at each iteration of the search for an optimal one and it will improve (in terms of the loss function) over time.  $\tau$ -ARGUS allows to stop CRP after the first RAPID is produced, but this solution is likely to be very far away from the optimal one.

### 2.13.3 Protection provided by controlled rounding

The protection provided by controlled rounding can be assessed by considering the uncertainty about the disclosive true values achieved releasing rounded values; that is the existence interval that an intruder can compute for the rounded value. We assume that also the values of the rounding base,  $b$ , and the number of steps allowed,  $K$ , are released together with the rounded table. Furthermore, we assume that it is known that the original values are frequencies (hence nonnegative integers).

#### Zero-restricted rounding

Given a rounded value,  $a$ , an intruder can compute the following existence intervals for the true value,  $z$ :

$$z \in [0, b - 1] \text{ if } a = 0$$

$$z \in [a - b + 1, a + b - 1] \text{ if } a \neq 0.$$

For example, if the rounding base is  $b=5$  and the rounded value is  $a=0$ , a user can determine that the original value is between 0 and 4. If the rounded value is not 0, then users can determine that the true value is between plus or minus 4 units from the published value.

---

<sup>14</sup> For further details see Salazar, Staggermeier and Bycroft (2005 [Controlled rounding implementation](#), UN-ECE Worksession on SDC, Geneva)

### K-step rounding

As mentioned before, it is assumed that the number of steps allowed is released together with the rounded table. Let  $K$  be the number of steps allowed, then an intruder can compute the following existence intervals for the true value  $z$ :

$$z \in [0, (K+1)b - 1] \text{ if } a < (K+1)b$$

$$z \in [a - (K+1)b + 1, a + (K+1)b - 1] \text{ if } a \geq (K+1)b.$$

For example, assume that for controlled rounding with  $b=5$  and  $K=1$ ,  $a=15$ , then a user can determine that  $z \in [6, 24]$ .

## 2.13.4 Choosing the parameters for Controlled Rounding

The parameters that can be chosen for rounding are the rounding base,  $b$ , and the number of steps allowed. If their value is released, users (including potential intruders) will be able to compute existence intervals for the true values according to the formulae given above. Then, the choice of the parameters' values depends on the protection required for the disclosive values. Of course, the larger the existence interval the greater the protection but also the *damage* caused to the data. The choice of the rounding base, then, should be made by the data protector considering the protection requirements and the damage caused to the data. A discussion on how existence intervals can be related to protection requirements can be found, for example, in Willenborg and de Waal (2001). Below we give some general considerations on the effect of different choices of the rounding base.

Frequencies are disclosive if their values are not larger than a chosen threshold, say  $f$ . In  $\tau$ -ARGUS the minimal rounding base is  $b=f$ . When this value is chosen, disclosive values can be rounded either to 0 or to  $b$ . Hence, an intruder would know that all published zeros are disclosive values, while he or she could not determine if a published value equal to  $b$  is a disclosive value or a larger, safe, one. In some cases this protection can be considered insufficient because it is required that the existence interval for values rounded to zero contains at least one safe value. Then the value of  $b$  must be chosen to be greater than  $f$  or the number of steps allowed must be greater than zero. It must be stressed, however, that the larger the base and the greater the damage inflicted to the data (including safe values). In some cases, data protector may be happy with a base that is less than the minimum frequency threshold. For example, it could be decided that the width of the existence interval must be not less than the minimum frequency. In this case, the base should be chosen to be the minimal integer not smaller than  $f/2$ . Using a smaller base than the minimum safe frequency can be achieved in  $\tau$ -ARGUS by lowering the threshold before computing the table. This "trick" is allowed in rounding because the procedure does not change if the disclosive cells are changed (unlike secondary suppression).

## 2.14 Audit

---

When a table is protected by cell suppression, by making use of the linear relation between published and suppressed cell values in a table (including its margins), it is always possible for any particular suppressed cell of a table to derive upper and

lower bounds for its true value. This holds for either tables with non-negative values, and those tables containing negative values as well, when it is assumed that instead of zero, some other (possibly tight) lower bound for any cell is available to data users in advance of publication. The interval given by these bounds is called the ‘*feasibility interval*’. The example below<sup>15</sup> illustrates the computation of the feasibility interval in the case of a simple two-dimensional table where all cells may only assume non-negative values:

Example 3

Example	1	2	Total
1	$X_{11}$	$X_{12}$	7
2	$X_{21}$	$X_{22}$	3
3	3	3	6
Total	9	7	16

For this table the following linear relations hold:

$$X_{11} + X_{12} = 7 \quad (R1)$$

$$X_{21} + X_{22} = 3 \quad (R2)$$

$$X_{11} + X_{21} = 6 \quad (C1)$$

$$X_{12} + X_{22} = 4 \quad (C2)$$

$$\text{with } X_{ij} \geq 0 \quad \text{for all } (i, j)$$

Using linear programming methodology, it is possible to derive systematically for any suppressed cell in a table a upper bound ( $X^{\max}$ ) and a lower bound ( $X^{\min}$ ) for the set of feasible values. In the example above, for cell (1,1) these bounds are ( $X^{\min}_{11} = 3$  and  $X^{\max}_{11} = 6$ ).

A general mathematical statement for the linear programming problem to compute upper and lower bounds for the suppressed entries of a table is given in Fischetti and Salazar (2000)<sup>16</sup>.

Note that in the current implementation the  $\tau$ -ARGUS audit routine computes upper and lower bounds (i.e. the feasibility intervals) for the suppressed entries of a hierarchical table considering the *full* set of table relations – even, if the table is a hierarchical table. After obtaining these feasibility intervals, they are compared to the protection intervals (c.f. subsection on protection levels in section 4.2.2 of the ESSNET SDC Handbook) and the result of this comparison will be reported to the user. When a table has been protected properly, the feasibility interval of each primary sensitive cell should cover the protection interval. These intervals will be shown by  $\tau$ -ARGUS.

#### Auditing a hierarchical table

It should be noted that secondary cell suppression algorithms like Modular and Hypercube relying on a backtracking procedure (c.f. the subsection on linked and

<sup>15</sup> (Geurts, 1992, Table 10, p 20)

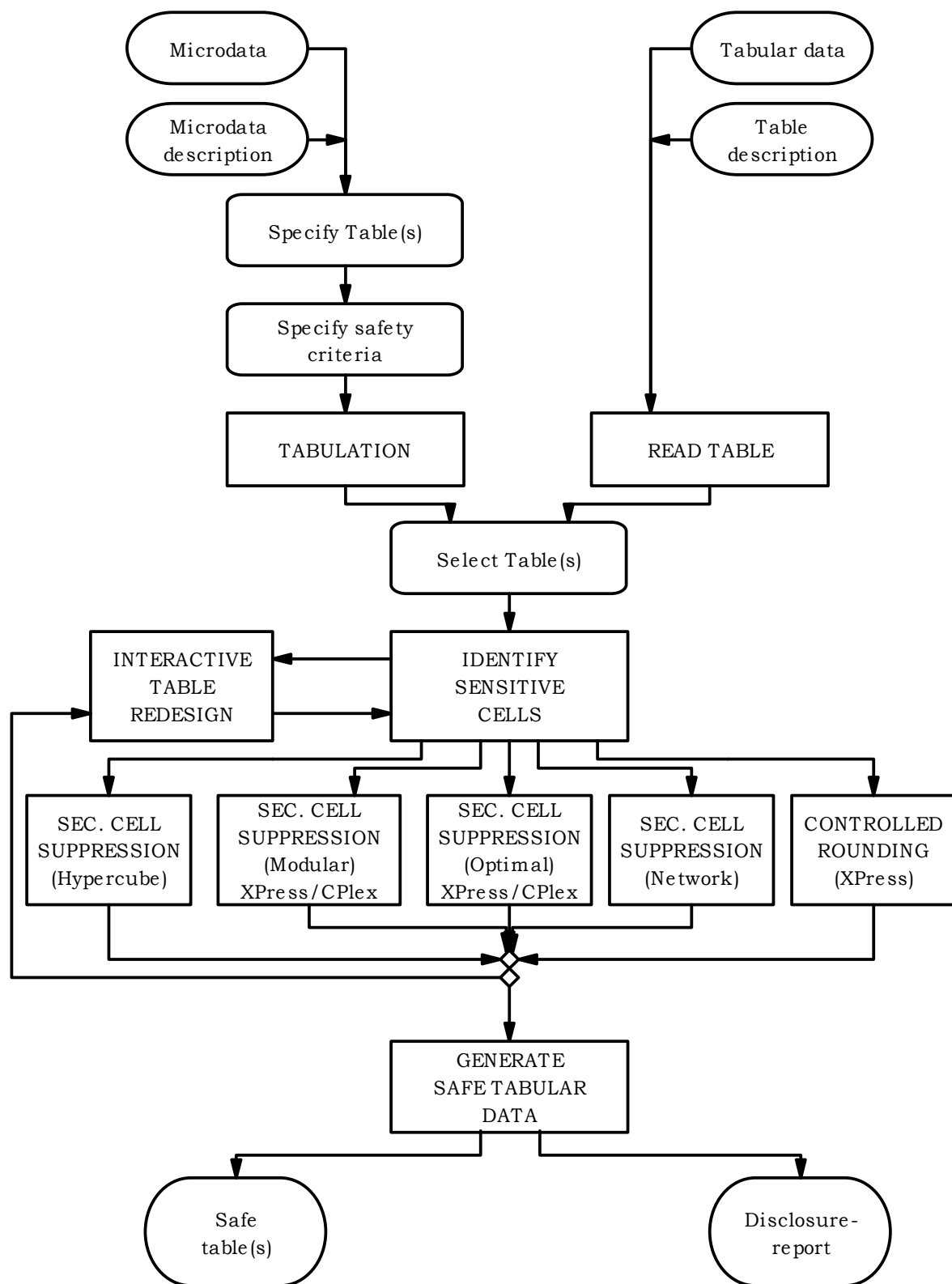
<sup>16</sup> Fischetti, M, Salazar Gonzales, J.J. (2000), ‘Models and Algorithms for Optimizing Cell Suppression Problem in Tabular Data with Linear Constraints’, in Journal of the American Statistical Association, Vol. 95, pp 916

hierarchical tables in section 4.3.1. of the ESSNET SDC Handbook) assign secondary suppressions considering only a part of the table relations at a time, e.g. those referring to the ‘current’ subtable. These methods are able to protect each subtable properly in the sense that the feasibility intervals of the sensitive cells indeed cover the protection intervals. But this holds only, if the feasibility intervals are computed considering only the table relations of the particular subtable. But for a hierarchical table, feasibility intervals computed on basis of the set of relations for the *full* table normally tend to be closer than those computed on basis of separate sets of relations corresponding to individual sub-tables. Hence, in a hierarchical table, it is not unlikely that the Audit routine discovers that some cells were not protected properly.

### **Discovering singleton problems**

Making use of the additional knowledge of a respondent, who is the single respondent to a cell (a so called ‘singleton’), it is possible to derive intervals that are much closer than without this knowledge. The audit routine could be used to identify problems in this respect in the following way: in advance of running the audit routine, set the status of a particular singleton cell from “unsafe” to “safe”.

## 2.15 Functional design of $\tau$ -ARGUS



### 3 A TOUR OF $\tau$ -ARGUS

---

In this chapter, we explain and display the key features of  $\tau$ -ARGUS.  $\tau$ -ARGUS is a menu driven program, and here we describe a number of menus which the user will follow in order to prepare a table for output in a ‘safe’ form. The aim of the tour is to guide the user through the basic features of the program without describing every feature in detail. The only pre-requisite knowledge is basic experience of the Windows environment. In Chapter 4 (Reference) a more systematic description of the different parts of  $\tau$ -ARGUS will be given. In Chapter 5Chapter 3 can be read as a standalone chapter as there is enough detail to enable the user to run the program. However, not every option is covered and the user is pointed in the direction of the Reference chapter in a number of instances. In addition, back references to the theory explained in Chapter 2 are also indicated. In this tour we will use the data in the file `tau_testW.asc`, which comes with the installation of  $\tau$ -ARGUS.

In this tour we will start with the fixed format data file `tau_testW.asc`, build a table from that file and go through the process of disclosure control and finish with a protected safe table.

The key windows for preparation of the data and the processes of disclosure control (depicted graphically in the figure in section 2.14) are explored in this tour, which are given below.

#### Preparation

- *First steps.* Before using  $\tau$ -ARGUS for the first time, some options should be set to make  $\tau$ -ARGUS better usable in your environment. Also you can select the solver you want to use in secondary cell suppression. See section 3.1.1
- *Open Microdata.* This involves declaring both the microdata and the associated metadata. See section 3.1.2
- *Specify Metafile.* This shows how the metafile can be entered when there is no metafile available, or can be edited after being read in but before any tables are being specified. This includes options such as declaring variables to be explanatory or response, and setting up the hierarchical structure of the data and the location of the variable in the file. See section 3.1.3
- *Specify Tables.* Declare the tables for which protection is required, along with the safety rule and minimum frequency rule on which the primary suppressions will be based. When this has been finished the tables will be computed or read in. See section 3.1.4

#### Process of Disclosure Control

- *View Table.* This is the heart of  $\tau$ -ARGUS. View the table will show the tables when they have been computed or read in and when all the safety rules for primary suppressions have been applied.

You can inspect the table; get information about the number of unsafe cells etc. It contains options to modify the table using global recoding. There are several options to make the table safe via secondary cell suppression and rounding.

Also an audit procedure is available to check quality of any secondary suppression pattern.

- *Save Table*. The user can save the ‘safe’ table in a number of formats as will be seen in section 3.2.2.

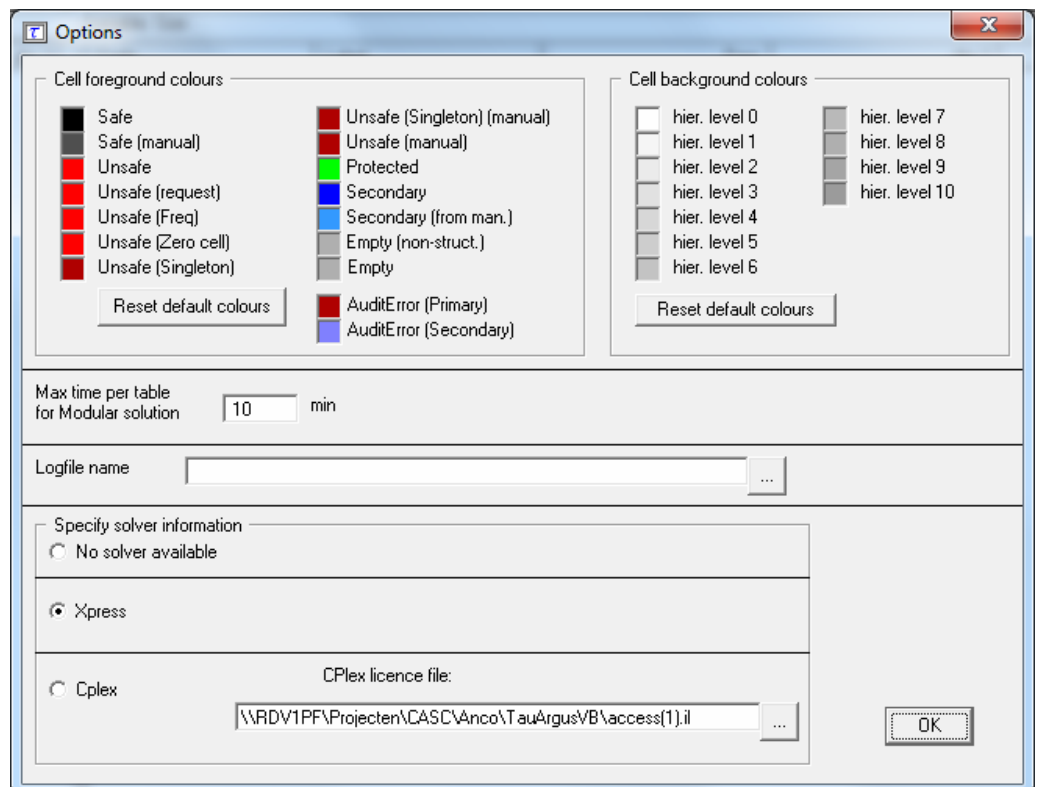
## 3.1 Preparation

### 3.1.1 First steps

Via (Help|Options) you can open the options window.

Before starting the process of protecting a table, you can customise  $\tau$ -ARGUS. Different colour setting can be selected here, but you can change them later as well. Some methods for secondary suppression (the modular and the optimal) require an external linear programming solver. For the complex problems of  $\tau$ -ARGUS we have decided to use high quality commercial solvers. Although  $\tau$ -ARGUS is freeware software these solvers are commercial packages and you have to acquire a licence for them separately. More information can be found on the CASC-website (<http://neon.vb.cbs.nl/casc/>) The choice of this solver must be decided before protecting a table. The choices are either Xpress or Cplex, the different LP\_solver supported by  $\tau$ -ARGUS. See also section 2.9 for more details.

You can select Xpress or Cplex. For Cplex the name of the licence file must specified. However also without a licence you can use  $\tau$ -ARGUS, but only a limited number of secondary suppression options are available.



Once this window has been opened details of the solver can be entered. Other

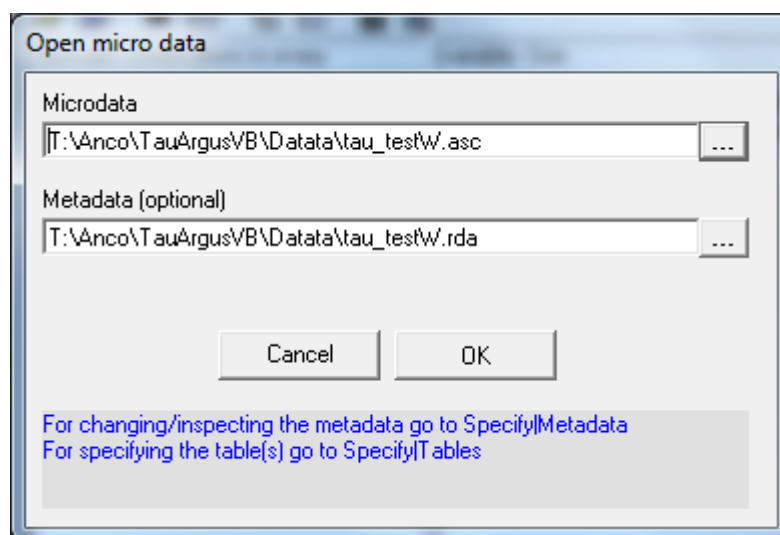
alterations to the default such as the colour of values in particular cells can also be made.

### 3.1.2 Open a microdata file

In this tour we only deal with how to open a fixed format microdata file (sections 3.1.2 to 3.1.4). If an already constructed table is to be used, then go to the Reference chapter (section 4.2.2). To start disclosure control with  $\tau$ -ARGUS there are two possible options:

1. Open a microdata file from which a table can be constructed,
2. Open an already-constructed table.

Both a microdata file and the metadata file describing this microdata file are required. The microdata file must be either a fixed format ASCII file or a free format file with a specified separator. By clicking (File|Open Microdata) you can specify both the name of the microdata file and the name of the file containing the metadata.



For most uses of  $\tau$ -ARGUS, the microdata and metadata file are stored in separate files. The simplest way to use the program is to use the extension .ASC for the (fixed format) datafile and .RDA for the metadata file. If the name of the metadata file is the same as the datafile, except for the extension, and it already exists in the same directory,  $\tau$ -ARGUS will fill in the name of this file automatically in the space under the metadata heading. If no metadata file is specified, the program has the facility to let you specify the metadata interactively via the menu option (Specify|Metafile). This is also the place to make changes to the metadata file. In subsection 3.1.3 we will give a description of the metadata file for  $\tau$ -ARGUS.

As an alternative to the ASCII files the microdata can be stored in an SPSS-system file. There is an option to read the system files.

### 3.1.3 Specify metafile

When you enter or change the metadata file interactively using  $\tau$ -ARGUS the option (Specify|Metafile) will bring you to the following screen:

The key elements of this window are the definitions for each variable. Most variables will be defined as one of the following.

- Explanatory Variable: a variable to be used as a categorical (spanning) variable when defining a table.
- Response Variable: a variable to be used as a cell item in a table.
- Weight variable: a variable containing the sampling weighting scheme.

More details on these variables along with the others can be found in the Reference chapter (subsection 4.3.1).

Other important features of this window are as follows.

- Codelist:  $\tau$ -ARGUS will automatically build the codelists for the explanatory variables, or you can specify a codelist file (a list-of-codes of the explanatory variables) as follows.
- Automatic: The codelist is created from the categories in the variable.
- Codelist file: The codes can be read in from an external file. Each category can contain a label. The codelist is only used for enhancing the presentation but always  $\tau$ -ARGUS will build a codelist from the datafile itself.
- Missing values: this gives information on the missing values which are attached to a codelist. Two distinct missing value indicators can be set (the reason for this is for the purposes of indicating different reasons for missing values: for example perhaps non-responses of different forms: maybe one code for the

response ‘*don't know*’, and another for ‘*refusal*’). Missing values however are not required.

- Hierarchical codes: The hierarchy can be derived from
- The digits of the individual codes in the data file or
- A specified file containing the hierarchical structure.

Examples are shown in the metafile information below.

### The Metafile

The metafile describes the variables in the microdata file, both the record layout and some additional information necessary to perform the SDC-process. Each variable is specified on one main line, followed by one or more option lines. An example is shown here. The leading spaces shown only serve only to make the file more readable; they have no other meaning.

```

Year 1  2 99
  <RECODEABLE>
IndustryCode 4  5 99999
  <RECODEABLE>
  <HIERARCHICAL>
  <HIERLEVELS>  3 1 1 0 0
Size 9  2 99
  <RECODEABLE>
Region 12  2 99
  <RECODEABLE>
  <CODELIST> Region.cdl
  <HIERCODELIST> Region2.hrc
  <HIERLEADSTRING> @
  <HIERARCHICAL>
Wgt 14  4 9999
  <NUMERIC>
  <DECIMALS>  1
  <WEIGHT>
Var1 19  9 9999999999
  <NUMERIC>
Var2 28  10 9999999999
  <NUMERIC>
  <DECIMALS>  2
.....
.....

```

### Details of the variables

‘*Year*’ : For this variable each record begins on position 1, is 2 characters long and missing values are represented by 99. It is also recodeable implicitly stating that it is an explanatory or spanning variable used to create the tables.

‘*IndustryCode*’: For this variable each record begins on position 4 and is 5 characters long. Missing values are represented by 99999. As well as being recodeable this variable is hierarchical and the hierarchy structure is specified. The first 3 characters are in the top hierarchy level, the 4<sup>th</sup> character in the second level and the 5<sup>th</sup> character in the lowest level. The two zeros at the end of this definition are redundant in this example.

‘Size’: For this variable each record begins on position 9 and is 2 characters long, and missing values are represented by 99. It is also recodeable.

‘Region’: For this variable each record begins on position 12 and is 2 characters long. Missing values are represented by 99. An example of a codelist file can be found in *region.cdl* and of a hierarchical codelist file in *region2.hrc*. Contents of these files are shown here.

The file *region.cdl*:

```
1,Groningen
2,Friesland
3,Drenthe
4,Overijssel
5,Flevoland
6,Gelderland
7,Utrecht
8,Noord-Holland
9,Zuid-Holland
10,Zeeland
11,Noord-Brabant
12,Limburg
Nr,North
Os,East
Ws,West
Zd,South
```

The file *region.hrc*:

```
Nr
@ 1
@ 2
@ 3
Os
@ 4
@ 5
@ 6
@ 7
Ws
@ 8
@ 9
@10
Zd
@11
@12
```

Additional details of these coding files can be found in the Reference chapter (section 4.2.1).

### 3.1.4 Specify tables

When the metadata file is ready, the tables to be protected can be specified. This is achieved via Specify|Specify Tables. A window to specify the tables is presented. In the example here we have a 2 dimensional table (2 explanatory variables) and a single response variable. A safety rule has been defined.

**Specify Tables**

**explanatory variables**

Year  
IndustryCode  
Size  
Region

**cell items**

Var1  
Var2  
Var3  
Var4  
Var5  
Var6  
Var7  
Var8  
<freq>

**response variable:**

Var2

**shadow variable:**

**cost variable**

☐ unity ☐ frequency ☒ variable ☐ distance function

lambda 1 MSC 20000

☐ Apply Weights  
☐ Missing = safe  
☐ Use holdings info

**Dom Rule**

	P	Q	N
Ind-1	15	100	1
Ind-2	0	100	0
Hold-1	0	100	0
Hold-2	0	100	0

☐ Minimum frequency  
Ind. 5 30 %  
Hold. 0 30 %

☐ Zero unsafe  
range: 0

☐ Manual safety range: 0 %

Expl. vars	rule	Resp. var	Shadow & Cost var
Size,Region	IND.: p= 15, q= 100, N= 1	Var2	Shadow=Default,Cost=Default

Cancel Compute tables

The key elements of this window are as follows.

### Explanatory variables

On the left is the listbox with the explanatory variables.

Click on '>' the selected variable to the next box in which the selected explanatory variables can be seen. From the box on the left hand side, containing explanatory variables, the variables that will be used in the row or the column of the table, in a 2-way table can be selected. Up to four explanatory variables can be selected to create a table.

### Cell items

The 'cell items' box contains the variables, which were declared as 'response variables' in the metafile. By using the '>' button they can be moved to the 'response variable' box to be used in the defined table.

### Response variable

Any variable in the cell items box can be chosen as the response variable. More than one response variable can be chosen.

### Shadow variable

The shadow variable is the variable which is used to apply the safety rule. By default this is the response variable. More detail on the Shadow variable can be found in section 4.3.4 in the Reference chapter.

### Cost variable

This variable describes the cost of each cell. These are the costs that are minimised when the secondary suppressed cells are calculated (See section 2.6 in the Theory

chapter for the further details). By default this is the response variable but other choices are possible. If the response or any other explicitly specified variable is used for this purpose, the circle next to 'variable' should be filled. Then, any variable name can be transferred from the cell items to the cost variable window. It is also possible to use the frequency of the cells as a cost-function. This will suppress cells with respect to number of contributors to each cell. A third option is that the number of cells to be suppressed is minimised, irrespective of the size of their contributions (unity option – cost variable is set to 1 for each cell). More details will be given in the Reference Chapter along with an example (section 4.3.4). Note that choice of the cost variable does not have any impact when using the hypercube method for secondary suppression.

### Weight

If the data file has a sample weight, specified in the metadata file, the table can be computed taking this weight into account. In this case, the 'apply weights' box should be ticked. More details will be given in the Reference Chapter along with an example (section 4.3.4).

### The safety rule

The concept of safety rules is explained in section 2.2 in the chapter on Theory. In this window the left side of the window allows the type of rule to be selected, this is usually either the dominance rule or p% rule, along with the necessary parameter values. Several rules together can be set for any particular table.

Additionally, the minimum number of contributors (threshold rule) can be chosen. In the window this is referred to as the 'Minimum Frequency'

*Now, brief summaries are provided to define the Dominance and p% rules.*

### Dominance Rule

This is sometimes referred to as the (n,k) rule where n is the number of contributors to a cell contributing more than k% of the total value of the cell (if the cell is to defined as unsafe for publication).

### p% rule

The p% rule says that if  $x_1$  can be determined to an accuracy of better than p% of the true value then it is disclosive where  $x_1$  is the largest contributor to a cell.

This rule can be written as:

$$\sum_{i=3}^c x_i \geq \frac{p}{100} x_1$$

for the cell to be non-disclosive where  $c$  is the total number of contributors to the cell and the intruder is a respondent in the cell.

It is important to know that when entering this rule in  $\tau$ -ARGUS the value of  $n$  refers to the number of intruders in coalition (who wish to group together to estimate the largest contributor).

A typical example would be that the sum of all reporting units excluding the largest two must be at least 10% of the value of the largest. Therefore, in  $\tau$ -ARGUS set  $p=10$  and  $n=1$  as there is just one intruder in the coalition, respondent  $x_2$ .

The choice of safety rule is specified by the user and the chosen parameters can then be entered. From these parameters symmetric safety ranges are computed automatically prior to the secondary suppressions.

For the minimum frequency rule, a safety range is calculated from the user given

range. This is usually a small positive value and is required to enable secondary suppression to be carried out.

A manual safety range is also required for cells that can be made unsafe by intervention of the user.

Other options such as the 'Request Rule' or the 'Holding Rule' will be looked at in more detail in the Reference chapter (section 4.3.4).

When everything has been filled in, click 'v' to transport all the specified parameters describing the table to the 'listwindow' on the bottom. As many tables as you want may be specified, only limited by the memory of the computer. If a table is to be modified press the '^' button.

### **Creating the Table**

Pressing the 'Compute tables' button will invoke  $\tau$ -ARGUS to actually compute the tables requested and the process to start disclosure control may be invoked.  $\tau$ -ARGUS will come back with the main window showing the number of unsafe cells per variable, per dimension, as explained in the next section 3.2.

## **3.2 The Process of disclosure control**

---

When the table(s) have been calculated, the main-window of  $\tau$ -ARGUS will be displayed again with an overview of all the unsafe cells per variable (over all the tables). An example is shown here. This window (underneath the main menu for  $\tau$ -ARGUS), shows the number of unsafe combinations per variable. For example there are no single unsafe cells in dimension one for either variable. (i.e. the one-way marginal total for different values of 'Size' and 'Region' are all not disclosive).

The right hand window gives the equivalent information for each level of the variable indicated on the left. For example, there are 12 unsafe cells in the two-way 'Size' x 'Region' table.

### **Size**

There are however 12 unsafe cells in the 2 way table *Size* by *Region* as can be seen by the right hand window which gives the equivalent information for each level of the variable indicated on the left. There are 5 unsafe cells where Size = 2, 6 unsafe cells where Size = 4 and a single unsafe cell where Size = 9.

**TauARGUS**

File Specify Modify Output Help

#unsafe combinations in every

Variable	dim 1	dim 2
Size	0	9
Region	0	9

variable: Region

Code	Label	Freq	dim 1	dim 2
	Total	42723	0	0
Nr	North	11395	0	2
·1	Groningen	6112	0	2
·2	Friesland	3798	0	0
·3	Drenthe	1485	0	0
Os	East	10227	0	1
·4	Overijssel	538	0	2
·5	Flevoland	1597	0	0
·6	Gelderland	5446	0	2
·7	Utrecht	2646	0	0
Ws	West	10054	0	0
·8	Noord-Holland	1182	0	0
·9	Zuid-Holland	8018	0	0
10	Zeeland	854	0	0
Zd	South	11047	0	0
11	Noord-Brabant	7511	0	0
12	Limburg	3536	0	0
99		0	0	0

Status 3/29/2011 2:01 PM

### Region

Two of the 'North' subtotal cells are disclosive. Within the 'North' region, 2 cells in Groningen are disclosive. Two 'East' subtotal cells are also unsafe. Within the 'East' region 2 cells in 'Overijssel' are unsafe along with 2 cells in 'Gelderland'. Finally there is 1 unsafe cell for the 'South' subtotal and within this region there is 1 unsafe cell for 'Noord-Brabant'

### 3.2.1 View table

The 'Modify/View table' option shows the calculated table plus some additional information. This is regarded as the key window in  $\tau$ -ARGUS.

The selected table is displayed in a spreadsheet view. Safe cells are shown in black, whilst cells failing the safety rule and/or minimum frequency rule are displayed in red. These are the default colours.

The user now has to decide whether to carry out secondary suppressions immediately or to perform some recoding first. There are other options such as changing the status of individual cells manually, this will be discussed further in the Reference chapter (see section 4.4.2).

## Cell information

Cells can be selected in the table by moving the cursor arrow. In each case, information about the selected cell is shown on the right.

The status of the cell can be one of the following. Some of the terms will be explained later in this section but others are expanded upon in the Reference section 4.4.2.

- Safe: Does not violate the safety rule
- Safe (from manual): manually made safe during this session
- Unsafe: According to the safety rule
- Unsafe (request): Unsafe according to the Request rule.
- Unsafe (frequency): Unsafe according to the minimum frequency rule.
- Unsafe (zero cell) Unsafe because the zero-cells are considered unsafe..
- Unsafe (from manual): manually made unsafe during this session.
- Protected: Cannot be selected as a candidate for secondary cell suppression.
- Secondary: Cell selected for secondary suppression.
- Secondary (from manual): Unsafe due to secondary suppression after primary suppressions carried out manually.
- Zero: Value is zero and cannot be suppressed.
- Empty: No records contributed to this cell and the cell cannot be suppressed.

## Change Status

The second pane ('Change Status') on the right will allow the user to change the

cell-status.

- Set to Safe: A cell, which has failed the safety rules is here declared safe by the user.
- Set to Unsafe: A cell, which has passed the safety rules is here declared to be unsafe by the user.
- Set to Protected: A safe cell is set so that it cannot be selected for secondary suppression.
- Set Cost: Change the value of the Cost-value for this cell
- Use 'a priori' information (see below).

### A Priori Info

This option is an *a priori* option to be mainly used for microdata which allows the user to feed  $\tau$ -ARGUS a list of cells where the status of the standard rules can be overruled i.e. the status of the cells is already specified. The associated file specifying this information is free format. The format will be:

Code of first spanning variable, Code of second spanning variable, Status of cell (u = unsafe, p = protected (not to be suppressed), s = safe).

Also the cost-function can be changed here for a cell. This will make the cell more likely to become secondary cell suppression, when the value is low or less likely when the value is high.

Nr, 4, u
Zd, 6, p
5, 5, c, 1

A full description of the aproiri file can be found in section 5.5

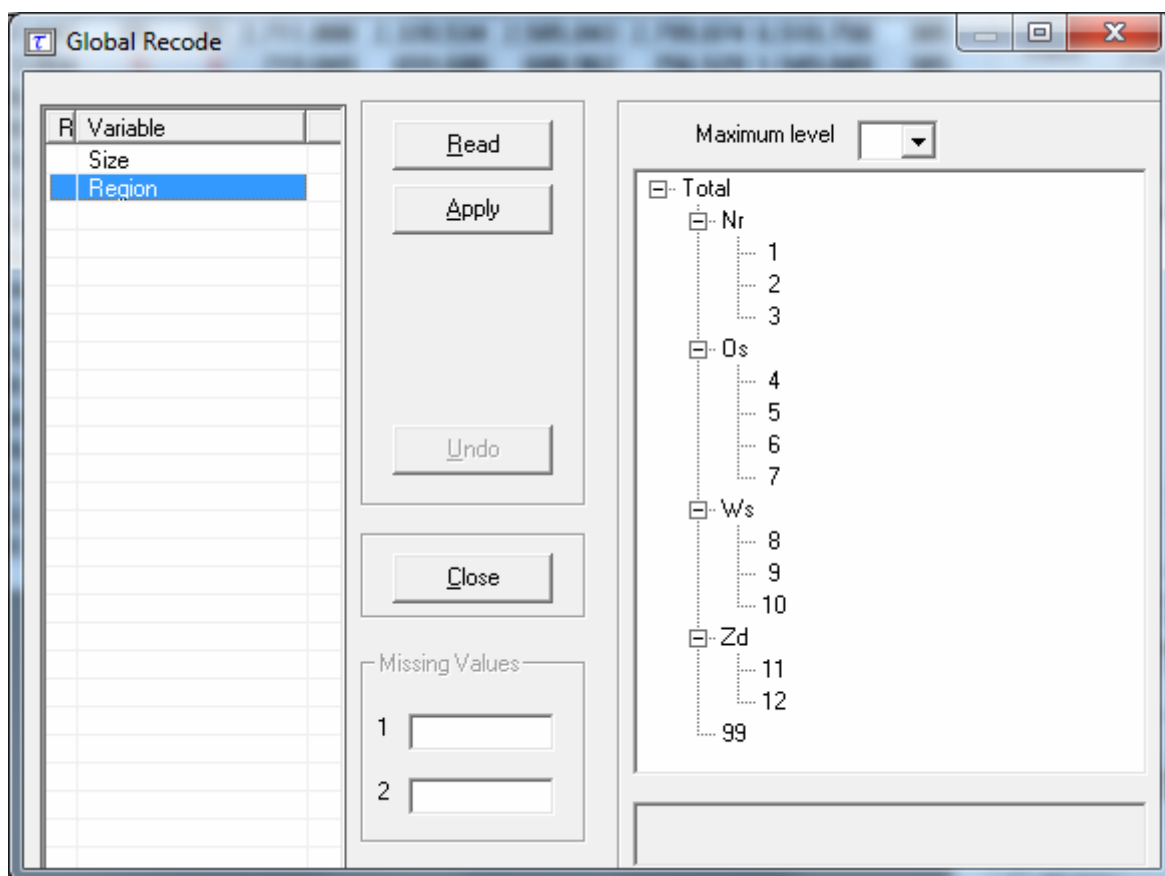
### Recode

The recode button will bring the user to the recoding system. Recoding is a very powerful method of protecting a table. Collapsed cells usually have more contributors and therefore tend to be much safer.

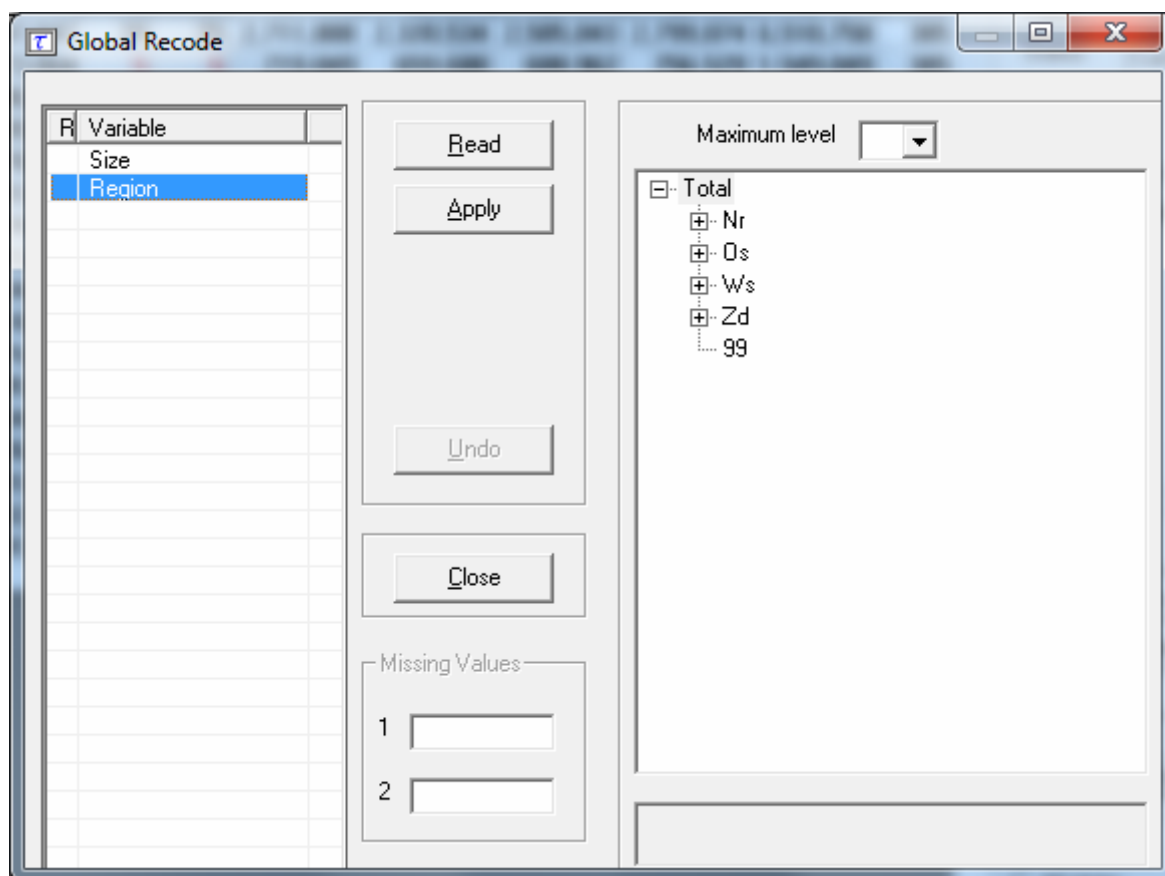
#### *Hierarchical Recoding*

The first window shows the codes awaiting recoding

In this example, the 'Region' variable to recode has been selected. As 'Region' is a hierarchical variable, the codes are shown in a hierarchical tree. The user can either fold or unfold the branches by clicking on the '+' or '-' boxes which results in showing or omitting codes from the table, or by choosing an overall maximum hierarchical level. (See the following windows for details). Pressing the 'Apply' button followed by 'Close' will actually apply the selected recoding. Press the undo-button – it is now possible to go back to the original recoding scheme. Below this there are two windows, one showing the recode window prior to applying the recoding for 'Region' and the second showing the table following recoding.



This window shows the new hierarchical codes after collapsing all second level categories



By clicking 'Apply', we obtain this window which shows the table after recoding

	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
.Nr	4,373,664	5	5	719,049	659,680	688,962	756,529	1,549,049	385
.Os	3,703,896	15	5	642,238	515,003	534,147	620,392	1,392,096	-
.Ws	4,576,116	-	-	648,972	543,570	663,897	775,132	1,944,545	-
.Zd	4,193,971	-	15	701,549	602,281	618,037	647,021	1,625,068	-
.99	-	-	-	-	-	-	-	-	-

Cell Information:

Value: 0

Status: Empty

Cost: 0

Shadow: 0

# contributions: 0

Top n of shadow: 0

☐ Holding level

Request: 0

Change status:

Set to Safe

Set to Unsafe

Set to Protected

Set Cost

A priori info

All Non-StructEmpty

Recode

Suppress:

☒ HyperCube

☐ Modular

☐ Network

☐ Optimal

☐ Marginal

☐ Rounding

☐ InverseWgt

Singleton

Undo Singleton

Suppress

Undo Suppress

Audit

3 dig. separator

Output View

Select Table

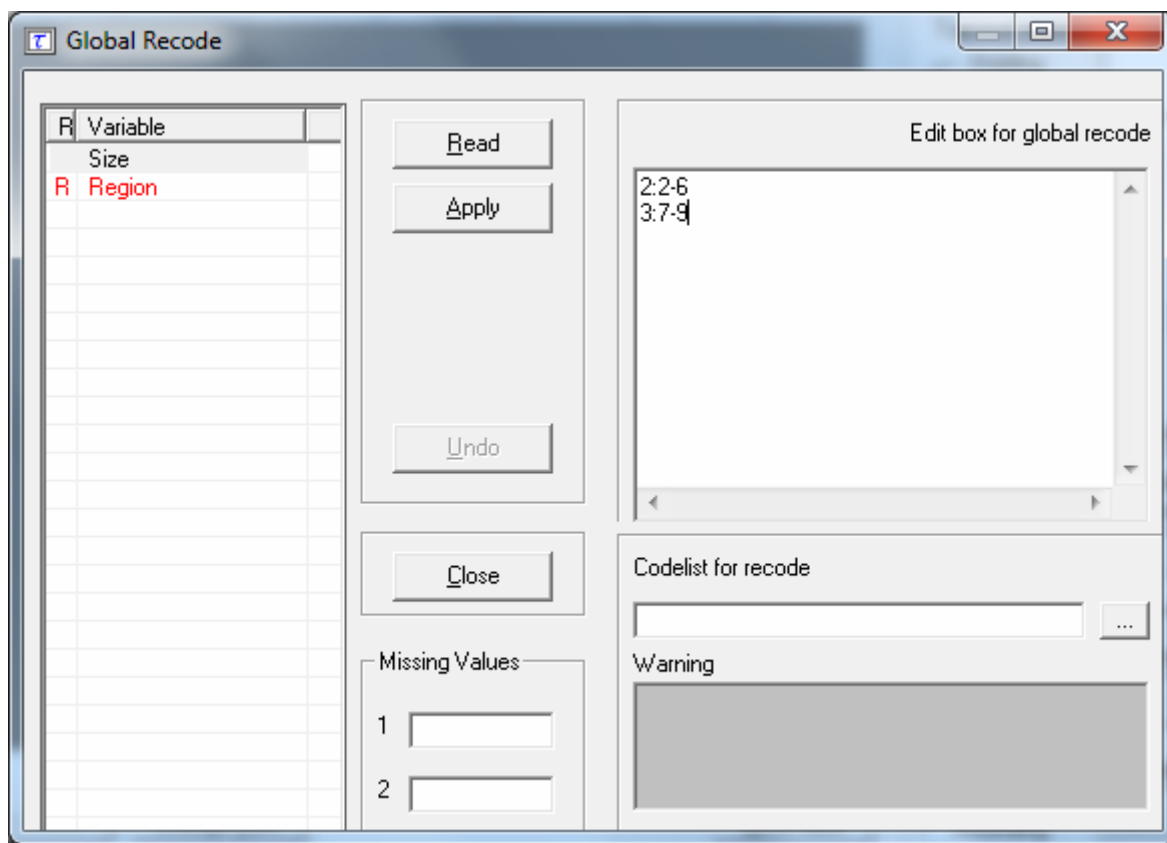
Change View

Table Summary

Write table

Close

### Non Hierarchical Recoding



In this example the non-hierarchical 'Size' variable has been selected to be recoded. The user can either write the required recodings in the edit box or import them from a previously written file. In the example the line 2:2-6 results that categories 2,3,4,5, and 6 will be recoded into a new category 2, whilst categories 7,8 and 9 will be recoded into the new category 3.

Once the recoding has been applied (both for hierarchical and non hierarchical data) the table can again be displayed. If there are now no cells, which fail the safety rules, the table can be saved as a protected table. However, if there are still a number of unsafe cells, secondary suppression needs to be carried out. This is necessary as the table is not yet safe. If only the cells failing the safety rules are suppressed, other cell values could be obtained by differencing.

### Secondary Suppression

The Suppress button is an important button. It will activate the modules for computing the necessary secondary suppressions as described above. There are a number of options here.

- Hypercube
- Modular
- Network
- Optimal

### Hypercube

This is also known as the GHMITER method. The approach builds on the fact that a suppressed cell in a simple n-dimensional table without substructure cannot be disclosed exactly if that cell is contained in a pattern of suppressed, nonzero cells, forming the corner points of a hypercube.

### **Modular**

This partial method will break the hierarchical table down to several non-hierarchical tables, protect them and compose a protected table from the smaller tables. As this method uses the optimisation routines, an LP-solver is required: this will be either XPRESS or CPLEX. The routine used can be specified in the Options window, this will be discussed later.

### **Optimal**

This method protects the hierarchical table as a single table without breaking it down into smaller tables. As this method uses the optimisation routines, an LP-solver is required: this will be either XPRESS or CPLEX. The routine used can be specified in the Options window, this will be discussed later.

### **Network**

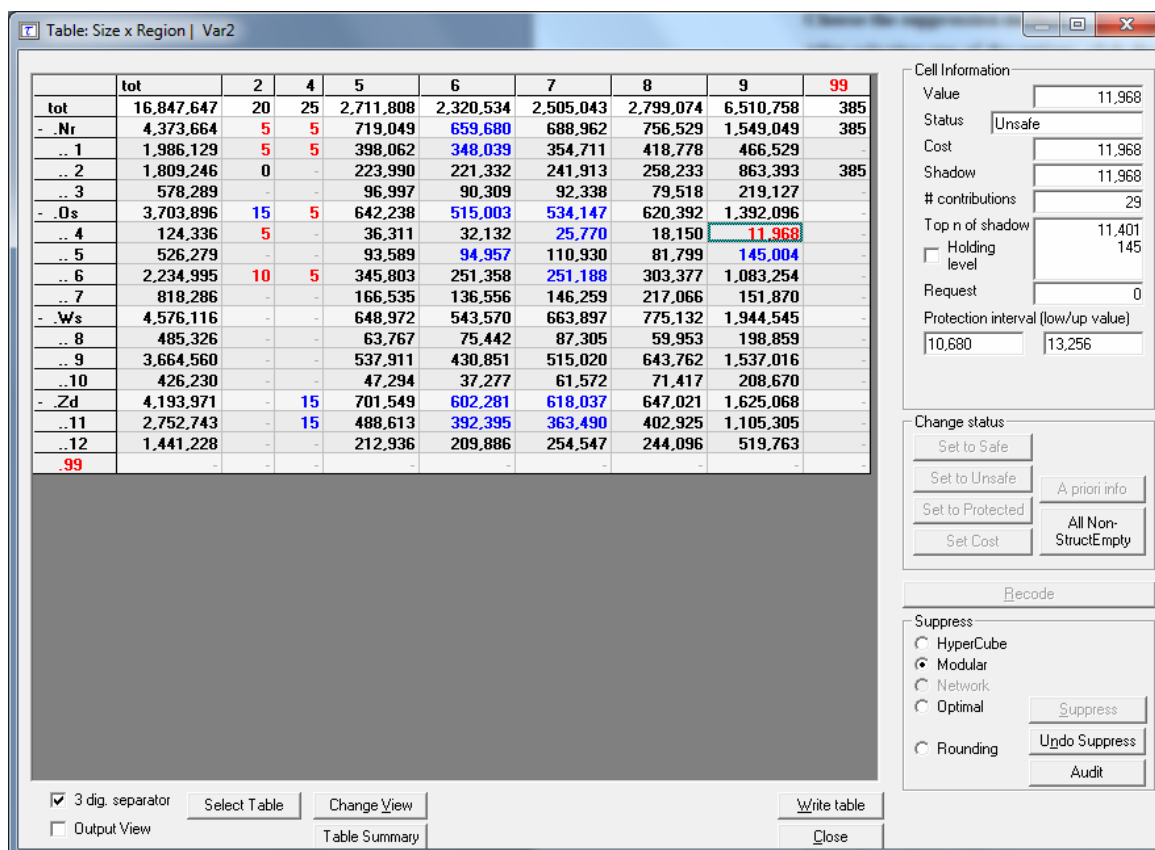
This is a Network Flow approach for large unstructured 2 dimensional tables or a 2 dimensional table with one hierarchy (the first variable specified). This method is also based on optimisation techniques, but does not require an external solver like XPress or Cplex.

### **Rounding**

The controlled rounding procedure can be applied. The user has to specify the rounding base. Note that this option requires the XPress solver.

### **Choose the suppression method**

The radio-buttons at the right lower part of the window allow selecting the desired suppression method. Clicking on the Suppress button will then start the process of calculating the secondary suppressions. When this process has finished the protected table will be displayed and also the user will be informed about the number of cells selected for secondary suppression and the time taken to perform the operation. The secondary suppressed cells will be shown in blue by default.



## Summary Window

By clicking on 'Table Summary', the summary window is obtained. The summary window gives an overview of the cells according to their status.

- Freq: The number of cells in each category
- # rec: The number of observations in each category
- Sum Resp: Total cell value in each category

SumCost: The sum of the cost variable. Here it is equal to the response variable.

By clicking on 'OK', we return to the table window. The table may now be written as an output file in the required format. Any cells which have been selected for suppression will be replaced by 'X'. The safe table can be saved by using the 'Write table' button in this window, or by using 'Output|Save table' on the main menu.

Summary for table no: 1

Explan. Var	# Codes	Status	Freq	# rec	Sum Resp	SumCost
Region	18	Safe	98	247673	98549517.04	98549517.04
Size	9	Safe (manual)	0	0	0	0
		Unsafe	9	43	12013	12013
		Unsafe (request)	0	0	0	0
		Unsafe (Freq)	0	0	0	0
		Unsafe (Zero cell)	0	0	0	0
		Unsafe (Singleton)	0	0	0	0
		Unsafe (Singleton) (m...	0	0	0	0
		Unsafe (manual)	0	0	0	0
		Protected	0	0	0	0
		Secondary	12	8622	2524351	2524351
		Secondary (from man.)	0	0	0	0
		Empty (non-struct.)	0	0	0	0
		Empty	43	0	0	0
		Total	162	256338	101085881.04	101085881.04

Protected by Modular

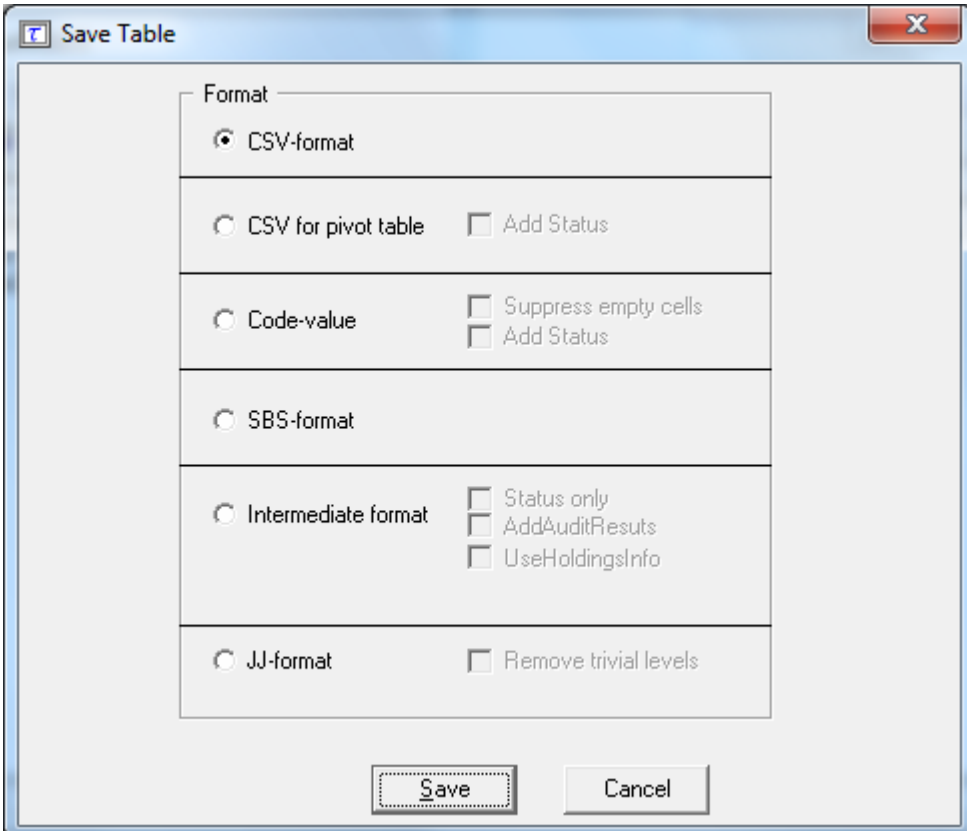
OK

### 3.2.2 Save the safe table

When the table is safe it may be written to the hard disk of the computer. The user has four options:

1. As a CSV file. This Comma separated file can easily be read into Excel. Please note that  $\tau$ -ARGUS uses the ',' as the field-separator in this CSV-file. This might influence opening the CSV file in Excel. A solution for this is to change the settings in the Windows control-panel. This is a typical tabular output maintaining the appearance of the table in  $\tau$ -ARGUS.
2. A CSV-file for a pivot table. This offers the opportunity to make use of the facilities of pivot table in Excel. The status of each cell can be added here as an option (Safe, Unsafe or Protected for example). The information for each cell is displayed on a single line unlike standard csv format.
3. A text file in the format code-value, this is separated by commas. Here, the cell status is again an option. Also empty cells can be suppressed from the output file if required. The information for each cell is displayed on a single line similar to the CSV file for a pivot table.
4. SBS format. This is a special format required for sending data to Eurostat.
5. A file in intermediate format for possible input into another program. This contains protection levels and external bounds for each cell. This table could even be read back into  $\tau$ -ARGUS.

Finally, a report will be generated to a user specified directory. This report will also be displayed on the screen when the table has been written. It will contain details such as table structure, safety rules (and number of cells failing), secondary suppression method (and number of cell failing) and details of any recodes. An example is shown in the Reference section 4.5.2. As this is an HTML-file it can be viewed easily later or printed.

A screenshot of a 'Save Table' dialog box. The dialog has a title bar with a small icon and the text 'Save Table', and a close button (X) in the top right corner. The main area is a light gray rectangle containing a 'Format' section. This section is a white box with a thin border, containing several radio buttons and checkboxes. The radio buttons are for 'CSV-format' (selected), 'CSV for pivot table', 'Code-value', 'SBS-format', 'Intermediate format', and 'JJ-format'. The checkboxes are for 'Add Status' (next to 'CSV for pivot table'), 'Suppress empty cells' and 'Add Status' (next to 'Code-value'), 'Status only', 'AddAuditResults', and 'UseHoldingsInfo' (next to 'Intermediate format'), and 'Remove trivial levels' (next to 'JJ-format'). At the bottom of the dialog are two buttons: 'Save' and 'Cancel'.

Save Table

Format

☒ CSV-format

☐ CSV for pivot table ☐ Add Status

☐ Code-value ☐ Suppress empty cells  
☐ Add Status

☐ SBS-format

☐ Intermediate format ☐ Status only  
☐ AddAuditResults  
☐ UseHoldingsInfo

☐ JJ-format ☐ Remove trivial levels

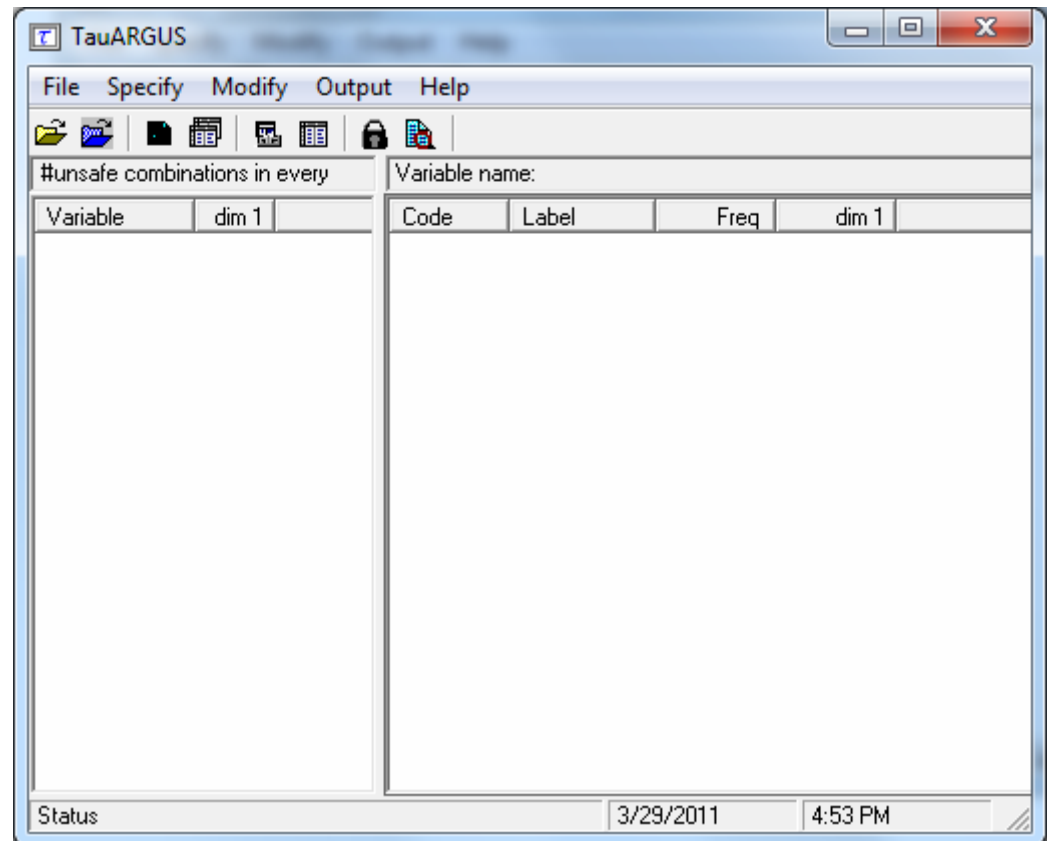
Save Cancel

## 4 REFERENCE SECTION - DESCRIPTION OF THE MENU ITEMS

---

Chapter 3 gave a brief introduction of the most frequently used options within  $\tau$ -ARGUS. In this section a more detailed description of the program by menu-item is presented. The information in this section is the same as the information shown when the help-facility of  $\tau$ -ARGUS is invoked.

In chapter 5 some general descriptions are given.



### 4.1 Main Window

---

There are five menu headings:

<b>File</b>	Under <b>F</b> ile either a microdata file or tabular data file can be opened together with the meta data file describing the data; in addition there is the option to open a Batch process file and to Exit.
<b>Specify</b>	<b>S</b> pecify allows the metadata to be entered or edited and the user can specify the tables to be protected along with primary sensitivity rules. When tabular data is the starting point the details about the table can be specified
<b>Modify</b>	Under <b>M</b> odify, the table can be selected, viewed and any secondary suppressions carried out. Also secondary suppression for linked tables can be performed.

<b>Output</b>	<b>Output</b> allows the suppressed table to be saved. In addition there is also view report and write a batchfile. Also a tool to generate a-priory information can be found here.
<b>Help</b>	Finally, there is a <b>Help</b> menu, with contents, news, options and about-box of the program.

Below is a list of the menu items which are shown under each of the menu headings. As some of the items are context specific they will not all be always available.

Overview of the menu-items

<a href="#"><u>File</u></a>	<a href="#"><u>Specify</u></a>	<a href="#"><u>Modify</u></a>	<a href="#"><u>Output</u></a>	<a href="#"><u>Help</u></a>
<a href="#"><u>Open Microdata</u></a>	<a href="#"><u>Metafile</u></a>	<a href="#"><u>Select Table</u></a>	<a href="#"><u>Save table</u></a>	<a href="#"><u>Contents</u></a>
<a href="#"><u>Open Table</u></a>	<a href="#"><u>Specify Tables</u></a>	<a href="#"><u>ViewTable</u></a>	<a href="#"><u>View Report</u></a>	<a href="#"><u>News</u></a>
<a href="#"><u>Open Table Set</u></a>		<a href="#"><u>Linked Tables</u></a>	<a href="#"><u>Generate Apriory</u></a>	<a href="#"><u>Options</u></a>
<a href="#"><u>Open Batch Process</u></a>			<a href="#"><u>Write Batch File</u></a>	<a href="#"><u>About</u></a>
<a href="#"><u>Exit</u></a>				

These menu items will be explained in detail in the following sections.

## 4.2 *The File menu*

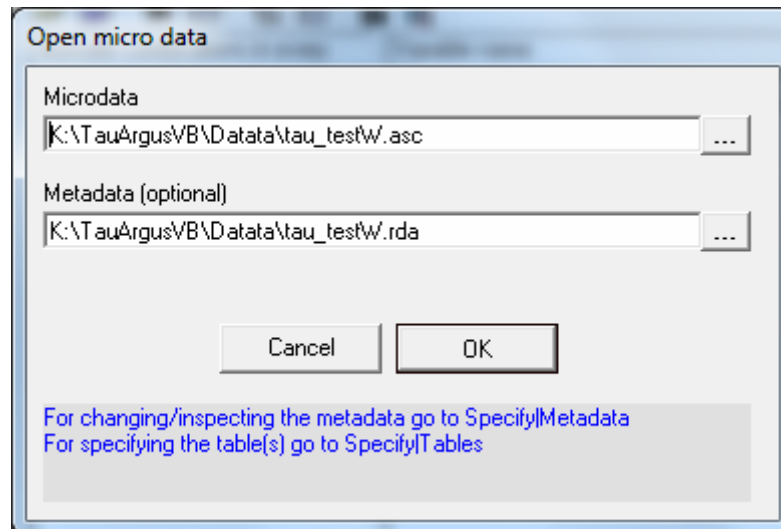
---

$\tau$ -ARGUS can read data in two ways. The first way is that  $\tau$ -ARGUS will read the data from a microdata file (fixed format, free format and a SPSS\_systemfile), which is explained in section 4.2.1. From this microdata  $\tau$ -ARGUS can then build a table and during this process compute all necessary additional information. This allows for most flexibility.

The second is input and treatment of a pre-formed tabulated data and is dealt with in section 4.2.2. Only one of these options can be used at one time, a table and a set of microdata cannot be read in  $\tau$ -ARGUS simultaneously.  $\tau$ -ARGUS can also be used in batch, see section 4.2.3


### 4.2.1 File | Open Microdata

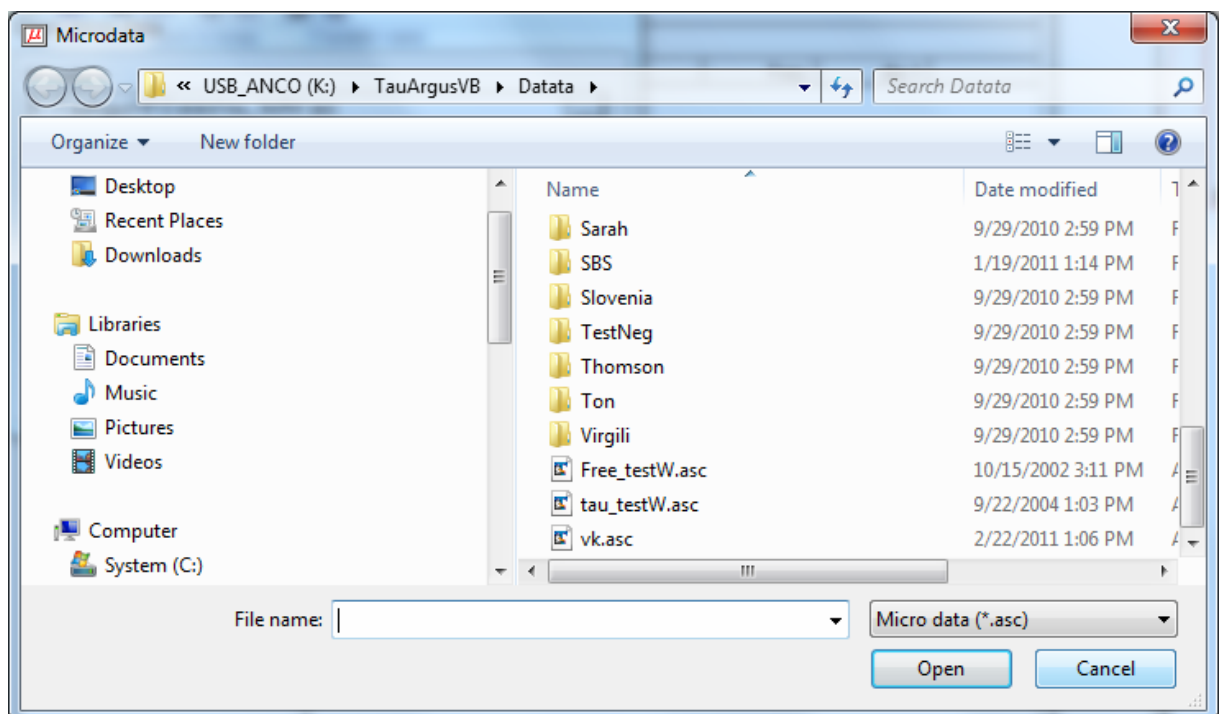
The File|Open microdata menu allows the user to specify the microdata file (both fixed and free format or a SPSS-system file) and optionally the metadata file.



In this dialog box the user can select the microdata-file or the SPSS system file and the corresponding metadata file

By default the microdata-file has extension .asc and the metafile .rda. (Note, the user may use any file extension, but is advised to use default names).

When the user clicks on  they get a traditional open file dialog box.



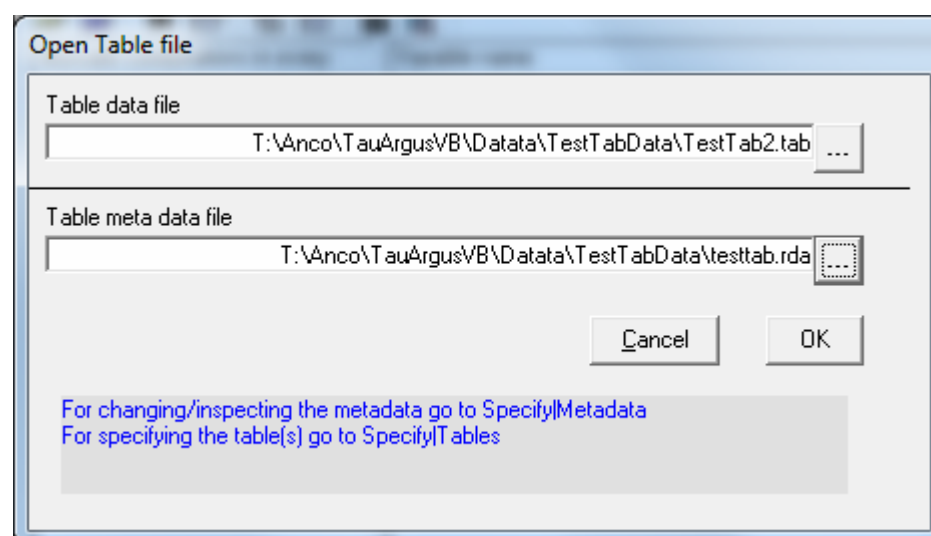
This box enables searching for the required files. Other file-extensions can be chosen when clicking on the files of type listbox. When the user has selected the microdata file a suggestion for the metafile (with the same name but with the extension .rda) is given but only when this file exists. Note, both files do not necessarily have to have the same name.

A full description of the metadata file can be found in section 5.1.

When the data file has been selected and optionally the meta data file. You can proceed to the menu options to edit/modify the meta data file and then specify the tables required. See section 4.3.1 and 4.3.4.

## 4.2.2 File | Open Table

This is the option allowing the input of tabular data into  $\tau$ -ARGUS. In this case, an already-constructed table is read in. This is reached by selecting 'Open a Table' on the main window of  $\tau$ -ARGUS.



The name of the datafile containing the table to be opened (in the format given below) needs to be specified in the top line. The name of the file containing the metadata is entered on line 2. Later on you will be offered the option of adapting the metadata or even enter the metadata from scratch.

There is a great flexibility with this option as it allows the status, the cell frequency, the top n values and the lower and upper protection levels to be entered for each cell. The more detail is given for each cell to more flexibility  $\tau$ -ARGUS offers in a later stage to apply sensitivity rules etc.

If only a cell status is provided,  $\tau$ -ARGUS can only use that and give each cell a fixed protection level of some percentage, if also the largest say 2 contributors are provided,  $\tau$ -ARGUS can apply most of the sensitivity rules.

Here, by clicking 'OK' this allows both re-specification of the metafile under the Specify|Metafile option and the setting safety rules using the 'Specify|Table Metadata' option.

**Format** (An example of a 2 dimensional table)

This (artificially generated) datafile shows 2 explanatory variables, cell value, cell frequency, the top 3 values in each cell. With this information  $\tau$ -ARGUS is still able to apply the primary sensitivity rules, like p% rule.

T,	T,	2940	,48,	200,200,200
T,	A,	745	,12,	200,100,100
T,	B,	810	,12,	200,100,100
T,	C,	685	,12,	200,100,100
T,	D,	700	,12,	200,100,100

1,	T,	795	,12,	200,100,100
1,	A,	350	,3,	200,100,50
1,	B,	190	,3,	100,50,40
1,	C,	150	,3,	100,40,10
1,	D,	115	,3,	50,40,25
2,	T,	670	,12,	200,100,100
2,	A,	115	,3,	50,40,25
2,	B,	340	,3,	200,100,40
2,	C,	115	,3,	50,40,25
2,	D,	120	,3,	100,10,10
3,	T,	785	,12,	200,100,100
3,	A,	190	,3,	100,50,40
3,	B,	115	,3,	50,40,25
3,	C,	325	,3,	200,100,25
3,	D,	165	,3,	100,40,25
4,	T,	690	,12,	200,100,100
4,	A,	100	,3,	50,25,25
4,	B,	175	,3,	100,50,25
4,	C,	115	,3,	50,40,25
4,	D,	310	,3,	200,100,10

Alternatively if only the status is given to  $\tau$ -ARGUS , there is no other option than to use the status and treat all unsafe cells as ‘manually’ unsafe.

T,	T,	2940	,u
T,	A,	745	,s
T,	B,	810	,s
T,	C,	685	,s
T,	D,	700	,s
1,	T,	795	,s
1,	A,	350	,s
1,	B,	190	,s
1,	C,	150	,s
1,	D,	115	,s
2,	T,	670	,s
2,	A,	115	,s
2,	B,	340	,s
2,	C,	115	,u
2,	D,	120	,u
3,	T,	785	,s
3,	A,	190	,s
3,	B,	115	,s
3,	C,	325	,s
3,	D,	165	,s
4,	T,	690	,s
4,	A,	100	,s
4,	B,	175	,s
4,	C,	115	,s
4,	D,	310	,s

For tables of dimension 3 or higher, additional columns for the explanatory variables would have to be added as well as additional rows to allow for the increased depth of the table.

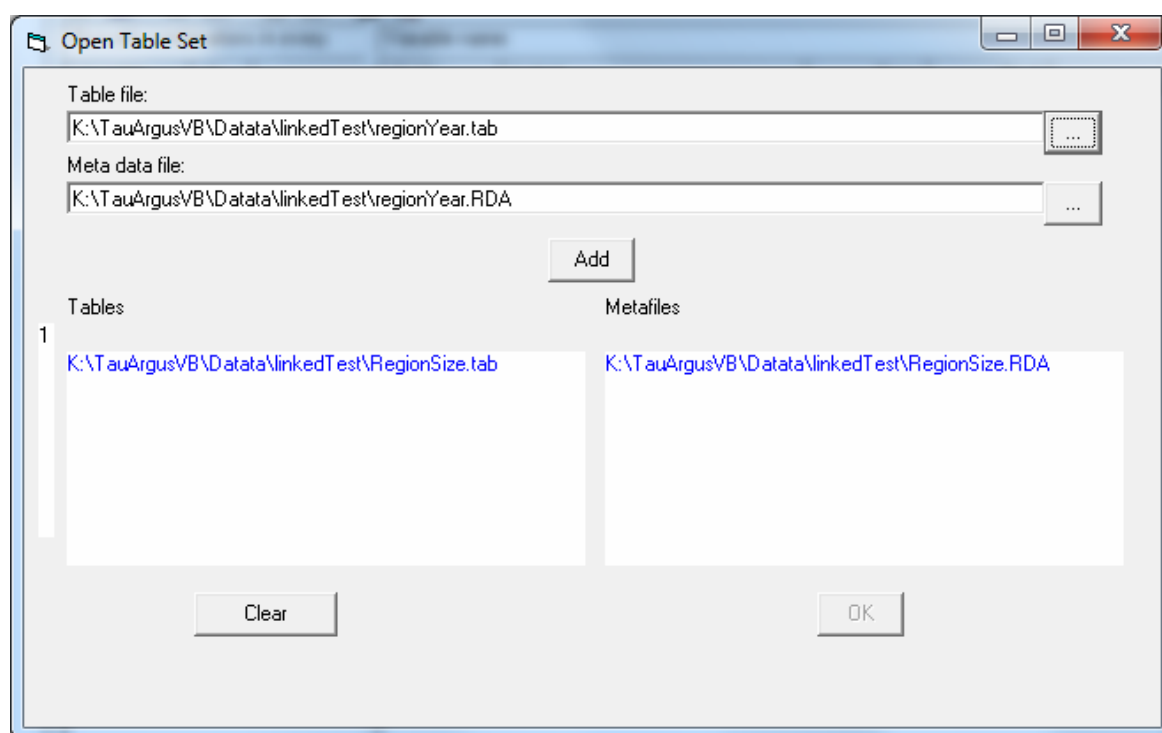
The next step will be to optionally edit the metadata and then read the table.

### 4.2.3 File | Open Table Set

When the linked tables procedure will be used in combination with tabular input, the option “Open Table Set” must be used to read a set of tables in  $\tau$ -ARGUS. The “Open Table” option described above (4.2.2) allows for only one single table.

In this option a set of tables with the corresponding meta data files (\*.rda) is specified.

When the set is complete, press the OK-button.



After pressing the OK-button, you will be guided automatically to the Specify Tables window for each table. This is described in section 4.3.5.

Please note that it is advisable to read each table in  $\tau$ -ARGUS before. This to be sure that the specification of the tables and the meta data is correct.

### 4.2.4 File | Open Batch Process

This option allows the user to run the commands in batch mode from opening the microdata and metadata, protecting the table and creating the output of the final table(s).

If the last line of the batch-file is <GOINTERACTIVE>  $\tau$ -ARGUS will first perform all the actions as specified in the batch-file and then open the main menu and giving the control to the user to continue the work in the interactive modus.

The lay-out of the batch-file is described in section 5.6.

Note that a log file is maintained of all actions. By default the log file is “Logbook.txt” in the temp-directory, but in the batch-file a different file can be chosen. Also from the command-line a log file name can be specified. See also section 5.6.

#### **4.2.5 File | Exit**

Exits the  $\tau$ -ARGUS-session.

### **4.3 *The Specify menu***

---

#### **4.3.1 Specify | Metafile [for microdata]**

Clicking on ‘Specify|metafile’ gives the user the opportunity to either edit the metafile already read in or to enter the metafile information directly at the terminal.

In this dialog box all attributes of the variables can be specified. This is a good alternative to editing the rda file outside  $\tau$ -ARGUS.  $\tau$ -ARGUS does a moderate checking of the rda-file, but no guarantee can be given for a proper functioning of a manually edited RDA-file. The rda file has been explained in detail in section 5.1.1. Here, editing of a rda file within  $\tau$ -ARGUS is looked at.

If under File|Open Microdata an *rda* file has been specified, this dialog box shows the contents of this file. If no *.rda* file has been specified the information can be specified in this dialog box after pushing the New button. As default "New" is substituted, but the user is expected to fill in a correct name. Apart from defining a new variable, an existing one can be modified or deleted.

An example of the metafile window is shown here.

In the left top field the file type (fixed or free format) can be specified.

The following attributes for each variable can be specified or edited:

- name of the variables
- its first position in the data file
- its field-length
- the number of decimals.

Furthermore, the kind of variable can be specified or edited (more detail on these can be seen in section 4.2.1):

- *explanatory variable*: This can be used as a spanning variable in the row or column of the table
- *response variable*: This can be used as a cell-item
- *weight variable*: This specifies the sampling-weight of the record and is based on the sampling design used.

The following are specialist variable types and have not been previously described. As they are specific to designating safety rules, more detail is given in section 4.3.4.

#### Holding Indicator

The Holding indicator: sometimes groups of records belong together. So it could be better to apply the confidentiality protection to businesses at a number of levels. This variable is the group identifier.  $\tau$ -ARGUS expects the records of a group to be together in the input datafile. An example is shown in section 4.3.4.

#### Request Protection

The Request protection option is used if the Request Rule under ‘Specify tables’ is to be applied. This variable indicates whether or not a records asked for protection. This is further explained in section 4.3.4. Additionally the codes specifying whether a respondent asked for asking protection is to be specified; two different codes are possible, corresponding to two different sets of parameters in the sensitivity rule. This rule is often used in Foreign Trade Statistics

#### Additional Specifications

Other attributes, which may be edited or specified are missing value options, (optional, not required) codelist files, (optional, not required) hierarchies.

Details on these options have been given in section 4.2.1.

In summary, for codelist the ‘automatic’ option simply generates the codes from the data. Specifying a codelist, allows the user to supply an additional file (usually .cdl) containing the labels attached to the codes. These labels are used to enhance the information by  $\tau$ -ARGUS on the screen. In both cases  $\tau$ -ARGUS will use the codes that it finds in the datafile.

Hierarchies can either be derived from the digits in the codes or from a file (usually .hrc)

#### The Rda file

Here is an example of a rda file for microdata. This has already been shown in section 4.2.1 and is shown here for completeness. (Note, the dots at the bottom just means that here a shortened version of the file is presented.)

```
YEAR 1 2 99
  <RECODEABLE>
IndustryCode 4 5 99999
  <RECODEABLE>
  <HIERARCHICAL>
  <HIERLEVELS> 3 1 1 0 0
Size 9 2 99
  <RECODEABLE>
Region 12 2 99
  <RECODEABLE>
  <CODELIST> Region.cdl
  <HIERCODELIST> Region2.hrc
  <HIERLEADSTRING> @
  <HIERARCHICAL>
Wgt 14 4 9999
```

```

<NUMERIC>
<DECIMALS> 1
<WEIGHT>
Var1 19 9 999999999
<NUMERIC>
Var2 28 10 9999999999
<NUMERIC>
<DECIMALS> 2
.....

```

See also section 5.1.1 for a more detailed description

T-ARGUS can also read free format data files. In that case there are slight differences. You select free format in the combo box in the left upper corner. And specify the separator used. The parameter starting position is no longer valid and will not be visible.

### 4.3.2 Specify | Metafile [SPSS System files]

When  $\tau$ -ARGUS works with a SPSS system file the specification of the meta data is twofold. First the variables of interest for building a table have to be specified. Secondly the meta data has to be filled in that could not be automatically retrieved from the system file. SPSS gives only the basic information like variable names, field length

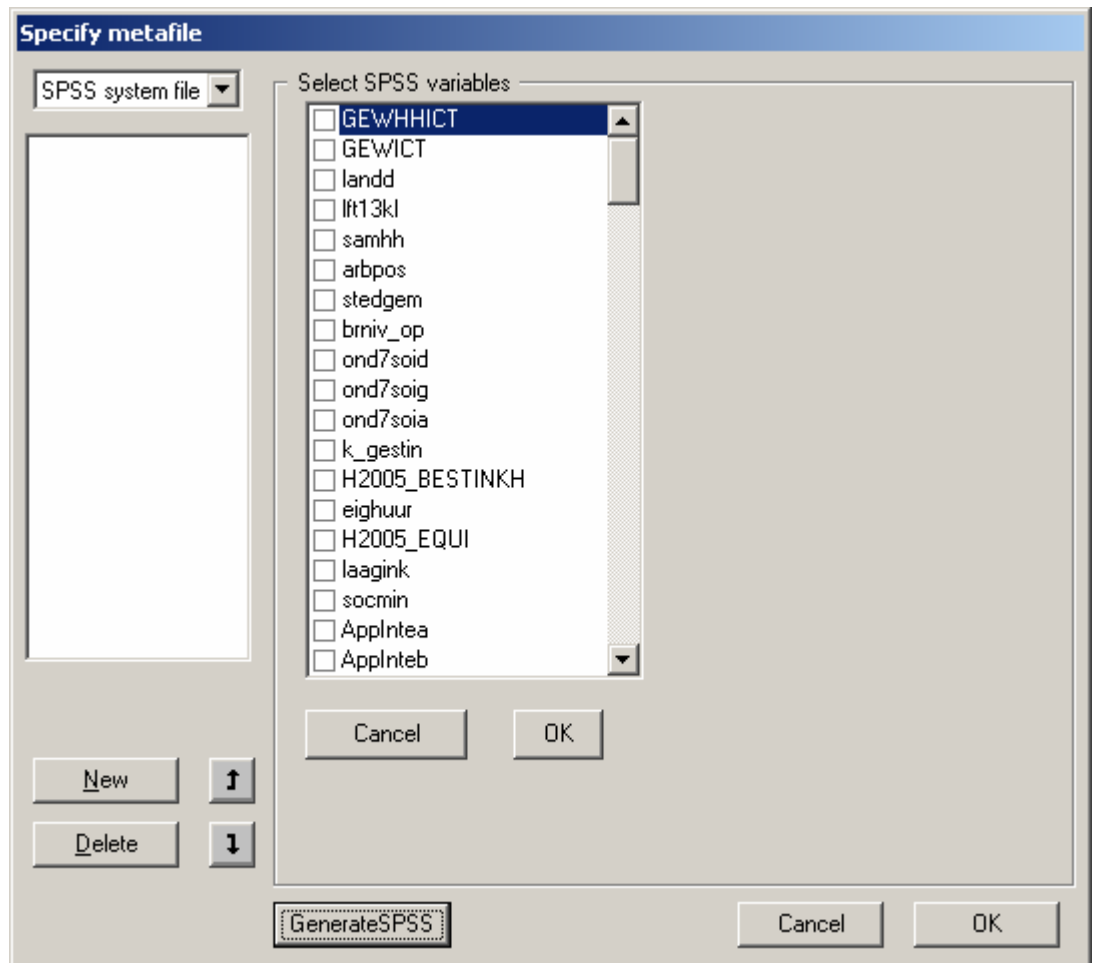
#### Selecting the variables.

By pressing the button 'Generate' a window will be shown, with all the variables available in the system file. You can now select the required variables.

#### Expanding the metadata

The working of  $\tau$ -ARGUS when using a SPSS system file is very similar to the fixed format version, However you will see that certain field cannot be changed as they are implied by SPSS.

Often SPSS cannot decide whether a variable is a spanning variable or a response variable (eg AGE recoded numerically in SPSS). Also the hierarchical information has to be added. Refer to section 4.3.1



### 4.3.3 Specify | Metafile [for tabular data]

When a tabular datafile has been selected, the metadata window will have a different form. Clicking on 'Specify|Metafile' gives the opportunity to either edit the metafile already read in or to enter the metafile information directly at the computer. In section 5.1.4 a detailed description of the metafile for tabular data can be found

Below is displayed the 'Specify metafile' window for tabular input data.

Above the list of variables the separator used to separate the variables in the datafile can be specified.

Here, the variables can be specified or edited as required.

The options are:

- 'Explanatory Variable' – The spanning variables used to produce the table.
- 'Response Variable' – The variable used to calculate the cell total.
- 'Shadow variable' – The variable is used as a shadow variable.
- 'Cost variable' – The variable is used as the cost-variable.
- 'Lower prot .Level' – The lower protection level

- ‘Upper prot. Level’ – The upper protection level
- ‘Frequency’ – This indicates the number of observations making up the cell total. If there is no frequency variable each cell is assumed to consist of a single observation.
- ‘topN variable’ – This shows if this variable is defined as one of the top N contributors to the cell. The pre-defined value for TopN is 1. The first variable declared as ‘topN’ will contain the largest values in each cell, the second variable so declared will contain the second largest values etc.
- ‘Status Indicator’ – allows a variable in the left-hand pane to be declared as a Status Indicator. Typically cells can be declared as Safe, Unsafe or Protected.

**Specify metafile**

Free format

Separator

Expvar1  
Expvar2  
Respvar  
Freq  
Top1  
Top2  
Top3  
Status

**Attributes**

name: Expvar1

length: 0

decimals: 0

☒ explanatory variable  
☐ response variable  
☐ shadow  
☐ cost var  
☐ lower prot. level  
☐ upper prot. level

☐ frequency  
☐ topN variable  
☐ status indicator

**Codelist**

☒ automatic  
☐ codelist filename

Code for Total T

Missings:  
1: X 2:

☐ hierarchical

☐ Levels from microdata

☐ Levels from file

Leading string .

New ↑

Delete ↓

Cancel OK

For explanatory variables the code for the total has to be specified. We recommend strongly that the user also provides the values for the totals himself, but if needed he can ask  $\tau$ -ARGUS to compute these totals. It should be noted that when the option to compute the totals by  $\tau$ -ARGUS you will loose vital information as the cell status. See also section 4.3.5 In any case,  $\tau$ -ARGUS needs these totals as they play an important role in the structure of a table and also are important for the suppression models.

### 4.3.4 Specify | Specify Tables [for microdata]

In this dialog box the user can specify the tables which require protection. In one run of  $\tau$ -ARGUS more than one table can be specified, but the tables will be protected separately unless they are linked (have at least one variable in common). In that case they can be protected simultaneously if required. In section 4.4.3 the idea of linked tables will be discussed.

Also, the user has to specify the parameters for the dominance rule or p% rule and the minimum number of contributors in a cell, etc. At present  $\tau$ -ARGUS allows up to 6-dimensional tables, but due to the capacities of the LP-solver used (either Xpress or Cplex depending on the user's license) and the complexity of the optimisations involved, tables of this complexity can only be protected by the hypercube method (see section 2.8 in the Theory chapter).

Below is a typical window obtained when specifying tables with the dominance rule applied.

**Specify Tables**

explanatory variables: Year, IndustryCode, Size, Region

cell items: Var1, Var2, Var3, Var4, Var5, Var6, Var7, Var8, <freq>

response variable:

shadow variable:

cost variable: ☐ unity ☐ frequency ☒ variable

lambda: 1

☐ Dominance rule

☒ P%-rule

☐ Request rule

	P	N
Ind-1	15	1
Ind-2	0	0
Hold-1	0	0
Hold-2	0	0

☐ Minimum frequency

	freq	range
Ind.	0	30 %
Hold.	0	30 %

☐ Zero unsafe

range: 0

Manual safety range: 0 %

Expl. vars	rule	Resp. var	Shadow & Cost var

Cancel Compute tables

In section 4.3.1 details of variable definitions in the metafile were explained. Now consider how the variables defined in the metafile are used to create a table along with an associated safety rule.

#### The explanatory (or spanning) variables

On the left is the listbox with the explanatory variables.

When the user clicks on ‘>’ or ‘<’ the selected variable is transported to the next box. From the left box with explanatory variables the user can select the variables that will be used as the spanning variables in the row or the column of the table.

### **Cell items**

Here, is a list of variables that can be used as response, shadow or cost variables in the disclosure control. By pressing the ‘>’ or ‘<’ they can be transferred to or from the windows on the right.

### **The response variable**

From the list of cell items the user can select a variable as a response variable. This is the variable for which the table to be protected is calculated. If a number of tables are required for the same explanatory variables then more than one response variable can be entered here. Each response variable is suppressed independently.

### **The shadow variable**

The shadow variable is the variable that is used to apply the safety rule. By default this is the response variable, but it is possible to select another variable. The safety rules are built on the principle of the characteristics of the largest contributors to a cell. If a variable other than the response variable is a better indicator this variable can be used here; e.g. the turnover (a proxy for the size of the enterprise) can be a suitable variable to apply the safety rule, although the table is constructed using another (response) variable.

### **The cost variable**

This variable describes the costs of suppressing each individual cell; these costs are used by the internal workings of the secondary suppression routines. Note that the choice of the cost variable does not have any effect when the hypercube method is used for secondary cell suppression. See 2.7.1 for information about how cell costs are determined during execution of the hypercube method. With exception of the hypercube method, these costs are minimised when the secondary suppressed cells are determined. By default, this is the response variable but two other choices are possible as well as the use of a different response variable.

Use the frequency of the cells as a cost-function: this will minimise the number of records contributing to the cells to be suppressed.

The number of cells to be suppressed is minimised, irrespective of the size of their contributions (Unity option).

A Box-Cox like-transformation can be applied to the individual values of the cost variable before minimisation of the cost function. The simplified Box Cox function used here is  $x^\lambda$  where  $x$  is the cost variable and  $\lambda$  is the transformation parameter. For example if  $\lambda = 0.5$  a square root transformation is used and if  $\lambda = 0$  a log transformation will be applied. Applying this to the unity-choice is rather meaningless.

### **Weight**

If the data file has a sample weight, specified in the metadata, the table can be computed taking these weights into account.

If the ‘Apply Weights’ box is ticked, the weights are applied to the cell entries as for the simple application of normal sampling weights in a survey. In addition these weights are used in applying the safety rules. When we have a sample the normal idea behind the sensitivity rules that the largest contributions can make a good estimate of each other is no longer valid. The solution is that we artificially create a

complete cell by assuming that each contribution is in fact as many contributions as its sample weight. And we apply the sensitivity rules on this cell. An example might help here.

For example if there is a cell with two contributions:

*100, weight 4*

*10, weight 7*

The cell value = (4 x 100) + (7 x 10) = 470. Without considering the weights there are only two contributors to the cell 100 and 10. However by taking account of the sampling weights the cell values are approximately 100, 100, 100, 100, 10, 10, 10, 10, 10, 10 and 10. The largest two contributors are now 100 and 100. These are regarded as the largest two values for application of the safety rules. If the weights are not integers, a simple extension is applied.

### **The safety rule**

The concept of safety rules is explained in section 2.2 On the left side of the window the type of rule that can be selected along with the value of the parameters is shown. The possible rules are:

- Dominance Rule
- P% Rule
- Request Rule (this rule is described in detail later in this section)

Additionally, the minimum number of contributors may be chosen (in the 'minimum frequency' box).

Two dominance rules and two P% rules can be applied to each table. When 2 rules are specified, for a cell to be declared non-disclosive, it must satisfy both rules.

### ***Dominance Rule***

This is sometimes referred to as the (n,k) rule where n is the number of contributors to a cell contributing more than k% of the total value of the cell (if the cell is to be defined as unsafe). A popular choice would be to set n equal to 3 and k equal to 75%. An example of the window when specifying a single dominance rule is shown at the start of this section.

### ***P% rule***

The P% rule says that if  $x_l$  can be determined to an accuracy of better than P% of the true value then it is disclosive where  $x_l$  is the largest contributor to a cell.

The rule can be written as:

$$\sum_{i=3}^c x_i \geq \frac{p}{100} x_1$$
 for the cell to be non-disclosive where  $c$  is the total number of contributors to the cell and the intruder is a respondent in the cell.

It is important to know that when entering this rule in  $\tau$ -ARGUS the value of  $N$  refers to the number of intruders in coalition (who wish to group together to estimate the largest contributor).

A typical example would be that the sum of all reporting units excluding the largest two must be at least 10% of the value of the largest. Therefore, in  $\tau$ -ARGUS set  $p=10$  and  $n=1$  as there is just one intruder in the coalition, respondent  $x_2$ .

For the dominance rule and the P%-rule the safety ranges required (as a result of applying the rule) can be derived automatically. The theory gives formulas for the

upper limit only, but for the lower limit there is a symmetric range. See e.g. Loeve (2001). (This is referenced in Section 2.2 (Theory))

### **Request Rule**

This is a special option applicable in certain countries relating to e.g. foreign trade statistics. Here, cells are protected only when the largest contributor represents over (for example) 70% of the total and that contributor asked for protection. Therefore, a variable indicating the request is required.

This option requires an additional variable in the data, with e.g. 0 representing no request for that particular business, and 1 representing a request where the particular cell value is  $> x\%$  of the cell total. In fact there is an option for two different thresholds. The min freq is interpreted such that if a cell has at least one request and the cell-freq is below the freq-threshold, that cell is considered to be unsafe as well. Even if the request is not the largest one. The idea is that in that case a large non requesting contributor could reveal the smaller requesting contributor.

Note that the 3 rules (dom. rule,  $p\%$  rule and request rule) do not make any sense if there are positive and negative contributions to a cell.

### **Minimum Frequency**

If this box is checked, a rule controlling the minimum number of contributors to a cell will be specified. If the number of contributors is less than this value, the cell is considered unsafe.

### **Freq**

Here, the minimum number of contributors can be stated. This is sometimes known as the threshold rule. It is also possible to specify no safety rule apart from a minimum frequency value.

### **Frequency-range**

As described above, for the dominance rule and the  $P\%$ -rule safety ranges can be derived automatically. However, the theory does not provide any safety range for the minimum frequency rule. Therefore, the user must provide a safety-range percentage required to allow secondary suppressions to be carried out. For example, if this value was set to equal 30%, it would mean an attacker would not be able to calculate an interval for this cell to within 30% of the actual value when looking at the safe output. Following this, the secondary suppressions may be carried out.

### **Manual Safety Range**

When a cell is set manually unsafe (an option to discussed later),  $\tau$ -ARGUS cannot calculate safety-ranges itself. Therefore, the user must supply a safety-percentage for this option for the same reasons as in the above section, to allow secondary suppressions to be applied.

### **Zero Unsafe**

If all contributions to a cell are zero, the cell value will be zero too. Applying sensitivity rules here has some problems. Is the sum of the largest 3 zeros larger than zero? Nevertheless all contributions to this cell can be easily disclosed. If cells with total contributions of zero are to be regarded as unsafe, this box has to be checked. A manual safety range will also be required, not as a percentage but as a value at the level of the cell-item.

### Missing = safe

If one of the spanning variables of a cell has a code missing, this cell is often no longer sensitive. The idea behind this is that the respondent in this cell is not identifiable. When this option is checked, all cells for which at least one spanning variable has a missing value is considered safe, whatever all the sensitivity rules say. If this option is not checked the normal procedures as for all other cells are applied.

### Holding Indicator

*This section on the Holding Indicator is best read after section 4.4.2*

In some countries, confidentiality protection is applied to businesses at different levels. For example, as in the U.K. a number of 'reporting units' (the lower level of unit) within a cell might belong to an 'enterprise group' (higher level). The level at which the confidentiality rule is applied clearly matters. The holding indicator allows such groupings to be defined and used in one or more of the safety rules.

This is now illustrated with an example looking at both the P% rule and the threshold rule at the same time.

The screenshot shows the 'Specify Tables' dialog box. The 'explanatory variables' list on the left contains 'Town', 'Size', 'NACE', and 'Employees'. The 'cell items' list on the right contains 'TurnOver' and '<freq>'. The 'P%-rule' tab is selected, showing a table with 'Ind-1', 'Ind-2', 'Hold-1', and 'Hold-2' rows, each with 'P' and 'N' columns. The 'Minimum frequency' section shows 'Ind.' with a value of 1 and a range of 10%. The 'Hold.' section shows a value of 0 and a range of 10%. The 'Zero unsafe' checkbox is unchecked, and the 'Manual safety range' is set to 0. The 'cost variable' section has 'variable' selected. The 'Missing = safe' checkbox is checked. The 'Compute tables' button is at the bottom right.

### Consider the following dataset

Cell Ref	Cell Ref	Cell value (reporting unit)	Enterprise group
800	20	599	1
800	20	344	1
800	20	244	1

800	30	355	1
800	20	644	2
800	30	433	2
800	30	323	3
800	30	343	3
900	20	23	4
900	20	43	5
900	20	34	5
900	20	53	5
900	30	700	6
900	30	200	6
900	30	60	7
900	30	40	8
900	30	10	9

Assume the following safety rules

- Threshold rule: At least 3 enterprise groups (higher level units) in a cell
- P% rule: The sum of all the reporting units (lower level units) excluding the largest 2 must be at least 10% of the value of the largest.

There are 4 cells in the table along with the margins. The cell we are interested in here is *Cellref 900,30*: 5 reporting units, 4 enterprise groups

At the reporting unit the values are 700,200,60,40,10

At the enterprise group the values are 900,60,40,10

This rule has been designed so that when the P% rule is applied to this cell:

With reporting units the cell is safe.  $10+60+40 = 110$ . This is greater than 10% of the largest value (70) so the cell is safe.

With enterprise groups the cell is unsafe.  $40+10 = 50$ . This is less than 10% of the largest value (90) so the cell is unsafe.

Apply the threshold rule to the enterprise groups (Hold. =3) and P% rule to the reporting units.

Once again a safety range percentage is required.

The output from the application of this rule is shown below. Two cells fail the threshold rule with the holding rule applied.

The threshold rule has been applied correctly using the holding indicator as the correct cells are safe (that would be unsafe if the holding indicator was not being used).

Table: Regio 2 x SBI 2 | Omzet

	tot	- a	aa	ab	99
tot	2,479	2,479	1,279	1,200	-
- .1	729	729	529	200	-
..11	209	209	209	-	-
..12	220	220	20	200	-
..13	300	300	300	-	-
- .2	1,750	1,750	750	1,000	-
..21	650	650	250	400	-
..22	500	500	500	-	-
..23	600	600	-	600	-
.99	-	-	-	-	-

Cell Information

Value: 0

Status: Empty

Cost: 0

Shadow: 0

# contributions: 0

Top n of shadow:

☐ Holding level

Request: 0

Change status

Set to Safe

Set to Unsafe

Set to Protected

Set Cost

A priori info

All Non-StructEmpty

Recode

Suppress

☐ HyperCube

☒ Modular

☐ Network

☐ Optimal

☐ Rounding

Suppress

Undo Suppress

Audit

☒ 3 dig. separator

Select Table

Change View

Write table

Table Summary

Close

### After all the options have been selected compute the table

When all the necessary information has been given, click 'v' to transport all the specified parameters to the 'listwindow' on the bottom. As many tables as required can be specified but as the size of the memory of a computer is restricted it is not advisable to select too many tables. To modify an already made table press the '^' button.

Click on 'Compute Tables' to compute the tables. In case of an SPSS system file SPSS will first be called to export the needed microdata automatically to a scratch fixed format ASCII file in the TEMP directory

When the table(s) have been computed, the main-window of  $\tau$ -ARGUS will be displayed again with an overview of all the unsafe cells per variable for every table calculated. An example is shown here for the *Size* by *Region* table. Firstly, the *Size* dimension is looked at and then secondly the *Region* dimension.

The window (underneath the main menu for  $\tau$ -ARGUS) shows the number of unsafe combinations per variable. For example, there are no unsafe cells in dimension one for either variable. (*i.e.* the one-way marginal total for different values of 'Size' and 'Region' are all not disclosive).

### Variable Size

There are however 12 unsafe cells in the 2 way table *Size* by *Region* as can be seen

by the right hand window which gives the equivalent information for each level of the variable indicated on the left. There are 5 unsafe cells where Size = 2, 6 unsafe cells where Size = 4 and a single unsafe cell where Size = 9.

**TauARGUS**

File Specify Modify Output Help

#unsafe combinations in every

Variable	dim 1	dim 2
Size	0	9
Region	0	9

variable: Size

Code	Label	Freq	dim 1	dim 2
	Total	42723	0	0
·2		9	0	4
·4		5	0	4
·5		20002	0	0
·6		8831	0	0
·7		5498	0	0
·8		4594	0	0
·9		3779	0	1
99		5	0	0

Status 5/2/2011 12:34 PM

### Variable region

The 12 unsafe cells when looked at by Region show that two of these are 'North' subtotal cells. Within the 'North' region, 2 cells in Groningen are disclosive.

Two 'East' subtotal cells are also unsafe. Within the 'East' region 2 cells in 'Overijssel' are unsafe along with 2 cells in 'Gelderland'

Finally there is 1 unsafe cell for the 'South' subtotal and within this region there is 1 unsafe cell for 'Noord-Brabant'

The screenshot shows the TauARGUS software interface. The 'Specify' menu is active, and a table displays the dimensions of variables. The table has columns for 'Variable', 'dim 1', and 'dim 2'. The variables listed are 'Size' and 'Region'. The 'Region' variable is expanded, showing a list of regions with their respective frequencies and dimensions. The 'Total' row is highlighted in blue.

Variable	dim 1	dim 2
Size	0	9
Region	0	9

Code	Label	Freq	dim 1	dim 2
	Total	42723	0	0
Nr	North	11395	0	2
·1	Groningen	6112	0	2
·2	Friesland	3798	0	0
·3	Drenthe	1485	0	0
Os	East	10227	0	1
·4	Overijssel	538	0	2
·5	Flevoland	1597	0	0
·6	Gelderland	5446	0	2
·7	Utrecht	2646	0	0
Ws	West	10054	0	0
·8	Noord-H...	1182	0	0
·9	Zuid-Holl...	8018	0	0
10	Zeeland	854	0	0
Zd	South	11047	0	0
11	Noord-Br...	7511	0	0
12	Limburg	3536	0	0
99		0	0	0

Status: 5/2/2011 11:39 AM

It should be noted that more than one response variable can be specified. This will produce tables for each of the Response variables using the Spanning variables specified.

#### 4.3.5 Specify | Specify tables [for tabular data]

When the 'Specify|Metafile' option is followed the 'Specify|Table metadata' option is also available and the window is displayed here. This will allow the application of safety rules such as the Dominance Rule and the P% rule. Section 4.3.4 (specifying tables from microdata) will explain these safety rules and other options in detail.

In the safety rule frame, the type of rule can be selected along with the value of the parameters. These are the dominance rule and P% rule. Additionally, the minimum number of contributors can be chosen (threshold rule), via ticking and filling-in the minimum frequency box. If both the status and some information to apply the sensitivity rules have been supplied, both options ‘use given status’ and ‘use safety rules’ are enabled and the user can choose which one to use.

Depending on the amount of detail in the table file some options will be disabled. If no top1 and top2 information is provided, the p%-rule cannot be used.

There is an option to calculate the possibly missing marginals and totals. This option should be used only as an emergency. It is always better to provide  $\tau$ -ARGUS with a full, complete table. When  $\tau$ -ARGUS has to compute these marginals all safety information will be ignored.

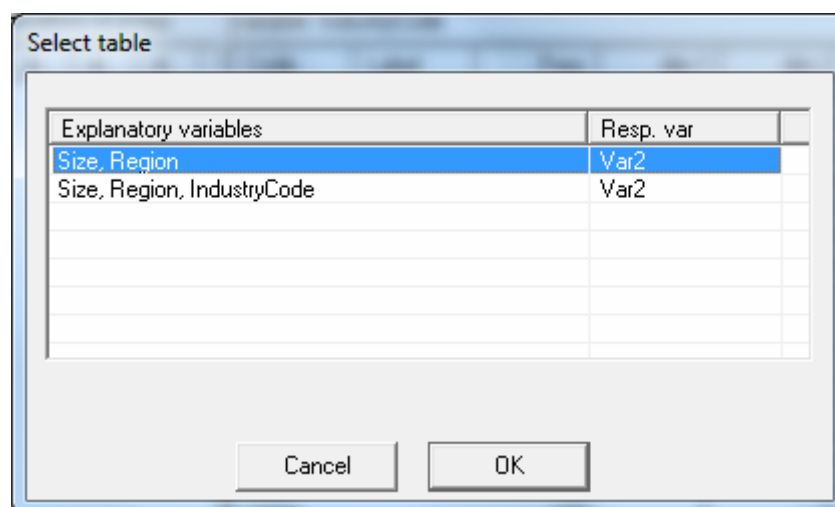
When all the options have been completed, pressing the ‘OK’ button will invoke  $\tau$ -ARGUS to actually compute the table requested. Now the process of disclosure control can begin.

## 4.4 The Modify menu

### 4.4.1 Modify | Select Table

This dialog box enables the user to select the table they want to see. If the user has specified only one table, this table will be selected automatically and this option cannot be accessed. In the example window shown here the first table is a 2

dimensional table (Size x Region) followed by a 3 dimensional table (Size x Region x IndustryCode). Select the table to be processed and press the OK-button.



## 4.4.2 Modify | View Table

This section is divided in four parts: a general description of the ‘View Table’ screen, global recoding, the secondary cell suppression and some other options

### 4.4.2.1 The View table screen

This window shows the table selected with Modify|View Table. On the left side the table itself is shown in a spreadsheet view. Safe cells are black, unsafe cells (those failing the primary suppression rule) are red. You can however adapt these colours in the options menu, see section 4.6.3. In this example there are 12 unsafe cells and by viewing the table the user can now see the actual cells that are unsafe.

Any secondary suppressed cells are shown in blue (there are none at this stage, in this example) and empty cells have a hyphen (-). The two check-boxes on the left-bottom give some control over the layout.

- If the 3-digit separator box is checked, the window will show the cell-values will be shown, using the 3 digits separator to give a more readable format.
- The Output view shows the table, with all the suppressed cells replaced by an ‘X’; this is how the safe table will be published, but without the colours distinguishing between primary and secondary suppressions.

Table: Size x Region | Var2

	tot	2	4	5	6	7
tot	16,847,647	20	25	2,711,808	2,320,534	2,505,043
- .Nr	4,373,664	5	5	719,049	659,680	688,962
.. 1	1,986,129	5	5	398,062	348,039	354,711
.. 2	1,809,246	0	-	223,990	221,332	241,913
.. 3	578,289	-	-	96,997	90,309	92,338
- .Os	3,703,896	15	5	642,238	515,003	534,147
.. 4	124,336	5	-	36,311	32,132	25,770
.. 5	526,279	-	-	93,589	94,957	110,930
.. 6	2,234,995	10	5	345,803	251,358	251,188
.. 7	818,286	-	-	166,535	136,556	146,259
- .Ws	4,576,116	-	-	648,972	543,570	663,897
.. 8	485,326	-	-	63,767	75,442	87,305
.. 9	3,664,560	-	-	537,911	430,851	515,020
.. 10	426,230	-	-	47,294	37,277	61,572
- .Zd	4,193,971	-	15	701,549	602,281	618,037
.. 11	2,752,743	-	15	488,613	392,395	363,490
.. 12	1,441,228	-	-	212,936	209,886	254,547
.99	-	-	-	-	-	-

Cell Information

Value: 10

Status: Unsafe

Cost: 10

Shadow: 10

# contributions: 2

Top n of shadow: 5

Holding level: 5

Request: 0

Protection interval (low/up value): 9 11

Change status

Set to Safe

Set to Unsafe

Set to Protected

Set Cost

A priori info

All Non-StructEmpty

Recode

Suppress

HyperCube

Modular

Network

Optimal

Rounding

Suppress

Undo Suppress

Audit

3 dig. separator

Output View

Select Table

Change View

Write table

Table Summary

Close

For some windows, the complete table cannot be seen on the screen. In these cases there will be scrollbars at the bottom and the right of the table above, which can be used to display the unseen columns.

For large tables one does not want to see the whole table on one screen, which is virtually impossible. Therefore  $\tau$ -ARGUS will show only the first two levels of the hierarchal structures. If you want to see more you can open and close certain parts of the table by clicking on the codes with a '+' or '-' sign. This works similar to the way you open and close certain parts in the Windows explorer. Via 'Change View' at the bottom of the screen you can also select the level of each hierarchy you want to see.

### Example of a 3D table

Variables in 3D tables are displayed at the top left of the typical 2 dimensional table. The arrow at the right of the box allows selection of the required level of this variable. The table shown will be the 2-dimensional table for the specified level(s) of the variables displayed at the top of the table.

Table: Size 2 x Region 2 x IndustryCode | Var2

IndustryCode  
tot

	tot	2	4	5	6	7
tot	16,847,647	20	25	2,711,808	2,320,534	2,505,043
- .Nr	4,373,664	5	5	719,049	659,680	688,962
.. 1	1,986,129	5	5	398,062	348,039	354,711
.. 2	1,809,246	0	-	223,990	221,332	241,913
.. 3	578,289	-	-	96,997	90,309	92,338
- .Os	3,703,896	15	5	642,238	515,003	534,147
.. 4	124,336	5	-	36,311	32,132	25,770
.. 5	526,279	-	-	93,589	94,957	110,930
.. 6	2,234,995	10	5	345,803	251,358	251,188
.. 7	818,286	-	-	166,535	136,556	146,259
- .Ws	4,576,116	-	-	648,972	543,570	663,897
.. 8	485,326	-	-	63,767	75,442	87,305
.. 9	3,664,560	-	-	537,911	430,851	515,020
..10	426,230	-	-	47,294	37,277	61,572
- .Zd	4,193,971	-	15	701,549	602,281	618,037
..11	2,752,743	-	15	488,613	392,395	363,490
..12	1,441,228	-	-	212,936	209,886	254,547
.99	-	-	-	-	-	-

Cell Information

Value: 10

Status: Unsafe

Cost: 10

Shadow: 10

# contributions: 2

Top n of shadow: 5

☐ Holding level: 5

Request: 0

Protection interval (low/up value): 9 11

Change status

Set to Safe

Set to Unsafe

Set to Protected

Set Cost

A priori info

All Non-StructEmpty

Recode

Suppress

☒ HyperCube

☐ Modular

☐ Network

☐ Optimal

☐ Rounding

Suppress

Undo Suppress

Audit

3 dig. separator

Output View

Select Table

Change View

Write table

Table Summary

Close

### Additional information in the View Table window

Clicking on a cell in the main body of the table makes information about this cell visible in the Cell Information pane.

Here, the following information can be seen:

1. the cell-value
2. the cell status
3. the total cost variable value for the cell
4. the total of the shadow variables for the cell
5. the number of contributors to a cell
6. the values of the shadow variable for the largest contributors.

In addition for a primary unsafe cell also the required lower and upper protection levels are shown. If you put your mouse over this value, also the lower and upper protection as a distance to the cell value is shown together with the same value as a percentage.

Information about the Holding level and the Request protection variable are also displayed here.

The status of the cell can be:

- Safe: Does not violate the safety rule
- Safe (from manual): manually made safe during this session
- Unsafe: According to the safety rule
- Unsafe (request): Unsafe according to the Request rule.
- Unsafe (frequency): Unsafe according to the minimum frequency rule.
- Unsafe (from manual): manually made unsafe during this session (see ‘Change Status’ below).
- Protected: Cannot be selected as a candidate for secondary cell suppression (see ‘Change Status’ below).
- Secondary: Cell selected for secondary suppression.
- Secondary (from manual): Unsafe due to secondary suppression after primary suppressions carried out manually (see ‘Change Status’ and ‘Secondary suppressions’ below).
- Zero: Value is zero and cannot be suppressed.
- Empty: No records contributed to this cell and the cell cannot be suppressed.

### **Change Status**

The second pane (‘Change Status’) on the right will allow the user to change the cell-status.

- Set to Safe: A cell, which has failed the safety rules is here declared safe by the user.
- Set to Unsafe: A cell, which has passed the safety rules is here declared to be unsafe by the user.
- Set to Protected: A safe cell is set so that it cannot be selected for secondary suppression.
- Set Cost: Change the cost function for a cell.

### **A priori info**

This option allows you to feed  $\tau$ -ARGUS a list of cells where the status of the standard rules can be overruled. E.g. a cell must be kept confidential or not for other reasons that just because of the sensitivity rules. By modifying the cost-function you can influence the selection of the secondaries. E.g. the cells suppressed last year can get a preference for the suppression this year by giving this cell a small value for the cost-function.

The option ‘trivial levels’ is important. Often in a table with hierarchies, some levels in a hierarchy break down in only one lower level. This implies that there are different cells in a table which are implicitly the same. Changing the status of one of them might lead to inconsistencies and serious problems. E.g. one of the two is unsafe and the other is protected, the solution is impossible. If you select the option ‘Expand for trivial levels’,  $\tau$ -ARGUS will always modify all cells that are the same if you modify one of them.

The format of the file is free format. The separator can be chosen.

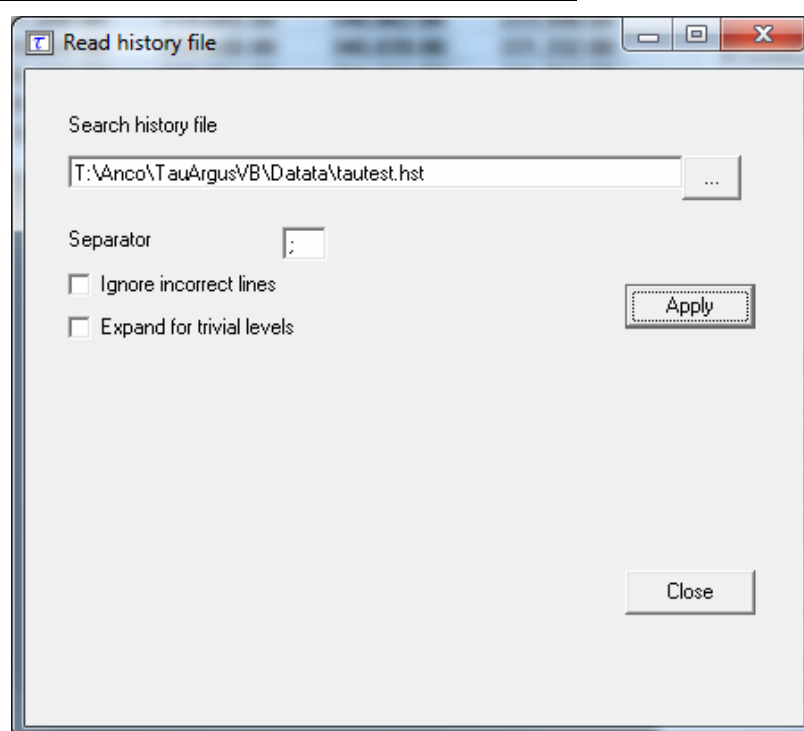
The format is:

Code of first spanning variable, Code of second spanning variable,..., Status of cell (u = unsafe, p = protected (not to be suppressed), s = safe).

Also the cost-function can be changed here for a cell. This will make the cell more likely to become secondary cell suppression, when the value is low, or less likely when the value is high.

Normally the sensitivity rules will also give the required protection levels for unsafe cells. But sometimes, e.g. in the case of ‘manual unsafe cells’ the user might want to specify the required protection level different for a standard percentage. After the keyword ‘pl’, the lower and upper protection levels can be given for a specific cell. Note that the protection levels will always have to be positive, as they are considered as distances from the cell-value. See also section 5.5

```
Nr, 4, u
Zd, 6, p
  5, 5, c, 1
Zd, 5, pl, 100, 200
```



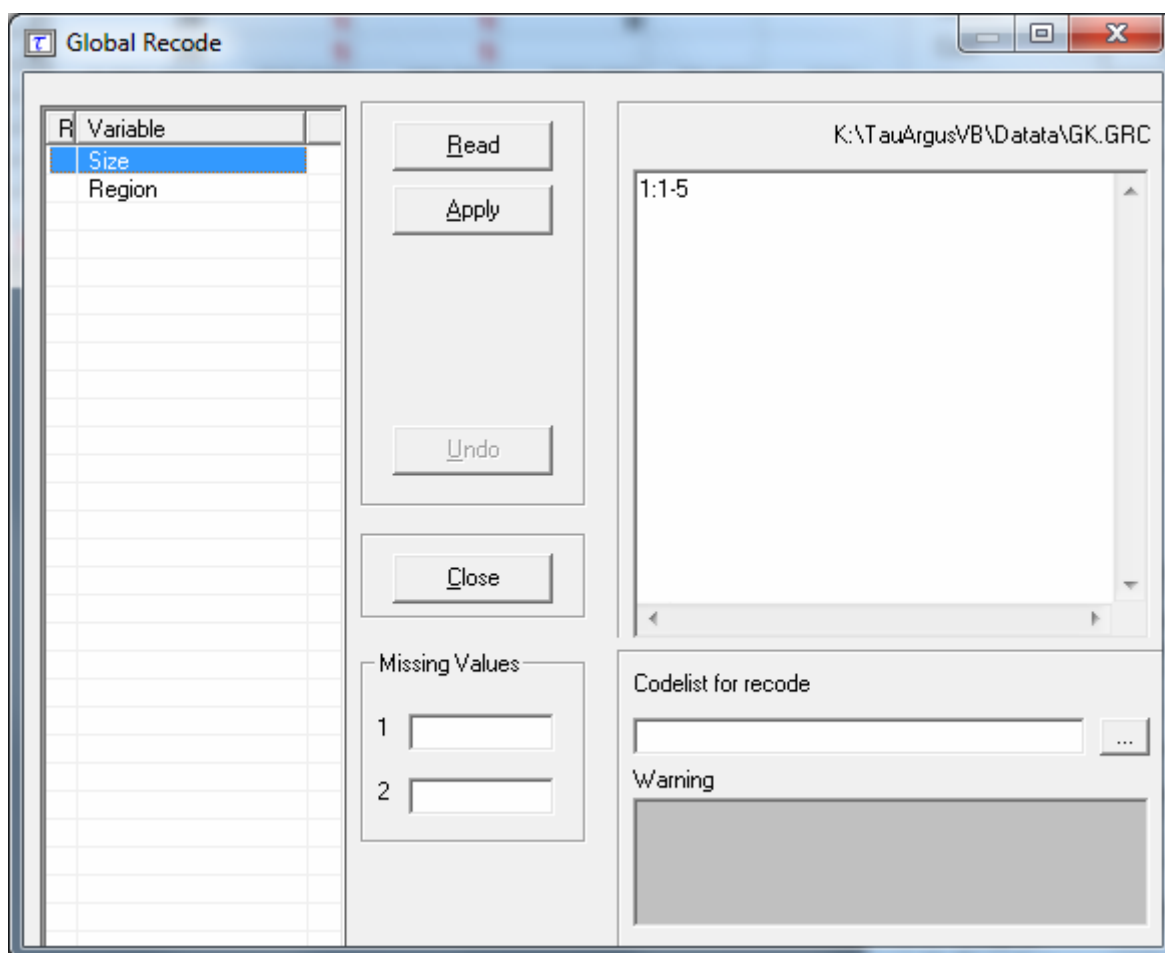
#### 4.4.2.2 Global recoding

The recode button will open the recoding options. Recoding is a very powerful method of protecting a table. Collapsed cells tend to have more contributors and therefore tend to be much safer.

##### **Recoding a non-hierarchical variable**

There is a clear difference in recoding a hierarchical variable compared to a non-hierarchical variable.

In the non-hierarchical case the user can specify a global recoding manually. Either enter the recoding described below manually or read it from a file. The default extension for this file is .GRC.



There are some standards about how to specify a recode scheme. All codelists are treated as alphanumeric codes. This means that codelists are not restricted to numerical codes only. However, this also implies that the codes '01' and '1' are considered different codes and also 'aaa' and 'AAA' are different. In a recoding scheme the user can specify individual codes separated by a comma (,) or ranges of codes separated by a hyphen (-). The range is determined by treating the codes as strings and using the standard string comparison. E.g. '0111' < '11' as the '0' precedes the '1' and 'ZZ' < 'a' as the uppercase 'Z' precedes the lowercase 'a'. Special attention should be paid when a range is given without a left or right value. This means every code less or greater than the given code. In the first example, the new category 1 will contain all the codes less than or equal to 49 and code 4 will contain everything larger than or equal to 150.

Example:

for a variable with the categories 1,...,182 a possible recode is then:

1:	-	49
2:	50	- 99
3:	100	- 149
4:	150	-

for a variable with the categories 01 till 10 a possible recode is:

```
1: 01 , 02
2: 03 , 04
3: 05 - 07
4: 08 , 09 , 10
```

An important point is not to forget the colon (:) if it is forgotten, the recode will not work.

Recoding 3: 05,06,07 can be shortened to 3: 05-07.

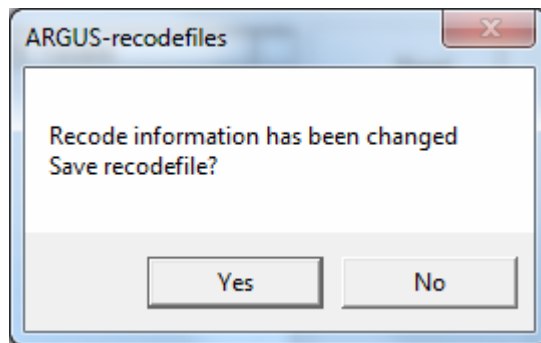
Additionally, changing the coding for the missing values can be performed by entering these codes in the relevant textboxes.

Also a new codelist with the labels for the new coding scheme can be specified. This is entered by means of a codelist file. An example is shown here. (note, there are no colons in this file)

```
1,Groningen
2,Friesland
3,Drenthe
4,Overijssel
5,Flevoland
6,Gelderland
7,Utrecht
8,Noord-Holland
9,Zuid-Holland
10,Zeeland
11,Noord-Brabant
12,Limburg
Nr,North
Os,East
Ws,West
Zd,South
```

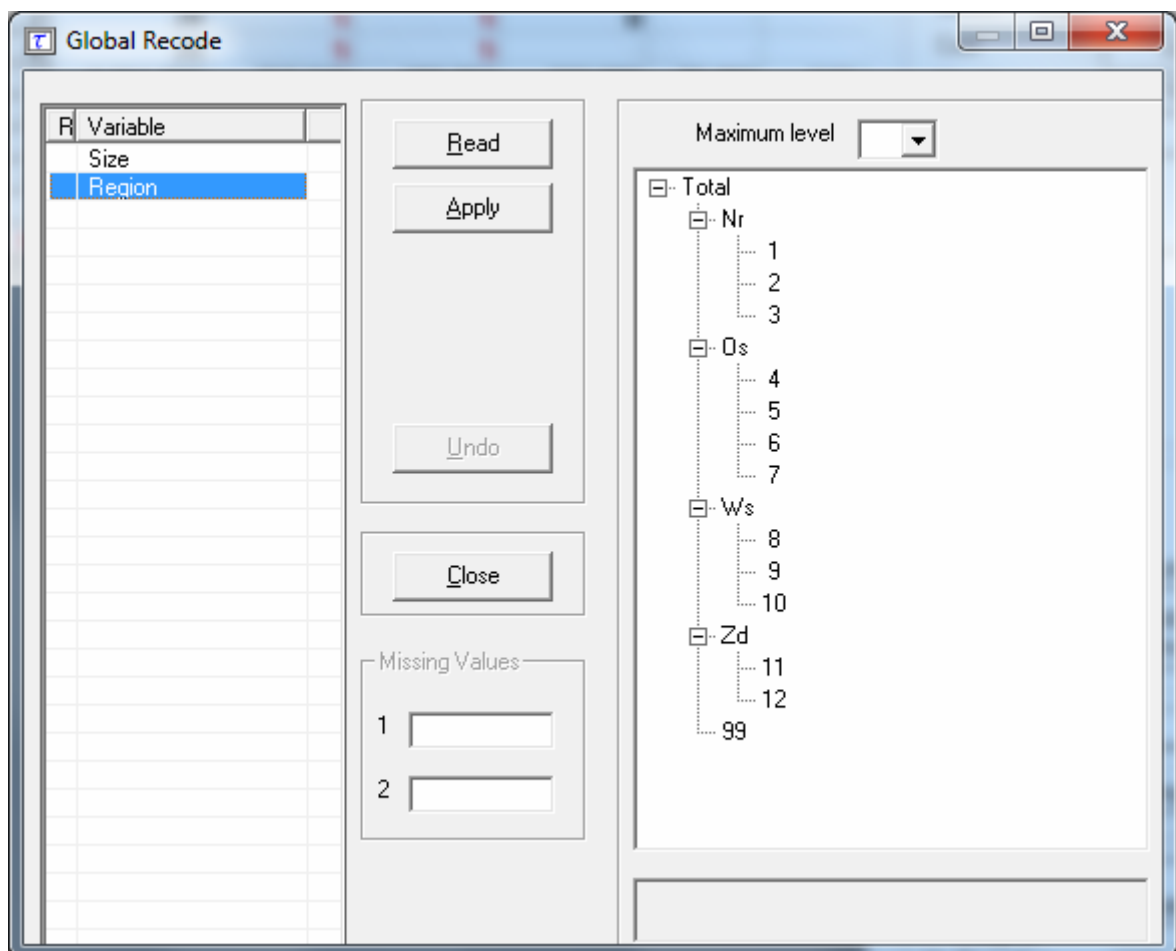
Pressing the 'Apply' button will actually restructure the table. If required, recoding can easily be undone by pressing 'undo recoding'. The window will return to the originally coding structure. If there is any error in the recoding such as certain codes not being found when pressing the 'Apply' button, an error message will be shown at the bottom of the screen. Alternatively, a warning could be issued; *e.g.* if the user did not recode all original codes,  $\tau$ -ARGUS will inform the user. This may have been the intention of the user, therefore the program allows it. In the above example a  $\tau$ -ARGUS message informs the user that 4 codes have not been changed.

At the end of the operation  $\tau$ -ARGUS will ask you whether or not a modified recoding scheme must be saved or not.



Once the 'Close' button has been pressed,  $\tau$ -ARGUS will present the table with the recoding applied.

### Recoding a hierarchical variable



In the hierarchical case the code scheme is typically a tree. To global recode a hierarchical variable requires a user to manipulate a tree structure. The standard Windows tree view is used to present a hierarchical code. Certain parts of a tree can be folded and unfolded with the standard Windows actions (clicking on '+' and '-').

The maximum level box at the top of the screen offers the opportunity to fold and unfold the tree to a certain level.

Additionally, the user can change the coding for the missing values by entering these codes in the relevant textboxes.

Pressing the 'Apply' button will actually restructure the table. If required, a recoding may always be undone.

### 4.4.2.3 Secondary suppression

The actions in the suppress pane in the table window (after selecting 'modify|table') are now looked at.

With suppress the table can be protected by causing additional cells to be suppressed. This is necessary to make a safe table.

### Suppression Options

There are a number of suppression options, which can be seen on the bottom right hand side of the window.

- Hypercube
- Modular
- Network
- Optimal

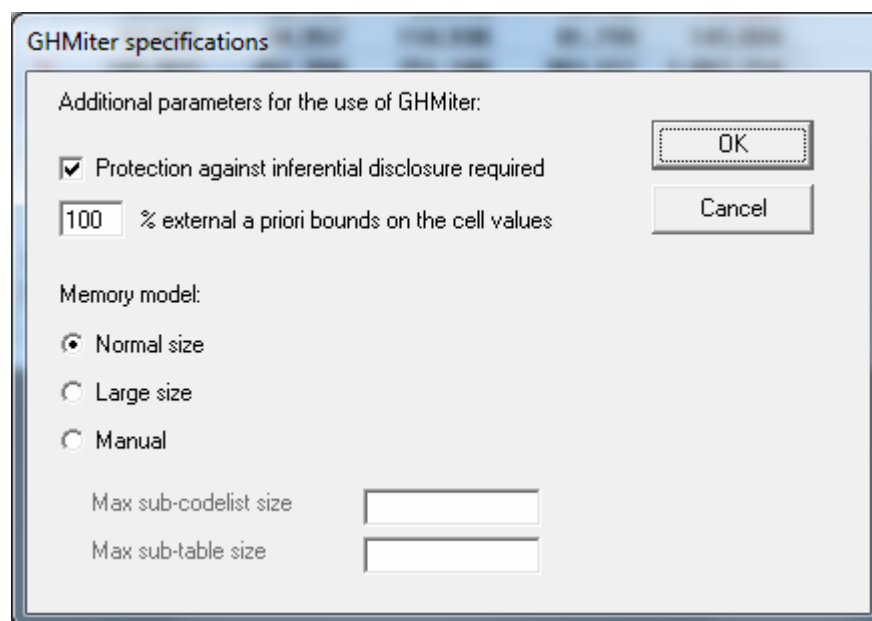
The screenshot shows the 'Table: Size x Region | Var2' window. The main table displays data for various categories (tot, Nr, 1, 2, 3, 0s, 4, 5, 6, 7, 8, 9, 10, 11, 12, .99) across columns 2 through 9. The cell at row 9, column 9 is highlighted with a red border and contains the value 11,968. To the right of the table is the 'Cell Information' panel, which shows details for the selected cell: Value 11,968, Status Unsafe, Cost 11,968, Shadow 11,968, # contributions 29, Top n of shadow 11,401, Holding level 145, Request 0, and Protection interval (low/up value) 10,680 / 13,256. Below this is the 'Change status' panel with buttons for 'Set to Safe', 'Set to Unsafe', 'Set to Protected', 'Set Cost', 'A priori info', 'All Non-StructEmpty', and 'Recode'. The 'Suppress' panel at the bottom right shows radio buttons for 'HyperCube' (selected), 'Modular', 'Network', 'Optimal', and 'Rounding', along with 'Suppress', 'Undo Suppress', and 'Audit' buttons. At the bottom of the window are checkboxes for '3 dig. separator' and 'Output View', and buttons for 'Select Table', 'Change View', 'Write table', 'Table Summary', and 'Close'.

Hypercube

This is also known as the GHMITER method. The approach builds on the fact that a suppressed cell in a simple n-dimensional table without substructure cannot be disclosed exactly if that cell is contained in a pattern of suppressed, nonzero cells, forming the corner points of a hypercube. Selecting the hypercube method will lead to the following window being showed by  $\tau$ -ARGUS.

GHMITER will select secondary suppressions that protect the sensitive cells properly against the risk of inferential disclosure, to some extent, if the user activates the option “*Protection against inferential disclosure required*”. If the option is inactivated, on the other hand, GHMITER will not check secondary suppressions to be sufficiently large.

For more explanation, and detailed information on the hypercube see section 2.8.



The lower part of the pop-up window above enables the user to affect the setting of two parameters, “Max sub-codelist size” and “Max sub-table size” that GHMITER uses for memory allocation.

If the option ‘normal size’ is active, the default values mentioned below will be used. Ticking the option ‘large size’ will lead to a setting of 250 and 25000, respectively.

“Max sub-codelist size” must exceed the largest maximum sub-codelist size of all explanatory variables of the table. The maximum sub-codelist size of a (hierarchical) variable is the largest number of categories on the same (hierarchical) level that contribute to the same category on the (hierarchical) level just above. The default value for “Max sub codelist size” is 200.

“Max sub-table size” must exceed the number of cells in the largest subtable, e.g. the product of the maximum sub-codelist sizes taken over all explanatory variables. The default value is 6000.

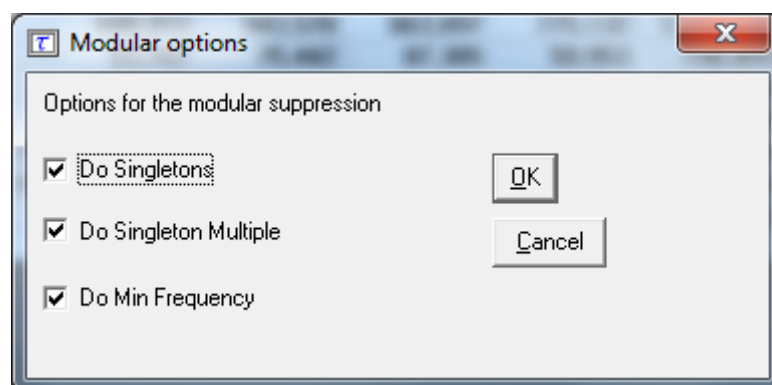
Note that we strongly recommend designing tables so that they fit the ‘normal’ setting, e.g. better think about restructuring the table rather than using the ‘large’ option. The better approach (instead of using the ‘large’ option) would be to introduce a (more detailed) hierarchical structure into the table, because in this way the table will provide more information to the user.

The Cancel button will bring you back to the Table View window, without protecting the table.

### Modular

This partial method will break down the hierarchical table into several non-hierarchical tables, protect them and compose a protected table from the smaller tables. As this method uses the optimisation routines, an LP-solver is required: this will be either XPRESS or CPLEX. The routine used can be specified in the Options box, this will be discussed later.

After starting the modular procedure a little window will be shown whether or not these three additional rules will be applied. At the end of section 2.10 more information on these three rules can be found.

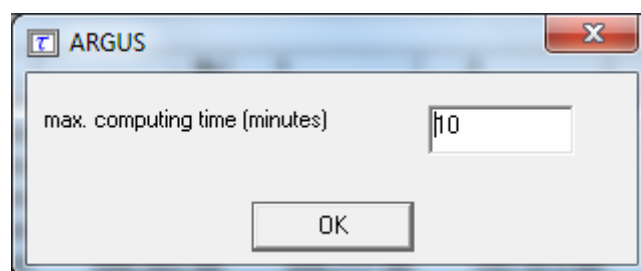


### Optimal

This method protects the hierarchical table as a single table without breaking it down into smaller tables. As this method uses the optimisation routines, an LP-solver is required: this will be either XPRESS or CPLEX. The routine used can be specified in the Options box, this will be discussed later.

It is the responsibility of the users of  $\tau$ -ARGUS to apply for a licence for one of these commercial packages themselves. Information on obtaining one of these licences will be found in a 'read.me' file supplied with the software or on the CASC website.

By choosing 'Suppress/Optimal' a further question is asked. The question is 'How much time do you allow the system to compute the optimal solution'.



When the specified time has been used  $\tau$ -ARGUS will ask you what to do. This can be twofold, you allow  $\tau$ -ARGUS to continue for a new amount of time, or not. The window below allows you to specify this.

Note that  $\tau$ -ARGUS will check only at a specific location in a cycle whether or not the time has elapsed.

**Time Check**

ARGUS has reached the time-limit

Lower limit optimisation 1000

Upper limit optimisation 1100

Difference (percentage) 9.09%

Number of suppressions 16

Time used so far 5 min

Do you want to proceed?

Time allowed next  min.

Yes No

## Network

This is a Network Flow approach for large unstructured 2 dimensional tables and 2 dimensional hierarchical tables with only one hierarchy (the first variable specified). The user has the option of selecting an optimisation method (PPRN and Dykstra). Both optimisation methods are available free of an additional licence. By default the Dykstra solution is advised.

**Parameters for Network solution**

Parameters for the hierarchical network flow solution

Solver type

☐ PPRN

☒ Dykstra

Order of the primaries

☒ Normal

☐ Ascending

☐ Descending

OK Cancel

As the network solution is an heuristic to find an approximation of the real optimal solution, it cannot be expected that always an optimal solution is found. Nevertheless it is guaranteed that at least a good feasible solution is found in a relatively short time. The order in which the primaries are provided to the network

algorithm could influence the solution found. Therefore three options are available to order the primaries.

### Choose the suppression method

After selecting one of the options, click the Suppress button.  $\tau$ -ARGUS will run and display a protected table after informing the user of the number of cells selected for secondary suppression and the time taken to perform the operation. The secondary suppressed cells will be shown in blue.

The screenshot shows the 'Table: Size x Region | Var2' window. The main table displays data for various regions (tot, Nr, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, .99) across different variables (tot, 2, 4, 5, 6, 7, 8, 9, .99). The table is color-coded: primary cells are black, secondary suppressed cells are blue, and cells selected for suppression are red. The right panel contains 'Cell Information' (Value: 11,968, Status: Unsafe, Cost: 11,968, Shadow: 11,968, # contributions: 29, Top n of shadow: 11,401, Holding level: 145, Request: 0, Protection interval: 10,680 to 13,256) and 'Change status' buttons (Set to Safe, Set to Unsafe, Set to Protected, Set Cost, A priori info, All Non-StructEmpty). The bottom panel has 'Suppress' options (HyperCube, Modular, Network, Optimal, Rounding) and buttons (Suppress, Undo Suppress, Audit). The bottom status bar includes '3 dig. separator', 'Output View', 'Select Table', 'Change View', 'Write table', 'Table Summary', and 'Close'.

When the user is satisfied with the table it can be saved (see section 4.5.1 for the possible formats). Press the write-table button. This is the same button as via the menu Output|Save table.

#### 4.4.2.4 Controlled rounding

The “Round” option in the View Table window is active only if the Xpress licence is selected in the Help|Option window. The reason for this is that for Xpress,  $\tau$ -ARGUS has access to the Mixed Integer Model (MIP), thanks to the cooperation of Dash Inc. This option allows to round the selected table with the Controlled Rounding Program (see Section 2.13 for details on this method).

The next figure shows the simple frequency table obtained from the test data using the variable Size and Region.

	tot	2	4	5	6	7	8	9	.99
tot	42,723	9	5	20,002	8,831	5,498	4,594	3,779	
- .Nr	11,395	6	1	5,137	2,471	1,487	1,426	862	
.. 1	6,112	1	1	2,856	1,334	770	744	406	
.. 2	3,798	5	-	1,657	802	497	527	305	
.. 3	1,485	-	-	624	335	220	155	151	
- .Os	10,227	3	1	4,808	2,041	1,380	1,101	893	
.. 4	538	1	-	293	111	58	46	29	
.. 5	1,597	-	-	676	370	283	152	116	
.. 6	5,446	2	1	2,648	1,014	658	572	551	
.. 7	2,646	-	-	1,191	546	381	331	197	
- .Ws	10,054	-	-	4,696	2,062	1,294	1,026	976	
.. 8	1,182	-	-	467	303	178	127	107	
.. 9	8,018	-	-	3,866	1,627	997	799	729	
..10	854	-	-	363	132	119	100	140	
- .Zd	11,047	-	3	5,361	2,257	1,337	1,041	1,048	
..11	7,511	-	3	3,791	1,521	822	666	708	
..12	3,536	-	-	1,570	736	515	375	340	
.99	-	-	-	-	-	-	-	-	

Simple frequency table obtained from the test data.

Rounding can be applied also to tables with no unsafe cells. The choice of the minimum threshold and whether zeros are safe or not has an effect on the minimal possible rounding base, as it will be explained in the *Option* paragraph.

When rounding has been chosen and the round-button has been pressed, the following window will be shown. You can enter a few parameters.

## Rounding Options

Rounding option window

The controlled rounding window allows to set the following parameters:

- **Rounding Base**  
Cell values will be changed to multiples of the base. The minimum rounding base is equal to the maximum between the minimum frequency threshold and twice the highest Protection Level set for an unsafe cell (with the Dominance or p-q rule). See the Section 2.2 for details on safety rules and section 2.6 protection levels. If no rule is specified the minimum base is 1. Rounding to base 1 can be used to round a table with fractional entries for “cosmetic” motives. Note that in some cases the choosing a small rounding base may lead to a nonfeasible problem; increasing the base may make the problem feasible;
- **Number of steps allowed**  
This value specifies the maximum number of *steps* allowed in order to find a feasible solution when a zero-restricted one does not exist. The default value is 0, *i.e.* zero-restricted. Higher values can be chosen by selecting the value from the drop-down menu. Note that the higher the number of steps allowed the lengthier is the search, hence the greater the risk of hitting the time constraint. At any rate, if a zero-restricted solution exists, this is the solution provided, whatever the number of steps allowed.
- **Max computing time**  
This value determines the time after which the user is prompted for a decision about continuing or stopping the search. The default value is 20 minutes. When the maximum time is hit the user is prompted to enter a new maximum time or to choose to terminate the search.

- **Partitions**  
This option enables the partitioning of the table into sub-tables corresponding to each category of the first spanning variable. This option is recommended for tables with more than approximately 150,000 cells. Partitioning can only be used in this version when the first variable is non-hierarchical. The first variable should be such that the sub-tables have maximum size of about 150,000 cells and also trying to keep their number low; performance may be improved by wisely choosing the partitioning variable. The number of subtables and the size of each subtable is shown in the window. See Section (rounding theory) for further details. When rounding, the progress of the partitioning process is shown at the bottom of the option window;
- **Stopping Rule**

These options allow to control the quality of the rounded solution. The user can choose:

- **First Rapid**  
The solution is obtained by rounding conventionally (to the closest multiple of the base) the internal cells and then the marginal values are obtained by addition. This solution is likely to present several values that have a large distance from the original values. This option should be used with extreme care and, likely, when everything else fails;
- **First feasible**  
The solution provided is the first rounded one that has the specified number of jumps, regardless of its optimality. This means that there could exist other solutions that have a lower overall distance from the original table. In many cases, when optimality is not crucial, this solution is quite close to the optimal one and it can be found in a shorter time;
- **Optimal**  
This option provides the fully optimal controlled rounded solution.

### **The rounded table**

The next figure shows the rounded table shows the values rounded to multiples of 5. Note that the values that were originally zero (hence empty cells denoted with a dash) are still shown as a dash while the values that have been rounded down to zero are shown as zeros.

Table: Size x Region   freq									
	tot	2	4	5	6	7	8	9	99
tot	42,725	10	5	20,000	8,830	5,500	4,595	3,780	5
- .Nr	11,395	5	0	5,135	2,470	1,490	1,430	860	5
.. 1	6,110	0	0	2,855	1,335	770	745	405	-
.. 2	3,800	5	-	1,655	800	500	530	305	5
.. 3	1,485	-	-	625	335	220	155	150	-
- .Os	10,230	5	0	4,810	2,040	1,380	1,100	895	-
.. 4	540	0	-	295	110	60	45	30	-
.. 5	1,600	-	-	675	370	285	155	115	-
.. 6	5,445	5	0	2,650	1,015	655	570	550	-
.. 7	2,645	-	-	1,190	545	380	330	200	-
- .Ws	10,055	-	-	4,695	2,065	1,295	1,025	975	-
.. 8	1,180	-	-	465	305	180	125	105	-
.. 9	8,020	-	-	3,865	1,630	995	800	730	-
..10	855	-	-	365	130	120	100	140	-
- .Zd	11,045	-	5	5,360	2,255	1,335	1,040	1,050	-
..11	7,510	-	5	3,790	1,520	820	665	710	-
..12	3,535	-	-	1,570	735	515	375	340	-
.99	-	-	-	-	-	-	-	-	-

Cell Information

Value: 42,725  
Status: Safe  
Cost: 42,723  
Shadow: 0  
# contributions: 42723  
Top n of shadow:  
☐ Holding level  
Request: 0

Change status

Set to Safe  
Set to Unsafe  
Set to Protected  
Set Cost  
A priori info  
All Non-StructEmpty

Recode

☐ HyperCube  
☐ Modular  
☐ Network  
☐ Optimal  
☒ Rounding

Round

Undo Rounding

Audit

☒ 3 dig. separator  
☐ Output View

Select Table

Change View

Write table

Close

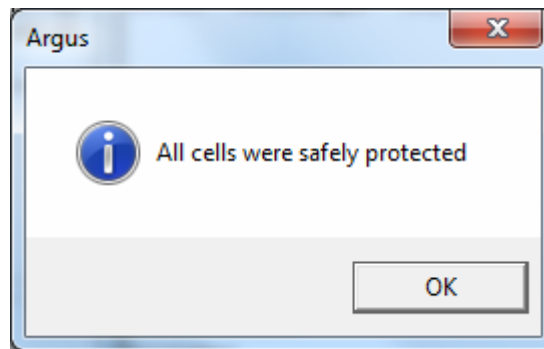
Table Summary

The rounded table

#### 4.4.2.5 The audit procedure

After the secondary cell suppression procedure has been carried out all cells should have been properly protected. Cell suppression guarantees that unsafe cells cannot be estimated to a narrower interval than the required protection interval. The realised upper and lower bounds can be computed by solving two linear programming problems for each unsafe cell. This can be rather an effort doing it all manually, but the audit procedure will do this. See also section 2.14.

The Audit option will only be active after secondary cell suppression. By activating the procedure all the linear programming problems for all unsafe (both primary and secondary) cell will be computed. When completed a message will be shown whether all cells were protected correctly.



If in the unfortunate case the protection was not optimal according to the audit procedure a list of problems will be shown.

For each unsafe cell the realised lower and upper bounds will be shown. If you put your mouse on the value also the distance to the real value and the corresponding percentage will be shown

The screenshot shows the main Argus application window. The title bar reads 'Table: Size x Region | Var2'. The main area contains a table with columns labeled 'tot', '2', '4', '5', '6', '7', '8', and '9'. The rows are labeled with various codes like 'tot', '.Nr', '1', '2', '3', '4', '5', '6', '7', '8', '9', '10', '11', '12', and '.99'. Some cells in the table are highlighted in red or blue. To the right of the table is a sidebar with 'Cell Information' showing details for a selected cell (Value: 11,968, Status: Unsafe, etc.). Below this are buttons for 'Change status' (Set to Safe, Set to Unsafe, Set to Protected, Set Cost) and 'Suppress' (HyperCube, Modular, Network, Optimal, Rounding). At the bottom of the window are checkboxes for '3 dig. separator' and 'Output View', and buttons for 'Select Table', 'Change View', 'Write table', 'Table Summary', and 'Close'.

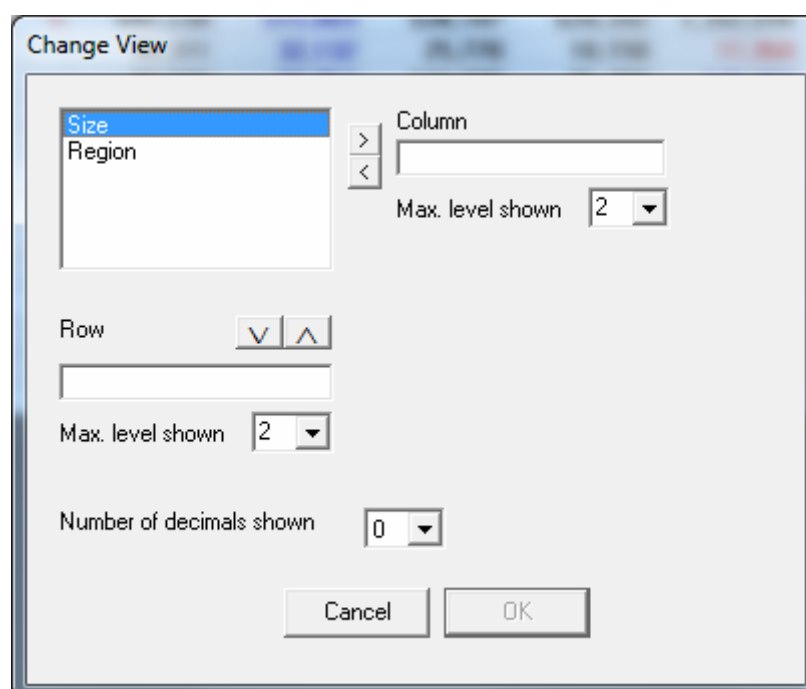
#### 4.4.2.6 The Options at the Bottom of the table

At the bottom of this window there are a few additional options. These options will be described here.

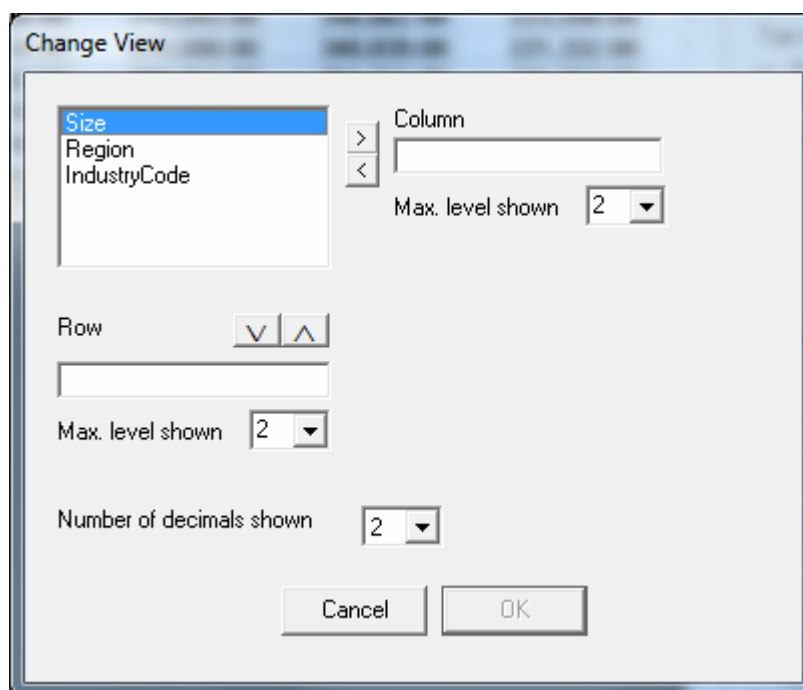
## Change View

By clicking on Change View in the Table window after clicking on Modify|ViewTable at an earlier stage the dialog box below pops up. The user can specify which variable is wanted in the row and the column. In the two-dimensional case, the table can only be transposed. In the higher dimensional case, the remaining variables will be in the layer. For these layer variables a combo-box will appear at the top of the table, where the user can select a code. This will show the corresponding slice of the table.

Also the level of detail for each (hierarchical) explanatory variable can be chosen. This will influence the size/level of the table shown. Of course after this, the level of detail in each spanning variable can be changed by folding and unfolding the hierarchy (clicking on the codes).



For a 3 dimensional table, this window is as follows:



### Table summary

Pressing 'table summary' provides a table summary giving an overview of the number of cells according to their status. The example shown here refers to the case after secondary suppression has been performed.

Summary for table no: 1

Explan. Var	# Codes	Status	Freq	# rec	Sum Resp	SumCost
Size	9	Safe	98	245051	98032974.04	98032974.04
Region	18	Safe (manual)	0	0	0	0
		Unsafe	9	43	12013	12013
		Unsafe (request)	0	0	0	0
		Unsafe (Freq)	3	9	45	45
		Unsafe (Zero cell)	0	0	0	0
		Unsafe (manual)	0	0	0	0
		Protected	0	0	0	0
		Secondary	9	11235	3040849	3040849
		Secondary (from man.)	0	0	0	0
		Empty (non-struct.)	0	0	0	0
		Empty	43	0	0	0
		Total	162	256338	101085881.04	101085881.04

Protected by Modular

The headings in the summary window are as follows:

*Freq*: The number of cells in each category

*# rec*: The number of observations in each category

*#Holding*: The number of holdings in each category (0 if holdings are not used for this table)

*Sum Resp*: Total cell value in each category

*SumCost*: The sum of the cost variable.

### **3 dig separator**

This removes or inserts the character separating the thousands for the values in the table.

### **Output View**

This option allows the table to be shown as it will be output, with suppressed cells (primary and secondary) replaced by a X.

## **4.4.3 Modify | Linked Tables**

This option is available when the tables specified have at least one explanatory or spanning variable in common and the same response variable.

When the tables are built from micro data, the tables can be specified using the screen below. See also section 4.3.1.

An example is shown.

Expl. vars	rule	Resp. var	Shadow & Cost var
Size,Region	IND.: p= 15, q= 100, N= 1, MinFreq ...	Var2	Shadow=Default, Cost=Default
Year,Region	IND.: p= 15, q= 100, N= 1, MinFreq ...	Var2	Shadow=Default, Cost=Default

When the tables are supplied to  $\tau$ -ARGUS as tabular input see section 4.2.3 (Open Table set). When supplying a set of ready made tables it should be clear to  $\tau$ -ARGUS which explanatory variables are in fact the same dimension. They should have the same name, even if the level of detail is different.

E.g. if a regional variable is an explanatory variable in two tables, but in one table it is at the level of province and in the other at the level of municipality, they should nevertheless have the same name. If not  $\tau$ -ARGUS will not recognise them as a link.

The set of linked tables can be protected using the hypercube (see section 2.8) and the extended modular approach (see section 2.11).

When protecting a set of linked tables the restriction is that all tables are a sub-set of a theoretical cover table. The cover table is formed by building a table using the longest code list for each dimension. The dimensions are decided by looking for different names of explanatory variables.

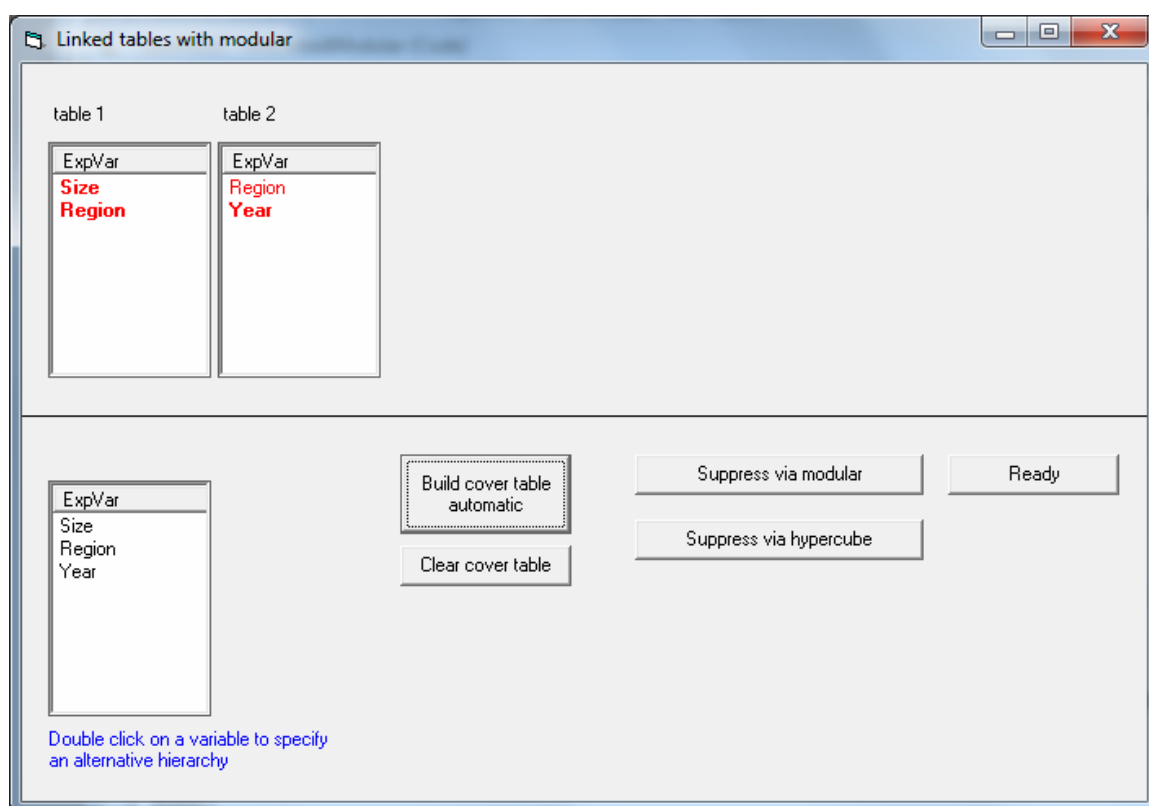
The cover table can be specified in two ways:

- “automatic”:  $\tau$ -ARGUS will look for the largest codelists per dimension and select the largest one for the cover table.
- The user can do it manually. By double clicking on a variable in one of the sub-tables. That variable is selected for the cover table. That variable will become red and bold, the other variables will become only red. When all variables are red the cover table has been specified completely.

As long as the cover table does not have more than 4 dimensions the linked tables approach is possible.

In the current implementation there is one restriction. For each of the spanning variables in the cover table the codelist and the hierarchy should be present in one of the linked tables. For all other tables the codelists and the hierarchy should be a subset of this cover hierarchy. And of course the set of linked tables should be consistent. The cells that are logical the same should have exactly the same value and status. If not the protection of the cover table will fail.

When tabular data is the starting point, it is the responsibility of the user that the tables are consistent. This means that the cell values of corresponding cells are the same and also the status. If not this is an inconsistent situation. The modular approach is very strict on complete additivity, as the optimisation routines behind modular require this. The hypercube is a bit more relaxed.



When the cover table has been defined press ‘Suppress via modular’ or ‘suppress via hypercube’.  $\tau$ -ARGUS will then start an automatic procedure.

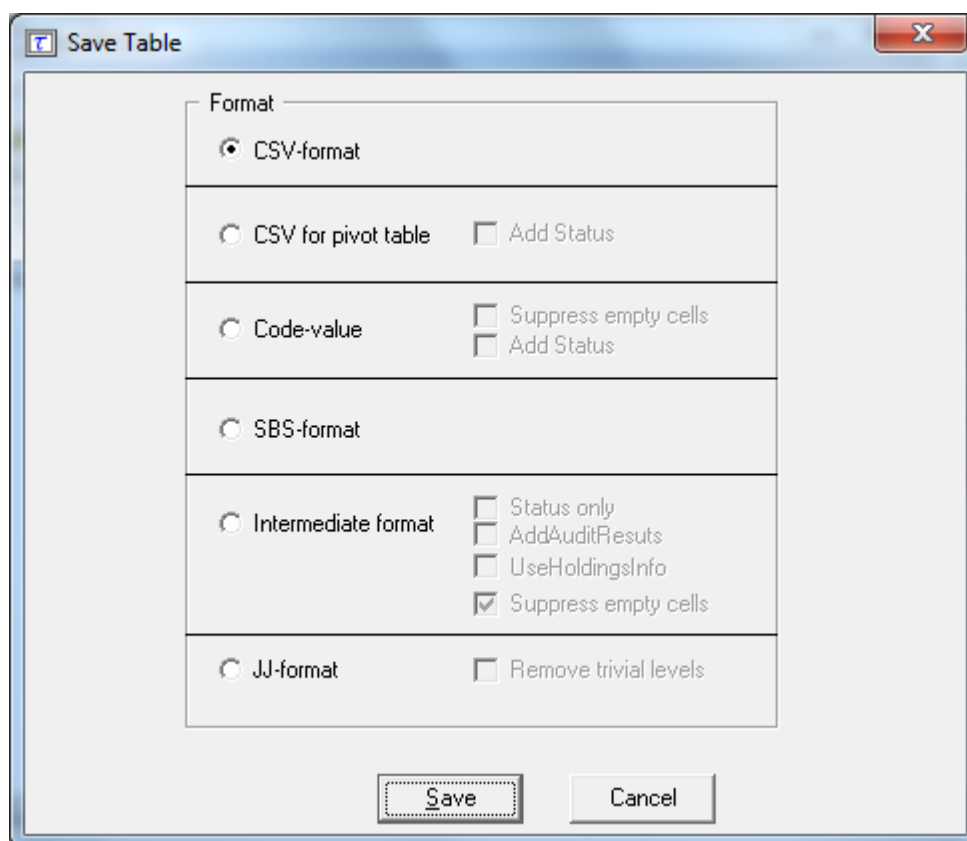
When the modular approach is selected, the subtables will be loaded in the cover table. The cover table will then be protected via an extra batch-run of  $\tau$ -ARGUS and in the end the results (suppression pattern) will be transferred to the original subtables. If this procedure might fail, information could be found in the log-file of  $\tau$ -ARGUS. See also section 5.7.

When the hypercube is selected, all the input files for the hypercube will be prepared and the linked table procedure of the hypercube will be started to protect the set of tables.

When the protection has been completed, the linked tables procedure can be closed and the individual protected subtables can be inspected and stored as normal tables.

### 4.5.1 Output | Save Table

There are five options of storing the tables



**As a CSV file.** This Comma Separated file can easily be read into Excel. Please note the Excel should interpret the comma as a separator. If your local settings are different you could use the Excel option 'Data|Text to Columns'. This is a typical tabular output maintaining the appearance of the table in  $\tau$ -ARGUS.

**CSV-file for a pivot table.** This offers the opportunity to make use of the facilities of pivot table in Excel. The status of each cell can be added here as an option (Safe, Unsafe or Protected for example). The information for each cell is displayed on a single line unlike standard csv format

**A text file** in the format code-value, this is separated by commas. Here, the cell status is again an option. Also empty cells can be suppressed from the output file if required. The information for each cell is displayed on a single line similar to the CSV file for a pivot table. There are two possibilities. Either the unsafe cells are shown as an 'x', as it should be in the final publication or the exact status can be printed in the output file in addition to the cell value. Optionally empty cells can be suppressed.

When the status is added to the output file  $\tau$ -ARGUS can use 14 different statuses. They can also be found in the report file.

Number	Status
1	Safe
2	Safe (manual)
3	Unsafe
4	Unsafe (request)
5	Unsafe (Freq)
6	Unsafe (Zero cell)
7	Unsafe (Singleton)
8	Unsafe (Singleton) (manual)
9	Unsafe (manual)
10	Protected
11	Secondary
12	Secondary (from man.)
13	Empty (non-struct.)
14	Empty

**A SBS-format file.** This file contains the information required by Eurostat for different surveys like the SBS-survey. Each line describes one cell in the table. First all the spanning variables, with the levels in the hierarchy, then the cell value, the cell frequency, the status and the dominance percentage. If the 2 largest contributors have been computed this percentage is the sum of the largest two, otherwise the largest one. It will be obvious that this output format is not possible if a table has been used as input, with only the status or maybe a cell frequency.

The cell status can be:

A	Frequency unsafe
B	Dominance unsafe with one contributor
C	Dominance unsafe with two contributors
D	Secondary unsafe
V	Safe

**A file in intermediate format** for possible input into another program. This contains protection levels and external bounds for each cell. This file could even be read back into  $\tau$ -ARGUS, using the read tables option

The options are:

- Write only the status
- Add the results of the audit procedure (realised lower and upper bounds)
- Write information at the holding level, like the frequency.

- Suppress the empty cells.

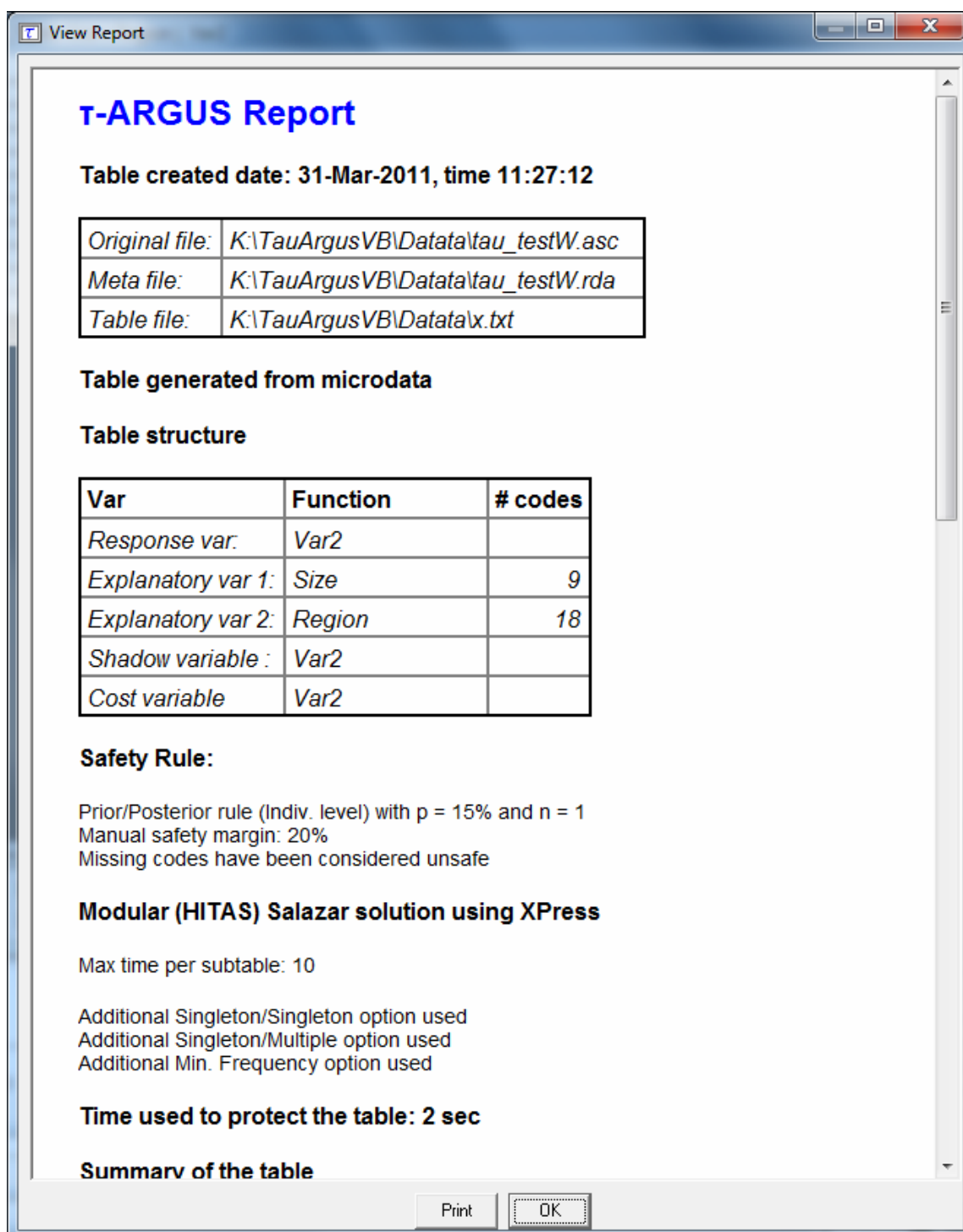
Of course certain options are only available if appropriate.

Finally, a report will be generated to a user specified directory. This report will be shown, when the table has been written. As this is an HTML-file it can be viewed easily later.

### **4.5.2 Output | View Report**

Views the report file which has been generated with Output|Save Table. An example of a part of the output HTML file is shown here.

As can be seen the essential information, for somebody other than the user, about which rules have been applied to make the data safe is displayed along with details of any recoding.



### 4.5.3 Output | GenerateApriory

In many situations it is desirable to coordinate the secondary suppressions between tables. This can be because of links between tables. Suppressions in one table should also be suppressions on the other table.

But also when protecting monthly tables it could be a good idea to coordinate suppressions between the different months. A secondary suppression in month 1 could be an ideal candidate for secondary suppression in month 2. This could be

achieved by changing the suppression weights for these cells.

The apriory option is the way to change the default suppression weights etc. But this leaves the task of generating the apriory file.

Via this option the protected file as generated by  $\tau$ -ARGUS can be converted into an apriory file. The table has to be saved in the format code-value with the 'Add Status'-option selected.

For each status the user can select which action in the apriory file has to be created. This can be a change of the suppression weight, give a new status, or nothing at all.

The user has to specify the protected file (written in the right format (saved as code/value plus status) and the apriory file to be generated.

Also the correspondence between the variables must be specified. It is not always the case that the first spanning variable is also the first spanning variable in the apriory file. Even the number of variable can be different. If not all variables of the safe file will be available in the newly to be protected file, only the score for the total will be used. This is often the case if the apriory file is generated for a linked tables problem.

The separator to be used in the apriory file must be specified as well; a comma is the default.

Pressing the 'Go'-button will generate the apriory file and the 'ready'-button will bring you back to the main menu of  $\tau$ -ARGUS.

**Make Apriory file**

Name of the safe file  
 ...

Separator

Name of the apriory file  
 ...

Dimension Output

1  ...  
 2  ...

Status	Omit	Safe	Unsafe	Protect	Weight	Weight Value
Safe	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="checkbox"/>	<input type="text" value="6"/>
Safe (manual)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="checkbox"/>	<input type="text" value="7"/>
Unsafe	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input checked="" type="checkbox"/>	<input type="text" value="5"/>
Unsafe (request)	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="checkbox"/>	<input type="text"/>
Unsafe (Freq)	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="checkbox"/>	<input type="text"/>
Unsafe (Zero cell)	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="checkbox"/>	<input type="text"/>
Unsafe (Singleton)	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="checkbox"/>	<input type="text"/>
Unsafe (Singleton) (manual)	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="checkbox"/>	<input type="text"/>
Unsafe (manual)	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="checkbox"/>	<input type="text"/>
Protected	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="checkbox"/>	<input type="text"/>
Secondary	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="checkbox"/>	<input type="text"/>
Secondary (from man.)	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="checkbox"/>	<input type="text"/>
Empty (non-struct.)	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="checkbox"/>	<input type="text"/>
Empty	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="checkbox"/>	<input type="text"/>

#### 4.5.4 Output | Write Batch File

The commands used in interactive mode can be saved into a file for future use. T-ARGUS will write a batch file containing the commands necessary to achieve the current situation of the  $\tau$ -ARGUS run so far. For more information on the batch facility see section 4.2.3

For example the following shows the dominance rule ( $n=3$ ,  $k=75$ ) applied to the Size by Region table with Var2 as the response variable. The threshold value = 5 with a safety range = 30%. Modular secondary suppression was applied. The last line indicates that  $\tau$ -ARGUS will not stop after these commands but become an

interactive program.

```
<OPENMICRODATA> "C:\Program  
Files\TauARGUS\data\tau testW.asc"  
<OPENMETADATA> "C:\Program  
Files\TauARGUS\data\tau testW.rda"  
<SPECIFYTABLE> "Size" "Region" | "Var2" | "" | ""  
<SAFETYRULE> NK(3,75) | NK(0,0) | FREQ(5,30) |  
<READMICRODATA>  
<SUPPRESS> MOD(1)  
<GOINTERACTIVE>
```

## 4.6 The Help menu

---

### 4.6.1 Help | Contents

This shows the contents page of the help file and from there makes the help available. This program has context-sensitive help.

### 4.6.2 Help | News

Information on the latest developments is shown. Old friends can see here which new extensions have been included in this version of  $\tau$ -ARGUS and information about bugs is shown here as well.

### 4.6.3 Help | Options

There are a number of options, which can be changed here. The colours indicating the status of a cell can be altered.

In order to make a hierarchical table more readable, the different levels of the hierarchy will be indicated with an increasing grey background. If you like different colours, you can adapt this.

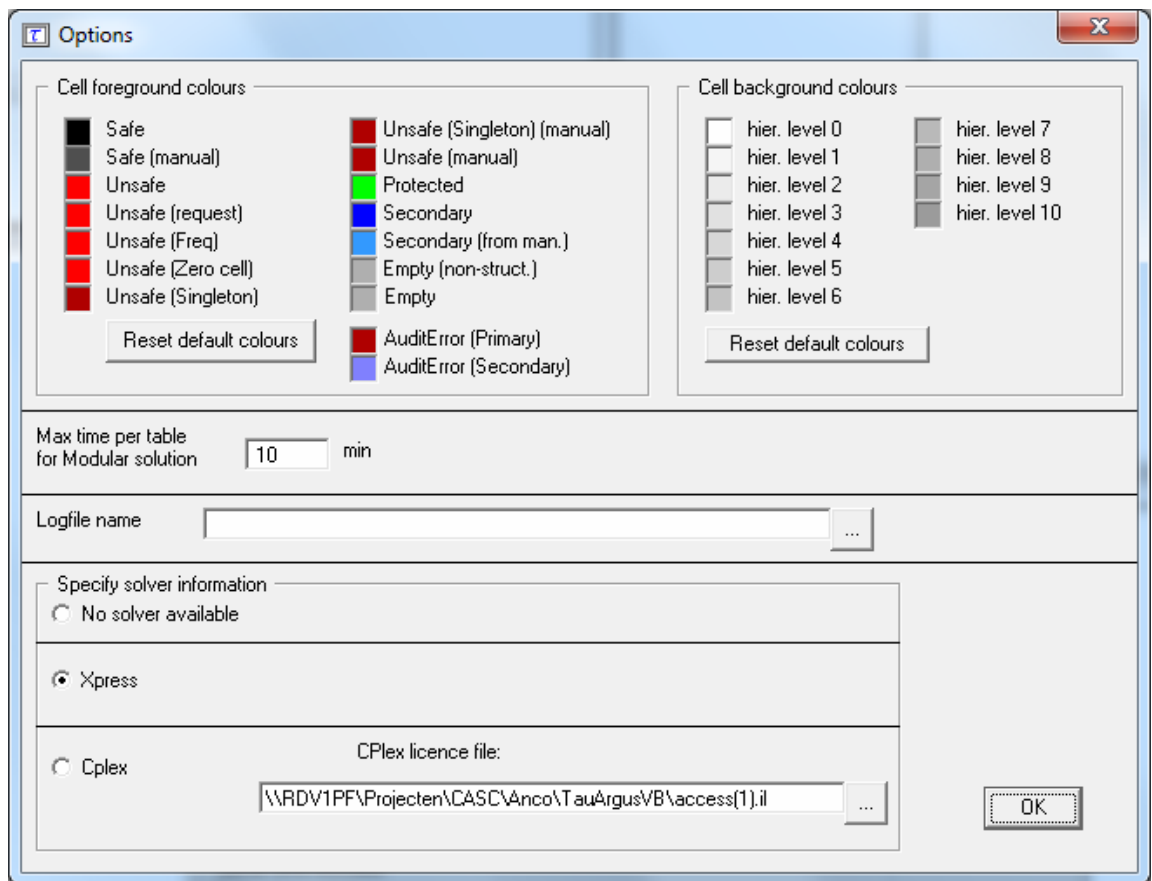
For the modular solution the maximum computing time per subtable can be specified. This could speed up the computations, but on the other hand might give a less optimal solution.

Also the name of the logfile (see section 5.7) can be changed here. By default it is Logbook.txt in the temp-directory.

Finally the solver for the optimisation routines must be specified. The options are: No Solver, CPLEX or XPRESS.  $\tau$ -ARGUS can work with both solvers.

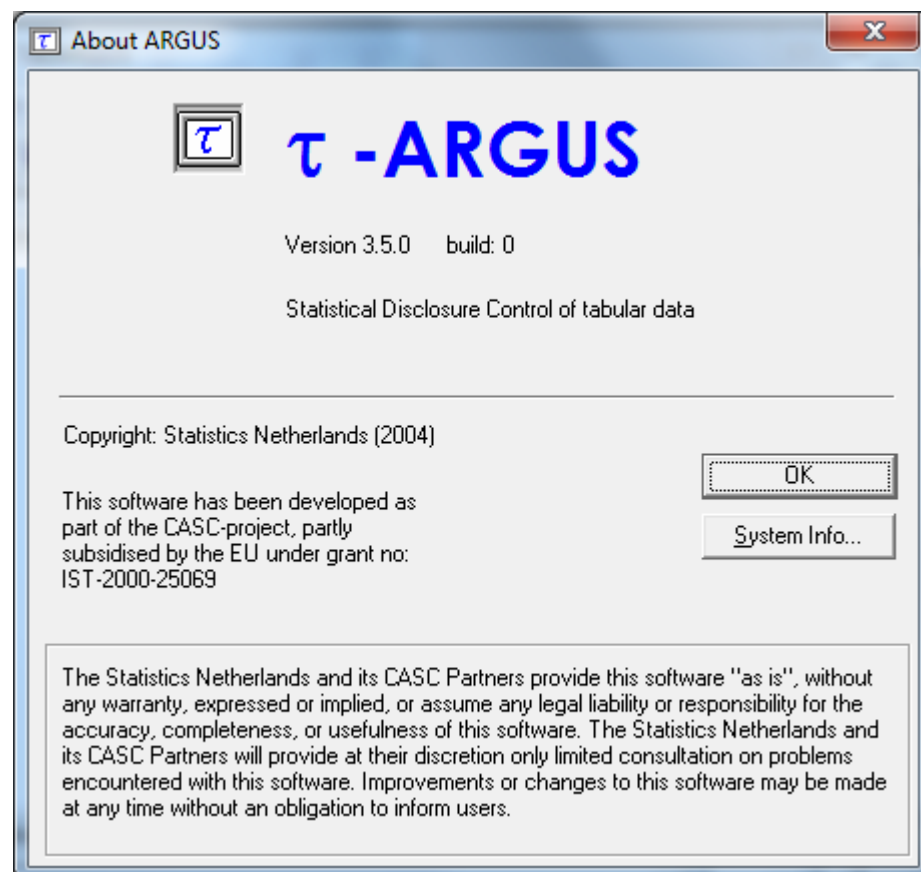
If the Cplex optimisation routine is being used, the location of the licence file can be specified here. For XPress the name of the licence file is prescribed and fixed (XPAUTH.XPR) The XPress licence file has to be stored in the  $\tau$ -ARGUS program directory.

$\tau$ -ARGUS will store the information of all these options in the registry and will use it in future runs. It is advisable but not necessary to open this window at the start of a  $\tau$ -ARGUS session to ensure the correct solver has been chosen.



#### 4.6.4 Help | About

Shows the about box.



## 5 FURTHER DESCRIPTIONS

---

### 5.1 *Meta data files*

---

The meta data plays a vital role in  $\tau$ -ARGUS. The meta data is always specified and stored in a separate file. As  $\tau$ -ARGUS can read both micro data as input as tabular data, the meta data descriptions will be different as well. Nevertheless there are many similarities, especially between meta data for fixed and free format micro data.

The meta data can always be changed/adapted/entered via the menu item Specify|Metadata, see sections 4.3.1, 4.3.2 and 4.3.2.

As no standard meta data system is available which is powerful enough to manage the complete metadata specification necessary for Statistical Disclosure Control we had to develop something specially for  $\tau$ -ARGUS.

The metadata file (default extension .RDA) has globally the same structure for all the different file types that can be handled by  $\tau$ -ARGUS. i.e fixed format/free format microdata, SPSS system files or tabular data.

For each variable the name is specified followed by its position in the file and possible missing values. Following this specification additional information can be specified. These specifications always start with a keyword enclosed by “<” and “>” followed by the specifications.

The metadata is always stored in a plain text file without any tabs or so. If you wish you could enter/modify the metadata file with e.g. Notepad, but not with Word. It is then your own risk that the metadata is syntactically correct. T-ARGUS will check the meta data file when it is read, but to a certain limit. The best way is to modify the metadata via the  $\tau$ -ARGUS program

We will first describe in section 5.1.1 the meta data file a fixed format micro data file. In the subsequent sections the special issues for the other file formats (free format and SPSS) will be described. In section 5.1.4 the meta data for tabular data files will be described.

#### 5.1.1 **Meta data for fixed format micro data**

For fixed format for each variable the starting position and the field length have to be specified. Also the possible missing values must be specified as well as the role that a variable can play in the SDC-analysis, like spanning variable (also known as explanatory variable), cell item, weight, etc. Additional extra specifications can be entered. as well, like codelists and hierarchical structures.

The metafile describes the variables in the microdata file, both the record layout and some additional information necessary to perform the SDC-process. Each variable is specified on one main line followed by one or more option lines.

The first line gives the name of the variable followed by the starting position for each record, the width of the field and optionally one or two missing value indicators for the record. Missing values are not required in  $\tau$ -ARGUS, but they

can play a role when deciding whether or not a cell is unsafe.

For fixed format microdata it is not necessary to specify all the variables in the file. Only the variables used in  $\tau$ -ARGUS have to be specified. When reading the data  $\tau$ -ARGUS will ignore the fields not described.

The following lines explain specific characteristics of the variable:

• <RECODEABLE>	This variable can be recoded and used as an explanatory variable in a table
• <CODELIST>	This explanatory (or spanning) variable can have an associated codelist which gives labels to the codes for this particular variable. The name of the codelist file follows this <CODELIST> command. The default extension is .CDL. See section 5.3
• <NUMERIC>	This numeric variable can be used as cell-item.
• <DECIMALS>	The number of decimal places specified for this variable
• <WEIGHT>	This variable contains the weighting scheme
• <HIERARCHICAL>	This variable is hierarchical. The codings are structured so that there is a top code such as Region (N,S,E,W) and within each of these are smaller more specific areas (and possibly sub-areas). Tables may be viewed at different levels of hierarchy.
• <HIERLEVELS>	The hierarchy is derived from the digits of the codes itself. The specification is followed by a list of integers denoting the width of each level. The sum of these integers should be the width of the total code. An example is shown beneath the rda file below.
• <HIERCODELIST>	The name of the file describing the hierarchical structure. Default extension .HRC. See section 5.2.
• <HIERLEADSTRING>	The string/character that is used to indicate the depth of a code in the hierarchy. See section 5.2
• <REQUEST>	This variable contains the status denoting whether or not a respondent asked for protection
• <HOLDING>	This variable contains the indication whether a group of records belong to the same group/holding

Here is an example of a rda file for microdata. (Note, the dots at the bottom just means that here a shortened version of the file is presented.)

```
YEAR 1 2 99
  <RECODEABLE>
IndustryCode 4 5 99999
  <RECODEABLE>
  <HIERARCHICAL>
  <HIERLEVELS> 3 1 1 0 0
Size 9 2 99
  <RECODEABLE>
Region 12 2 99
  <RECODEABLE>
  <CODELIST> Region.cdl
  <HIERCODELIST> Region2.hrc
  <HIERLEADSTRING> @
  <HIERARCHICAL>
```

```

Wgt 14 4
  <NUMERIC>
  <DECIMALS> 1
  <WEIGHT>
Var1 19 9 999999999
  <NUMERIC>
Var2 28 10 9999999999
  <NUMERIC>
  <DECIMALS> 2
.....

```

#### Explanation of the details of the variables

‘*Year*’: For this explanatory/spanning variable each record begins on position 1, is 2 characters long and missing values are represented by 99. It is also recodeable implicitly stating that it is an explanatory or spanning variable used to create the tables.

‘*IndustryCode*’: For this variable each record begins on position 4 and is 5 characters long. Missing values are represented by 99999. As well as being recodeable this variable is hierarchical and the hierarchy structure is specified. The first 3 characters are in the top hierarchy level, the 4<sup>th</sup> character in the second level and the 5<sup>th</sup> character in the lowest level. As ‘*Industry*’ is a 5 digit variable there are 5 digits specified for the hierarchical structure. This is the reason for the 2 zeros at the end.

‘*Size*’: For this variable each record begins on position 9 and is 2 characters long, and missing values are represented by 99. It is also recodeable.

‘*Region*’: For this variable each record begins on position 12 and is 2 characters long. Missing values are represented by 99. Region has a codelist. See section 5.3. Region is also a hierarchical variable. As the hierarchical structure cannot be derived from the structure of the coding scheme itself the hierarchical structure is described in a special .HRC file. See section 5.2. The hierarchical structure is described with an indentation structure. Therefore the indentation character (HIERLEADSTRING) has to be specified. Here an @ was chose.

‘*Wgt*’: For this variable each record begins on position 14 and is 4 characters in length. There is 1 decimal place for these values and the variable is defined as a weight. A missing value is not allowed here.

Two numeric variables are also shown in the above rda file. These numeric variables (not defined as weights) are those to be used as cell items *i.e.* response variables used in creating the table.

‘*Var1*’: This variable begins on position 19 and is 9 characters long. Missing values are represented by 999999999 and it is numeric. However the missing values for numerical variables will be ignored. The missing values problem should have been solved by e.g. imputation techniques, but is outside of the scope of  $\tau$ -ARGUS

‘*Var2*’: This variable begins on position 28 and is 10 characters long. Missing values are represented by 9999999999 and it is numeric. This variable has 2 decimal places.

*The representation in an rda file for the Request rule and Holding Indicator are shown here for completeness.*

*Request rule*

```
Request 99 1
<REQUEST> "1" "2"
```

Here the request indicator is in column 99 and is one character long. Individuals (or companies) wishing to make use of this rule are represented by 1 or 2. Any other value will be interpreted as 'no request'. Two different parameters-sets for the request rule can be specified, the first set will be applied to the companies where the first code has been specified, the second set to the companies with the second code. The request rule is further explained in section 4.3.4.

This rule is used in foreign trade statistics and based on a special regulation..

#### *Holding Indicator*

```
entgroup 101 4
<HOLDING>
```

Here the variable 'entgroup' is in column 101 and is four characters long. This variable is to act as the holding indicator (see section 4.3.1 for further explanation). The records of a holding should be grouped together in the input datafile.  $\tau$ -ARGUS will not search through the whole file to try to find all records for a holding. Before applying the sensitivity rules all records of one holding are grouped together and treated as one contribution.

### 5.1.2 Meta data for free format micro data

For a free-format datafile the RDA is a little bit different. Notably the first line specifies the separator used. This indicates to  $\tau$ -ARGUS that the record description is for a free-format file. And for each variable the starting position is no longer specified, as this is meaningless in a free-format datafile. For the rest there are no differences compared to the fixed format version. The example given above for a fixed format file will now look as:

```
<SEPARATOR> ", "
YEAR 2 99
<RECODEABLE>
Sbi 5 99999
<RECODEABLE>
<HIERARCHICAL>
<HIERLEVELS> 3 1 1 0 0
GK 2 99
<RECODEABLE>
Regio 2 99
<RECODEABLE>
<CODELIST> REGION.CDL
<HIERARCHICAL>
<HIERCODELIST> region2.hrc
<HIERLEADSTRING> @
Wgt 4 9999
<NUMERIC>
<DECIMALS> 1
<WEIGHT>
Var1 9 999999999
<NUMERIC>
Var2 10 9999999999
<NUMERIC>
<DECIMALS> 2
```

```
.....
.....
```

### 5.1.3 Meta data for SPSS system files

When the microdata is stored in a SPSS System file  $\tau$ -ARGUS can also read this data. However some special rules have to be taken into account. It is assumed that a valid license for SPSS is available on the computer, because  $\tau$ -ARGUS will call SPSS to read the data in the systemfile. Also part of the metadata will be retrieved from SPSS. However not all meta data needed for  $\tau$ -ARGUS is available in SPSS, so the user has to enter the additional metadata himself. See section 4.3.2

In fact  $\tau$ -ARGUS will call SPSS to export the data and will create a fixed format scratch file in the temp-directory. After that  $\tau$ -ARGUS will work similar to working with fixed format ASCII files.

The first time you open a SPSS systemfile, no metadata file can and has to be specified.

After opening the SPSS system file in this menu option SPSS will be called and the meta data (Variable names, field length, missing values) available in SPSS will be read. This is a process that takes a bit of time and should not be interrupted by pressing any key or so. However no progress information can be showed on the screen.

If you reopen an SPSS system file with a meta data file,  $\tau$ -ARGUS will check whether all the variables in the RDA file are really available in the system file.

The RDA file is very similar to the RDA file for a fixed format ASCII file. One exception is that the first line will read

```
<SPSS>
```

### 5.1.4 Meta data for tabular data files

When a tabular datafile has been selected, the metadata file will have a different structure. Clicking on 'Specify|Metafile' gives the opportunity to either edit the metafile already read in or to enter the metafile information directly at the computer.

As tabular input is always expected to be free format, first the separator has to be specified.

The variables can have the following role:

• <RECODEABLE>	The spanning variables used to produce the table. The same as for microdata input files, like hierarchical structures and codelist
• <TOTCODE>	Code for the total of a codelist
• <NUMERIC>	Response Variable – The variable used to calculate the cell total.
• <NUMERIC> <SHADOW>	Shadow variable – The variable is used as a shadow variable.

• <NUMERIC> <COST>.	Cost variable – The variable is used as the cost-variable
• <NUMERIC> <LOWERPL>	Lower protection level – The lower protection level
• <NUMERIC> <UPPERPL>	Upper protection level – The upper protection level
• <FREQUENCY>.	Frequency – This indicates the number of observations making up the cell total. If there is no frequency variable each cell is assumed to consist of a single observation
• <MAXSCORE>	‘topN variable’ – This shows if this variable is defined as one of the top N contributors to the cell. The pre-defined value for TopN is 1. The first variable declared as ‘topN’ will contain the largest values in each cell, the second variable so declared will contain the second largest values etc.
• <STATUS>	>‘Status Indicator’ – allows a variable in the left-hand pane to be declared as a Status Indicator. Typically cells can be declared as Safe, Unsafe or Protected
• <SAFE>	The code used for indicating that a cell is safe
• <UNSAFE>	The code used for indicating that a cell is unsafe
• <PROTECT>	The code used for indicating that a cell is protected and cannot be used for secondary suppression

For explanatory variables the code for the total has to be specified. We recommend strongly that the user also provides the values for the totals himself, but if needed he can ask  $\tau$ -ARGUS to compute these totals. In any case,  $\tau$ -ARGUS needs these totals as they play an important role in the structure of a table and also are important for the suppression models.

```

<SEPARATOR>      ", "
<SAFE> s
<UNSAFE> u
<PROTECT> p
expvar1
  <RECODABLE>
  <TOTCODE> T
expvar2
  <RECODABLE>
  <TOTCODE> T
respvar
  <NUMERIC>
freq

```

```

    <FREQUENCY>
top1
    <MAXSCORE>
top2
    <MAXSCORE>
top3
    <MAXSCORE>
stat
    <STATUS>

```

## 5.2 *Hierarchy file*

---

Hierarchical structures play an important role in  $\tau$ -ARGUS. The hierarchical structures can often be derived from the code itself. E.g. the NACE classification is an example of this. In other situations the structure is not so clear. In that case the whole structure has to be specified. A hierarchical structure is in fact a tree. And a tree can be described easily by indentation. In  $\tau$ -ARGUS a hierarchical structure can be described in a simple text-file, using Notepad or something similar. The default extension is .HRC.

One level deeper means a new sub-node in the tree. In the example given below only two levels are shown, but many more levels are allowed. The indentation (an @ in this example) character has to be specified separately in the RDA file.

Note that the total code is never specified in these .HRC files, as  $\tau$ -ARGUS always assumes that the total will be computed.

Note also that in this situation the codes 1 to 9 in a fixed format file have a leading space. This space should be used in the HRC-file as well.

region2.hrc

```

Nr
@ 1
@ 2
@ 3
Os
@ 4
@ 5
@ 6
@ 7
Ws
@ 8
@ 9
@10
Zd
@11
@12

```

### 5.3 Codelist file

---

Codelists can be specified for explanatory variables. The codes are stored in a separate file (default extension .CDL).

However the codes are only used to enhance certain windows during the processing. T-ARGUS itself will create the coding schemes for the variables used during the processing of the datafile. So a code not specified in the .CDL-file will not cause any problem, only the label is not available. Also codes specified but not found in the data file will be ignored.

The structure of the file is simple. Each line contains a code and a label separated by a “,”

*region.cdl*

```
1,Groningen
2,Friesland
3,Drenthe
4,Overijssel
5,Flevoland
6,Gelderland
7,Utrecht
8,Noord-Holland
9,Zuid-Holland
10,Zeeland
11,Noord-Brabant
12,Limburg
Nr,North
Os,East
Ws,West
Zd,South
```

### 5.4 Global recode file

---

Global recoding is a powerful method to reduce the number of primary unsafe cells. It reduces the size of the table, but the advantage is also that the number of primary unsafe cells is reduced. It is a classical balance to decide how far you should go when applying global recodes, but often the resulting table contains much more information, compared to a table with many, but suppressed cells.

For a hierarchical coding scheme  $\tau$ -ARGUS allows recoding via collapsing the tree structure of the hierarchy. But for non-structured codelists the global recode must be specified manually

The structure is always: A new code is assigned to a set of old codes. So all the old codes are collapsed into the new code. A set can be either a list of individual codes, separated by a comma, or an interval indicated by a lower code, dash upper code. If the upper or lower code is not specified an open interval is assumed.

Examples:

For a variable with the categories 1,...,182 a possible recode is then:

```

1: - 49
2: 50 - 99
3: 100 - 149
4: 150 -

```

This implies that every code below 49 will be recoded into the new code 1, all codes between 50 and 99 will be the new code 2 etc.

For a variable with the categories 01 till 10 a possible recode is:

```

1: 01 , 02
2: 03 , 04
3: 05 - 07
4: 08 , 09 , 10

```

An important point is not to forget the colon (:) if it is forgotten, the recode will not work.

Recoding 3: 05,06,07 can be shortened to 3: 05-07.

And the two different schemes can be combined as well

```

1: 02 - 06, 09

```

is a valid recode as well.

## 5.5 *The apriori file*

---

The apriori file can be used to modify the characteristics of a cell before the secondary cell suppression routines are called. You can modify the following characteristics:

- Cell status
- Cost-function
- Apriori bounds
- Protection levels

The apriori file is a simple text-file that can be created with notepad and similar programs. The layout of the apriori file is simple. First the codes of the spanning variables are given, separated by a semicolon (“;”), then the code indicating the change requested and the depending on the code some additional parameters

Code	Parameters	Description
S	-	Status becomes safe
U	-	Status become (manually) unsafe
P	-	Status becomes protected

C	New cost value	<p>A low cost-value will make it more likely that this cell becomes a candidate for secondary suppressions. A high value will decrease this chance.</p> <p>This can be used to coordinate suppression patterns between successive years of a certain table</p>
AB	New apriori bound (two values for the lower and upper bound)	If external knowledge of a cell is available and the range of the cell is smaller than the normal external bounds of a cell the new value can be given here
PL	New protection level	If smaller or larger protection is required, this can be indicated here

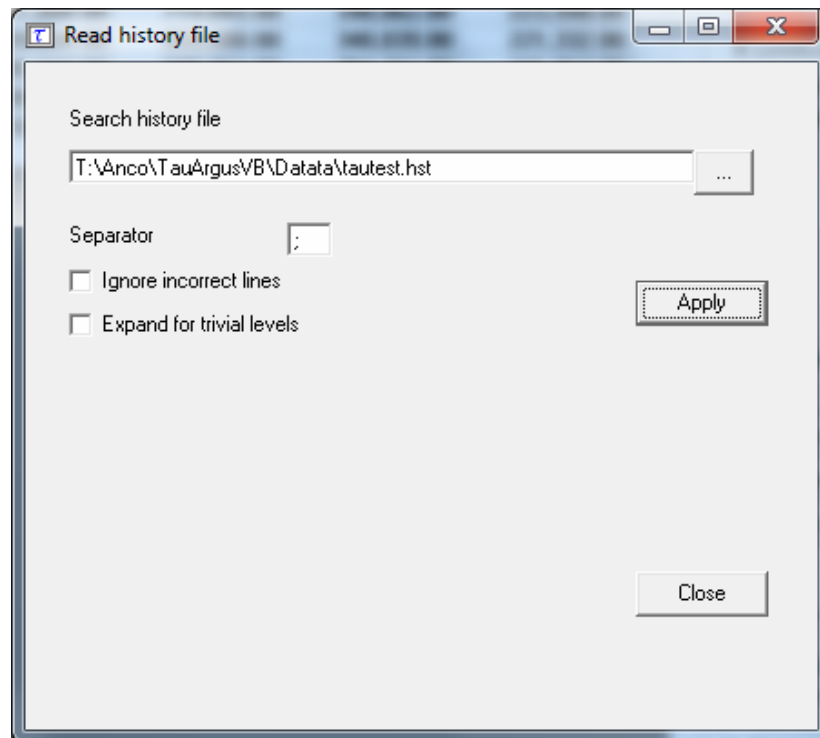
Note: changing the status of a cell is of course limited. E.g. a primary unsafe cell can not become protected, nor can a protected cell become unsafe.

The cost function must always be positive.

It is recommended to restrict the use of setting a cell status to protected. If you want to prevent that a cell will become a secondary suppression, give it a high cost value. If this cell is nevertheless suppressed, there will be a good reason for this. Putting a cell to protected, might lead to infeasible problems with all the consequences of that.

An example:

```
Nr, 4, u
Zd, 6, p
  5, 5, c, 1
Zd, 5, pl, 100, 200
```



The apriory file allows you to feed  $\tau$ -ARGUS a list of cells where the status of the standard rules can be overruled. E.g. a cell must be kept confidential or not for other reasons that just because of the sensitivity rules. By modifying the cost-function you can influence the selection of the secondaries. E.g. the cells suppressed last year can get a preference for the suppression this year by giving this cell a small value for the cost-function.

The option ‘trivial levels’ is important. Often in a table with hierarchies, some levels in a hierarchy break down in only one lower level. This implies that there are different cells in a table which are implicitly the same. Changing the status of one of them might lead to inconsistencies and serious problems. E.g. one of the two is unsafe and the other is protected, the solution is impossible. If you select the option ‘Expand for trivial levels’,  $\tau$ -ARGUS will always modify all cells that are the same if you modify one of them.

## 5.6 The Batch command file

---

$\tau$ -ARGUS has originally been designed as an interactive program. A complete menu-driven design guides you through all steps of the process. However a growing need for a batch version emerged after that. Since then  $\tau$ -ARGUS has been extended with a batch version. The batch commands are stored in a separate text-file. These commands can be executed from the command line or via the menu (File|Open Batch process. See section 4.2.3.

Alternatively the batch file can be used in a real batch environment as well. Just invoke  $\tau$ -ARGUS with the command

```
Taupath\TAUARGUS param1 param2 param3
```

where Taupath is the name of the directory where you installed  $\tau$ -ARGUS, param1

is the name of the above described file with batch commands. Param2 is optional, and is the name of the logfile. If omitted  $\tau$ -ARGUS will write a logbook in the file LOGBOOK.TXT in the temp-directory. See also section 5. Param3 is the parameter specifying the temp-directory. If omitted the default temp-directory will be used.

When using  $\tau$ -ARGUS interactively a batch file can be generated via the menu Output|Write Batch file. See section 4.5.4. But we advise you to inspect the results of this action before using this generated batch-file.

### Layout of the batch-file

A file can be written in a text editor and called from this command. Lines starting with “//” will be considered as comment and will be ignored.

The possible commands are shown here.

Command	Parameters
LOGBOOK	Name of the logbook file; If not specified the default logbook file will be used.
OPENMICRODATA	Data file name with microdata
OPENTABLEDATA	File name containing tabular data
OPENMETADATA	Metadata file name
SPECIFYTABLE	<p>“ExpVar1”“ExpVar2”“ExpVar3” RespVar ShadowVar Costvar, Lambda</p> <p>Shadow and cost variables are optional. If not specified then they equal the Response Variable.</p> <p>If the cost variable is specified either a numerical variable is specified or ‘-1’ is chosen for frequency or ‘-2’ for unity</p> <p>For lambda the default is 1.</p> <p>(See section 4.3.4 for the explanation for the use of lambda)</p>
CLEAR	Clears all and start a new session.
SAFETYRULE	<p>This command is used for primary suppression.</p> <p>All these parameters are described in the section Specify tables, 4.3.4</p> <p>A set of safety rule specifications separated by a “ ”</p> <p>Each safety spec starts with "P", "NK" "ZERO", "FREQ", "REQ", "WGT", "MIS" or "MAN" and between brackets the parameters</p> <p><b>P:</b> (p,n) with the n optional. (default = 1). So (20,3)→. A p% rule with p=20% and n=3</p> <p><b>NK:</b> (n,k). A n,k-dominance rule with</p> <p><b>ZERO:</b> (ZeroSafetyRange)</p> <p><b>FREQ:</b>(MinFreq, FrequencySafetyRange)</p> <p><b>REQ:</b> (Percent1, Percent2, SafetyMargin)</p> <p>All rules can appear several times,</p> <p>The first two <b>P</b>, <b>NK</b> are for the individual level; the following two for the holding level,</p> <p>The first <b>FREQ</b> and <b>REQ</b> are at the individual level the second one is for the holding.</p> <p><b>ZERO</b>, the zero safety range parameter, can be given only once for each safety rule.</p> <p><b>MIS:</b> 0 = cells with a missing code are unsafe if the safety-rules are violated; 1 = these cells are always safe. (Default = 0.)</p>

	<p><b>WGT:</b> 0 no weights are used, 1 = apply weights for computing the tables and in the safety rules Default = 0</p> <p><b>MAN:</b> (Manual safety margin). This margin is used e.g. of a table with only the status is read, or if via the apriori option a cell is set to manually unsafe. The default value = 20.</p>	
READMICRODATA	Just reads the microdata file and calculates the table(s), no parameters are required	
READTABLE	Just reads the tabular inputfile; If the only parameter = 1 the “compute missing totals” procedure will be used. Default = do not compute this.	
APRIORI	This reads an a-priori file The parameters are: Filename, Table number and the separator “Filename”, TabNo, Separator	
SUPPRESS	<p>This command applies the secondary suppression. The possible options are: GH: Hypercube MOD: Modular OPT: Optimal NET: Network RND: Controlled rounding The parameters are a few parameters between brackets; The first parameter is always the table number.</p> <p><b>GH</b>(TabNo, A priori Bounds Percentage, ModelSize) Model size 0 = normal, 1 indicates the large model. <b>MOD</b>(TabNo, MaxTimePerSubtable) <b>OPT</b>(TabNo, MaxComputingTime) <b>NET</b>(TabNo) <b>RND</b>(TabNo, RoundingBase, Steps, 1, Time, Partitions, StopRule) - Steps: number of steps allowed, normally 0 (default) - Time: Max computing time (20 = default) - Partitions: 0, 1 (0 = no partitioning (default), 1 = apply the partitioning procedure) - StopRule: 1 = Rapid only, 2 = First feasible solution, 3 = optimal solution (3 = default) Note that the fourth parameter is always 1. It is there only for future extensions.</p>	
SOLVER	Indicate whether you will be using CPLEX or XPress. Only needed when the type of solver has not yet been specified on the computer during a previous interactive session of $\tau$ -ARGUS. The only parameter allowed is CPLEX or XPRESS.	
WRITETABLE	P1 (OutputType)	P2 (parameter)

	1. CVS-file 2. CSV file for pivot table 3. Code, value file  4. SBS-output format 5. Intermediate file	Not used 1 = AddStatus 0 = not 1 = AddStatus 2 = suppress empty cells 3 = both options 0 = none Not used 0 = Status only 1 = also Top-n scores
GOINTERACTIVE	Should be the last command. If omitted the program will stop. If specified $\tau$ -ARGUS will go on as an interactive program. This command works only if the batch-file is invoked from the menu. If the batch-version is started from the command-line, this option will be ignored.	

A typical batch file would look like this: (note that everything after a // will be treated as comment)

#### A batch file using micro data

```
//datafile
<OPENMICRODATA> "C:\Program Files\TauARGUS\datatau_testW.asc"
//metafile
<OPENMETADATA> "C:\Program Files\TauARGUS\datatau_testW.rda"
//Exp|resp|shadow|cost -1= unit -2 = freq
<SPECIFYTABLE> "Size" "region" | "var2" | "var3" | "var3"
<SAFETYRULE> P(20,3) | FREQ(3,30) | ZERO(20)
<SPECIFYTABLE> "Size" "Year" | "var2" | "var3" | "var3"
<SAFETYRULE> NK(3,70) | FREQ(3,30) | ZERO(20)
<READMICRODATA>
<SUPPRESS> GH(1,75)
<WRITETABLE> (1,1,1,"C:\Program Files\TauARGUS\datax1.csv")
<SUPPRESS> GH(2,75)
<WRITETABLE> (2,2,1,"C:\Program Files\TauARGUS\datay11.csv")
<SUPPRESS> MOD(1)
<WRITETABLE> (1,3,0,"C:\Program Files\TauARGUS\datax20.txt")
<SUPPRESS> MOD(2)
<WRITETABLE> (2,4,0,"C:\Program Files\TauARGUS\datay20.tab")
<SUPPRESS> OPT(1,5)
<WRITETABLE> (1,1,1,"C:\Program Files\TauARGUS\datax3.csv")
<GOINTERACTIVE>
```

#### A batch file using tabular data

```
<OPENTABLEDATA> "E:\TauArgusVB\Datata\ \Nace3Size.tab"
<OPENMETADATA> E:\TauArgusVB\Datata\Nace3Size.RDA"
<SPECIFYTABLE> "IndustryCode" "Size" | "Var2" | "Var2" | "Var2"
//<SAFETYRULE>
<READTABLE>
<SUPPRESS> MOD(1)
<WRITETABLE> (1,3,3,"E:\TauArgusVB\Datata\Nace3SizeSafe.txt")
<GOINTERACTIVE>
```

In the above example the <SAFETYRULE> command was disabled as in this example it is assumed that that table already contained the status of each cell. However if the tabular input contains more information (frequency, TopN) the safety rule command could easily be used here as well.

If more than one table has to be processed, the >CLEAR> command could make a new start in a session. However it is also possible to use the standard Windows batch-language to combine several commands in one run.

```
start /wait taupath\tauargus.exe Tab1.arb Log1.txt
start /wait taupath\tauargus.exe Tab2.arb Log2.txt
start /wait taupath\tauargus.exe Tab3.arb Log3.txt
```

## 5.7 Log file

---

$\tau$ -ARGUS will write a log-file. This describes among others the commands used during the runs of  $\tau$ -ARGUS. It gives a log of the use of  $\tau$ -ARGUS. Especially for the batch process this file could give some information about the progress of the process. Below is given a small example. Please note that new information is always added to this file. So from time to time the user should erase this file to clean his computer.

By default the logfile is the file LOGBOOK.TXT in the temp-directory. In the options window the name of the logfile can be changed for the remainder of the current session.

The Temp directory is normally something like

C:\Documents and Settings\USER\Local Settings\Temp

where USER is the name of the current user. Of course in specific circumstances the network administrator might have chosen a different location.

When running in batch-mode it is possible to change the name of the log-file with a batch command in the batch-file or as the second parameter on the commandline. See section 4.2.3

```
08-feb-2007 15:45:19 : Start preparing for tabulation
08-feb-2007 15:45:19 : Table 1: Size x Region | Var2
08-feb-2007 15:45:19 : Start explore file: C:\Program
                        Files\TauARGUS\data\tau_testW.asc
08-feb-2007 15:45:20 : Start computing tables
08-feb-2007 15:45:20 : Compute tables completed
08-feb-2007 15:45:25 : Start Modular optimisation
08-feb-2007 15:45:26 : MODULAR finished with a problem
08-feb-2007 15:45:39 : Start Modular optimisation
08-feb-2007 15:45:40 : MODULAR finished successfully
08-feb-2007 15:45:50 : Table no: 1 has been saved
                        FileName: C:\Program Files\TauARGUS\data\x.sbs
                        as Save table in SBS-format
```

## 6 INDEX

---

A		
Apriory .....	95, 108	
B		
Batch .....	113	
C		
Cplex .....	7, 12, 19, 31, 45, 60, 79, 97, 111	
D		
dominance rule ..	10, 11, 12, 37, 60, 62, 68, 69, 96, 110	
G		
global recoding .....	11, 12, 30, 70, 74, 76, 107	
H		
Holding .....	38, 56, 64, 72, 89, 101, 103, 110	
M		
magnitude table .....	10	
	meta data .....	100
	missing .....	33, 34, 56, 63, 69, 75, 77, 100, 102, 110
N		
negative values .....	17, 21	
P		
p % rule .....	10, 11, 12, 37, 60, 62, 65, 68, 69, 110	
R		
Request rule .....	38, 40, 56, 63, 101, 102	
S		
Sensitive cell .....	10, 12, 27, 78	
Singleton .....	20, 28	
X		
Xpress .....	6, 7, 12, 19, 31, 45, 60, 79, 81, 111	