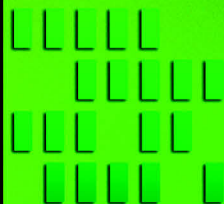


# Statistical Disclosure Control



Anco Hundepool • Josep Domingo-Ferrer  
Luisa Franconi • Sarah Giessing • Eric Schulte Nordholt  
Keith Spicer • Peter-Paul de Wolf









# Statistical Disclosure Control



## WILEY SERIES IN SURVEY METHODOLOGY

Established in Part by WALTER A. SHEWHART and SAMUEL S. WILKS

Editors: Mick P. Couper, Graham Kalton, Lars Lyberg, J.N.K. Rao,  
Norbert Schwarz, Christopher Skinner

A complete list of the titles in this series appears at the end of this volume.



# Statistical Disclosure Control

**Anco Hundepool**

*Statistics Netherlands, The Netherlands*

**Josep Domingo-Ferrer**

*Universitat Rovira i Virgili, Catalonia, Spain*

**Luisa Franconi**

*Italian National Institute of Statistics, Italy*

**Sarah Giessing**

*Federal Statistical Office of Germany, Germany*

**Eric Schulte Nordholt**

*Statistics Netherlands, The Netherlands*

**Keith Spicer**

*Office for National Statistics, UK*

**Peter-Paul de Wolf**

*Statistics Netherlands, The Netherlands*



A John Wiley & Sons, Ltd., Publication



This edition first published 2012  
© 2012 John Wiley & Sons, Ltd

*Registered office*

John Wiley & Sons Ltd, The Atrium, Southern Gate, Chichester, West Sussex, PO19 8SQ, United Kingdom

For details of our global editorial offices, for customer services and for information about how to apply for permission to reuse the copyright material in this book please see our website at [www.wiley.com](http://www.wiley.com).

The right of the author to be identified as the author of this work has been asserted in accordance with the Copyright, Designs and Patents Act 1988.

All rights reserved. No part of this publication may be reproduced, stored in a retrieval system, or transmitted, in any form or by any means, electronic, mechanical, photocopying, recording or otherwise, except as permitted by the UK Copyright, Designs and Patents Act 1988, without the prior permission of the publisher.

Wiley also publishes its books in a variety of electronic formats. Some content that appears in print may not be available in electronic books.

Designations used by companies to distinguish their products are often claimed as trademarks. All brand names and product names used in this book are trade names, service marks, trademarks or registered trademarks of their respective owners. The publisher is not associated with any product or vendor mentioned in this book. This publication is designed to provide accurate and authoritative information in regard to the subject matter covered. It is sold on the understanding that the publisher is not engaged in rendering professional services. If professional advice or other expert assistance is required, the services of a competent professional should be sought.

*Library of Congress Cataloging-in-Publication Data*

Hundepool, Anco.

Statistical disclosure control / Anco Hundepool [and six others].  
pages cm. – (Wiley series in survey methodology)

Includes bibliographical references and index.

ISBN 978-1-119-97815-2

1. Confidential communications—Statistical services. I. Title.

HA34.H86 2012

352.7/52384—dc23

2012015785

A catalogue record for this book is available from the British Library.

ISBN 978-1-119-97815-2

Typeset in 10/12pt Times by Aptara Inc., New Delhi, India



# Contents

<b>Preface</b>	<b>xi</b>
<b>Acknowledgements</b>	<b>xv</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Concepts and definitions	2
1.1.1 Disclosure	2
1.1.2 Statistical disclosure control	3
1.1.3 Tabular data	3
1.1.4 Microdata	3
1.1.5 Risk and utility	4
1.2 An approach to Statistical Disclosure Control	7
1.2.1 Why is confidentiality protection needed?	7
1.2.2 What are the key characteristics and uses of the data?	8
1.2.3 What disclosure risks need to be protected against?	8
1.2.4 Disclosure control methods	8
1.2.5 Implementation	9
1.3 The chapters of the handbook	9
<b>2 Ethics, principles, guidelines and regulations – a general background</b>	<b>10</b>
2.1 Introduction	10
2.2 Ethical codes and the new ISI code	11
2.2.1 ISI Declaration on Professional Ethics	11
2.2.2 New ISI Declaration on Professional Ethics	12
2.2.3 European Statistics Code of Practice	15
2.3 UNECE principles and guidelines	16
2.3.1 UNECE Principles and Guidelines on Confidentiality Aspects of Data Integration	18
2.3.2 Future activities on the UNECE principles and guidelines	19
2.4 Laws	19
2.4.1 Committee on Statistical Confidentiality	20
2.4.2 European Statistical System Committee	20



<b>3</b>	<b>Microdata</b>	<b>23</b>
3.1	Introduction	23
3.2	Microdata concepts	24
3.2.1	Stage 1: Assess need for confidentiality protection	24
3.2.2	Stage 2: Key characteristics and use of microdata	27
3.2.3	Stage 3: Disclosure risk	30
3.2.4	Stage 4: Disclosure control methods	32
3.2.5	Stage 5: Implementation	34
3.3	Definitions of disclosure	36
3.3.1	Definitions of disclosure scenarios	37
3.4	Definitions of disclosure risk	38
3.4.1	Disclosure risk for categorical quasi-identifiers	39
3.4.2	Notation and assumptions	40
3.4.3	Disclosure risk for continuous quasi-identifiers	41
3.5	Estimating re-identification risk	43
3.5.1	Individual risk based on the sample: Threshold rule	44
3.5.2	Estimating individual risk using sampling weights	44
3.5.3	Estimating individual risk by Poisson model	47
3.5.4	Further models that borrow information from other sources	48
3.5.5	Estimating per record risk via heuristics	49
3.5.6	Assessing risk via record linkage	50
3.6	Non-perturbative microdata masking	51
3.6.1	Sampling	51
3.6.2	Global recoding	52
3.6.3	Top and bottom coding	53
3.6.4	Local suppression	53
3.7	Perturbative microdata masking	53
3.7.1	Additive noise masking	54
3.7.2	Multiplicative noise masking	57
3.7.3	Microaggregation	60
3.7.4	Data swapping and rank swapping	72
3.7.5	Data shuffling	73
3.7.6	Rounding	73
3.7.7	Re-sampling	74
3.7.8	PRAM	74
3.7.9	MASSC	78
3.8	Synthetic and hybrid data	78
3.8.1	Fully synthetic data	79
3.8.2	Partially synthetic data	84
3.8.3	Hybrid data	86
3.8.4	Pros and cons of synthetic and hybrid data	98
3.9	Information loss in microdata	100
3.9.1	Information loss measures for continuous data	101
3.9.2	Information loss measures for categorical data	108



3.10	Release of multiple files from the same microdata set	110
3.11	Software	111
3.11.1	$\mu$ -ARGUS	111
3.11.2	<i>sdcMicro</i>	113
3.11.3	<i>IVEware</i>	115
3.12	Case studies	116
3.12.1	Microdata files at Statistics Netherlands	116
3.12.2	The European Labour Force Survey microdata for research purposes	118
3.12.3	The European Structure of Earnings Survey microdata for research purposes	121
3.12.4	NHIS-linked mortality data public use file, USA	128
3.12.5	Other real case instances	130
<b>4</b>	<b>Magnitude tabular data</b>	<b>131</b>
4.1	Introduction	131
4.1.1	Magnitude tabular data: Basic terminology	131
4.1.2	Complex tabular data structures: Hierarchical and linked tables	132
4.1.3	Risk concepts	134
4.1.4	Protection concepts	137
4.1.5	Information loss concepts	137
4.1.6	Implementation: Software, guidelines and case study	138
4.2	Disclosure risk assessment I: Primary sensitive cells	138
4.2.1	Intruder scenarios	138
4.2.2	Sensitivity rules	140
4.3	Disclosure risk assessment II: Secondary risk assessment	152
4.3.1	Feasibility interval	152
4.3.2	Protection level	154
4.3.3	Singleton and multi cell disclosure	155
4.3.4	Risk models for hierarchical and linked tables	155
4.4	Non-perturbative protection methods	157
4.4.1	Global recoding	157
4.4.2	The concept of cell suppression	157
4.4.3	Algorithms for secondary cell suppression	158
4.4.4	Secondary cell suppression in hierarchical and linked tables	161
4.5	Perturbative protection methods	163
4.5.1	A pre-tabular method: Multiplicative noise	165
4.5.2	A post-tabular method: Controlled tabular adjustment	165
4.6	Information loss measures for tabular data	166
4.6.1	Cell costs for cell suppression	166
4.6.2	Cell costs for CTA	167
4.6.3	Information loss measures to evaluate the outcome of table protection	167



4.7	Software for tabular data protection	168
4.7.1	Empirical comparison of cell suppression algorithms	169
4.8	Guidelines: Setting up an efficient table model systematically	173
4.8.1	Defining spanning variables	174
4.8.2	Response variables and mapping rules	175
4.9	Case studies	178
4.9.1	Response variables and mapping rules of the case study	178
4.9.2	Spanning variables of the case study	179
4.9.3	Analysing the tables of the case study	179
4.9.4	Software issues of the case study	181
<b>5</b>	<b>Frequency tables</b>	<b>183</b>
5.1	Introduction	183
5.2	Disclosure risks	184
5.2.1	Individual attribute disclosure	185
5.2.2	Group attribute disclosure	186
5.2.3	Disclosure by differencing	187
5.2.4	Perception of disclosure risk	190
5.3	Methods	191
5.3.1	Pre-tabular	191
5.3.2	Table re-design	192
5.3.3	Post-tabular	193
5.4	Post-tabular methods	193
5.4.1	Cell suppression	193
5.4.2	ABS cell perturbation	193
5.4.3	Rounding	194
5.5	Information loss	199
5.6	Software	201
5.6.1	Introduction	201
5.6.2	Optimal, first feasible and RAPID solutions	202
5.6.3	Protection provided by controlled rounding	203
5.7	Case studies	204
5.7.1	UK Census	204
5.7.2	Australian and New Zealand Censuses	205
<b>6</b>	<b>Data access issues</b>	<b>208</b>
6.1	Introduction	208
6.2	Research data centres	209
6.3	Remote execution	209
6.4	Remote access	210
6.5	Licensing	211
6.6	Guidelines on output checking	211
6.6.1	Introduction	211
6.6.2	General approach	212
6.6.3	Rules for output checking	215



6.6.4	Organisational/procedural aspects of output checking	224
6.6.5	Researcher training	233
6.7	Additional issues concerning data access	236
6.7.1	Examples of disclaimers	236
6.7.2	Output description	236
6.8	Case studies	237
6.8.1	The US Census Bureau Microdata Analysis System	237
6.8.2	Remote access at Statistics Netherlands	239
<b>Glossary</b>		<b>243</b>
<b>References</b>		<b>261</b>
<b>Author index</b>		<b>279</b>
<b>Subject index</b>		<b>282</b>







# Preface

In the last 20 years, work in official statistics has changed drastically due to new possibilities of information technology. While in the old days, the output of statistical offices mainly consisted of a set of tables, limited by the mere size of paper publications, nowadays new information technology makes it possible to publish much more detailed tabular data. Moreover, on the input side, the use of computers has increased the amount and detail of the data collected. The use of external registers, now available for a number of statistical offices, has led to an enormous amount of very detailed microdata. Of course, this is a very positive development as the mission of the statistical offices is to describe the society with as much detail as possible. For example, this makes it possible for policy makers to make well-informed decisions.

On the other hand, publishing a vast amount of detailed information has the risk of disclosing sensitive information both on individual persons as well as on economic entities. Privacy issues have become more and more of a concern for people. The growth of the internet has made people aware of consequences of the concept ‘big brother is watching you’. Statistical Disclosure Control (SDC) is thus a rapidly evolving field: disclosure control thinking has to keep pace with increase in computing power, developments in matching software and the proliferation of public and private databases. Statistical offices need to find the right balance between the need to inform society as well as possible, on the one hand, and the need to safeguard the privacy of the respondents on the other.

Though the main audience we target consists of employees of statistical offices that need to apply the methods discussed in this book, it may be of interest to other readers as well. For example, data-archiving institutes and health data-collecting institutes have to deal with similar problems concerning confidentiality. Moreover, users of statistical output should be aware of the reasoning and methodology behind statistical disclosure control. As the issue of confidentiality grows in society, the new generation of graduates at universities should be exposed to the concepts and the methods described in this book. Last but not the least, computer scientists, database experts, data miners, practitioners dealing with medical data, etc., may find the contents of this book useful: indeed, the methods used for data anonymisation and privacy-preserving data mining are in essence the same as used in SDC.

There are several reasons to take privacy protection seriously. Firstly, there are legal frameworks that regulate what is allowed and what is not allowed with regard to publication of private information. However, there are other reasons for statistical



offices to take confidentiality protection seriously as well. For example, offices need to maintain good relationships with respondents. After all, they are an essential source of information on which statistical offices build their statistics. Without respondents, there are no statistics. The respondents must be able to trust that their private and often sensitive information is safe in the hands of statistical offices. In Chapter 2, we give an overview of the major international guidelines, regulations and principles. Eric Schulte Nordholt is the author mainly responsible for this chapter.

Although disclosure control of tabular data is the oldest part of SDC, we will first focus on confidentiality issues with microdata. Josep Domingo-Ferrer and Luisa Franconi are the main contributors here. Since powerful computers and powerful statistical software are available to standard researchers, the analysis of statistical information has changed from analysing tabular data towards analysing individual data (microdata). As statistical institutes nowadays have very large databases available to produce their statistical output, these databases could ideally have a second life as the basis for various statistical research projects, e.g. at universities. Indeed, it would be a waste of time and money if researchers had to collect the information themselves again; collaboration is much more efficient. However, before sensitive statistical databases can be made available to universities for research, confidentiality must be guaranteed. Several methods exist here and tools have been developed to implement this. Various aspects of the protection of microdata are described in Chapter 3.

There is a much longer tradition of publishing tabular data. In cases where data are aggregated into tables, it is a misunderstanding to believe that there would be no privacy issues. Indeed, tabular data can disclose individual information as well. For example, if a cell in a table has only one contributor, there is an obvious risk of disclosure. It was as late as the 1970s that the first rules were proposed to assess whether tabular data could be disclosive. Since then, considerable developments have taken place. Sarah Giessing and Peter-Paul de Wolf are the main contributors to Chapter 4 on magnitude tables. This chapter deals mainly with two issues. The first step is to decide which cells in a table are disclosive (the easier part) and the second, how to protect disclosive tables adequately (the more difficult part).

Magnitude tables have attracted the most attention when solving disclosure control issues, but frequency tables can also disclose individual sensitive data. There is still a lot of work to be done to make people aware of the problems here. The approaches in the case of frequency tables are very different from those for magnitude tables. Keith Spicer has contributed to Chapter 5 that deals with frequency tables.

The protection of microdata is the subject of Chapter 3. But when data cannot be adequately protected without adversely affecting the ability to answer key research questions, alternative ways to provide access to the data must be sought. A discussion on this can be found in Chapter 6. The first step is to make moderately protected microdata. This data can be made available to serious, trustworthy researchers. Such moderate protection together with a strict contract may be enough to meet both the needs of the researchers and the need for privacy protection. When this is not an option, offices have opened special secure environments on their premises, where researchers can analyse the microdata, while the data remains under the control of the institute. In case researchers want to take results home, the results have to be



checked on disclosure before they are released. For that, a series of guidelines are proposed.

When different research groups cooperate but work in different institutes, a common understanding of the terminology is needed. To facilitate this, we have added a glossary of statistical terms used in SDC. The glossary is based on a glossary proposed in 2005 at the UNECE Work Session on Statistical Disclosure Control.

Anco Hundepool  
Josep Domingo-Ferrer  
Luisa Franconi  
Sarah Giessing  
Eric Schulte Nordholt  
Keith Spicer  
Peter-Paul de Wolf







# Acknowledgements

This book is the outcome of work on Statistical Disclosure Control (SDC) that has been carried out in Europe over the past years. Smaller teams at Statistical Offices and universities have cooperated intensively. This very fruitful cooperation has been supported financially by several European projects. It started with the SDC project from 1996 to 1998 in the Fourth Framework Programme of the EU, followed by the CASC project (2000–2003) in the Fifth Framework Programme of the EU. Eurostat, the European statistical office, has since then supported our work via various projects. The authors are grateful for the highly appreciated cooperation with their colleagues from both the statistical offices and the universities.







# 1

## Introduction

National Statistical Institutes (NSIs) publish a wide range of trusted, high-quality statistical outputs. To achieve their objective of supplying society with rich statistical information, these outputs are as detailed as possible. However, this objective conflicts with the obligation NSIs have to protect the confidentiality of the information provided by the respondents. Statistical Disclosure Control (SDC) or Statistical Disclosure Limitation (SDL) seeks to protect statistical data in such a way that they can be released without giving away confidential information that can be linked to specific individuals or entities.

In addition to official statistics, there are several other areas of application of SDC techniques, including:

- *Health information.* This is one of the most sensitive areas regarding privacy.
- *E-commerce.* Electronic commerce results in the automated collection of large amounts of consumer data. This wealth of information is very useful to companies, which are often interested in sharing it with their subsidiaries or partners. Such consumer information transfer is subject to strict regulations.

This handbook aims to provide technical guidance on SDC for NSIs on how to approach this problem of balancing the need to provide users with statistical outputs and the need to protect the confidentiality of respondents. SDC should be combined with other tools (administrative, legal and IT) in order to define a proper data-dissemination strategy based on a risk-management approach.

A data-dissemination strategy offers many different statistical outputs covering a range of different topics for many types of users. Different outputs require different approaches to SDC and different mixtures of tools.



- *Tabular data protection.* Tabular data protection is the oldest and best established part of SDC, because tabular data have been the traditional output of NSIs. The goal here is to publish static aggregate information, i.e. tables, in such a way that no confidential information on specific individuals among those to which the table refers can be inferred. In the majority of cases, confidentiality protection is achieved only by statistical tools due to the absence of legal and IT restrictions.
- *Dynamic databases.* The scenario here is a database to which the user can submit statistical queries (sums, averages, etc.). The aggregate information obtained as a result of successive queries should not allow him to infer information on specific individuals. The mixture of tools here may vary according to the setting and the data provided.
- *Microdata protection.* In recent years, with the widespread use of personal computers and the public demand for data, microdata (that is data sets containing for each respondent the scores on a number of variables) are being disseminated to users in universities, research institutes and interest groups (Trewin 2007). Microdata protection is the youngest sub-discipline and has experienced continuous evolution in the last years. If microdata are freely disseminated then SDL methods will be very severe to protect confidentiality of respondents; if, on the other hand, legal restrictions are in place (such as Commission Regulation 831/2002; see Section 2.4.2) a different amount of information may be released.
- *Protection of output of statistical analyses.* The need to allow access to microdata has encouraged the creation of Microdata Laboratories (Safe Centres) in many NSIs. Due to an IT-protected environment, legal and administrative restrictions users may analyse detailed microdata. Checking the output of these analyses to avoid confidentiality breaches is another field which is developing in SDC research. This handbook provides guidance on how to protect confidentiality for all of these types of output using statistical methods.

This first chapter provides a brief introduction to some of the key concepts and definitions involved with this field of work as well as a high-level overview of how to approach problems associated with confidentiality.

## 1.1 Concepts and definitions

### 1.1.1 Disclosure

A disclosure occurs when a person or an organisation recognises or learns something that they did not know already about another person or organisation, via released data. There are two types of disclosure risk: (1) identity disclosure and (2) attribute disclosure. Identity disclosure occurs with the association of a respondents' identity with a disseminated data record containing confidential information (see Duncan *et al.* 2001).



Attribute disclosure occurs with the association of either an attribute value in the disseminated data or an estimated attribute value based on the disseminated data with the respondent (see Duncan *et al.* 2001).

Some NSIs may also be concerned with the perception of disclosure risk. For example, if small values appear in tabular output users may perceive that no (or insufficient) protection has been applied. More emphasis has been placed on this type of disclosure risk in recent years because of declining response rates and decreasing data quality.

### 1.1.2 Statistical disclosure control

SDC techniques can be defined as the set of methods to reduce the risk of disclosing information on individuals, businesses or other organisations. SDC methods minimise the risk of disclosure to an acceptable level while releasing as much information as possible. There are two types of SDC methods; perturbative and non-perturbative methods. Perturbative methods falsify the data before publication by introducing an element of error purposely for confidentiality reasons. Non-perturbative methods reduce the amount of information released by suppression or aggregation of data. A wide range of different SDC methods are available for different types of outputs.

### 1.1.3 Tabular data

There are two types of tabular output:

1. *Magnitude tables.* In a magnitude table, each cell value represents the sum of a particular response, across all respondents that belong to that cell. Magnitude tables are commonly used for business or economic data providing, for example turnover of all businesses of a particular industry within a region.
2. *Frequency tables.* In a frequency table, each cell value represents the number of respondents that fall into that cell. Frequency tables are commonly used for Census or social data providing, for example the number of individuals within a region who are unemployed.

### 1.1.4 Microdata

A microdata set  $\mathbf{V}$  can be viewed as a file with  $n$  records, where each record contains  $m$  variables (also called *attributes*) on an individual respondent, who can be a person or an organisation (e.g. a company). Microdata are the form from which all other data outputs are derived and they are the primary form that data are stored in. While in the past, NSIs simply derived outputs of other forms, more and more, microdata are becoming a key output by themselves.

Depending on their sensitivity, the variables in an original unprotected microdata set can be classified into four categories which are not necessarily disjoint:

1. *Identifiers.* These are variables that *unambiguously* identify the respondent. Examples are passport number, social security number, full name, etc. Since



the objective of SDC is to prevent confidential information from being linked to specific respondents, SDC normally assumes that identifiers in  $\mathbf{V}$  have been removed/encrypted in a pre-processing step.

2. *Quasi-identifiers or key variables.* Borrowing the definition from Dalenius (1986), Samarati (2001), a quasi-identifier is a set of variables in  $\mathbf{V}$  that, in combination, can be linked with external information to re-identify (some of) the respondents to whom (some of) the records in  $\mathbf{V}$  refer. Unlike identifiers, quasi-identifiers cannot be removed from  $\mathbf{V}$ . The reason is that any variable in  $\mathbf{V}$  potentially belongs to a quasi-identifier (depending on the external data sources available to the user of  $\mathbf{V}$ ). Thus, one would need to remove all variables (!) to make sure that the data set no longer contains quasi-identifiers.
3. *Confidential outcome variables.* These are variables which contain sensitive information on the respondent. Examples are salary, religion, political affiliation, health condition, etc.
4. *Non-confidential outcome variables.* Those variables which contain non-sensitive information on the respondent. Examples are town and country of residence, etc. Note that variables of this kind cannot be neglected when protecting a data set, because they can be part of a quasi-identifier. For instance, if 'Job' and 'Town of residence' can be considered non-confidential outcome variables, but their combination can be a quasi-identifier, because everyone knows who is the doctor in a small village.

Depending on their data type, the variables in a microdata set can be classified as:

- *Continuous.* A variable is considered continuous if it is numerical and arithmetical operations can be performed on it. Examples are income and age.
- *Categorical.* A variable is considered categorical when it takes values over a finite set and standard arithmetical operations do not make sense. Two main types of categorical variables can be distinguished:
  - *Ordinal.* An ordinal variable takes values in an ordered range of categories. Thus, the  $\leq$ , max and min operators are meaningful with ordinal data. The instruction level and the political preferences (left-right) are examples of ordinal variables.
  - *Nominal.* A nominal variable takes values in an unordered range of categories. The only possible operator is comparison for equality. The eye color and the address of an individual are examples of nominal variables.

### 1.1.5 Risk and utility

SDC seeks to optimise the trade-off between the disclosure risk and the utility of the protected released data. Yet, SDC is a discipline which was born from daily statistical practice and any theory trying to bestow a unified scientific standing to it should be flexible enough to deal with the risk-utility trade-off in a rather vast number of situations (different data structures, different contexts of previously released information



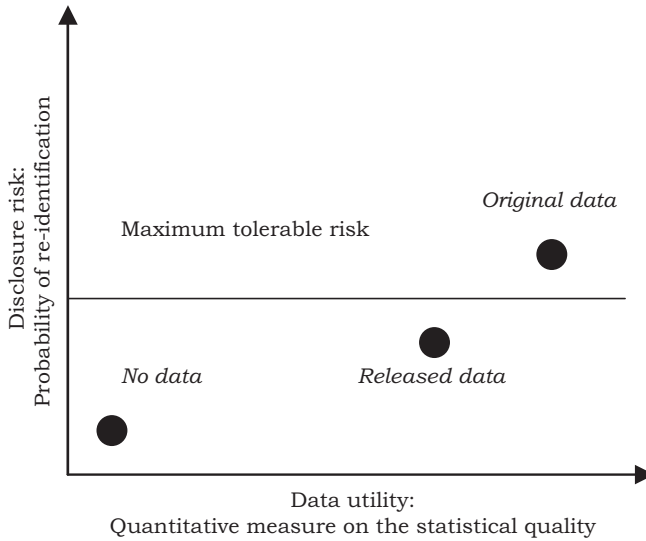


Figure 1.1 R-U confidentiality map.

and intruders' side knowledge, etc.). How to accurately measure risk and utility, how to distinguish legitimate users from intruders and how far to implement transparency are some of the open challenges SDC faces today, most of which are insightfully highlighted in Cox *et al.* (2011) and discussed in Domingo-Ferrer (2011).

We next review briefly some proposed models to deal with the risk-utility tension. One of them is merely descriptive (R-U maps depicting risk and utility), while the other two ( $k$ -anonymity and differential privacy) adopt a minimax approach: subject to a minimum guaranteed privacy, utility is to be maximised.

### R-U maps

NSIs should aim to determine optimal SDC methods and solutions that minimise disclosure risk while maximising the utility of the data. Figure 1.1 contains an R-U confidentiality map developed by Duncan *et al.* (2001), where R is a quantitative measure of disclosure risk and U is a quantitative measure of data utility.

In the lower left hand quadrant of the graph, low disclosure risk is achieved but also low utility, where no data is released at all. In the upper right hand quadrant of the graph, high disclosure risk is realised but also high utility, represented by the point where the original data is released. The NSI must set the maximum tolerable disclosure risk based on standards, policies and guidelines. The goal in this disclosure risk, data utility decision problem is then to find the balance in maintaining the utility of the data but reducing the risk below the maximum tolerable risk threshold.

### $k$ -anonymity

$k$ -anonymity is a concept that was proposed by Samarati (2001); Samarati and Sweeney (1998); Sweeney (2002a, 2002b) as a different approach to face the conflict between information loss and disclosure risk.



A data set is said to satisfy  $k$ -anonymity for  $k > 1$  if, for each combination of values of quasi-identifiers (e.g. name, address, age, gender, etc.), at least  $k$  records exist in the data set sharing that combination.

Note that, if a protected data set  $V'$  satisfies  $k$ -anonymity, an intruder trying to link  $V'$  with an external non-anonymous data source will find at least  $k$  records in  $V'$  that match any value of the quasi-identifier the intruder uses for record linkage. Thus re-identification, i.e. mapping a record in  $V'$  to a non-anonymous record in the external data source, is not possible; the best the intruder can hope for is to map groups of  $k$  records in  $V'$  to each non-anonymous external record.

If, for a given  $k$ ,  $k$ -anonymity is assumed to be enough protection, one can concentrate on minimising information loss with the only constraint that  $k$ -anonymity should be satisfied. This is a clean way of solving the tension between data protection and data utility.

In Samarati (2001) and Sweeney (2002b) the approach suggested to reach  $k$ -anonymity is to combine generalisation and local suppression; in Domingo-Ferrer and Torra (2005), the use of microaggregation was proposed as an alternative; see Section 3.6 for a description of those methods.

$k$ -anonymity has been criticised as being necessary but sometimes not sufficient to guarantee privacy. For example, if a group of  $k$  records sharing a quasi-identifier combination also shares the same value for a confidential variable (e.g. AIDS='YES'), it suffices for an intruder to know that John Smith's record is one of  $k$  records in the group to learn that John Smith suffers from AIDS. To remedy this, other privacy properties evolving from  $k$ -anonymity have been proposed; see Domingo-Ferrer and Torra (2008) for a survey of critiques and evolved models.

### Differential privacy

In general, SDC methods take the data to be published as input and modify them with the aim of reducing the risk of disclosure by removing/changing combinations of variables likely to be re-identifiable. Differential privacy (Dwork 2006, 2011) tackles privacy preservation from a different perspective. On one side, differential privacy is not based on the precise understanding that certain combinations of variables might be re-identifiable. Instead, it seeks to guarantee, in a probabilistic sense, that after the addition of a new record to a database any information extracted from the database will remain close to what it was before. On the other side, differential privacy does not focus on a specific database. It seeks to protect all the possible databases that may arise from record addition.

The following formal definition of differential privacy can be found in Dwork (2006). A randomised function  $\kappa$  gives  $\varepsilon$ -differential privacy if, for all data sets  $D$  and  $D'$  differing in at most one row, and all  $S \subset \text{Range}(\kappa)$

$$P[\kappa(D) \in S] \leq e^\varepsilon P[\kappa(D') \in S]. \quad (1.1)$$

Similarly to what happened with  $k$ -anonymity, the idea is that, among the randomisations offering differential privacy for a certain parameter value  $\varepsilon$ , one should choose the one causing minimal information loss.

Computationally, differential privacy is achieved by output perturbation; the responses are computed on the real data and the result is perturbed before release. In



most of the literature on differential privacy, the noise added for perturbation is taken to follow a Laplace distribution. However, it has been shown in Soria-Comas and Domingo-Ferrer (2011) that better noise distributions exist, in the sense that, given a parameter  $\varepsilon$ , they guarantee  $\varepsilon$ -differential privacy and have a smaller variance than the Laplace distribution (which implies less distortion for the data and hence less information loss).

However elegant, differential privacy has been criticised as not caring about data utility. In particular, if the original variables have bounded ranges,  $\varepsilon$ -differentially private data are likely to go off-range if  $\varepsilon$  is small (high protection); see Sarathy and Muralidhar (2011) and Soria-Comas and Domingo-Ferrer (2011) for details.

## 1.2 An approach to Statistical Disclosure Control

This section describes the approach that a data provider within an NSI should take in order to meet data users needs while managing confidentiality risks. A general framework for addressing the question of confidentiality protection for different statistical outputs is proposed based on the following five key stages, and we outline how the handbook provides guidance on the different aspects of this process:

1. Why is confidentiality protection needed?
2. What are the key characteristics and uses of the data?
3. What disclosure risks need to be protected against?
4. Disclosure control methods.
5. Implementation.

### 1.2.1 Why is confidentiality protection needed?

There are three main reasons why confidentiality protection is needed for statistical outputs:

1. It is a fundamental principle for Official Statistics that the statistical records of individual persons, businesses or events used to produce Official Statistics are strictly confidential, and are to be used only for statistical purposes. Principle 6 of the UN Economic Commission report 'Fundamental Principles for Official Statistics', April 1992 states: 'Individual data collected by statistical agencies for statistical compilation, whether they refer to natural or legal persons, are to be strictly confidential and used exclusively for statistical purposes'. The disclosure control methods applied for the outputs from an NSI should meet the requirements of this principle.
2. There may be legislation that places a legal obligation on an NSI to protect individual business and personal data. In addition, where public statements are made about the protection of confidentiality or pledges are made to respondents of business or social surveys these place a duty of confidence on the NSI that the NSI must legally comply with.



3. One of the reasons why the data collected by NSIs is of such high quality is that data suppliers or respondents have confidence and trust in the NSI to preserve the confidentiality of individual information. It is essential that this confidence and trust is maintained and that identifiable information is held securely, only used for statistical purposes and not revealed in published outputs.

More information on regulations and legislation is provided in Chapter 2.

### **1.2.2 What are the key characteristics and uses of the data?**

When considering confidentiality protection of a statistical output it is important to understand the key characteristics of the data since all of these factors influence both disclosure risks and appropriate disclosure control methods. This includes knowing the type of data, e.g. full population or sample survey; sample design, an assessment of quality, e.g. the level of non-response and coverage of the data; variables and whether they are categorical or continuous; and type of outputs, e.g. microdata, magnitude or frequency tables. Producers of statistics should design publications according to the needs of users, as a first priority. It is, therefore, vital to identify the main users of the statistics, and understand why they need the figures and how they will use them. This is necessary to ensure that the design of the output is relevant and the amount of disclosure protection used has the least possible adverse impact on the usefulness of the statistics. Section 3.2 addresses some examples on how to carry out this initial analysis.

### **1.2.3 What disclosure risks need to be protected against?**

Disclosure risk assessment combines the understanding gained above with a method to identify situations where there is a likelihood of disclosure. Risk is a function of likelihood (related to the design of the output), and impact of disclosure (related to the nature of the underlying data). In order to be explicit about the disclosure risks to be managed, one should consider a range of potentially disclosive situations or scenarios and take action to prevent them. A disclosure scenario describes (i) which information is potentially available to an intruder and (ii) how the intruder would use the information to identify an individual. A range of intruder scenarios should be determined for different outputs to provide an explicit statement of what the disclosure risks are, and what elements of the output pose an unacceptable risk of disclosure. Issues in developing disclosure scenarios are provided in Section 3.3.1. Risk-assessment methods for microdata are covered in Section 3.4 and different rules applied to assess the risk of magnitude and frequency tables are described in Chapters 4 and 5, respectively.

### **1.2.4 Disclosure control methods**

Once an assessment of risk has been undertaken, an NSI must then take steps to manage any identified risks. The risk within the data is not entirely eliminated but is reduced to an acceptable level, this can be achieved either through the application



of SDC methods or through the controlled use of outputs, or through a combination of both. Several factors must be balanced through the choice of approach. Some measure of information loss and impact on main uses of the data can be used to compare alternatives. Any method must be implemented within a given production system so available software and efficiency within demanding production timetables must be considered. SDC methods used to reduce the risk of microdata, magnitude tables and frequency tables are covered in Chapters 3–5, respectively. Chapter 6 provides information on how disclosure risk can be managed by restricting access.

### **1.2.5 Implementation**

The final stage in this approach to a disclosure control problem is implementation of the methods and dissemination of the statistics. This will include identification of the software to be used along with any options and parameters. The proposed guidance will allow data providers to set disclosure control rules and select appropriate disclosure control methods to protect different types of outputs. The most important consideration is maintaining confidentiality but these decisions will also accommodate the need for clear, consistent and practical solutions that can be implemented within a reasonable time and using available resources. The methods used will balance the loss of information against the likelihood of individuals information being disclosed. Data providers should be open and transparent in this process and document their decisions and the whole risk-assessment process so that these can be reviewed. Users should be aware that a data set has been assessed for disclosure risk, and whether methods of protection have been applied. For quality purposes, users of a data set should be provided with an indication of the nature and extent of any modification due to the application of disclosure control methods. Any technique(s) used may be specified, but the level of detail made available should not be sufficient to allow the user to recover disclosive data. Each chapter of the handbook provides details of software that can be used to assess and manage disclosure risk of the different statistical outputs.

## **1.3 The chapters of the handbook**

This book starts with an overview of regulations describing the legal underpinning of SDC in Chapter 2. Microdata are covered in Chapter 3, magnitude tables are addressed in Chapter 4 and Chapter 5 provides guidance for frequency tables. Chapter 6 describes the confidentiality problems associated with microdata access issues. Within each chapter, different approaches to assessing and managing disclosure risks are described and the advantages and disadvantages of different SDC methods are discussed. Where appropriate recommendations are made for best practice. In Chapter 7, a glossary of statistical terms used in SDC has been included.



# Ethics, principles, guidelines and regulations – a general background

## 2.1 Introduction

Information from statistics becomes available for the public in tabular and microdata form. Historically, only tabular data were available and National Statistical Institutes (NSIs) had a monopoly on the microdata. Since the 1980s, the PC revolution led to the end of this monopoly. Now also other users of statistics have the possibility of using microdata. These microdata can be conveyed with CD-ROMs, USB sticks and other means. Recently, also other possibilities of getting statistical information have become more popular: remote access and remote execution. With these techniques, researchers can get access to data that remain in a statistical office or can execute set-ups without having the data on their own PC. For very sensitive information, some NSIs have the possibility to let bona fide researchers work on-site within the premises of the NSI.

The task of statistical offices is to produce and publish statistical information about society. The data collected are ultimately released in a suitable form to policy makers, researchers and the general public for statistical purposes. The release of such information may have the undesirable effect that information on individual entities instead of on sufficiently large groups of individuals is disclosed. The question then arises how the information available can be modified in such a way that the data released can be considered statistically useful and do not jeopardise the privacy of the entities concerned. The Statistical Disclosure Control theory is used to solve the problem of how to publish and release as much detail in these data as possible



without disclosing individual information (see, e.g. Willenborg and de Waal 1996, Willenborg and de Waal 2001, Lenz 2010 or Duncan *et al.* 2011).

In the current chapter, the ethical codes (Section 2.2), UNECE principles and guidelines (Section 2.3) and laws (Section 2.4) will be described.

## 2.2 Ethical codes and the new ISI code

Many countries have an ethical code that forms the basis of the production of official statistics. An internationally recognised ethical code is the declaration on professional ethics by the International Statistical Institute (ISI).

### 2.2.1 ISI Declaration on Professional Ethics

After an intense preparation process taking place from 1979 to 1985, ISI adopted the ISI Declaration on Professional Ethics in 1985 (see ISI 1985). For statistical disclosure control, clauses 4.5 and 4.6 of the declaration are of importance, and therefore, they are cited below.

#### **4.5 Maintaining confidentiality of records**

*Statistical data are unconcerned with individual identities. They are collected to answer questions such as ‘how many?’ or ‘what proportion?’, not ‘who?’. The identities and records of cooperating (or non-cooperating) subjects should, therefore, be kept confidential, whether or not confidentiality has been explicitly pledged.*

#### **4.6 Inhibiting disclosure of identities**

*Statisticians should take appropriate measures to prevent their data from being published or otherwise released in a form that would allow any subject’s identity to be disclosed or inferred.*

There can be no absolute safeguards against breaches of confidentiality, that is the disclosure of identified or identifiable data in contravention of an implicit or explicit obligation to the source. Many methods exist for lessening the likelihood of such breaches, the most common and potentially secure of which is anonymity. Its virtue as a security system is that it helps to prevent unwitting breaches of confidentiality. As long as data travel incognito, they are more difficult to attach to individuals or organisations.

There is a powerful case for identifiable statistical data to be granted ‘privileged’ status in law so that access to them by third parties is legally blocked in the absence of the permission of the responsible statistician (or his or her subjects). However, even without such legal protection it is the statistician’s responsibility to ensure that the identities of subjects are protected.



Anonymity alone is by no means a guarantee of confidentiality. A particular configuration of attributes can, like a fingerprint, frequently identify its owner beyond reasonable doubt. So statisticians need to counteract the opportunities for others to infer identities from their data. They may decide to group data in such a way as to disguise identities or to employ a variety of available measures that seek to impede the detection of identities without inflicting very serious damage to the aggregate data set. Some damage to analysis possibilities is unavoidable in these circumstances, but it needs to be weighted against the potential damage to the sources of data in the absence of such action.

The widespread use of computers is often regarded as a threat to individuals and organisations because it provides new methods of disclosing and linking identified records. On the other hand, the statistician should attempt to exploit the impressive capacity of computers to disguise identities and to enhance data security.

### **2.2.2 New ISI Declaration on Professional Ethics**

A declaration on professional ethics has to be renewed from time to time. After a number of years of preparation, the new declaration on professional ethics was adopted by the ISI Council on 22 and 23 July 2010 in Reykjavik (Iceland) (see ISI 2010).

The New ISI Declaration on Professional Ethics consists of a statement of shared professional values and a set of ethical principles that derive from these values.

For the purposes of the new declaration, the definition of who is a statistician goes well beyond those with formal degrees in the field, to include a wide array of creators and users of statistical data and tools. Statisticians work within a variety of economic, cultural, legal and political settings, each of which influences the emphasis and focus of statistical inquiry. They also work within one of several different branches of their discipline, each involving its own techniques and procedures and, possibly, its own ethical approach.

The aim of the new declaration is to enable the statistician's individual ethical judgements and decisions to be informed by shared values and experience, rather than by rigid rules imposed by the profession.

The declaration recognises that the operation of one principle may impede the operation of another. Statisticians thus have competing obligations not all of which can be fulfilled simultaneously. Thus, statisticians will sometimes have to make choices between principles. The declaration does not attempt to resolve these choices or to establish priorities among the principles. Instead, it offers a framework within which the conscientious statistician should be able to work comfortably. It is urged that departures from the framework of principles be the result of deliberation rather than of ignorance.

The declaration's first intention is to be informative and descriptive rather than authoritarian or prescriptive. Second, it is designed to be applicable as far as possible to the wide and changing areas of statistical methodology and application. For this reason, its provisions are drawn quite broadly. Third, although the principles are framed so as to have wider application to decisions than to the issues it specifically