

---

# **Statistical Disclosure Control: A Theory Guide**

**Thijs Benschop, Cathrine Machingauta, Matthew Welch**

**Jul 26, 2021**



---

## Table of Content

---

<b>1</b>	<b>Introduction</b>	<b>3</b>
1.1	Building a knowledge base . . . . .	4
1.2	Using this guide . . . . .	5
1.3	Statistical software tools for technical implementation . . . . .	5
1.4	Outline of this guide . . . . .	5
<b>2</b>	<b>Glossary and list of acronyms</b>	<b>7</b>
2.1	List of Acronyms . . . . .	7
2.2	Glossary . . . . .	8
<b>3</b>	<b>Statistical Disclosure Control (SDC): An Introduction</b>	<b>13</b>
3.1	Need for SDC . . . . .	13
3.2	The risk-utility trade-off in the SDC process . . . . .	15
<b>4</b>	<b>Release Types</b>	<b>17</b>
4.1	Conditions for PUFs . . . . .	18
4.2	Conditions for SUFs . . . . .	19
4.3	Conditions for microdata available in a controlled research data center . . . . .	20
<b>5</b>	<b>Measuring Risk</b>	<b>23</b>
5.1	Types of disclosure . . . . .	23
5.2	Classification of variables . . . . .	24
5.3	Disclosure scenarios . . . . .	25
5.4	Levels of risk . . . . .	26
5.5	Individual risk . . . . .	27
5.6	Special Uniques Detection Algorithm (SUDA) . . . . .	31
5.7	Risk measures for continuous variables . . . . .	34
5.8	Global risk . . . . .	36
5.9	Household risk . . . . .	37
<b>6</b>	<b>Anonymization Methods</b>	<b>39</b>
6.1	Classification of SDC methods . . . . .	39
6.2	Non-perturbative methods . . . . .	40
6.3	Perturbative methods . . . . .	50
6.4	Anonymization of geospatial variables . . . . .	60
6.5	Anonymization of the quasi-identifier household size . . . . .	61
6.6	Special case: census data . . . . .	61

<b>7</b>	<b>Measuring Utility and Information Loss</b>	<b>63</b>
7.1	General utility measures for continuous and categorical variables . . . . .	64
7.2	Utility measures based on the end user's needs . . . . .	66
7.3	Regression . . . . .	66
7.4	Assessing data utility with the help of data visualizations . . . . .	67
7.5	Choice of utility measure . . . . .	71
<b>8</b>	<b>The SDC Process</b>	<b>75</b>
8.1	Step 1: Need for confidentiality protection . . . . .	75
8.2	Step 2: Data preparation and exploring data characteristics . . . . .	76
8.3	Step 3: Type of release . . . . .	77
8.4	Step 4: Intruder scenarios and choice of key variables . . . . .	78
8.5	Step 5: Data key uses and selection of utility measures . . . . .	79
8.6	Step 6: Assessing disclosure risk . . . . .	80
8.7	Step 7: Assessing utility measures . . . . .	80
8.8	Step 8: Choice and application of SDC methods . . . . .	80
8.9	Step 9: Re-measure risk . . . . .	81
8.10	Step 10: Re-measure utility . . . . .	81
8.11	Step 11: Audit and Reporting . . . . .	82
8.12	Step 12: Data release . . . . .	82
<b>9</b>	<b>Case Studies (Illustrating the SDC Process)</b>	<b>85</b>
9.1	Case study 1- SUF . . . . .	85
9.2	Case study 2 - PUF . . . . .	111
<b>10</b>	<b>Appendices</b>	<b>135</b>
10.1	Appendix A: Overview of Case Study Variables . . . . .	135
10.2	Appendix B: Example of Blanket Agreement for SUF . . . . .	136
10.3	Appendix C: Internal and External Reports for Case Studies . . . . .	138
10.4	<b>Case study 1 - Internal report</b> . . . . .	138
10.5	<b>Case study 1 - External report</b> . . . . .	145
10.6	<b>Case study 2 - Internal report</b> . . . . .	146
10.7	<b>Case study 2- External report</b> . . . . .	151
10.8	Appendix D: Execution Times for Multiple Scenarios Tested using Selected Sample Data . . . . .	151
<b>11</b>	<b>Acknowledgements</b>	<b>155</b>
	<b>Bibliography</b>	<b>157</b>

This guide provides an introduction to the theory of Statistical Disclosure Control (SDC) for microdata. It includes an overview of the most commonly applied methods in SDC, a step-by-step overview of the complete SDC process and many examples from practice in National Statistics Offices (NSOs).

For guidance on the technical implementation of the theory mentioned in the guide, please refer to our guides:

- (i) for guidance on the application of methods and on using *sdcMicro* from the command-line
- (ii) for guidance on the application of methods and on using the GUI *sdcApp* available for *sdcMicro*

**Authors of this guide:** Thijs Benschop and Matthew Welch, The World Bank

**Acknowledgments:** The authors thank Olivier Dupriez and Cathrine Machingauta (The World Bank) for their technical comments and inputs throughout the process.

The production of this guide was made possible through a World Bank Knowledge for Change II Grant: KCP II - A microdata dissemination challenge: Balancing data protection and data utility. Grant number: TF 015043, Project Number P094376. As well as from United Kingdom - DFID funding to the World Bank Multi-Donor Trust Fund - International Household Survey and Accelerated Data Program – TF071804/TF011722/TF0A7461.



# CHAPTER 1

---

## Introduction

---

National statistics agencies are mandated to collect microdata<sup>1</sup> from surveys and censuses to inform and measure policy effectiveness. In almost all countries, statistics acts and privacy laws govern these activities. These laws require that agencies protect the identity of respondents, but may also require that agencies disseminate the results and, in appropriate cases, the microdata. Data producers who are not part of national statistics agencies are also often subject to restrictions, through privacy laws or strict codes of conduct and ethics that require a similar commitment to privacy protection. This has to be balanced against the increasing requirement from funders that data produced using donor funds be made publically available.

This tension between complying with confidentiality requirements while at the same time requiring that microdata be released means that a demand exists for practical solutions for applying Statistical Disclosure Control (SDC), also known as microdata anonymization. The provision of adequate solutions and technical support has the potential to “unlock” a large number of datasets.

The [International Household Survey Network \(IHSN\)](#) and the World Bank have contributed to successful programs that have generated tools, resources and guidelines for the curation, preservation and dissemination of microdata and resulted in the documentation of thousands of surveys by countries and agencies across the world. While these programs have ensured substantial improvements in the preservation of data and dissemination of good quality metadata, many agencies are still reluctant to allow access to the microdata. The reasons are technical, legal, ethical and political, and sometimes involve a fear of being criticized for not providing perfect data. When combined with the tools and guidelines already developed by the IHSN/World Bank for the curation, preservation and dissemination of microdata, tools and guidelines for the anonymization of microdata should further reduce or remove some of these obstacles.

Working with the IHSN, PARIS21 (OECD), Statistics Austria and the Vienna University of Technology, the World Bank has contributed to the development of an open source software package for SDC, called *sdcMicro*. The package was developed for use with the open source *R* statistical software, available from the Comprehensive R Archive Network (CRAN) at <http://cran.us.r-project.org>. The package includes numerous methods for the assessment and reduction of disclosure risk in microdata.

Ensuring that a free open source solution is available to agencies was an important step forward, but not a sufficient one. There is still limited consolidated and reported knowledge on the impact of disclosure risk reduction methods

---

<sup>1</sup> Microdata are unit-level data obtained from sample surveys, censuses and administrative systems. They provide information about characteristics of individual people or entities such as households, business enterprises, facilities, farms or even geographical areas such as villages or towns. They allow in-depth understanding of socio-economic issues by studying relationships and interactions among phenomena. Microdata are thus key to designing projects and formulating policies, targeting interventions and monitoring and measuring the impact and results of projects, interventions and policies.

on data utility. This limited access to knowledge combined with a lack of experience in using the tools and methods makes it difficult for many agencies to implement optimal solutions, i.e., solutions that meet their obligations towards both privacy protection and the release of data useful for policy monitoring and evaluation. This guide attempts to fill this critical gap by:

- (i) consolidating knowledge gained at the World Bank through experiments conducted during a large-scale evaluation of anonymization techniques
- (ii) translating the experience and key results into practical guidelines

It should be stressed that SDC is only one part of the data release process, and its application must be considered within the complete data release framework. The level and methods of SDC depend on the laws of the country, the sensitivity of the data and the access policy (i.e., who will gain access) considered for release. Agencies that are currently releasing data are already using many of the methods described in this guide and applying appropriate access policies to their data before release. The primary objective of this guide is to provide a primer to those new to the process who are looking for guidance on both theory and practical implementation. This guide is not intended to prescribe or advocate for changes in methods that specific data producers are already using and which they have designed to fit and comply with their existing data release policies.

The guide seeks to provide practical steps to those agencies that want to unlock access to their data in a safe way and ensure that the data remain fit for purpose.

## 1.1 Building a knowledge base

The release of data is important, as it allows researchers and policymakers to replicate officially published results, generate new insights into issues, avoid duplication of surveys and provide greater returns to the investment in the survey process.

Both the production of reports, with aggregate tables of indicators and statistics, and the release of microdata result in privacy challenges to the producer. In the past, for many agencies, the only requirement was to release a report and some key indicators. The recent movement around Open Data, Open Government and transparency means that agencies are under greater pressure to release their microdata to allow broader use of data collected through public and donor funds. This guide focuses on the methods and processes for the release of microdata.

Releasing data in a safe way is required to protect the integrity of the statistical system, by ensuring agencies honor their commitment to respondents to protect their identity. Agencies do not widely share, in substantial detail, their knowledge and experience using SDC and the processes for creating safe data with other agencies. This makes it difficult for agencies new to the process to implement solutions. To fill this experience and knowledge gap, we evaluated the use of a broad suite of SDC methods on a range of survey microdata covering important development topics related to health, labor, education, poverty and inequality. The data we used were all previously treated to make them safe for release. Given that their producers had already treated these data, it was not possible, nor was it our goal, to pass any judgment on the safety of these data, many of which are in the public domain. The focus was rather on measuring the effects that various methods would have on the risk-utility trade-off for microdata produced to measure common development indicators. We used the experience from this large-scale experimentation to inform our discussion of the processes and methods in this guide.

---

### Important

At no point was any attempt made to re-identify, through matching or any other method, any respondents in the surveys we used in building our knowledge base. All risk assessments were based on frequencies and probabilities.

---



## 1.2 Using this guide

The methods discussed in this guide originate from a large body of literature on SDC. The processes underlying many of the methods are the subject of extensive academic research and many, if not all, of them are used extensively by agencies experienced in preparing microdata for release.

Where possible, for each method and topic, we provide elaborate examples, references to the original or seminal work describing the methods and algorithms in detail and recommended readings. This, when combined with the discussion of the method and practical considerations in this guide, should allow the reader to understand the methods and their strengths and weaknesses. It should also provide enough detail for readers to use an existing software solution to implement the methods or program the methods in statistical software of their choice.

## 1.3 Statistical software tools for technical implementation

There are several software solutions for the technical implementation of SDC. The two are widely used in statistics offices in the European Union and implement many of the methods discussed in this guide are *R* package *sdcMicro*<sup>2</sup> and the application  $\mu$ -ARGUS<sup>3</sup> developed by Statistics Netherlands.

*sdcMicro* is an add-on package to the statistical software *R*. The package was developed and is maintained by Matthias Templ, Alexander Kowarik and Bernhard Meindl. The statistical software *R* and the *sdcMicro* package, as well as any other packages needed for the SDC process, are freely available from the Comprehensive R Archive Network (CRAN) mirrors (<http://cran.r-project.org/>). The software is available for Linux, Windows and Macintosh operating systems. We chose to use *R* and *sdcMicro* because it is freely available, accommodates all main data formats and is easy to adapt by the user. The World Bank, through the IHSN, has also provided funding towards the development of the *sdcMicro* package to ensure it meets the requirements of the agencies we support.

*sdcMicro* can be used from *R* command-line or by using the integrated GUI *sdcApp* for those not familiar with *R*. This guide is accompanied by two technical practice guides both available as Read the Docs documentations:

- (i) [Statistical Disclosure Control for Microdata: A Practice Guide](#) for using *sdcMicro* from command-line
- (ii) [sdcApp manual](#) for using the GUI *sdcApp*

## 1.4 Outline of this guide

This guide is divided into the following main sections:

- (i) the Section [Statistical Disclosure Control \(SDC\): An Introduction](#) is a primer on SDC.
- (ii) the Section [Release Types](#) gives an introduction to different release types for microdata.
- (iii) the Sections [Anonymization Methods](#), [Measuring Risk](#) and [Measuring Utility and Information Loss](#) cover SDC methods, risk and utility measurement. The goal here is to provide knowledge that allows the reader to independently apply and execute the SDC process. This section is enriched with real examples. The interested reader can also find more information in the references and recommended readings at the end of each section.
- (iv) the Section [The SDC Process](#) provides a step-by-step guide to disclosure control, which draws upon the knowledge presented in the previous sections.
- (v) the Section [Case Studies \(Illustrating the SDC Process\)](#) presents a number of detailed case studies that demonstrate the use of the methods, their implementation in *sdcMicro* and the process that should be followed to reach the optimal risk-utility solution.

<sup>2</sup> See <http://cran.r-project.org/web/packages/sdcMicro/index.html> and the GitHub <https://github.com/sdcTools/sdcMicro> of the developers. The GitHub repository can also be used to submit bugs found in the package.

<sup>3</sup>  $\mu$ -ARGUS is available at: <http://neon.vb.cbs.nl/casc/mu.htm>. The software was recently ported to open source.



---

### Glossary and list of acronyms

---

#### 2.1 List of Acronyms

**AFR** Sub Saharan Africa

**COICOP** Classification of Individual Consumption by Purpose

**CRAN** Comprehensive R Archive Network

**CTBIL** Contingency Table-Based Information Loss

**DHS** Demographic and Health Surveys

**DIS** Data Intrusion Simulation

**EAP** East Asia and the Pacific

**ECA** Europe and Central Asia

**EU** European Union

**GIS** Geographical Information System

**GPS** Global Positioning System

**GUI** Graphical User Interface

**HIV/AIDS** Human Immunodeficiency Virus/Acquired Immune Deficiency Syndrome

**I2D2** International Income Distribution Database

**IHSN** International Household Survey Network

**LAC** Latin America and the Caribbean

**LSMS** Living Standards Measurement Survey

**MDAV** Maximum Distance Average Vector

**MDG** Millennium Development Goal

**MENA** Middle East and North America

**MICS** Multiple Indicator Cluster Survey  
**MME** Mean Monthly Expenditures  
**MMI** Mean Monthly Income  
**MSU** Minimal Sample Uniques  
**NSI** National Statistical Institute  
**NSO** National Statistical Office  
**OECD** Organization for Economic Cooperation and Development  
**PARIS21** Partnership in Statistics for Development in the 21<sup>st</sup> century  
**PRAM** Post Randomization Method  
**PC** Principal Component  
**PUF** Public Use File  
**SA** South Asia  
**SDC** Statistical Disclosure Control  
**SSE** Sum of Squared Errors  
**SHIP** Survey-based Harmonized Indicators Program  
**SUDA** Special Uniques Detection Algorithm  
**SUF** Scientific Use File  
**UNICEF** United Nations Children's Fund

## 2.2 Glossary

**Administrative data** Data collected for administrative purposes by government agencies. Typically, administrative data require specific SDC methods.

**Anonymization** Use of techniques that convert confidential data into anonymized data/ removal or masking of identifying information from datasets.

**Attribute disclosure** Attribute disclosure occurs if an intruder is able to determine new characteristics of an individual or organization based on the information available in the released data.

**Categorical variable** A variable that takes values over a finite set, e.g., gender. Also called factor in *R*.

**Confidentiality** Data confidentiality is a property of data, usually resulting from legislative measures, which prevents it from unauthorized disclosure.<sup>2</sup>

**Confidential data** Data that will allow identification of an individual or organization, either directly or indirectly.<sup>1</sup>

**Continuous variable** A variable with which numerical and arithmetic operations can be performed, e.g., income.

**Data protection** Data protection refers to the set of privacy-motivated laws, policies and procedures that aim to minimize intrusion into respondents' privacy caused by the collection, storage and dissemination of personal data.<sup>2</sup>

**Deterministic methods** Anonymization methods that follow a certain algorithm and produce the same results if applied repeatedly to the same data with the same set of parameters.

---

<sup>2</sup> OECD, <http://stats.oecd.org/glossary>

<sup>1</sup> Australian Bureau of Statistics, <http://www.nss.gov.au/nss/home.nsf/pages/Confidentiality--+Glossary>

- Direct identifier** A variable that reveals directly and unambiguously the identity of a respondent, e.g., names, social identity numbers.
- Disclosure** Disclosure occurs when a person or an organization recognizes or learns something that they did not already know about another person or organization through released data.<sup>1</sup> See also Identity disclosure, Attribute disclosure and Inferential disclosure.
- Disclosure risk** A disclosure risk occurs if an unacceptably narrow estimation of a respondent's confidential information is possible or if exact disclosure is possible with a high level of confidence.<sup>2</sup> Disclosure risk also refers to the probability that successful disclosure could occur.
- End user** The user of the released microdata file after anonymization. Who is the end user depends on the release type.
- Factor variable** Factor variables are one way to classify categorical variables in *R*.
- Hierarchical structure** Data is made up of collections of records that are interconnected through links, e.g., individuals belonging to groups/households or employees belonging to companies.
- Identifier** An identifier is a variable/ information that can be used to establish identity of an individual or organization. Identifiers can lead to direct or indirect identification.
- Identity disclosure** Identity disclosure occurs if an intruder associates a known individual or organization with a released data record.
- Indirect identification** Indirect identification occurs when the identity of an individual or organization is disclosed, not using direct identifiers but through a combination of unique characteristics in key variables.<sup>1</sup>
- Inferential disclosure** Inferential disclosure occurs if an intruder is able to determine the value of some characteristic of an individual or organization more accurately with the released data than otherwise would have been possible.
- Information loss** Information loss refers to the reduction of the information content in the released data relative to the information content in the raw data. Information loss is often measured with respect to common analytical measures, such as regressions and indicators. See also Utility.
- Interval** A set of numbers between two designated endpoints that may or may not be included. Brackets (e.g., [0, 1]) denote a closed interval, which includes the endpoints 0 and 1. Parentheses (e.g., (0, 1) denote an open interval, which does not include the endpoints.
- Intruder** A user who misuses released data by trying to disclose information about an individual or organization, using a set of characteristics known to the user.
- k*-anonymity** The risk measure *k*-anonymity is based on the principle that the number of individuals in a sample sharing the same combination of values (key) of categorical key variables should be higher than a specified threshold *k*.
- Key** A combination or pattern of key variables/quasi-identifiers.
- Key variables** A set of variables that, in combination, can be linked to external information to re-identify respondents in the released dataset. Key variables are also called "quasi-identifiers" or "implicit identifiers".
- Microaggregation** Anonymization method that is based on replacing values for a certain variable with a common value for a group of records. The grouping of records is based on a proximity measure of variables of interest. The groups of records are also used to calculate the replacement value.
- Microdata** A set of records containing information on individual respondents or on economic entities. Such records may contain responses to a survey questionnaire or administrative forms.
- Noise addition** Anonymization method based on adding or multiplying a stochastic or randomized number to the original values to protect data from exact matching with external files. Noise addition is typically applied to continuous variables.

**Non-perturbative methods** Anonymization methods that reduce the detail in the data or suppress certain values (masking) without distorting the data structure.

**Observation** A set of data derived from an object/unit of experiment, e.g., an individual (in individual-level data), a household (in household-level data) or a company (in company data). Observations are also called “records”.

**Original data** The data before SDC/anonymization methods were applied. Also called “raw data” or “untreated data”.

**Outlier** An unusual value that is correctly reported but is not typical of the rest of the population. Outliers can also be observations with an unusual combination of values for variables, such as 20-year-old widow. On their own age, 20 and widow are not unusual values, but their combination may be.<sup>1</sup>

**Perturbative methods** Anonymization methods that alter values slightly to limit disclosure risk by creating uncertainty around the true values, while retaining as much content and structure as possible, e.g. microaggregation and noise addition.

**Population unique** The only record in the population with a particular set of characteristics, such that the individual or organization can be distinguished from other units in the population based on that set of characteristics.

**Post Randomization Method (PRAM)** Anonymization method for microdata in which the scores of a categorical variable are altered according to certain probabilities. It is thus intentional misclassification with known misclassification probabilities.<sup>1</sup>

**Probabilistic methods** Anonymization methods that depend on a probability mechanism or a random number-generating mechanism. Every time a probabilistic method is used, a different outcome is generated.

**Privacy** Privacy is a concept that applies to data subjects while confidentiality applies to data. The concept is defined as follows: “It is the status accorded to data which has been agreed upon between the person or organization furnishing the data and the organization receiving it and which describes the degree of protection which will be provided.”<sup>2</sup>

**Public Use File (PUF)** Type of release of microdata file, which is freely available to any user, for example on the internet.

**Quasi-identifiers** A set of variables that, in combination, can be linked to external information to re-identify respondents in the released dataset. Quasi-identifiers are also called “key variables” or “implicit identifiers”.

**Raw data** The data before SDC/anonymization methods were applied. Also called “original data” or “untreated data”.

**Recoding** Anonymization method for microdata in which groups of existing categories/values are replaced with new values, e.g. the values ‘protestant’, and ‘catholic’ are replaced with ‘Christian’. Recoding reduces the detail in the data. Recoding of continuous variables leads to a transformation from continuous to categorical, e.g. creating income bands.

**Record** A set of data derived from an object/unit of experiment, e.g., an individual (in individual-level data), a household (in household-level data) or a company (in company data). Records are also called “observations”.

**Regression** A statistical process of measuring the relation between the mean value of one variable and corresponding values of other variables.

**Re-identification risk** See Disclosure risk

**Release** Dissemination – the release to users of information obtained through a statistical activity.<sup>2</sup>

**Respondents** Individuals or units of observation whose information/responses to a survey make up the data file.

**Sample unique** The only record in the sample with a particular set of characteristics, such that the individual or organization can be distinguished from other units in the sample based on that set of characteristics.

**Scientific Use File (SUF)** Type of release of microdata file, which is only available to selected researchers under contract. Also known as “licensed file”, “microdata under contract” or “research file”.

***sdcMicro*** An *R* based package authored by Templ, M., Kowarik, A. and Meindl, B. with tools for the anonymization of microdata, i.e. for the creation of public- and scientific-use files.

***sdcMicroGUI*** A GUI for the *R* based *sdcMicro* package, which allows users to use the *sdcMicro* tools without *R* knowledge.

**Sensitive variables** Sensitive or confidential variables are those whose values must not be discovered for any respondent in the dataset. The determination of sensitive variables is often subject to legal and ethical concerns.

**Statistical Disclosure Control (SDC)** Statistical Disclosure Control techniques can be defined as the set of methods to reduce the risk of disclosing information on individuals, businesses or other organizations. Such methods are only related to the dissemination step and are usually based on restricting the amount of or modifying the data released.<sup>2</sup>

**Suppression** Data suppression involves not releasing information that is considered unsafe because it fails confidentiality rules being applied. Sometimes this is done by replacing values signifying individual attributes with missing values. In the context of this guide, usually to achieve a desired level of *k*-anonymity.

**Threshold** An established level, value, margin or point at which values that fall above or below it will deem the data safe or unsafe. If unsafe, further action will need to be taken to reduce the risk of identification.

**Utility** Data utility describes the value of data as an analytical resource, comprising analytical completeness and analytical validity.

**Untreated data** The data before SDC/anonymization methods were applied. Also called “raw data” or “original data”.

**Variable** Any characteristic, number or quantity that can be measured or counted for each unit of observation.





---

## Statistical Disclosure Control (SDC): An Introduction

---

### 3.1 Need for SDC

A large part of the data collected by statistical agencies cannot be published directly due to privacy and confidentiality concerns. These concerns are both of legal and ethical nature. SDC seeks to treat and alter the data so that the data can be published or released without revealing the confidential information it contains, while, at the same time, limit information loss due to the anonymization of the data. In this guide, we discuss only disclosure control for microdata.<sup>1</sup> Microdata are datasets that provide information on a set of variables for each individual respondent. Respondents can be natural persons, but also legal entities such as companies.

The aim of anonymizing microdata is to transform the datasets to achieve an “acceptable level” of disclosure risk. The level of acceptability of disclosure risk and the need for anonymization are usually at the discretion of the data producer and guided by legislation. These are formulated in the dissemination policies and programs of the data providers and based on considerations including “[...] the costs and expertise involved; questions of data quality, potential misuse and misunderstanding of data by users; legal and ethical matters; and maintaining the trust and support of respondents” (*DuBo10*). There is a moral, ethical and legal obligation for the data producers to ensure that data provided by the respondents are used only for statistical purposes.

In some cases, the dissemination of microdata is a legal obligation, but, in most cases, the legislation will formulate restrictions. Thus, a country’s legislative framework will shape its microdata dissemination policy. It is crucial for data producers to “ensure there is a sound legal and ethical base (as well as the technical and methodological tools) for protecting confidentiality. This legal and ethical base requires a balanced assessment between the public good of confidentiality protection on the one hand, and the public benefits of research on the other. A decision on whether or not to provide access might depend on the merits of specific research proposals and the credibility of the researcher, and there should be some allowance for this in the legal arrangements.” (*DuBo10*).

“Data access arrangements should respect the legal rights and legitimate interests of all stakeholders in the public research enterprise. Access to, and use of, certain research data will necessarily be limited by various types of legal requirements, which may include restrictions for reasons of:

- National security: data pertaining to intelligence, military activities, or political decision making may be classified and therefore subject to restricted access.

---

<sup>1</sup> There is another strand of literature on the anonymization of tabular data, see e.g., *HDFG12*.

- Privacy and confidentiality: data on human subjects and other personal data are subject to restricted access under national laws and policies to protect confidentiality and privacy. However, anonymization or confidentiality procedures that ensure a satisfactory level of confidentiality should be considered by custodians of such data to preserve as much data utility as possible for researchers.
- Trade secrets and intellectual property rights: data on, or from, businesses or other parties that contain confidential information may not be accessible for research. (...)” (*DuBo10*).

Box 1, extracted from *DuBo10*, provides several examples of statistical legislation on microdata release.

---

### Info-box - Examples of statistical legislation on microdata release

1. The US Bureau of the Census operates under Title 13-Census of the US Code. “Title 13, U.S.C., Section 9 prohibits the publication or release of any information that would permit identification of any particular establishment, individual, or household. Disclosure assurance involves the steps taken to ensure that Title 13 data prepared for public release will not result in wrongful disclosure. This includes both the use of disclosure limitation methods and the review process to ensure that the disclosure limitation techniques used provide adequate protection to the information.”

Source: <https://www.census.gov/srd/sdc/wendy.drb.faq.pdf>

2. In Canada the Statistics Act states that “no person who has been sworn under section 6 shall disclose or knowingly cause to be disclosed, by any means, any information obtained under this Act in such a manner that it is possible from the disclosure to relate the particulars obtained from any individual return to any identifiable individual person, business or organization.”

Source: <http://laws-lois.justice.gc.ca/eng/acts/S-19/page-2.html>

Statistics Canada does release microdata files. Its microdata release policy states that Statistics Canada will authorise the release of microdata files for public use when:

- (a) the release substantially enhances the analytical value of the data collected; and
  - (b) the Agency is satisfied all reasonable steps have been taken to prevent the identification of particular survey units.
3. In Thailand the Act states: “Personal information obtained under this act shall be strictly considered confidential. A person who performs his or her duty hereunder or a person who has the duty of maintaining such information cannot disclose it to anyone who doesn’t have a duty hereunder except in the case that:
    - (1) Such disclosure is for the purpose of any investigation or legal proceedings in a case relating to an offense hereunder.
    - (2) Such disclosure is for the use of agencies in the preparation, analysis or research of statistics provided that such disclosure does not cause damage to the information owner and does not identify or disclose the data owner.”

Source: [http://web.nso.go.th/eng/en/about/stat\\_act2007.pdf](http://web.nso.go.th/eng/en/about/stat_act2007.pdf) section 15.

---

Besides the legal and ethical concerns and codes of conducts of agencies producing statistics, SDC is important because it guarantees data quality and response rates in future surveys. If respondents feel that data producers are not protecting their privacy, they might not be willing to participate in future surveys. “[...] one incident, particularly if it receives strong media attention, could have a significant impact on respondent cooperation and therefore on the quality of official statistics” (*DuBo10*). At the same time, if data users are unable to gain enough utility from the data due to excessive or inappropriate SDC protection, or are unable to access the data, then the large investment in producing the data will be lost.

## 3.2 The risk-utility trade-off in the SDC process

SDC is characterized by the trade-off between risk of disclosure and utility of the data for end users. The risk–utility scale extends between two extremes; (i) no data is released (zero risk of disclosure) and thus users gain no utility from the data, to (ii) data is released without any treatment, and thus with maximum risk of disclosure, but also maximum utility to the user (i.e., no information loss). The goal of a well-implemented SDC process is to find the optimal point where utility for end users is maximized at an acceptable level of risk. Fig. 3.1 illustrates this trade-off. The triangle corresponds to the raw data. The raw data have no information loss, but generally have a disclosure risk higher than the acceptable level. The other extreme is the square, which corresponds to no data release. In that case there is no disclosure risk, but also no utility from the data for the users. The points in-between correspond to different choices of SDC methods and/or parameters for these methods applied to different variables. The SDC process looks for the SDC methods and the parameters for those methods and applies these in a way that reduces the risk sufficiently, while minimizing the information loss.

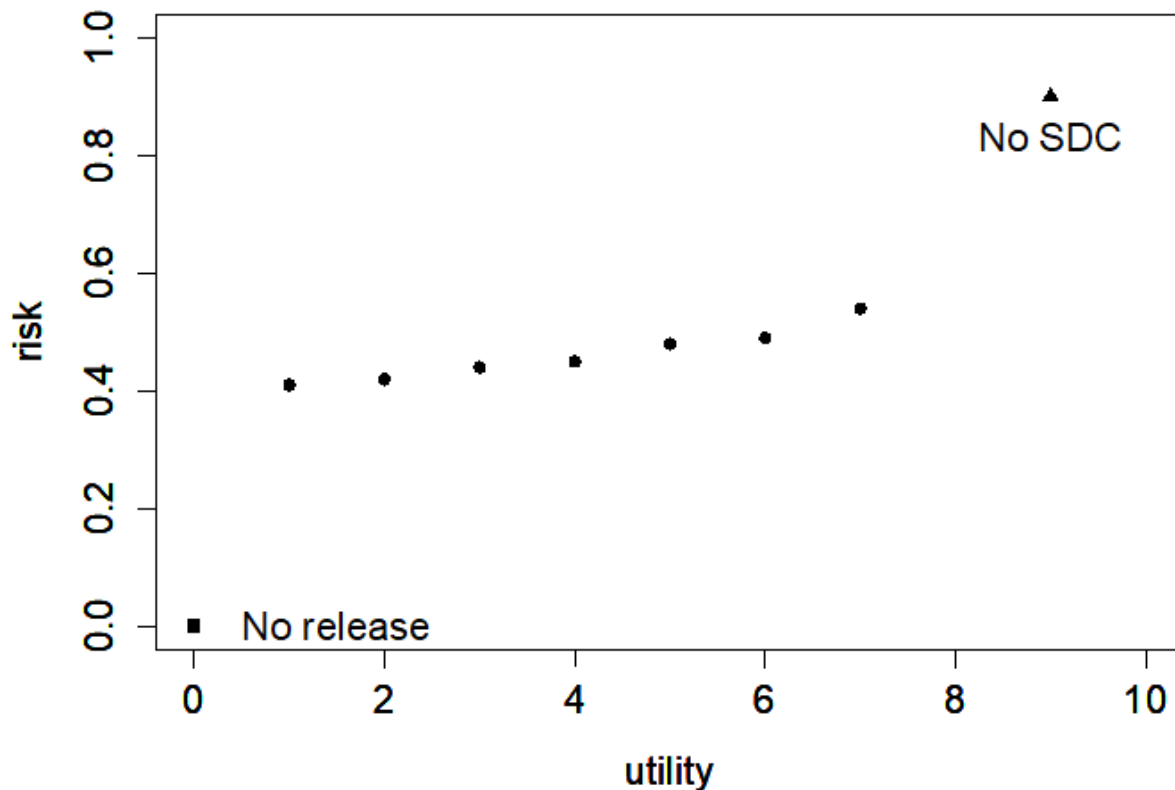


Fig. 3.1: Risk-utility trade-off

SDC cannot achieve total risk elimination, but can reduce the risk to an acceptable level. Any application of SDC methods will suppress or alter values in the data and as such decrease the utility (i.e., result in information loss) when compared to the original data. A common thread that will be emphasized throughout this guide will be that the process of SDC should prioritize the goal of protecting respondents, while at the same time keeping the data users in mind to limit information loss. In general, the lower the disclosure risk, the higher the information loss and the lower the data utility for end-users.

In practice, choosing SDC methods is partially trial and error: after applying methods, disclosure risk and data utility are re-measured and compared to the results of other choices of methods and parameters. If the result is satisfactory,

the data can be released. We will see that often the first attempt will not be the optimal one. The risk may not be sufficiently reduced or the information loss may be too high and the process has to be repeated with different methods or parameters until a satisfactory solution is found. Disclosure risk, data utility and information loss in the SDC context and how to measure them are discussed in subsequent chapters of this guide.

Again, it must be stressed that the level of SDC and methods applied depend to a large extent on the entire data release framework. For example, a key consideration is to whom and under what conditions the data are to be released (see also the Section [Release types](#)). If data are to be released as public use data, then the level of SDC applied will necessarily need to be higher than in the cases where data are released under license conditions to trusted users after careful vetting. With careful preparation, data may be released under both public and licensed versions. We discuss how this might be achieved later in the guide.

## References

---

### Release Types

---

This section discusses data release. Rather than rewriting work that has already been conducted through the World Bank and its partners at the IHSN, this section extracts from an excellent guide published by *DuBo10*.

The trade-off between risk and utility in the anonymization process depends greatly on who the users are<sup>1</sup> and under what conditions a microdata file is released. Generally, three types of data release methods are practiced and apply to different target groups.

- **Public Use File (PUF)**: the data “are available to anyone agreeing to respect a core set of easy-to-meet conditions. Such conditions relate to what cannot be done with the data (e.g. the data cannot be sold), upon gaining access to the data. In some cases PUFs are disseminated with no conditions; often being made available on-line, [e.g. on the website of the statistical agency]. These data are made easily accessible because the risk of identifying individual respondents is considered minimal. Minimising the risk of disclosure involves eliminating all content that can identify respondents directly—for instance, names, addresses and telephone numbers. In addition this requires purging relevant indirect identifiers from the microdata file. These vary across survey designs, but commonly-suppressed indirect identifiers include geographical information below the sub-national level at which the sample is representative. Occasionally, certain records may be suppressed also from PUFs, as might variables characterised by extremely skewed distribution or outliers. However, in lieu of deleting entire records or variables from microdata files, alternative SDC methods can minimise the risk of disclosure while maximizing information content. Such methods include top-and-bottom coding, local suppression or using data perturbation techniques [(see the Section [Anonymization methods](#) for an overview of anonymization methods)]. PUFs are typically generated from census data files using a sub-set [or sample] of records rather than the entire file and [from sample surveys, such as] household surveys.” (*DuBo10*).
- **Scientific Use File (SUF)** (also known as a licensed file, microdata under contract or research file): the “dissemination is restricted to users who have received authorization to access them after submitting a documented application and signing an agreement governing the data’s use. While typically licensed files are also anonymised to ensure the risk of identifying individuals is minimised when used in isolation, they may still [potentially] contain identifiable data if linked with other data files. Direct identifiers such as respondents’ names must be removed from a licensed dataset. The data files may, however, still contain indirect variables that could identify respondents by matching them to other data files such as voter lists, land registers or school records. When disseminating licensed files, the recommendation is to establish and sign an agreement between the data producer and external bona fide users – trustworthy users with legitimate need to access the data. Such an agreement

---

<sup>1</sup> See Section 5 in *DuBo10* as to who the users of microdata are and to whom microdata should be made available.

should govern access and use of such microdata files<sup>2</sup>. Sometimes, licensing agreements are only entered into with users affiliated to an appropriate sponsoring institution. i.e., research centers, universities or development partners. It is further recommended that, before entering into a data access and use agreement, the data producer asks potential users to complete an application form to demonstrate the need to use a licensed file (instead of the PUF version, if available) for a stated statistical or research purpose” (*DuBo10*). This also allows the data producer to learn which characteristics of the data are important for the users, which is valuable information for optimizing future anonymization processes.

- **Microdata available in a controlled research data center** (also known as data enclave): “Some files may be offered to users under strict conditions in a data enclave. This is a facility [(often on the premises of the data provider)] equipped with computers not linked to the internet or an external network and from which no information can be downloaded via USB ports, CD-DVD or other drives. Data enclaves contain data that are particularly sensitive or allow direct or easy identification of respondents. Examples include complete population census datasets, enterprise surveys and certain health related datasets containing highly-confidential information. Users interested in accessing a data enclave will not necessarily have access to the full dataset – only to the particular data subset they require. They will be asked to complete an application form demonstrating a legitimate need to access these data to fulfill a stated statistical or research purpose [...] The outputs generated must be scrutinised by way of a full disclosure review before release. Operating a data enclave may be expensive – it requires special premises and computer equipment. It also demands staff with the skills and time to review outputs before their removal from the data enclave in order to ensure there is no risk of disclosure. Such staff must be familiar with data analysis and be able to review the request process and manage file servers. Because of the substantial operating costs and technical skills required, some statistical agencies or other official data producers opt to collaborate with academic institutions or research centres to establish and manage data enclaves.”

There are other data access possibilities besides these, such as teaching files, files for other specific purposes, remote execution or remote access. Obviously, the required level of protection depends on the type of release; a PUF file must be protected to a much larger extent than a SUF file, which in turn has to be protected more than a file which is only available in an on-site facility. The Section [Step 3: Type of release](#) gives more guidance on the choice of the release type and its implications for the anonymization process. The same microdata set can be released in different ways for different users, e.g., as SUF and teaching file. The Section [Step 3: Type of release](#) discusses the particular issues of multiple releases of one dataset.

The first step for any agency that wants to release data would be formulation of clear data dissemination policies for the release of microdata. We will see later that deciding on the level of anonymization needed will depend partly on knowing under what conditions the data will be released. Access policies and conditions provide the framework for the whole release process.

The following sections further specify the conditions under which microdata should be provided under different release types.

## 4.1 Conditions for PUFs

“Generally, data regarded as public are open to anyone with access to an [National Statistical Office] (NSO) website. It is, however, normally good practice to include statements defining suitable uses for and precautions to be adopted in using the data. While these may not be legally binding, they serve to sensitise the user. Prohibitions such as attempts to link the data to other sources can be part of the ‘use statement’ to which the user must agree, on-line, before the data can be downloaded. [...] Dissemination of microdata files necessarily involves the application of rules or principles. [The info-box] below [taken from *DuBo10*] shows basic principles normally applying to PUFs.” (*DuBo10*).

---

### Info-box - Conditions for accessing and using PUFs

1. Data and other material provided by the NSO will not be redistributed or sold to other individuals, institutions or organisations without the NSO’s written agreement.

---

<sup>2</sup> Appendix B provides an example of a blanket agreement.

2. Data will be used for statistical and scientific research purposes only. They will be employed solely for reporting aggregated information, including modelling, and not for investigating specific individuals or organisations.
3. No attempt will be made to re-identify respondents, and there will be no use of the identity of any person or establishment discovered inadvertently. Any such discovery will be reported immediately to the NSO.
4. No attempt will be made to produce links between datasets provided by the NSO or between NSO data and other datasets that could identify individuals or organisations.
5. Any books, articles, conference papers, theses, dissertations, reports or other publications employing data obtained from the NSO will cite the source, in line with the citation requirement provided with the dataset.
6. An electronic copy of all publications based on the requested data will be sent to the NSO.
7. The original collector of the data, the NSO, and the relevant funding agencies bear no responsibility for the data's use or interpretation or inferences based upon it.

Note: Items 3 and 6 in the list require that users be provided with an easy way to communicate with the data provider. It is good practice to provide a contact number, an email address, and possibly an on-line "feedback provision" system.

Source: *DuBo10*

---

## 4.2 Conditions for SUFs

"For [SUFs], terms and conditions must include the basic common principles plus some additional ones applying to the researcher's organisation. There are two options: firstly, data are provided to a researcher or a team for a specific purpose; secondly, data are provided to an organization under a blanket agreement for internal use, e.g., to an international body or research agency. In both cases, the researcher's organisation must be identified, as must suitable representatives to sign the licence" (*DuBo10*).

*Access to a researcher or research team for a specific purpose*

"If data are provided for an individual research project, the research team must be identified. This is covered by requiring interested users to complete a formal request to access the data (a model of such a request form is provided in Appendix 1 [in *DuBo10*]). The conditions to obtain the data (see example in the info-box below) will specify that the files will not be shared outside the organisation and that data will be stored securely. To the possible extent, the intended use of the data – including a list of expected outputs and the organisation's dissemination policy – must be identified. Access to licensed datasets is only granted when there is a legally-registered sponsoring agency, e.g., government ministry, university, research centre or national or international organization" (*DuBo10*).

---

### Info-box - Conditions for accessing and using SUFs

Note: Items 1 to 8 below are similar to the conditions for use of public use files in the info-box above. Items 9 and 10 would have to be adapted in the case of a blanket agreement.

1. Data and other material provided by the NSO will not be redistributed or sold to other individuals, institutions or organisations without the NSO's written agreement.
2. Data will be used for statistical and scientific research purposes only. They will be employed solely for reporting aggregated information, including modelling, and not for investigating specific individuals or organisations.
3. No attempt will be made to re-identify respondents, and there will be no use of the identity of any person or establishment discovered inadvertently. Any such discovery will be reported immediately to the NSO.
4. No attempt will be made to produce links between datasets provided by the NSO or between NSO data and other datasets that could identify individuals or organisations.

5. Any books, articles, conference papers, theses, dissertations, reports or other publications employing data obtained from the NSO will cite the source, in line with the citation requirement provided with the dataset.
6. An electronic copy of all publications based on the requested data will be sent to the NSO.
7. The NSO and the relevant funding agencies bear no responsibility the data's use or for interpretation or inferences based upon it.
8. An electronic copy of all publications based on the requested data will be sent to the NSO.
9. The researcher's organisation must be identified, as must the principal and other researchers involved in using the data must be identified. The principal researcher must sign the licence on behalf of the organization. If the principal researcher is not authorized to sign on behalf of the receiving organization, a suitable representative must be identified.
10. The intended use of the data, including a list of expected outputs and the organisation's dissemination policy must be identified.

(Conditions 9 to 11 may be waived in the case of educational institutions)

Source: *DuBo10*

---

### *Blanket agreement to an organization*

"In the case of a blanket agreement, where it is agreed the data can be used widely but securely within the receiving organisation, the licence should ensure compliance, with a named individual formally assuming responsibility for this. Each additional user must be made aware of the terms and conditions that apply to data files: this can be achieved by having to sign an affidavit. Where such an agreement exists, with security in place, it is not necessary for users to destroy the data after use" (*DuBo10*). [Appendix B](#) provides an example of the formulation of such an agreement.

## 4.3 Conditions for microdata available in a controlled research data center

Access to microdata in research data centers is "used for particularly sensitive data or for more detailed data for which sufficient anonymisation to release them outside the NSO premises is not possible. These can be referred to also as data laboratories or research data centres. A [research data centre] may be located at the NSO headquarters or in major centres such as universities close to the research community. They are used to give researchers access to complete data files but without the risk of releasing confidential data. In a typical [research data centre], NSO staff supervise access and use of the data; the computers must not be able to communicate outside the [research data centre]; and the results obtained by the researchers must be screened for confidentiality by an NSO analyst before taken outside. A model of a data enclave access policy is provided in Appendix 2 [in *DuBo10*], and a model of a data enclave access request form is in Appendix 3 [in *DuBo10*]" (*DuBo10*).

Research data centers "have the advantage of providing access to detailed microdata but the disadvantage of requiring researchers to work at a different location. And they are expensive to set up and operate. It is, however, quite likely that many countries have used on-site researchers as a way of providing access to microdata. These researchers are sworn in under the statistics' acts in the same way as regular NSO employees. This approach tends to favour researchers who live near NSO headquarters." (*DuBo10*)

---

### Recommended Reading Material on Release Types

Dupriez, O., & Boyko, E. (2010). *Dissemination of Microdata Files; Principles, Procedures and Practices*. International Household Survey Network (IHSN).

---



## References



## 5.1 Types of disclosure

Measuring disclosure risk is an important part of the SDC process: risk measures are used to judge whether a data file is safe enough for release. Before measuring disclosure risk, we have to define what type of disclosure is relevant for the data at hand. The literature commonly defines three types of disclosure; we take these directly from *Lamb93* (see also *HDFG12*).

- **Identity disclosure**, which occurs if the intruder associates a known individual with a released data record. For example, the intruder links a released data record with external information, or identifies a respondent with extreme data values. In this case, an intruder can exploit a small subset of variables to make the linkage, and once the linkage is successful, the intruder has access to all other information in the released data related to the specific respondent.
- **Attribute disclosure**, which occurs if the intruder is able to determine some new characteristics of an individual based on the information available in the released data. Attribute disclosure occurs if a respondent is correctly re-identified and the dataset contains variables containing information that was previously unknown to the intruder. Attribute disclosure can also occur without identity disclosure. For example, if a hospital publishes data showing that all female patients aged 56 to 60 have cancer, an intruder then knows the medical condition of any female patient aged 56 to 60 in the dataset without having to identify the specific individual.
- **Inferential disclosure**, which occurs if the intruder is able to determine the value of some characteristic of an individual more accurately with the released data than would otherwise have been possible. For example, with a highly predictive regression model, an intruder may be able to infer a respondent's sensitive income information using attributes recorded in the data, leading to inferential disclosure.

SDC methods for microdata are intended to prevent identity and attribute disclosure. Inferential disclosure is generally not addressed in SDC in the microdata setting, since microdata is distributed precisely so that researchers can make statistical inference and understand relationships between variables. In that sense, inference cannot be likened to disclosure. Also, inferences are designed to predict aggregate, not individual, behavior, and are therefore usually poor predictors of individual data values.

## 5.2 Classification of variables

For the purpose of the SDC process, we use the classifications of variables described in the following paragraphs (see Fig. 5.1 for an overview). The initial classification of variables into identifying and non-identifying variables depends on the way the variables can be used by intruders for re-identification (*HDFG12*, *TeMK14*):

- **Identifying variables:** these contain information that can lead to the identification of respondents and can be further categorized as:
  - **Direct identifiers** reveal directly and unambiguously the identity of the respondent. Examples are names, passport numbers, social identity numbers and addresses. Direct identifiers should be removed from the dataset prior to release. Removal of direct identifiers is a straightforward process and always the first step in producing a safe microdata set for release. Removal of direct identifiers, however, is often not sufficient.
  - **Quasi-identifiers (or key variables)** contain information that, when combined with other quasi-identifiers in the dataset, can lead to re-identification of respondents. This is especially the case when they can be used to match the information with other external information or data. Examples of quasi-identifiers are race, birth date, sex and ZIP/postal codes, which might be easily combined or linked to publically available external information and make identification possible. The combinations of values of several quasi-identifiers are called keys (see also the section *Levels of Risk*). The values of quasi-identifiers themselves often do not lead to identification (e.g., male/female), but a combination of several values of quasi-identifier can render a record unique (e.g. male, 14 years, married) and hence identifiable. It is not generally advisable to simply remove quasi-identifiers from the data to solve the problem. In many cases, they will be important variables for any sensible analysis. In practice, any variable in the dataset could potentially be used as a quasi-identifier. SDC addresses this by identifying variables as quasi-identifiers and anonymizing them while still maintaining the information in the dataset for release.
- **Non-identifying** variables are variables that cannot be used for re-identification of respondents. This could be because these variables are not contained in any other data files or other external sources and are not observable to an intruder. Non-identifying variables are nevertheless important in the SDC process, since they may contain confidential/sensitive information, which may prove damaging should disclosure occur as a result of identity disclosure based on identifying variables.

These classifications of variables depend partially on the availability of external datasets that might contain information that, when combined with the current data, could lead to disclosure. The identification and classification of variables as quasi-identifiers depends, amongst others, on the availability of information in external datasets. An important step in the SDC process is to define a list of possible disclosure scenarios based on how the quasi-identifiers might be combined with each other and information in external datasets and then treating the data to prevent disclosure. We discuss disclosure scenarios in more detail in the section *Disclosure scenarios*.

For the SDC process, it is also useful to further classify the quasi-identifiers into **categorical**, **continuous** and **semi-continuous** variables. This classification is important for determining the appropriate SDC methods for that variable, as well as the validity of risk measures.

- **Categorical** variables take values over a finite set, and any arithmetic operations using them are generally not meaningful or not allowed. Examples of categorical variables are gender, region and education level.
- **Continuous** variables can take on an infinite number of values in a given set. Examples are income, body height and size of land plot. Continuous variables can be transformed into categorical variables by constructing intervals (such as income bands).<sup>1</sup>
- **Semi-continuous** variables are continuous variables that take on values that are limited to a finite set. An example is age measured in years, which could take on values in the set  $\{0, 1, \dots, 100\}$ . The finite nature of the

---

<sup>1</sup> Recoding a continuous variable is sometimes useful in cases where the data contains only a few continuous variables. We will see in the Section *Individual risk* that many methods used for risk calculation depend on whether the variables are categorical. We will also see that it is easier for the measurement of risk if the data contains only categorical or only continuous variables.

values for these variables means that they can be treated as categorical variables for the purpose of SDC.<sup>2</sup>

Apart from these classifications of variables, the SDC process further classifies variables according to their sensitivity or confidentiality. Both quasi-identifiers and non-identifying variables can be classified as **sensitive** (or confidential) or **non-sensitive** (or non-confidential). This distinction is not important for direct identifiers, since direct identifiers are removed from the released data.

- **Sensitive** variables contain confidential information that should not be disclosed without suitable treatment using SDC methods to reduce disclosure risk. Examples are income, religion, political affiliation and variables concerning health. Whether a variable is sensitive depends on the context and country: a certain variable can be considered sensitive in one country and non-sensitive in another.
- **Non-sensitive** variables contain non-confidential information on the respondent, such as place of residence or rural/urban residence. The classification of a variable as non-sensitive, however, does not mean that it does not need to be considered in the SDC process. Non-sensitive variables may still serve as quasi-identifiers when combined with other variables or other external data.

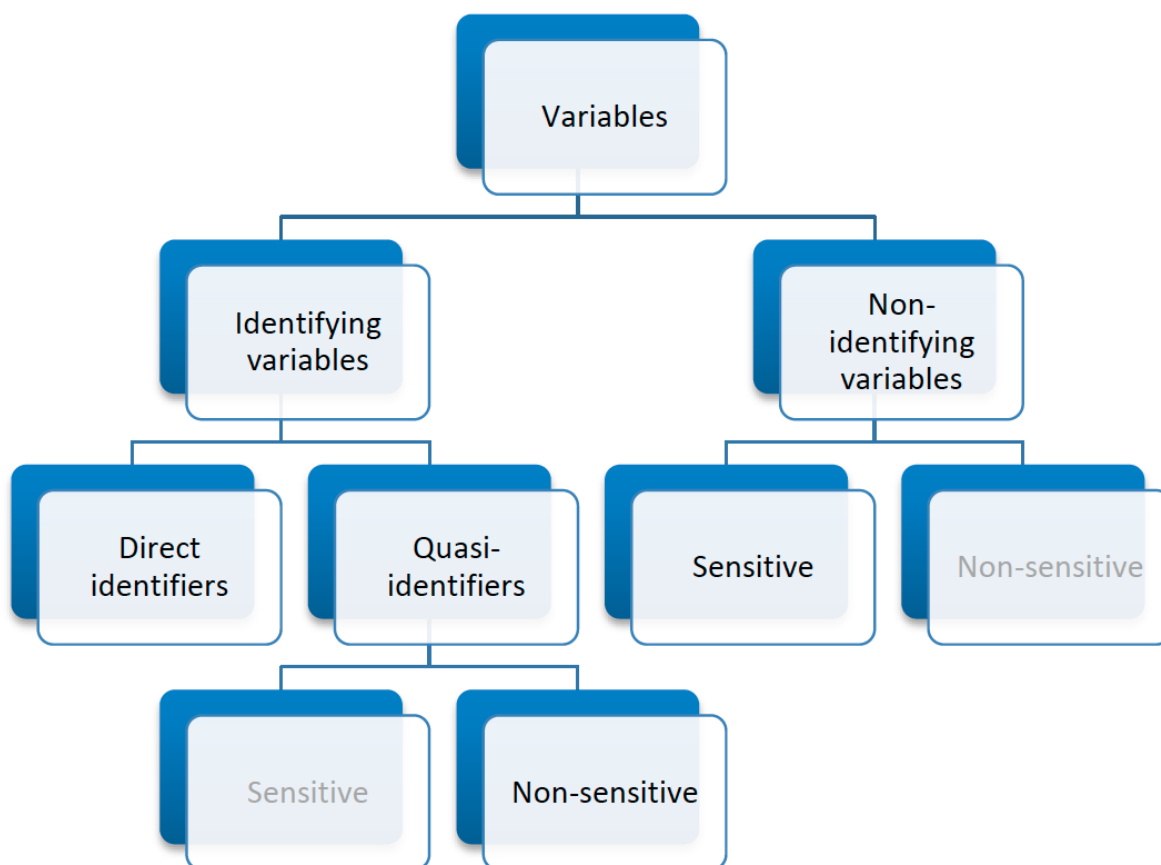


Fig. 5.1: Classification of variables

### 5.3 Disclosure scenarios

Evaluation of disclosure risk is carried out with reference to the available data sources in the environment where the dataset is to be released. In this setting, disclosure risk is the possibility of correctly re-identifying an individual

<sup>2</sup> This is discussed in greater detail in the following sections. In cases where the number of possible values is large, recoding the variable, or parts of the set it takes values on, to obtain fewer distinct values is recommended.

in the released microdata file by matching their data to an external file based on a set of quasi-identifiers. The risk assessment is done by identifying so-called disclosure or intrusion scenarios. A disclosure scenario describes the information potentially available to the intruder (e.g., census data, electoral rolls, population registers or data collected by private firms) to identify respondents and the ways such information can be combined with the microdata set to be released and used for re-identification of records in the dataset. Typically, these external datasets include direct identifiers. In that case, the re-identification of records in the released dataset leads to identity and, possibly, attribute disclosure. The main outcome of the evaluation of disclosure scenarios is the identification of a set of quasi-identifiers (i.e., key variables) that need to be treated during the SDC process (see *ELMP10*).

An example of a disclosure scenario could be the spontaneous recognition of a respondent by a researcher. For instance, while going through the data, the researcher recognizes a person with an unusual combination of the variables age and marital status. Of course, this can only happen if the person is well-known or is known to the researcher. Another example of a disclosure scenario for a publicly available file would be if variables in the data could be linked to a publically available electoral register. An intruder might try matching the entire dataset with individuals in the register. However, this might be difficult and take specialized expertise, or software, and other conditions have to be fulfilled. Examples are that the point in time the datasets were collected should approximately match and the content of the variables should be (nearly) identical. If these conditions are not fulfilled, exact matching is much less likely.

---

**Note:** Not all external data is necessarily in the public domain. Also privately owned datasets or datasets which are not released should be taken into consideration for determining the suitable disclosure scenario.

---

---

### Info-box - Disclosure scenarios and different release types

A dataset can have more than one disclosure scenario. Disclosure scenarios also differ depending on the data access type that the data will be released under; for example, Public Use Files (PUF) or Scientific Use Files (SUF, also known as licensed) or in a data enclave. The required level of protection, the potential avenues of disclosure as well as the availability of other external data sources differ according to the access type under which the data will be released. For example, the user of a Scientific Use File (SUF) might be contractually restricted by an agreement as to what they are allowed to do with the data, whereas a Public Use File (PUF) might be freely available on the internet under a much looser set of conditions. PUFs will in general require more protection than SUFs and SUFs will require more protection than those files only released in an data enclave. Disclosure scenarios should be developed with all of this in mind.

---

The evaluation of disclosure risk is based on the quasi-identifiers, which are identified in the analysis of disclosure risk scenarios. The disclosure risk directly depends on the inclusion or exclusion of variables in the set of quasi-identifiers chosen. This step in the SDC process (making the choice of quasi-identifiers) should therefore be approached with great thought and care. We will see later, as we discuss the steps in the SDC process in more detail, that the first step for any agency is to undertake an exercise in which an inventory is compiled of all datasets available in the country. Both datasets released by the national statistical office and from other sources are considered and their availability to intruders as well as the variables included in these datasets is analyzed. It is this information that will serve as a key metric when deciding which variables to choose as potential identifiers, as well as dictate the level of SDC and methods needed.

## 5.4 Levels of risk

With microdata from surveys and censuses, we often have to be concerned about disclosure at the individual or unit level, i.e., identifying individual respondents. Individual respondents are generally natural persons, but can also be units, such as companies, schools, health facilities, etc. Microdata files often have a hierarchical structure where individual units belong to groups, e.g., people belong to households. The most common hierarchical structure in microdata is the household structure in household survey data. Therefore, in this guide, we sometimes call disclosure risk for data with a hierarchical structure “household risk”. The concepts, however, apply equally to establishment

data and other data with hierarchical structures, such as school data with pupils and teachers or company data with employees.

We will see that this hierarchical structure is important to take into consideration when measuring disclosure risk. For hierarchical data, information collected at the higher hierarchical level (e.g., household level) would be the same for all individuals in the group belonging to that higher hierarchical level (e.g., household).<sup>3</sup> Some typical examples of variables that would have the same values for all members of the same higher hierarchical unit are, in the case of households, those relating to housing and household income. These variables differ from survey to survey and from country to country.<sup>4</sup> This hierarchical structure creates a further level of disclosure risk for two reasons:

1. if one individual in the household is re-identified, the household structure allows for re-identification of the other household members in the same household,
2. values of variables for other household members that are common for all household members can be used for re-identification of another individual of the same household. This is discussed in more detail in the Section *Household Risk*.

Next, we first discuss risk measures used to evaluate disclosure risk in the absence of a hierarchical structure. This includes risk measures that seek to aggregate the individual risk for all individuals in the microdata file; the objective is to quantify a global disclosure risk measure for the file. We then discuss how risk measures change when taking the hierarchical structure of the data into account.

We will also discuss how risk measures differ for categorical and continuous key variables. For categorical variables, we will use the concept of uniqueness of combinations of values of quasi-identifiers (so-called “keys”) used to identify individuals at risk. The concept of uniqueness, however, is not useful for continuous variables, since it is likely that all or many individuals will have unique values for that variable, by definition of a continuous variable. Risk measures for categorical variables are generally a priori measures, i.e., they can be evaluated before applying anonymization methods since they are based on the principle of uniqueness. Risk measures for continuous variables are a posteriori measures; they are based on comparing the microdata before and after anonymization and are, for example, based on the proximity of observations between the original and the treated (anonymized) datasets.

Files that are limited to only categorical or only continuous key variables are easiest for risk measurement. We will see in later sections that, in cases where both types of variables are present, recoding of continuous variables into categories is one approach to use to simplify the SDC process, but we will also see that from a utility perspective this may not be desirable. An example might be the use of income quintiles instead of the actual income variables. We will see that measuring the risk of disclosure based on the categorical and continuous variables separately is generally not a valid approach.

The risk measures discussed in the next section are based on several assumptions. In general, these measures rely on quite restrictive assumptions and will often lead to conservative risk estimates. These conservative risk measures may overstate the risk as they assume a worst-case scenario. Two assumptions should, however, be fulfilled for the risk measures to be valid and meaningful; the microdata should be a sample of a larger population (no census) and the sampling weights should be available. The Section *Special case: census data* briefly discusses how to deal with census data.

## 5.5 Individual risk

### 5.5.1 Categorical key variables and frequency counts

The main focus of risk measurement for categorical quasi-identifiers is on identity disclosure. Measuring disclosure risk is based on the evaluation of the probability of correct re-identification of individuals in the released data. We

<sup>3</sup> Besides variables collected at the higher hierarchical level, also variables collected at the lower level but with no (or little) variation within the groups formed by the hierarchical structure should be treated as higher level variables. An example could be mother tongue, where most households are monolingual, but the variable is collected at the individual level.

<sup>4</sup> Religion, for example, can be shared by all household members in some countries, whereas in other countries this variable is measured at the individual level and mixed-religion households exist.

use measures based on the actual microdata to be released. In general, the rarer a combination of values of the quasi-identifiers (i.e., key) of an observation in the sample, the higher the risk of identity disclosure. An intruder that tries to match an individual who has a relatively rare key within the sample data with an external dataset in which the same key exists will have a higher probability of finding a correct match than when a larger number of individuals share the same key. This can be illustrated with the following example that is illustrated in Table 5.1.

Table 5.1 shows values for 10 respondents for the quasi-identifiers “residence”, “gender”, “education level” and “labor status”. In the data, we find seven unique combinations of values of quasi-identifiers (i.e., patterns or keys) of the four quasi-identifiers. Examples of keys are {‘urban’, ‘female’, ‘secondary incomplete’, ‘employed’} and {‘urban’, ‘female’, ‘primary incomplete’, ‘non-LF’}. Let  $f_k$  be the sample frequency of the  $k^{\text{th}}$  key, i.e., the number of individuals in the sample with values of the quasi-identifiers that coincide with the  $k^{\text{th}}$  key. This would be 2 for the key {urban, female, secondary incomplete, employed}, since this key is shared by individuals 1 and 2 and 1 for the key {‘urban’, ‘female’, ‘primary incomplete’, ‘non-LF’}, which is unique to individual 3. By definition,  $f_k$  is the same for each record sharing a particular key.

The fewer the individuals with whom an individual shares his or her combination of quasi-identifiers, the more likely the individual is to be correctly matched in another dataset that contains these quasi-identifiers. Even when direct identifiers are removed from the dataset, that individual has a higher disclosure risk than others, assuming that their sample weights are the same. Table 5.1 reports the sample frequencies  $f_k$  of the keys for all individuals. Individuals with the same keys have the same sample frequency. If  $f_k = 1$ , this individual has a unique combination of values of quasi-identifiers and is called “sample unique”. The dataset in Table 5.1 contains four sample uniques. Risk measures are based on this sample frequency.

Table 5.1: Example dataset showing sample frequencies, population frequencies and individual disclosure risk

No	Residence	Gender	Education level	Labor status	Weight	$f_k$	$F_k$	risk
1	Urban	Female	Secondary incomplete	Employed	180	2	360	0.0054
2	Urban	Female	Secondary incomplete	Employed	180	2	360	0.0054
3	Urban	Female	Primary incomplete	Non-LF	215	1	215	0.0251
4	Urban	Male	Secondary complete	Employed	76	2	152	0.0126
5	Rural	Female	Secondary complete	Unemployed	186	1	186	0.0282
6	Urban	Male	Secondary complete	Employed	76	2	152	0.0126
7	Urban	Female	Primary complete	Non-LF	180	1	180	0.0290
8	Urban	Male	Post-secondary	Unemployed	215	1	215	0.0251
9	Urban	Female	Secondary incomplete	Non-LF	186	2	262	0.0074
10	Urban	Female	Secondary incomplete	Non-LF	76	2	262	0.0074

For sample data, it is more interesting to look at  $F_k$ , the population frequency of a combination of quasi-identifiers (key)  $k$ , which is the number of individuals in the population with the key that corresponds to key  $k$ . The population frequency is unknown if the microdata is a sample and not a census. Under certain assumptions, the expected value of the population frequencies can be computed using the sample design weight  $w_i$  (in a simple sample, this is the inverse of the inclusion probability) for each individual  $i$

$$F_k = \sum_{i | \text{key of individual } i \text{ corresponds to key } k} w_i$$

$F_k$  is the sum of the sample weights of all records with the same key  $k$ . Hence, like  $f_k$ ,  $F_k$  is the same for each record with key  $k$ . The risk of correct re-identification is the probability that the key is matched to the correct individual in the population. Since every individual in the sample with key  $k$  corresponds to  $F_k$  individuals in the population, the probability of correct re-identification is  $1/F_k$ . This is the probability of re-identification in the worst-case scenario and can be interpreted as disclosure risk. Individuals with the same key have the same frequencies, i.e., the frequency of the key.

If  $F_k = 1$ , the key  $k$  is both a sample and a population unique and the disclosure risk would be 1. Population uniques are an important factor to consider when evaluating risk, and deserve special attention. Table 5.1 also shows  $F_k$  for



the example dataset. This is further discussed in the case studies the Section [Case Studies](#).

In practice, this approach leads to conservative risk estimates, as it does not adequately take the sampling methods into account. In this case, the estimates of re-identification risk may be estimated too high. If this overestimated risk is used, the data may be overprotected (i.e., information loss will be higher than was necessary) when applying SDC measures. Instead, a Bayesian approach to risk measurement is recommended, where the posterior distribution of  $F_k$  is used (see e.g., *HDFG12*) to estimate an individual risk measure  $r_k$  for each key  $k$ .

ADD: more info on the exact computation of  $r_k$

The risk measure  $r_k$  is, as  $f_k$  and  $F_k$ , the same for all individuals sharing the same pattern of values of key variables and is referred to as individual risk. The values  $r_k$  can also be interpreted as the probability of disclosure for the individuals or as the probability for a successful match with individuals chosen at random from an external data file with the same values of the key variables. This risk measure is based on certain assumptions<sup>5</sup>, which are strict and may lead to a relatively conservative risk measure. The risk measures are also presented in [Table 5.1](#).

The main factors influencing the individual risk are the sample frequencies  $f_k$  and the sampling design weights  $w_i$ . If an individual is at relatively high risk of disclosure, in our example this would be individuals 3, 5, 7 and 8 in [Table 5.1](#), the probability that a potential intruder correctly matches these individuals with an external data file is high **relative to the other individuals in the released data**. In our example, the reason for the high risk is the fact that these individuals are sample uniques ( $f_k = 1$ ). This risk is the worst-case scenario risk and does not imply that the person will be re-identified with certainty with this probability. For instance, if an individual included in the microdata is not included in the external data file, the probability for a correct match is zero. Nevertheless, the risk measure computed based on the frequencies will be positive.

## 5.5.2 $k$ -anonymity

The risk measure  $k$ -anonymity is based on the principle that, in a safe dataset, the number of individuals sharing the same combination of values (keys) of categorical quasi-identifiers should be higher than a specified threshold  $k$ .  $k$ -anonymity is a risk measure based on the microdata to be released, since it only takes the sample into account. An individual violates  $k$ -anonymity if the sample frequency count  $f_k$  for the key  $k$  is smaller than the specified threshold  $k$ . For example, if an individual has the same combination of quasi-identifiers as two other individuals in the sample, these individuals satisfy 3-anonymity but violate 4-anonymity. In the dataset in [Table 5.1](#), six individuals satisfy 2-anonymity and four violate 2-anonymity. The individuals that violate 2-anonymity are sample uniques. The risk measure is the number of observations that violates  $k$ -anonymity for a certain value of  $k$ , which is

$$\sum_i I(f_k < k),$$

where  $I$  is the indicator function and  $i$  refers to the  $i^{\text{th}}$  record. This is simply a count of the number of individuals with a sample frequency of their key lower than  $k$ . The count is higher for larger  $k$ , since if a record violates  $k$ -anonymity, it also violates  $(k + 1)$ -anonymity. The risk measure  $k$ -anonymity does not consider the sample weights, but it is important to consider the sample weights when determining the required level of  $k$ -anonymity. If the sample weights are large, one individual in the dataset represents more individuals in the target population, the probability of a correct match is smaller, and hence the required threshold can be lower. Large sample weights go together with smaller datasets. In a smaller dataset, the probability to find another record with the same key is smaller than in a larger dataset. This probability is related to the number of records in the population with a particular key through the sample weights.

<sup>5</sup> The assumptions for this risk measure are strict and the risk is estimated in many cases higher than the actual risk. Among other assumptions, it is assumed that all individuals in the sample are also included in the external file used by the intruder to match against. If this is not the case, the risk is much lower; if the individual in the released file is not contained in the external file, the probability of a correct match is zero. Other assumptions are that the files contain no errors and that both sets of data were collected simultaneously, i.e. they contain the same information. These assumptions will often not hold generally, but are necessary for computation of a measure. An example of a violation of the last assumptions is could occur if datasets are collected at different points in time and records have changed. This could happen when people move or change jobs and makes correct matching impossible. The assumptions assume a worst-case scenario.

Assuming that the example dataset in Table 5.1 represents the full sample, we find that four observations violate 2-anonymity ( $f_k < 2$ ) and all 10 observations violate 3-anonymity ( $f_k < 3$ ). The relative number of 2-anonymity and 3-anonymity violators are resp. 40% and 100%. For other levels of  $k$ -anonymity, it is possible to compute the number of violating individuals by using the sample frequency counts.  $k$  can be replaced with any required threshold. The choice of the required threshold that all individuals in the microdata file should satisfy depends on many factors and is discussed further in the Section Local suppression on local suppression. In many institutions, typically required thresholds for  $k$ -anonymity are 3 and 5.

It is important to note that missing values are treated as if they were any other value. Two individuals with keys {‘Male’, missing, ‘Employed’} and {‘Male’, ‘Secondary complete’, ‘Employed’} share the same key, and similarly, {‘Male’, missing, ‘Employed’} and {‘Male’, ‘Secondary incomplete’, ‘Employed’} also share the same key. Therefore, the missing value in the first key is first interpreted as ‘Secondary complete’, and then as ‘Secondary incomplete’. This is illustrated in Table 5.2.

**Note:** The sample frequency of the third record is 3, since it is regarded to share its key both with the first and second record.

This principle is used when applying local suppression to achieve a certain level of  $k$ -anonymity (see the Section Local suppression) and is based on the fact that the value NA could replace any value.

ADD: parameter alpha and treating missing values

Table 5.2: Example dataset to illustrate the effect of missing values on  $k$ -anonymity

No	Gender	Education level	Labor status	$f_k$
1	Male	Secondary complete	Employed	2
2	Male	Secondary incomplete	Employed	2
3	Male	NA	Employed	3

If a dataset satisfies  $k$ -anonymity, an intruder will always find at least  $k$  individuals with the same combination of quasi-identifiers.  $k$ -anonymity is often a necessary requirement for anonymization for a dataset before release, but is not necessarily a sufficient requirement. The  $k$ -anonymity measure is only based on frequency counts and does not take (differences in) sample weights into account. Often  $k$ -anonymity is achieved by first applying recoding and subsequently applying local suppression, and in some cases by microaggregation, before using other risk measures and disclosure methods to further reduce disclosure risk. These methods are discussed in the Section Anonymization methods.

### 5.5.3 $l$ -diversity

$k$ -anonymity has been criticized for not being restrictive enough. Sensitive information might be disclosed even if the data satisfies  $k$ -anonymity. This might occur in cases where the data contains sensitive (non-identifying) categorical variables that have the same value for all individuals that share the same key. Examples of such sensitive variables are those containing information on an individual’s health status. Table 5.3 illustrates this problem by using the same data as previously used, but adding a sensitive variable, “health”. The first two individuals satisfy 2-anonymity for the key variables “residence”, “gender”, “education level” and “labor status”. This means that an intruder will find at least two individuals when matching to the released microdata set based on those four quasi-identifiers. Nevertheless, if the intruder knows that someone belongs to the sample and has the key {‘Urban’, ‘Female’, ‘Secondary incomplete’ and ‘Employed’}, with certainty the health status is disclosed (‘yes’), because both observations with this key have the same value. This information is thus disclosed without the necessity to match exactly to the individual. This is not the case for the individuals with the key {‘Urban’, ‘Male’, ‘Secondary complete’, ‘Employed’}. Individuals 4 and 6 have different values (‘yes’ and ‘no’) for “health”, and thus the intruder would not gain information about the health status from this dataset by matching an individual to one of these individuals.

Table 5.3:  $l$ -diversity illustration

No	Residence	Gender	Education level	Labor status	Health	$f_k$	$F_k$	$l$ -diversity
1	Urban	Female	Secondary incomplete	Employed	yes	2	360	1
2	Urban	Female	Secondary incomplete	Employed	yes	2	360	1
3	Urban	Female	Primary incomplete	Non-LF	yes	1	215	1
4	Urban	Male	Secondary complete	Employed	yes	2	152	2
5	Rural	Female	Secondary complete	Unemployed	yes	1	186	1
6	Urban	Male	Secondary complete	Employed	no	2	152	2
7	Urban	Female	Primary complete	Non-LF	no	1	180	1
8	Urban	Male	Post-secondary	Unemployed	yes	1	215	1
9	Urban	Female	Secondary incomplete	Non-LF	no	2	262	2
10	Urban	Female	Secondary incomplete	Non-LF	yes	2	262	2

The concept of (distinct)  $l$ -diversity addresses this shortcoming of  $k$ -anonymity (see *MKGV07*). A dataset satisfies  $l$ -diversity if for every key  $k$  there are at least  $l$  different values for each of the sensitive variables. In the example, the first two individuals satisfy only 1-diversity, individuals 4 and 6 satisfy 2-diversity. The required level of  $l$ -diversity depends on the number of possible values the sensitive variable can take. If the sensitive variable is a binary variable, the highest level of  $l$ -diversity that can be achieved is 2. A sample unique will always only satisfy 1-diversity.

ADD: recursive  $l$ -diversity and parameter `l_rekurs_c` in `ldiversity()`

$l$ -diversity is useful if the data contains categorical sensitive variables that are not quasi-identifiers themselves. It is not possible to select quasi-identifiers to calculate the  $l$ -diversity.  $l$ -diversity has to be calculated for each sensitive variable separately.

## 5.6 Special Uniques Detection Algorithm (SUDA)

The previously discussed risk measures depend on identifying key variables for which there may be information available from other sources or other datasets, and which, when combined with the current data, may lead to re-identification. In practice, however, it might not always be possible to conduct an inventory of all available datasets and their variables and thus assess all known external linkages and risks.

To overcome this, an alternative heuristic measure based on special uniques has been developed to determine the riskiness of a record, which leads to a SUDA metric or score (see *ELMF02*). These measures are based on the search for special uniques. To find these special uniques, algorithms, called SUDA (Special Uniqueness Detection Algorithm), have been developed. SUDA algorithms are based on the concept of special uniqueness, which is introduced in the next subsection. Since this is a heuristic approach, its performance is only tested in actual datasets, which is done in *ELMF02* for UK census data. These tests have shown that the performance of the algorithm leads to good risk estimates for these test datasets.

### 5.6.1 Sample unique vs. special unique

The previous measures of risk focused on the uniqueness of the key of a record in the dataset. Table 5.4 reproduces the data from Table 5.1. The sample dataset has 10 records and four pre-determined quasi-identifiers {"Residence", "Gender", "Education level" and "Labor status"}. Given the four quasi-identifiers, we have seven distinct patterns in those key variables, or keys (e.g., {'Urban', 'Female', 'Secondary incomplete', 'Employed'}). The sample frequency counts of the first and second records equal 2, because the two records share the same pattern (i.e., {'Urban', 'Female', 'Secondary incomplete', 'Employed'}). Record 3 is a sample unique because it is the only individual in the sample who is a female living in an urban area who is employed without completing primary school. Similarly, records 5, 7 and 8 are sample uniques, because they possess distinct patterns with respect to the four key variables.

Table 5.4: Sample uniques and special uniques

No	Residence	Gender	Education level	Labor status	Weight	$f_k$	$F_k$	risk
1	Urban	Female	Secondary incomplete	Employed	180	2	360	0.0054
2	Urban	Female	Secondary incomplete	Employed	180	2	360	0.0054
3	Urban	Female	Primary incomplete	Non-LF	215	1	215	0.0251
4	Urban	Male	Secondary complete	Employed	76	2	152	0.0126
5	Rural	Female	Secondary complete	Unemployed	186	1	186	0.0282
6	Urban	Male	Secondary complete	Employed	76	2	152	0.0126
7	Urban	Female	Primary complete	Non-LF	180	1	180	0.0290
8	Urban	Male	Post-secondary	Unemployed	215	1	215	0.0251
9	Urban	Female	Secondary incomplete	Non-LF	186	2	262	0.0074
10	Urban	Female	Secondary incomplete	Non-LF	76	2	262	0.0074

In addition to the records 3, 5, 7 and 8 in Table 5.4 being sample uniques with respect to the key variable set {"Residence", "Gender", "Education level", "Labor status"}, we can find unique patterns in these records without even having to consider the complete set of key variables. For instance, a unique pattern can be found in record 5 when we look only at the variables "Education level" and "Labor status" ({'Secondary complete', 'Unemployed'}). While the values {'Secondary complete'} and {'Unemployed'} are not unique in the sample, the combination of them, {'Secondary complete', 'Unemployed'} makes record 5 unique. This variable subset is referred to as the Minimal Sample Unique (MSU) as any smaller subset of this set of variables is not unique (in this case {'Secondary complete'} and {'Unemployed'}). It is an MSU of size 2. This holds as well for three other combinations in record 5, i.e., {'Female', 'Unemployed'} and {'Female', 'Secondary Complete'}, which are also MSUs of size 2 and {'Rural'} of size 1. In total, record 5 has four MSUs<sup>7</sup>. To determine if a set is an MSU of size  $k$ , we check whether it fulfills the minimal requirement. It suffices to check whether all subsets of size  $k - 1$  of the MSU are unique. If any of these subsets are also unique in the sample, the set found may be a sample unique, but violates the minimal requirement and is hence not an MSU. The unique subset of size  $k - 1$  could be a MSU. In our example, to determine if the MSU {'Secondary complete', 'Unemployed'} is a MSU, we checked as to whether its subsets {'Secondary complete'} and {'Unemployed'} were not unique in the sample. By definition, only sample uniques can be special uniques.

The SUDA algorithm identifies all the MSUs in the sample, which in turn are used to assign a SUDA score to each record. This score indicates how "risky" a record is. The potential risk of the records is determined based on two observations:

- The smaller the size of the MSU within a record (i.e., the fewer variables are needed to reach uniqueness), the greater the risk of the record
- The larger the number of MSUs possessed by a record, the greater the risk of the record

A record is defined as a special unique if it is a sample unique both on the complete set of quasi-identifiers (e.g., in the data in Table 5.4, the variables "Residence", "Gender", "Education level" and "Labor status") and simultaneously has at least one MSU (EISD98). Special uniques can be classified according to the number and size of subsets that are MSUs. Research has shown that special uniques are more likely to be population uniques than random uniques (EIMF02) and are thus relevant for risk assessment.

## 5.6.2 Calculating SUDA scores

The SUDA algorithm is used to search for MSUs in the data among the sample uniques to determine which sample uniques are also special uniques i.e., have subsets that are also unique (see Elliot et al., 2005). First the SUDA algorithm is used to identify the MSUs for each sample unique. To simplify the search and because smaller subsets are more important for disclosure risk, the search is limited to a maximum subset size. Subsequently, a score is assigned to each individual, which ranks the individuals according to their level of risk.

<sup>7</sup> There are more combinations of quasi-identifiers that make record 5 unique (e.g., {'Rural', 'Female'} and {'Female', 'Secondary Complete', 'Unemployed'}). These combinations, however, are not considered MSUs because they do not fulfill the **minimal** subset requirement. They contain subsets that are MSUs.

For each MSU of size  $k$  contained in a given record, a score is computed by  $\prod_{i=k}^M (ATT - i)$ , where  $M$  is the user-specified maximum size of MSUs<sup>8</sup>, and  $ATT$  is the total number of attributes or variables in the dataset. By definition, the smaller the size  $k$  of the MSU, the larger the score for the MSU, which reflects greater risk (see *EMMG05*). The final SUDA score for each record is computed by adding the scores for each MSU in the record. In this way, records with more MSUs are assigned a higher SUDA score, which also reflects the higher risk. The SUDA score ranks the individuals according to their level of risk. The higher the SUDA score, the riskier the sample unique.

#### Calculating SUDA scores – a simplified example

To illustrate how SUDA scores are calculated, we compute the SUDA scores for the sample uniques in the data in Table 5.5, which replicates the data from Table 5.5. Record 5 contains four MSUs: {Rural} of size 1, and {‘Secondary Complete’, ‘Unemployed’}, {‘Female’, ‘Unemployed’} and {Female, Secondary Complete} of size 2. Suppose the maximum size of MSUs we search for in the data,  $M$ , is set at 3. Knowing that,  $ATT$ , the number of selected key variables in the dataset, is 4; the score assigned to {Rural} is computed by  $\prod_{i=1}^3 (4 - i) = 3 * 2 * 1 = 6$ ; and the score assigned to {Secondary complete, Unemployed}, {Female, Unemployed} and {Female, Secondary Complete} is  $\prod_{i=2}^3 (4 - i) = 2 * 1 = 2$ . The SUDA score for the fifth record in Table 5.5 is then  $6 + 2 + 2 + 2 = 12$ , which is the sum of these four scores per MSU. The SUDA scores for the other sample uniques are computed accordingly<sup>9</sup>. The values that are in the MSUs in the sample uniques are shaded in Table 5.5. Records that are not sample uniques ( $f_k > 1$ ) cannot be special uniques and are assigned the score 0.

Table 5.5: Illustrating the calculation of SUDA and DIS-SUDA scores

No	Residence	Gender	Education level	Labor status	Weight	$f_k$	SUDA score	DIS-SUDA
1	Urban	Female	Secondary incomplete	Employed	180	2	0	0.0000
2	Urban	Female	Secondary incomplete	Employed	180	2	0	0.0000
3	Urban	Female	Primary incomplete	Non-LF	215	1	6	0.0051
4	Urban	Male	Secondary complete	Employed	76	2	0	0.0000
5	Rural	Female	Secondary complete	Unemployed	186	1	12	0.0107
6	Urban	Male	Secondary complete	Employed	76	2	0	0.0000
7	Urban	Female	Primary complete	Non-LF	180	1	6	0.0051
8	Urban	Male	Post-secondary	Unemployed	215	1	10	0.0088
9	Urban	Female	Secondary incomplete	Non-LF	186	2	0	0.0000
10	Urban	Female	Secondary incomplete	Non-LF	76	2	0	0.0000

ADD: compare with values in code block 4.7

To estimate record-level disclosure risks, SUDA scores can be used in combination with the Data Intrusion Simulation (DIS) metric (*ElMa03*), a method for assessing disclosure risks for the entire dataset (i.e., file-level disclosure risks). Roughly speaking, the DIS-SUDA method distributes the file-level risk measure generated by the DIS metric between records according to the SUDA scores of each record. This way, SUDA scores are calibrated against a consistent measure to produce the DIS-SUDA scores, which provide the record-level disclosure risk. These scores are used to compute the conditional probability that a unique match found by an intruder between the sample unique in the released microdata and an external data source is also a correct match, and hence a successful disclosure. The DIS-SUDA measure can be computed in *sdcMicro*. Since the DIS score is a probability, its values are in the interval  $[0, 1]$ . A full description of the DIS-SUDA method is provided by *ElMa03*.

<sup>8</sup> OECD, <http://stats.oecd.org/glossary>

<sup>9</sup> The third record has one MSU, {‘Primary incomplete’}; the seventh record has one MSU, {‘Primary complete’}; and the eighth record has three MSUs, {‘Urban, Unemployed’}, {‘Male, Unemployed’} and {‘Post-secondary’}.

Note that unlike the risk methods discussed earlier, the DIS-SUDA score does not fully account for the sampling weights. Risk measures based on the previous methods (i.e., negative binomial models) will in general have lower risks for those records with greater sampling weight, given the same sample frequency count, than those measured using DIS-SUDA. Therefore, instead of replacing the risk measures introduced in the previous section, the SUDA scores and DIS-SUDA approach should be used as a complementary method. As mentioned earlier, DIS-SUDA is particularly useful in situations where taking an inventory of all already available datasets and their variables is difficult.

ADD: why use SUDA: (Typically, after applying SDC methods, one would recalculate the SUDA scores and compare them to the original values. One way to quickly see the differences would be to rerun these visualizations and compare them to the base for risk changes.)

ADD: reference to fig2

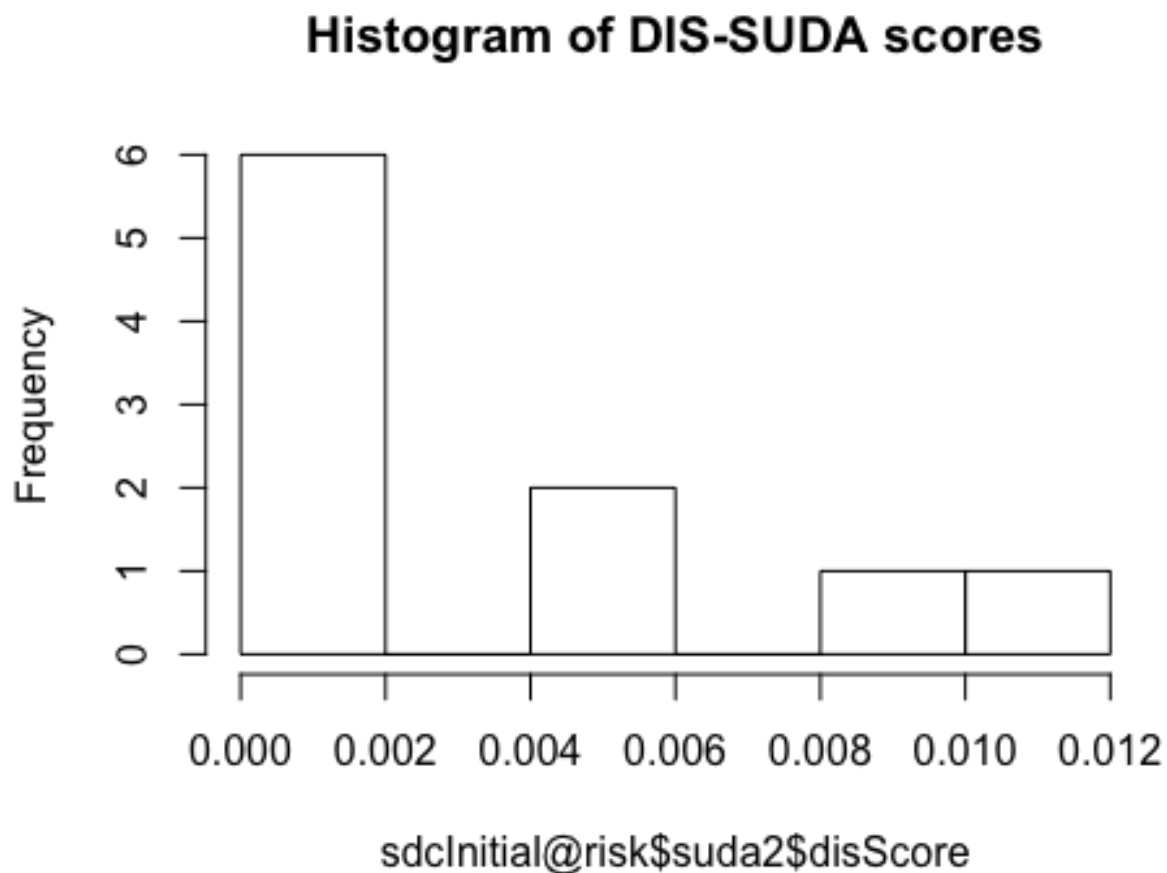


Fig. 5.2: Visualizations of DIS-SUDA scores

## 5.7 Risk measures for continuous variables

The principle of rareness or uniqueness of combinations of quasi-identifiers (keys) is not useful for continuous variables, because it is likely that all or many individuals will have unique keys. Therefore, other approaches are exploited for measuring the disclosure risk of continuous variables. These methods are based on uniqueness of the values in the



neighborhood of the original values. The uniqueness is defined in different ways: in absolute terms (interval measure) or relative terms (record linkage). Most measures are a posteriori measures: they are evaluated after anonymization of the raw data, compare the treated data with the raw data and evaluate for each individual the distance between the values in the raw and the treated data. This means that these methods are not useful for identifying individuals at risk within the raw data, but rather show the distance/difference between the dataset before and after anonymization and can therefore be interpreted as evaluation of the anonymization method. For that reason, they resemble the information loss measures discussed in the Section [Measuring utility and information loss](#). Finally, risk measures for continuous quasi-identifiers are also based on outlier detection. Outliers play an important role in the re-identification of these records.

### 5.7.1 Record linkage

Record linkage is an a posteriori method that evaluates the number of correct linkages when linking the perturbed values with the original values. The linking algorithm is based on the distance between the original and the perturbed values (i.e., distance-based record linkage). The perturbed values are matched with the closest individual. It is important to note that this method does not give information on the initial risk, but is rather a measure to evaluate the perturbation algorithm (i.e., it is designed to indicate the level of uncertainty introduced into the variable by counting the number of records that could be correctly matched).

Record linkage algorithms differ with respect to which distance measure is used. When a variable has very different scaling than other continuous variables in the dataset, rescaling the variables before using record linkage is recommended. Very different scales may lead to undesired results when measuring the multivariate distance between records based on several continuous variables. Since these methods are based on both the raw data and treated data, examples of their applications require the introduction of SDC methods and are therefore postponed to the case studies in the Section [Case Studies](#).

Besides distance-based record linkage, another method for linking is probabilistic record linkage (see *DoTo03*). The literature shows, however, that results from distance-based record linkage are better than the results from probabilistic record linkage. Individuals in the treated data that are linked to the correct individuals in the raw data are considered at risk of disclosure.

### 5.7.2 Interval measure

Successful application of an SDC method should result in perturbed values that are considered not too close to their initial values; if the value is relatively close, re-identification may be relatively easy. In the application of interval measures, intervals are created around each perturbed value and then a determination is made as to whether the original value of that perturbed observation is contained in this interval. Values that are within the interval around the initial value after perturbation are considered too close to the initial value and hence unsafe and need more perturbation. Values that are outside of the intervals are considered safe. The size of the intervals is based on the standard deviation of the observations and a scaling parameter  $k$ .

ADD: size of interval

The size of the intervals is  $k$  times the standard deviation. The larger  $k$ , the larger the intervals are, and hence the larger the number of observations within the interval constructed around their original values and the higher the risk measure.

For most values, this is a satisfactory approach. It is not a sufficient measure for outliers, however. After perturbation, outliers will stay outliers and are easily re-identifiable, even if they are sufficiently far from their initial values. Therefore, outliers should be treated with caution.

### 5.7.3 Outlier detection

Outliers are important for measuring re-identification risk in continuous microdata. Continuous data are often skewed, especially right-skewed. This means that there are a few outliers with very high values relative to the other observations of the same variable. Examples are income in household data, where only few individuals/households may have very high incomes, or turnover data for firms that are much larger than other firms in the sample are. In cases like these, even if these values are perturbed, it may still be easy to identify these outliers, since they will stay the largest values even after perturbation. (The perturbation will have created uncertainty as to the exact value, but because the value started out so much further away from other observations, it may still be easy to link to the high-income individual or very large firm.) Examples would be the only doctor in a geographical area with a high income or one single large firm in one industry type. Therefore, identifying outliers in continuous data is an important step when identifying individuals at high risk. In practice, identifying the values of a continuous variable that are larger than a predetermined  $p\%$ -percentile might help identify outliers, and thus units at greater risk of identification. The value of  $p$  depends on the skewness of the data.

A second approach for outlier detection is a posteriori measure comparing the treated and raw data. An interval is constructed around the perturbed values as described in the previous section. If the original values fall into the interval around the perturbed values, the perturbed values are considered unsafe since they are too close to the original values. There are different ways to construct such intervals, such as rank-based intervals and standard deviation-based intervals. *TeMe08* propose a robust alternative for these intervals. They construct the intervals based on the squared Robust Mahalanobis Distance (RMD) of the individual values. The intervals are scaled by the RMD such that outliers obtain larger intervals and hence need to have a larger perturbation in order to be considered safe than values that are not outliers. This method is illustrated in the Section [Case Studies](#).

## 5.8 Global risk

To construct one aggregate risk measure at the global level for the complete dataset, we can aggregate the measures for risk at the individual level in several ways. Global risk measures should be used with caution: behind an acceptable global risk can hide some very high-risk records that are compensated by many low risk records.

### 5.8.1 Mean of individual risk measures

A straightforward way of aggregating the individual risk measures is taking the mean of all individuals in the sample, which is equal to summing over all keys in the sample if multiplied by the sample frequencies of these keys and dividing by the sample size  $n$ :

$$R_1 = \frac{1}{n} \sum_i r_k = \frac{1}{n} \sum_k f_k r_k$$

$r_k$  is the individual risk of key  $k$  that the  $i^{\text{th}}$  individual shares (see the Section [Categorical key variables and frequency counts](#)).

The global risk in the example data in [Table 5.1](#) is 0.01582, which is the expected proportion of all individuals in the sample that could be re-identified by an intruder. Another way of expressing the global risk is the number of expected re-identifications,  $n * R_1$ , which is in the example  $10 * 0.01582$ .

---

**Note:** This global risk measure should be used with caution. The average risk can be relatively low, but a few individuals could have a very high probability of re-identification.

---

An easy way to check for this is to look at the distribution of the individual risk values or the number of individuals with risk values above a certain threshold, as shown in the next section.



### 5.8.2 Count of individuals with risks larger than a certain threshold

All individuals belonging to the same key have the same individual risk,  $r_k$ . Another way of expressing the total risk in the sample is the total number of observations that exceed a certain threshold of individual risk. Setting the threshold can be absolute (e.g., all those individuals who have a disclosure risk higher than 0.05 or 5%) or relative (e.g., all those individuals with risks higher than the upper quartile of individual risk). In the example, no individual has a higher disclosure risk than 0.05.

These calculations can then be used to treat data for individuals whose risk values are above a predetermined threshold. For example, one could suppress values of certain key variables for those individuals with risk above a specified threshold. This is explained further in the Section [Local suppression](#).

## 5.9 Household risk

In many social surveys, the data have a hierarchical structure where an individual belongs to a higher-level entity (see the Section [Levels of risk](#)). Typical examples are households in social surveys or pupils in schools. Re-identification of one household member can lead to re-identification of the other household members, too. It is therefore easy to see that if we take the household structure into account, the re-identification risk is the risk that at least one of the household members is re-identified.

$$r^h = P(A_1 \cup A_2 \cup \dots \cup A_J) = 1 - \prod_{j=1}^J 1 - P(A_j),$$

where  $A_j$  is the event that the  $j^{\text{th}}$  member of the household is re-identified and  $P(A_j) = r_k$  is the individual disclosure risk of the  $j^{\text{th}}$  member. For example, if a household member has three members with individual disclosure risks based on their respective keys 0.02, 0.03 and 0.03, respectively, the household risk is

$$1 - (1 - 0.02)(1 - 0.03)(1 - 0.03) = 0.078$$

The hierarchical or household risk cannot be lower than the individual risk, and the household risk is always the same for all household members. The household risk should be used in cases where the data contain a hierarchical structure, i.e., where a household structure is present in the data.

---

**Note:** The size of a household is an important identifier itself, especially for large households. Suppression of the actual size variable (e.g., number of household members), however, does not suffice to remove this information from the dataset, as a simple count of the household members for a particular household will allow reconstructing this variable as long as a household ID is in the data, which allows assigning individuals to households. We flag this for the reader's attention as it is important. Further discussion on approaches to the SDC process that take into account the household structure where it exists can be found in the Section [Anonymization of the quasi-identifier household size](#)

---

### Recommended Reading Material on Risk Measurement

Elliot, Mark J, Anna Manning, Ken Mayes, John Gurd, and Michael Bane. 2005. "SUDA: A Program for Detecting Special Uniques." *Joint UNECE/Eurostat Work Session on Statistical Data Confidentiality*. Geneva.

Hundepool, Anco, Josep Domingo-Ferrer, Luisa Franconi, Sarah Giessing, Eric Schulte Nordholt, Keith Spicer, and Peter Paul de Wolf. 2012. *Statistical Disclosure Control*. Chichester: John Wiley & Sons Ltd. doi:10.1002/9781118348239.

Lambert, Diane. 1993. "Measures of Disclosure Risk and Harm." *Journal of Official Statistics* 9(2) : 313-331.

Machanavajjhala, Ashwin, Daniel Kifer, Johannes Gehrke, and Muthuramakrishnan Venkatasubramanian. 2007. "L-diversity: Privacy Beyond K-anonymity." *ACM Trans. Knowl. Discov. Data* 1 (Article 3) (1556-4681). doi:10.1145/1217299.1217302. <https://ptolemy.berkeley.edu/projects/truststc/pubs/465/L%20Diversity%20Privacy.pdf>. Accessed July 6, 2018.

Templ, Matthias, Bernhard Meindl, Alexander Kowarik, and Shuang Chen. 2014. "Introduction to Statistical Disclosure Control (SDC)." <http://www.ihsn.org/home/sites/default/files/resources/ihsn-working-paper-007-Oct27.pdf>. August 1. Accessed July 6, 2018.

---

## References

---

## Anonymization Methods

---

This Section describes the SDC methods most commonly used. We discuss for every method for what type of data the method is suitable, both in terms of data characteristics and type of data. Furthermore, options such as specific parameters for each method are discussed as well as their impacts. These findings are meant as guidance but should be used with caution, since every dataset has different characteristics and our findings may not always address your particular dataset. The last three sections are on the anonymization of variables and datasets with particular characteristics that deserve special attention. The Section *Anonymization of geospatial variables* deals with for anonymizing geographical data, such as GPS coordinates, the Section *Anonymization of the quasi-identifier household size* discusses the anonymization of data with a hierarchical structure (household structure) and the Section *Special case: census data* describes the peculiarities of dealing with and releasing census microdata.

To determine which anonymization methods are suitable for specific variables and/or datasets, we begin by presenting some classifications of SDC methods.

### 6.1 Classification of SDC methods

SDC methods can be classified as **non-perturbative** and **perturbative** (see *HDFG12*).

- **Non-perturbative methods** reduce the detail in the data by generalization or suppression of certain values (i.e., masking) without distorting the data structure.
- **Perturbative methods** do not suppress values in the dataset but perturb (i.e., alter) values to limit disclosure risk by creating uncertainty around the true values.

Both non-perturbative and perturbative methods can be used for categorical and continuous variables.

We also distinguish between **probabilistic** and **deterministic** SDC methods.

- **Probabilistic methods** depend on a probability mechanism or a random number-generating mechanism. Every time a probabilistic method is used, a different outcome is generated. For these methods it is often recommended that a seed be set for the random number generator if you want to produce replicable results.
- **Deterministic methods** follow a certain algorithm and produce the same results if applied repeatedly to the same data with the same set of parameters.

SDC methods for microdata intend to prevent identity and attribute disclosure. Different SDC methods are used for each type of disclosure control. Methods such as recoding and suppression are applied to quasi-identifiers to prevent identity disclosure, whereas top coding a quasi-identifier (e.g., income) or perturbing a sensitive variable prevent attribute disclosure.

In this guide we discuss the most commonly applied methods from the literature and used in most agencies experienced in using these methods. All discussed methods are implemented in the *sdcMicro* package. [Table 6.1](#) gives an overview of the SDC methods discussed in this guide, their classification and types of data to which they are applicable.

Table 6.1: Overview of SDC methods

Method	Classification of SDC method	Data Type
Global recoding	non-perturbative, deterministic	continuous and categorical
Top and bottom coding	non-perturbative, deterministic	continuous and categorical
Local suppression	non-perturbative, deterministic	categorical
PRAM	perturbative, probabilistic	categorical
Micro aggregation	perturbative, probabilistic	continuous
Noise addition	perturbative, probabilistic	continuous
Shuffling	perturbative, probabilistic	continuous
Rank swapping	perturbative, probabilistic	continuous

## 6.2 Non-perturbative methods

### 6.2.1 Recoding

Recoding is a deterministic method used to decrease the number of distinct categories or values for a variable. This is done by combining or grouping categories for categorical variables or constructing intervals for continuous variables. Recoding is applied to all observations of a certain variable and not only to those at risk of disclosure. There are two general types of recoding: global recoding and top and bottom coding.

#### Global recoding

Global recoding combines several categories of a categorical variable or constructs intervals for continuous variables. This reduces the number of categories available in the data and potentially the disclosure risk, especially for categories with few observations, but also, importantly, it reduces the level of detail of information available to the analyst. To illustrate recoding, we use the following example. Assume that we have five regions in our dataset. Some regions are very small and when combined with other key variables in the dataset, produce high re-identification risk for some individuals in those regions. One way to reduce risk would be to combine some of the regions by recoding them. We could, for example, make three groups out of the five, call them ‘North’, ‘Central’ and ‘South’ and re-label the values accordingly. This reduces the number of categories in the variable region from five to three.

---

**Note:** Any grouping should be some logical grouping and not a random joining of categories.

---

Examples would be grouping districts into provinces, municipalities into districts, or clean water categories together. Grouping all small regions without geographical proximity together is not necessarily the best option from the utility perspective. [Table 6.2](#) illustrates this with a very simplified example dataset. Before recoding, three individuals have distinct keys, whereas after recoding (grouping ‘Region 1’ and ‘Region 2’ into ‘North’, ‘Region 3’ into ‘Central’ and ‘Region 4’ and ‘Region 5’ into ‘South’), the number of distinct keys is reduced to four and the frequency of every key is at least two, based on the three selected quasi-identifiers. The frequency counts of the keys  $f_k$  are shown in the last column of [Table 6.2](#). An intruder would find at least two individuals for each key and cannot distinguish any more

between individuals 1 – 3, individuals 4 and 6, individuals 5 and 7 and individuals 8 – 10, based on the selected key variables.

Table 6.2: Illustration of effect of recoding on frequency counts of keys

.	Before recoding				After recoding			
Individual	Region	Gender	Religion	$f_k$	Region	Gender	Religion	$f_k$
1	Region 1	Female	Catholic	1	North	Female	Catholic	3
2	Region 2	Female	Catholic	2	North	Female	Catholic	3
3	Region 2	Female	Catholic	2	North	Female	Catholic	3
4	Region 3	Female	Protestant	2	Central	Female	Protestant	2
5	Region 3	Male	Protestant	1	Central	Male	Protestant	2
6	Region 3	Female	Protestant	2	Central	Female	Protestant	2
7	Region 3	Male	Protestant	2	Central	Male	Protestant	2
8	Region 4	Male	Muslim	2	South	Male	Muslim	3
9	Region 4	Male	Muslim	2	South	Male	Muslim	3
10	Region 5	Male	Muslim	1	South	Male	Muslim	3

Recoding is commonly the first step in an anonymization process. It can be used to reduce the number of unique combinations of values of key variables. This generally increases the frequency counts for most keys and reduces the risk of disclosure. The reduction in the number of possible combinations is illustrated in Table 6.3 with the quasi-identifiers “region”, “marital status” and “age”. Table 6.3 shows the number of categories of each variable and the number of theoretically possible combinations, which is the product of the number of categories of each quasi-identifier, before and after recoding. “Age” is interpreted as a semi-continuous variable and treated as a categorical variable. The number of possible combinations and hence the risk for re-identification are reduced greatly by recoding. One should bear in mind that the number of possible combinations is a theoretical number; in practice, these may include very unlikely combinations such as age = 3 and marital status = widow and the actual number of combinations in a dataset may be lower.

Table 6.3: Illustration of the effect of recoding on the theoretically possible number of combinations in a dataset

Number of categories	Region	Marital status	Age	Possible combinations
before recoding	20	8	100	16,000
after recoding	6	6	15	540

The main parameters for global recoding are the size of the new groups, as well as defining which values are grouped together in new categories.

**Note:** Care should be taken to choose new categories in line with the data use of the end users and to minimize information loss as a result of recoding.

We illustrate this with three examples:

- *Age variable:* The categories of age should be chosen so that they still allow data users to make calculations relevant for the subject being studied. For example, if indicators need to be calculated for children of school going ages 6 – 11 and 12 – 17, and age needs to be grouped to reduce risk, then care should be taken to create age intervals that still allow the calculations to be made. A satisfactory grouping could be, for example, 0 – 5, 6 – 11, 12 – 17, etc., whereas a grouping 0 – 10, 11 – 15, 16 – 18 would destroy the data utility for these users. While it is common practice to create intervals (groups) of equal width (size), it is also possible (if data users require this) to recode only part of the variables and leave some values as they were originally. This could be done, for example, by recoding all ages above 20, but leaving those below 20 as they are. If SDC methods other than recoding will be used later or in a next step, then care should be taken when applying recoding to only part of the distribution, as this might increase the information loss due to the other methods, since the

grouping does not protect the ungrouped variables. Partial recoding followed by suppression methods such as local suppression may, for instance, lead to a higher number of suppressions than desired or necessary in case the recoding is done for the entire value range (see the next section on local suppression). In the example above, the number of suppressions of values below 20 will likely be higher than for values in the recoded range. The disproportionately high number of suppressions in this range of values that are not recoded can lead to higher utility loss for these groups.

- *Geographic variables*: If the original data specify administrative level information in detail, e.g., down to municipality level, then potentially those lower levels could be recoded or aggregated into higher administrative levels, e.g., province, to reduce risk. In doing so, the following should be noted: Grouping municipalities into abstract levels that intersect different provinces would make data analysis at the municipal or provincial level challenging. Care should be taken to understand what the user requires and the intention of the study. If a key component of the survey is to conduct analysis at the municipal level, then aggregating up to provincial level could damage the utility of the data for the user. Recoding should be applied if the level of detail in the data is not necessary for most data users and to avoid an extensive number of suppressions when using other SDC methods subsequently. If the users need information at a more detailed level, other methods such as perturbative methods might provide a better solution than recoding.
- *Toilet facility*: An example of a situation where a high level of detail might not be necessary and recoding may do very little harm to utility is the case of a detailed household toilet facility variable that lists responses for 20 types of toilets. Researchers may only need to distinguish between improved and unimproved toilet facilities and may not require the exact classification of up to 20 types. Detailed information of toilet types can be used to re-identify households, while recoding to two categories – improved and unimproved facilities – reduces the re-identification risk and in this context, hardly reduces data utility. This approach can be applied to any variable with many categories where data users are not interested in detail, but rather in some aggregate categories. Recoding addresses aggregation for the data users and at the same time protects the microdata. Important is to take stock of the aggregations used by data users.

Recoding should be applied only if removing the detailed information in the data will not harm most data users. If the users need information at a more detailed level, then recoding is not appropriate and other methods such as perturbative methods might work better.

### Examples of global recoding

In this section, we illustrate global recoding with two examples, one categorical variable and one continuous variable. Assume that the variable “sizeRes”, size of the residence area, has four different categories: ‘capital, large city’, ‘small city’, town’, and ‘countryside’). The first three are recoded or regrouped as ‘urban’ and the category ‘countryside’ is renamed ‘rural’. Fig. 6.1 illustrates the effect of recoding the variable “sizeRes” and show respectively the frequency counts before and after recoding. We see that the number of categories has reduced from 4 to 2 and the small categories (‘small city’ and ‘town’) have disappeared.

Global recoding of a numerical (continuous) variables changes it into a categorical variable. The intervals should cover the entire value range of the variable. Fig. 6.2 shows the effect of recoding the variable “age”, age in years, in ten-year intervals.

Instead of creating intervals of equal width, we can also create intervals of unequal width. This is illustrated in code53, where we use the age groups 1-5, 6-11, 12-17, 18-21, 22-25, 26-49, 50-64 and 65+. In this example, this is a useful step, since even after recoding in 10-year intervals, the categories with high age values have low frequencies. We chose the intervals by respecting relevant school age and employment age values (e.g., retirement age is 65 in this example) such that the data can still be used for common research on education and employment. Fig. 6.3 shows the effect of recoding the variable “age” after adjusting the intervals.

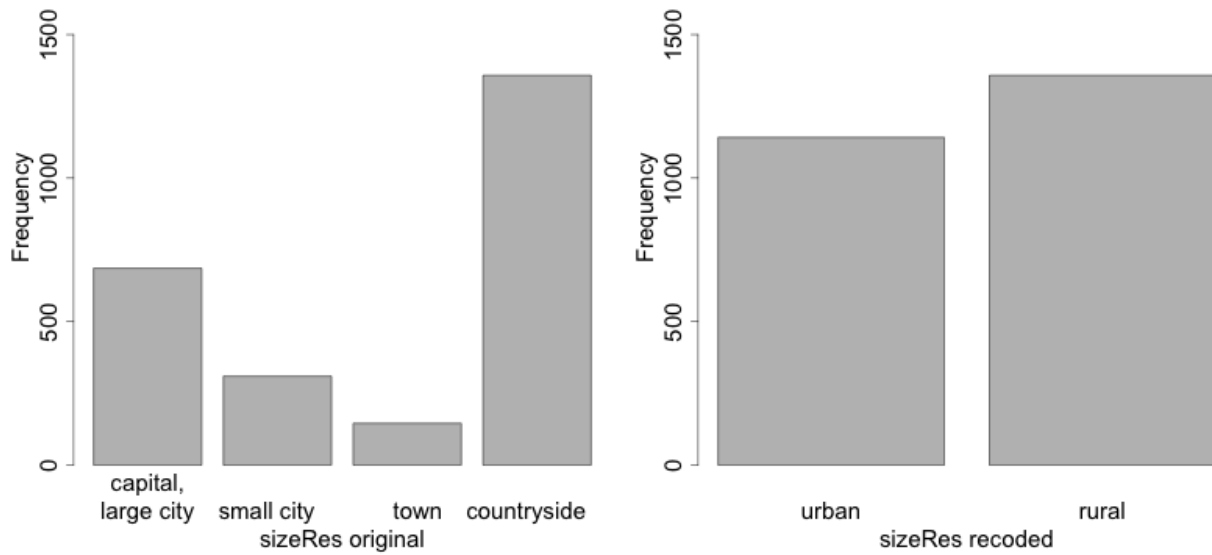


Fig. 6.1: Effect of recoding – frequency counts before and after recoding

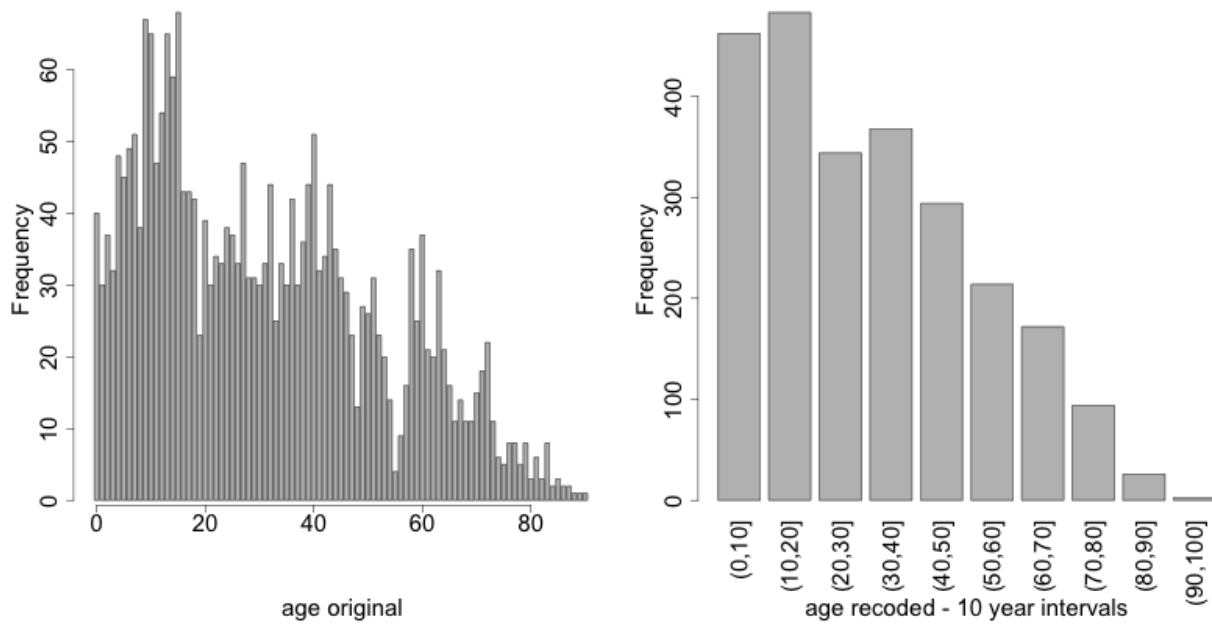


Fig. 6.2: Age variable before and after recoding

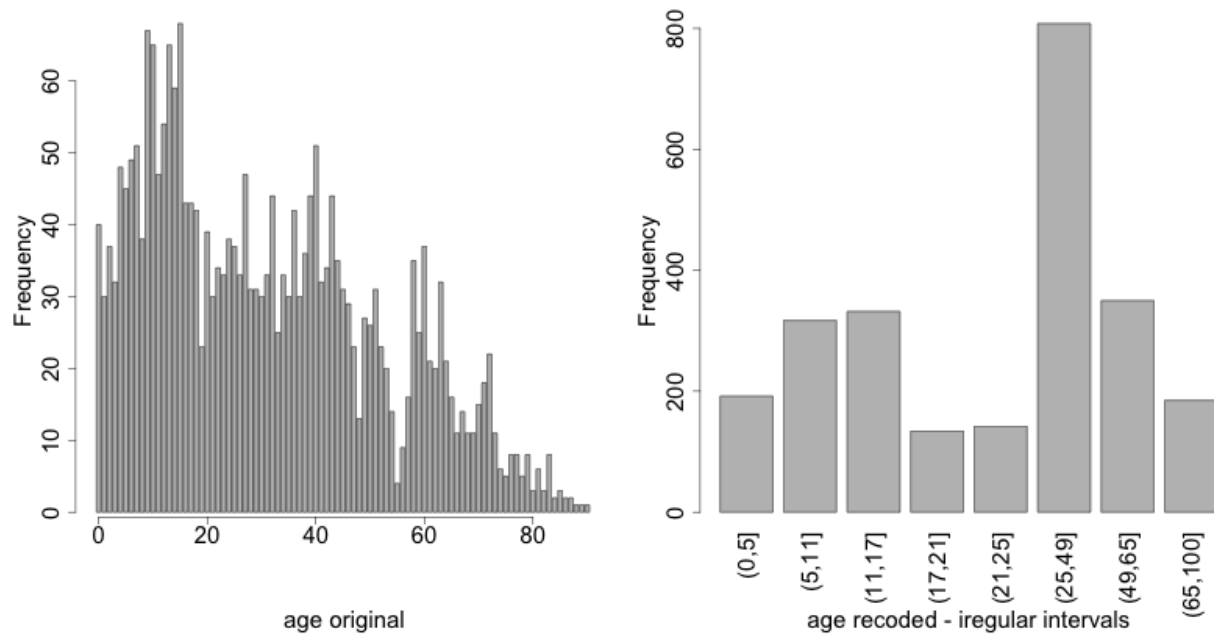


Fig. 6.3: Age variable before and after recoding

### Top and bottom coding

Top and bottom coding are similar to global recoding, but instead of recoding all values, only the top and/or bottom values of the distribution or categories are recoded. This can be applied only to ordinal categorical variables and (semi-)continuous variables, since the values have to be at least ordered. Top and bottom coding is especially useful if the bulk of the values lies in the center of the distribution with the peripheral categories having only few observations (outliers). Examples are age and income; for these variables, there will often be only a few observations above certain thresholds, typically at the tails of the distribution. The fewer the observations within a category, the higher the identification risk. One solution could be grouping the values at the tails of the distribution into one category. This reduces the risk for those observations, and, importantly, does so without reducing the data utility for the other observations in the distribution.

Deciding where to apply the threshold and what observations should be grouped requires:

- Reviewing the overall distribution of the variable to identify at which point the frequencies drop below the desired number of observations and identify outliers in the distribution. Fig. 6.4 shows the distribution of the age variable and suggests 65 (red vertical line) for the top code age.
- Taking into account the intended use of the data and the purpose for which the survey was conducted. For example, if the data are typically used to measure labor force participation for those aged 15 to 64, then top and bottom coding should not interfere with the categories 15 to 64. Otherwise the analyst would find it impossible to create the desired measures for which the data were intended. In the example, we consider this and code all age higher than 64.

### Rounding

Rounding is similar to grouping, but used for continuous variables. Rounding is useful to prevent exact matching with external data sources. In addition, it can be used to reduce the level of detail in the data. Examples are removing decimal figures or rounding to the nearest 1,000.

The next section discusses the method local suppression. Recoding is often used before local suppression to reduce the number of necessary suppressions.



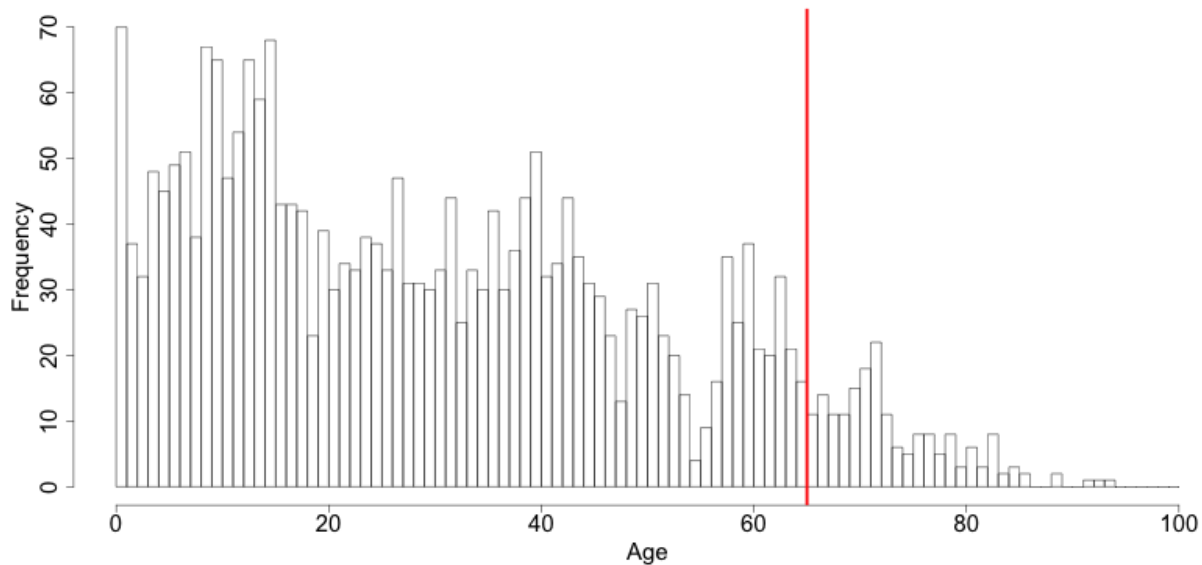


Fig. 6.4: Utilizing the frequency distribution of variable age to determine threshold for top coding

### Recommended Reading Material on Recoding

Hundepool, Anco, Josep Domingo-Ferrer, Luisa Franconi, Sarah Giessing, Rainer Lenz, Jane Naylor, Eric Schulte Nordholt, Giovanni Seri, and Peter Paul de Wolf. 2006. *Handbook on Statistical Disclosure Control*. ESSNet SDC. <http://neon.vb.cbs.nl/casc/handbook.htm>.

Hundepool, Anco, Josep Domingo-Ferrer, Luisa Franconi, Sarah Giessing, Eric Schulte Nordholt, Keith Spicer, and Peter Paul de Wolf. 2012. *Statistical Disclosure Control*. Chichester: John Wiley & Sons Ltd. doi:10.1002/9781118348239.

Templ, Matthias, Bernhard Meindl, Alexander Kowarik, and Shuang Chen. 2014. Statistical Disclosure Control (SDCMicro). <http://www.ihsn.org/home/software/disclosure-control-toolbox>. (accessed June 9, 2018).

De Waal, A.G., and Willenborg, L.C.R.J. 1999. *Information loss through global recoding and local suppression*. Netherlands Official Statistics, 14:17-20, 1999. Special issue on SDC

## 6.2.2 Local suppression

It is common in surveys to encounter values for certain variables or combinations of quasi-identifiers (keys) that are shared by very few individuals. When this occurs, the risk of re-identification for those respondents is higher than the rest of the respondents (see the Section *k-anonymity*). Often local suppression is used after reducing the number of keys in the data by recoding the appropriate variables. Recoding reduces the number of necessary suppressions as well as the computation time needed for suppression. Suppression of values means that values of a variable are replaced by a missing value. The Section *k-anonymity* discusses how missing values influence frequency counts and *k-anonymity*. It is important to note that not all values for all individuals of a certain variable are suppressed, which would be the case when removing a direct identifier, such as “name”; only certain values for a particular variable and a particular respondent or set of respondents are suppressed. This is illustrated in the following example and Table 6.4.

Table 6.4 presents a dataset with seven respondents and three quasi-identifiers. The combination {‘female’, ‘rural’, ‘higher’} for the variables “gender”, “region” and “education” is an unsafe combination, since it is unique in the

sample. By suppressing either the value ‘female’ or ‘higher’, the respondent cannot be distinguished from the other respondents anymore, since that respondent shares the same combination of key variables with at least three other respondents. Only the value in the unsafe combination of the single respondent at risk is suppressed, not the values for the same variable of the other respondents. The freedom to choose which value to suppress can be used to minimize the total number of suppressions and hence the information loss. In addition, if one variable is very important to the user, we can choose not to suppress values of this variable, unless strictly necessary. In the example, we can choose between suppressing the value ‘female’ or ‘higher’ to achieve a safe data file; we chose to suppress ‘higher’. This choice should be made taking into account the needs of data users. In this example we find “gender” more important than “education”.

Table 6.4: Local suppression illustration - sample data before and after suppression

Variable	Before local suppression			After local suppression		
ID	Gender	Region	Education	Gender	Region	Education
1	female	rural	higher	female	rural	<b>missing</b>
2	male	rural	higher	male	rural	higher
3	male	rural	higher	male	rural	higher
4	male	rural	higher	male	rural	higher
5	female	rural	lower	female	rural	lower
6	female	rural	lower	female	rural	lower
7	female	rural	lower	female	rural	lower

Since continuous variables have a high number of unique values (e.g., income in dollars or age in years),  $k$ -anonymity and local suppression are not suitable for continuous variables or variables with a very high number of categories. A possible solution in those cases might be to first recode to produce fewer categories (e.g., recoding age in 10-year intervals or income in quintiles). Always keep in mind, though, what effect any recoding will have on the utility of the data.

Several different algorithms can be used to determine which values to suppress. One common algorithm determines an optimal suppression pattern to achieve on a specified set of quasi-identifiers a certain level of  $k$ -anonymity for these quasi-identifiers. This algorithm used seeks to minimize the total number of suppressions while achieving the required  $k$ -anonymity threshold. By default, this algorithm is more likely to suppress values of variables with many different categories or values, and less likely to suppress variables with fewer categories. For example, the values of a geographical variable, with 12 different areas, are more likely to be suppressed than the values of the variable “gender”, which has typically only two categories. If variables with many different values are important for data utility and suppression is not desired for them, one can rank variables by importance and thus specify the order in which the algorithm will seek to suppress values within quasi-identifiers to achieve  $k$ -anonymity. The algorithm seeks to apply fewer suppressions to variables of high importance than to variables with lower importance. Nevertheless, suppressions in the variables with high importance might be inevitable to achieve the required level of  $k$ -anonymity.

### Example of local suppression

In this example local suppression is applied to achieve the  $k$ -anonymity threshold of 5 on the quasi-identifiers “gender”, “region”, “religion”, “age” and “ethnicity”. Without ranking the importance of the variables, the value of the variable “age” is more likely to be suppressed, since this is the variable with most categories. The variable “age” has 10 categories after recoding. The variable “gender” is least likely to be suppressed, since it has only two different values: ‘male’ and ‘female’. The other variables have 4 (“sizeRes”), 2 (“region”), and 8 (“ethnicity”) categories. The standard local suppression algorithm suppresses most values in the variable “age” (80). In fact, only the variable “ethnicity” of the other variables also needed suppressions (8) to achieve the  $k$ -anonymity threshold of 5. The variable “ethnicity” is the variable with the second highest number of suppressions.

The variable “age” is typically an important variable. Therefore, if possible, we would like to reduce the number of suppressions on “age” by specifying the order of importance of the variables and giving high importance (little

suppression) to the quasi-identifier “age”. We also assign importance to the variable “gender”. The effect is clear: there are no suppressions in the variables “age” and “gender”. For that, the other variables, especially “sizeRes” (87) and “ethnicity” (62), received many more suppressions. The total number of suppressed values has increased from 88 to 166.

---

**Note:** Fewer suppressions in one variable increase the number of necessary suppressions in other variables.

---

Generally, the total number of suppressed values needed to achieve the required level of  $k$ -anonymity increases when specifying the order of importance, since this prevents to use the optimal suppression pattern. The importance of variables should be specified only in cases where the variables with many categories play an important role in data utility for the data users<sup>1</sup>.

Fig. 6.5 demonstrates the effect of the required  $k$ -anonymity threshold and the order of importance on the data utility by using several labor market-related indicators from an I2D2<sup>2</sup> dataset before and after anonymization. Fig. 6.5 displays the relative changes as a percentage of the initial value after re-computing the indicators with the data to which local suppression was applied. The indicators are the proportion of active females and males, and the number of females and males of working age. The values computed from the raw data were, respectively, 68%, 12%, 8,943 and 9,702. The vertical line at 0 is the benchmark of no change. The numbers indicate the required  $k$ -anonymity threshold (3 or 5) and the colors indicate the importance vector: red (no symbol) is no importance vector, blue (with \* symbol) is high importance on the variable with the employment status information and dark green (with + symbol) is high importance on the age variable.

A higher  $k$ -anonymity threshold leads to greater information loss (i.e., larger deviations from the original values of the indicators, the 5's are further away from the benchmark of no change than the corresponding 3's) caused by local suppression. Reducing the number of suppressions on the employment status variable by specifying an importance vector does not improve the indicators. Instead, reducing the number of suppressions on age greatly reduces the information loss. Since specific age groups have a large influence on the computation of these indicators (the rare cases are in the extremes and will be suppressed), high suppression rates on age distort the indicators. It is generally useful to compare utility measures (see the Section [Measuring Utility and Information Loss](#)) to specify the importance vector, since the effects can be unpredictable.

The threshold of  $k$ -anonymity to be set depends on several factors, which are amongst others: 1) the legal requirements for a safe data file; 2) other methods that will be applied to the data; 3) the number of suppressions and related information loss resulting from higher thresholds; 4) the type of variable; 5) the sample weights and sample size; and 6) the release type (see the Section [Release Types](#)). Commonly applied levels for the  $k$ -anonymity threshold are 3 and 5.

Table 6.5 illustrates the influence of specifying the order of importance and the  $k$ -anonymity threshold on the global risk after suppression and total number of suppressions required to achieve this  $k$ -anonymity threshold. The dataset contains about 63,000 individuals. The higher the  $k$ -anonymity threshold, the more suppressions are needed and the lower the risk after local suppression (expected number of re-identifications). In this particular example, the computation time is shorter for higher thresholds. This is due the higher number of necessary suppressions, which reduces the difficulty of the search for an optimal suppression pattern.

The age variable is recoded in five-year intervals and has 20 age categories. This is the variable with the highest number of categories. Prioritizing the suppression of other variables leads to a higher total number of suppressions and a longer computation time.

---

<sup>1</sup> This can be assessed with utility measures.

<sup>2</sup> I2D2 is a dataset with data related to the labor market.

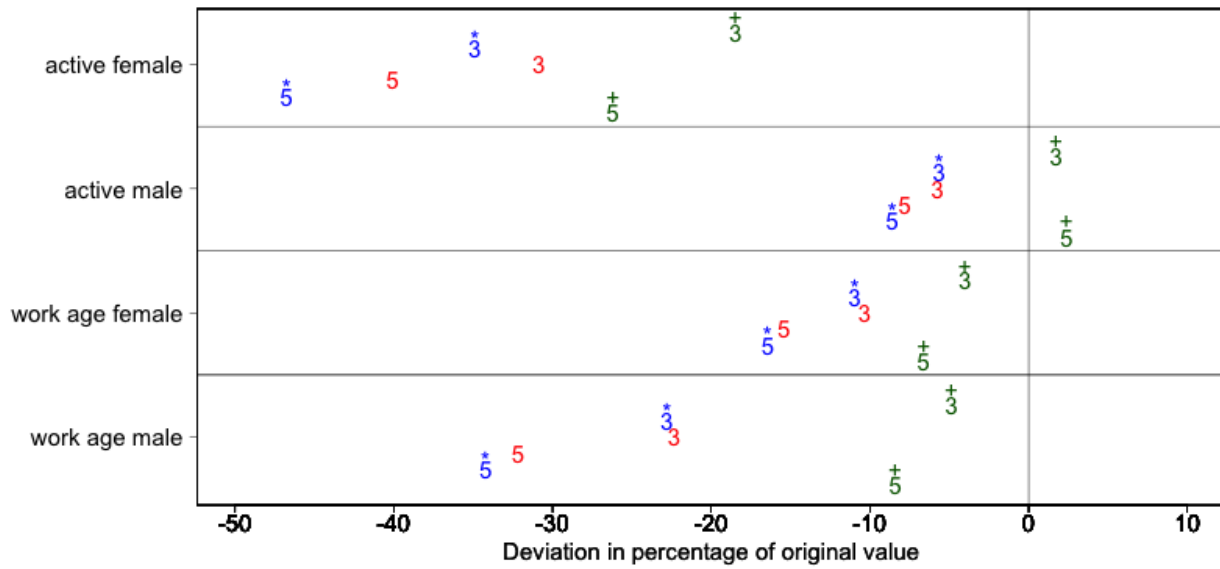


Fig. 6.5: Changes in labor market indicators after anonymization of I2D2 data

Table 6.5: How the order of importance and the  $k$ -anonymity threshold affect total number of suppressions

Threshold	Importance	Total number of	Threshold	Total number of
k-anonymity	vector	suppressions	k-anonymity	suppressions
3	none (default)	6,676	5,387	11.8
3	employment status	7,254	5,512	13.1
3	age variable	8,175	60	4.5
5	none (default)	9,971	7,894	8.5
5	employment status	11,668	8,469	10.2
5	age variable	13,368	58	3.8

In cases where there are a large number of quasi-identifiers and the variables have many categories, the number of possible combinations increases rapidly (see  $k$ -anonymity). If the number of variables and categories is very large, the computation time of the local suppression algorithms can be very long. Therefore, reducing the number of quasi-identifiers and/or categories before applying local suppression is recommended. This can be done by recoding variables or selecting some variables for other (perturbative) methods, such as PRAM. This is to ensure that the number of suppressions is limited and hence the loss of data is limited to only those values that pose most risk.

### All- $m$ approach

In some datasets, it might prove difficult to reduce the number of quasi-identifiers and even after reducing the number of categories by recoding, the local suppression algorithm takes a long time to compute the required suppressions. A solution in such cases can be the so-called ‘all- $m$  approach’ (see *Wolf15*). The all- $m$  approach consists of applying the local suppression algorithm as described above to all possible subsets of size  $m$  of the total set of quasi-identifiers. The advantage of this approach is that the partial problems are easier to solve and computation time will be slower. Caution should be applied since this method does not necessarily lead to  $k$ -anonymity in the complete set of quasi-identifiers. There are two possibilities to reach the same level of protection: 1) to choose a higher threshold for  $k$  or 2) to re-apply the local suppression algorithm on the complete set of quasi-identifiers after using the all- $m$  approach to achieve the

required threshold. In the second case, the all- $m$  approach leads to a shorter computation time at the cost of a higher total number of suppressions.

**Note:** The required level is not achieved automatically on the entire set of quasi-identifiers if the all- $m$  approach is used.

Therefore, it is important to evaluate the risk measures carefully after using the all- $m$  approach.

Table 6.6 presents the results of using the all- $m$  approach of a test dataset with 9 key variables and 4,000 records. The table shows the parameters 'k' and 'combs', which are respectively the  $k$ -anonymity threshold and the size of the subsets, the number of  $k$ -anonymity violators for different levels of  $k$  as well as the total number of suppressions. We observe that the different combinations do not always lead to the required level of  $k$ -anonymity. For example, when setting  $k = 3$ , and combs 3 and 7, there are still 15 records in the dataset (with a total of 9 quasi-identifiers) that violate 3-anonymity after local suppression. Due to the smaller sample size, the gains in running time are not yet apparent in this example, since the rerunning algorithm several times takes up time. A larger dataset would benefit more from the all- $m$  approach, as the algorithm would take longer in the first place.

Table 6.6: Effect of the all- $m$  approach on  $k$ -anonymity

Arguments		Number of violators for different levels of $k$ -anonymity on complete set			Total number of suppressions
k	combs	k = 2	k = 3	k = 5	
Before local suppression		2,464	3,324	3,877	0
3	.	0	0	1,766	2,264
5	.	0	0	0	3,318
3	3	2,226	3,202	3,819	3,873
3	3, 7	15	108	1,831	6,164
3	3, 9	0	0	1,794	5,982
3	5, 9	0	0	1,734	6,144
5	3	2,047	3,043	3,769	3,966
5	3, 7	0	6	86	7,112
5	3, 9	0	0	0	7,049
5	5, 9	0	0	0	7,129
5, 3	3, 7	11	108	1,859	6,140
5, 3	3, 9	0	0	1,766	2,264
5, 3	5, 9	0	0	0	3,318

Often the dataset contains variables that are related to the key variables used for local suppression. Examples are rural/urban to regions in case regions are completely rural or urban or variables that are only answered for specific categories (e.g., sector for those working, schooling related variables for certain age ranges). In those cases, the variables rural/urban or sector might not be quasi-identifiers themselves, but could allow the intruder to reconstruct suppressed values in the quasi-identifiers region or employment status. For example, if region 1 is completely urban, and all other regions are only semi-urban or rural, a suppression in the variable region for a record in region 1 can be simply reconstructed by the rural/urban variable. Therefore, it is useful to suppress the values corresponding to the suppressions in those linked variables.

Another simpler alternative for the local suppression algorithm described above is to suppress values of certain key variables of individuals with risks above a certain threshold. In this case, all values of the specified variable for respondents with a risk higher than the specified threshold will be suppressed. The risk measure used is the individual risk (see the Section [Individual risk](#)). This is useful if one variable has sensitive values that should not be released for individuals with high risks of re-identification. What is considered high re-identification probability depends on legal requirements.

## 6.3 Perturbative methods

Perturbative methods do not suppress values in the dataset, but perturb (alter) values to limit disclosure risk by creating uncertainty around the true values. An intruder is uncertain whether a match between the microdata and an external file is correct or not. Most perturbative methods are based on the principle of matrix masking, i.e., the altered dataset  $Z$  is computed as

$$Z = AXB + C$$

where  $X$  is the original data,  $A$  is a matrix used to transform the records,  $B$  is a matrix to transform the variables and  $C$  is a matrix with additive noise.

---

**Note:** Risk measures based on frequency counts of keys are no longer valid after applying perturbative methods.

---

This can be seen in [Table 6.7](#), which displays the same data before and after swapping some values. The swapped values are in italics. Both before and after perturbing the data, all observations violate  $k$ -anonymity at the level 3 (i.e., each key does not appear more than twice in the dataset). Nevertheless, the risk of **correct** re-identification of the records is reduced and hence information contained in other (sensitive) variables possibly not disclosed. With a certain probability, a match of the microdata with an external data file will be wrong. For example, an intruder would find one individual with the combination {‘male’, ‘urban’, ‘higher’}, which is a sample unique. However, this match is not correct, since the original dataset did not contain any individual with these characteristics and hence the matched individual cannot be a correct match. The intruder cannot know with certainty whether the information disclosed from other variables for that record is correct.

Table 6.7: Sample data before and after perturbation

Variable	Original data			After perturbing the data		
	Gender	Region	Education	Gender	Region	Education
1	female	rural	higher	female	rural	higher
2	female	rural	higher	female	rural	<i>lower</i>
3	male	rural	lower	male	rural	lower
4	male	rural	lower	<i>female</i>	rural	lower
5	female	urban	lower	<i>male</i>	urban	<i>higher</i>
6	female	urban	lower	female	urban	lower

One advantage of perturbative methods is that the information loss is reduced, since no values will be suppressed, depending on the level of perturbation. One disadvantage is that data users might have the impression that the data was not anonymized before release and will be less willing to participate in future surveys. Therefore, there is a need for reporting both for internal and external use (see the [Section Step 11: Audit and Reporting](#)).

An alternative to perturbative methods is the generation of synthetic data files with the same characteristics as the original data files. Synthetic data files are not discussed in these guidelines. For more information and an overview of the use of synthetic data as SDC method, we refer to *Drec11* and Section 3.8 in *HDFG12*. We discuss here five perturbative methods: Post Randomization Method (PRAM), microaggregation, noise addition, shuffling and rank swapping.

### 6.3.1 PRAM (Post RAndomization Method)

PRAM is a perturbative method for categorical data. This method reclassifies the values of one or more variables, such that intruders that attempt to re-identify individuals in the data do so, but with positive probability, the re-identification made is with the wrong individual. This means that the intruder might be able to match several individuals between external files and the released data files, but cannot be sure whether these matches are to the correct individual.

PRAM is defined by the transition matrix  $P$ , which specifies the transition probabilities, i.e., the probability that a value of a certain variable stays unchanged or is changed to any of the other  $k - 1$  values.  $k$  is the number of categories or factor levels within the variable to be PRAMmed. For example, if the variable region has 10 different regions,  $k$  equals 10. In case of PRAM for a single variable, the transition matrix is size  $k * k$ . We illustrate PRAM with an example of the variable “region”, which has three different values: ‘capital’, ‘rural1’ and ‘rural2’. The transition matrix for applying PRAM to this variable is size 3\*3:

$$P = \begin{bmatrix} 1 & 0 & 0 \\ 0.05 & 0.8 & 0.15 \\ 0.05 & 0.15 & 0.8 \end{bmatrix}$$

The values on the diagonal are the probabilities that a value in the corresponding category is not changed. The value 1 at position (1,1) in the matrix means that all values ‘capital’ stay ‘capital’; this might be a useful decision, since most individuals live in the capital and no protection is needed. The value 0.8 at position (2,2) means that an individual with value ‘rural1’ will stay with probability 0.8 ‘rural1’. The values 0.05 and 0.15 in the second row of the matrix indicate that the value ‘rural1’ will be changed to ‘capital’ or ‘rural2’ with respectively probability 0.05 and 0.15. If in the initial file we had 5,000 individuals with value ‘capital’ and resp. 500 and 400 with values ‘rural1’ and ‘rural2’, we expect after applying PRAM to have 5,045 individuals with capital, 460 with rural1 and 395 with rural2<sup>3</sup>. The recoding is done independently for each individual. We see that the tabulation of the variable “region” yields different results before and after PRAM, which are shown in Table 6.8. The deviation from the expectation is due to the fact that PRAM is a probabilistic method, i.e., the results depend on a probability-generating mechanism; consequently, the results can differ every time we apply PRAM to the same variables of a dataset.

**Note:** The number of changed values is larger than one might think when inspecting the tabulations in Table 6.8. Not all 5,052 individuals with value capital after PRAM had this value before PRAM and the 457 individuals in rural1 after PRAM are not all included in the 500 individuals before PRAM. The number of changes is larger than the differences in the tabulation (cf. transition matrix).

Given that the transition matrix is known to the end users, there are several ways to correct statistical analysis of the data for the distortions introduced by PRAM.

Table 6.8: Tabulation of variable “region” before and after PRAM

Value	Tabulation before PRAM	Tabulation after PRAM
capital	5,000	5,052
rural1	500	457
rural2	400	391

One way to guarantee consistency between the tabulations before and after PRAM is to choose the transition matrix so that, in expectation, the tabulations before and after applying PRAM are the same for all variables. This condition is fulfilled if the vector with the tabulation of the absolute frequencies of the different categories in the original data is an eigenvector of the transition matrix that corresponds to the unit eigenvalue. PRAM using such transition matrix is called invariant PRAM.

**Note:** Invariant does not guarantee that cross-tabulations of variables (unlike univariate tabulations) stay the same.

PRAM is a probabilistic method and the results can differ every time we apply PRAM to the same variables of a dataset. To overcome this and make the results reproducible, it is good practice to set a seed for the random number generator, so the same random numbers will be generated every time.

Table 6.9 shows the tabulation of the variable after applying invariant PRAM. We can see that the deviations from the initial tabulations, which are in expectation 0, are smaller than with the transition matrix that does not fulfill the invariance property. The remaining deviations are due to the randomness.

<sup>3</sup> The 5,045 is the expectation computed as  $5,000 * 1 + 500 * 0.05 + 400 * 0.05$ .



Table 6.9: Tabulation of variable “region” before and after (invariant) PRAM

Value	Tabulation before PRAM	Tabulation after PRAM	Tabulation after invariant PRAM
capital	5,000	5,052	4,998
rural1	500	457	499
rural2	400	391	403

Table 6.10 presents the cross-tabulations with the variable gender. Before applying invariant PRAM, the share of males in the city is much higher than the share of females (about 60%). This property is not maintained after invariant PRAM (the shares of males and females in the city are roughly equal), although the univariate tabulations are maintained. One solution is to apply PRAM separately for the males and females in this example<sup>4</sup>.

Table 6.10: Cross-tabulation of variable “region” and variable “gender” before and after invariant PRAM

.	Tabulation before PRAM		Tabulation after invariant PRAM	
Value	male	female	male	female
capital	3,056	1,944	2,623	2,375
rural1	157	343	225	274
rural2	113	287	187	216

PRAM is especially useful when a dataset contains many variables and applying other anonymization methods, such as recoding and local suppression, would lead to significant information loss. Checks on risk and utility are important after PRAM.

To do statistical inference on variables to which PRAM was applied, the researcher needs knowledge about the PRAM method as well as about the transition matrix. The transition matrix, together with the random number seed, can, however, lead to disclosure through reconstruction of the non-perturbed values. Therefore, publishing the transition matrix but not the random seed is recommended.

A disadvantage of using PRAM is that very unlikely combinations can be generated, such as a 63-year-old who goes to school. Therefore, the PRAMmed variables need to be audited to prevent such combinations from happening in the released data file. In principal, the transition matrix can be designed in such a way that certain transitions are not possible (probability 0). For instance, for those that go to school, the age must range within 6 to 18 years and only such changes are allowed.

This is illustrated in the following example. Assume, we have two variables “toilet” and “region” and the variable “toilet” needs to be PRAMmed. By applying PRAM to the variable “toilet” within the strata generated by the “region” variable, we prevent changes in the variable “toilet”, where toilet types in a particular region are exchanged with those in other regions. For instance, in the capital region certain types of unimproved toilet types are not in use and therefore these combinations should not occur after PRAMming. Values are only changed with those that are available in the same strata. Strata can be formed by any categorical variable, e.g., gender, age groups, education level.

---

### Recommended Reading Material on PRAM

Gouweleeuw, J. M, P Kooiman, L.C.R.J Willenborg, and P.P de Wolf. “Post Randomization for Statistical Disclosure Control: Theory and Implementation.” *Journal of Official Statistics* 14, no. 4 (1998a): 463-478. Available at <http://www.jos.nu/articles/abstract.asp?article=144463>

Gouweleeuw, J. M, P Kooiman, L.C.R.J Willenborg, and Peter Paul de Wolf. “The Post Randomization Method for Protecting Microdata.” *Qüestió, Quaderns d’Estadística i Investigació Operativa* 22, no. 1 (1998b): 145-156. Available at <http://www.raco.cat/index.php/Questiio/issue/view/2250>

---

<sup>4</sup> This can also be achieved with multidimensional transition matrices. In that case, the probability is not specified for ‘male’ -> ‘female’, but for ‘male’ + ‘rural’ -> ‘female’ + ‘rural’ and for ‘male’ + ‘urban’ -> ‘female’ + ‘urban’.



Marés, Jordi, and Vicenç Torra. 2010. "PRAM Optimization Using an Evolutionary Algorithm." In *Privacy in Statistical Databases*, by Josep Domingo-Ferrer and Emmanouil Magkos, 97-106. Corfú, Greece: Springer.

Warner, S.L. "Randomized Response: A Survey Technique for Eliminating Evasive Answer Bias." *Journal of American Statistical Association* 57 (1965): 622-627.

### 6.3.2 Microaggregation

Microaggregation is most suitable for continuous variables, but can be extended in some cases to categorical variables.<sup>5</sup> It is most useful where confidentiality rules have been predetermined (e.g., a certain threshold for  $k$ -anonymity has been set) that permit the release of data only if combinations of variables are shared by more than a predetermined threshold number of respondents ( $k$ ). The first step in microaggregation is the formation of small groups of individuals that are homogeneous with respect to the values of selected variables, such as groups with similar income or age. Subsequently, the values of the selected variables of all group members are replaced with a common value, e.g., the mean of that group. Microaggregation methods differ with respect to (i) how the homogeneity of groups is defined, (ii) the algorithms used to find homogeneous groups, and (iii) the determination of replacement values. In practice, microaggregation works best when the values of the variables in the groups are more homogeneous. When this is the case, then the information loss due to replacing values with common values for the group will be smaller than in cases where groups are less homogeneous.

In the univariate case, and also for ordinal categorical variables, formation of homogeneous groups is straightforward: groups are formed by first ordering the values of the variable and then creating  $g$  groups of size  $n_i$  for all groups  $i$  in  $1, \dots, g$ . This maximizes the within-group homogeneity, which is measured by the within-groups sum of squares ( $SSE$ )

$$SSE = \sum_{i=1}^g \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)^T (x_{ij} - \bar{x}_i)$$

The lower the  $SSE$ , the higher the within-group homogeneity. The group sizes can differ amongst groups, but often groups of equal size are used to simplify the search<sup>6</sup>.

Choice of group size depends on the homogeneity within the groups and the required level of protection. In general it holds that the larger the group, the higher the protection. A disadvantage of groups of equal sizes is that the data might be unsuitable for this. For instance, if two individuals have a low income (e.g., 832 and 966) and four individuals have a high income (e.g., 3,313, 3,211, 2,987, 3,088), the mean of two groups of size three (e.g.,  $(832 + 966 + 2,987) / 3 = 1,595$  and  $(3,088 + 3,211 + 3,313) / 3 = 3,204$ ) would represent neither the low nor the high income.

Often, values are replaced by the group mean. An alternative, more robust approach is to replace group values with the median. In cases where the median is chosen, one individual in every group keeps the same value if groups have odd sizes. In cases where there is a high degree of heterogeneity within the groups (this is often the case for larger groups), the median is preferred to preserve the information in the data. An example is income, where one outlier can lead to multiple outliers being created when using microaggregation. This is illustrated in Table 6.11. If we choose the mean as replacement for all values, which are grouped with the outlier (6,045 in group 2), these records will be assigned values far from their original values. If we chose the median, the incomes of individuals 1 and 2 are not perturbed, but no value is an outlier. Of course, this might in itself present problems.

**Note:** If microaggregation alters outlying values, this can have a significant impact on the computation of some measures sensitive to outliers, such as the GINI index.

<sup>5</sup> Microaggregation can also be used for categorical data, as long as there is a possibility to form groups and an aggregate replacement for the values in the group can be calculated. This is the case for ordinal variables.

<sup>6</sup> Here all groups can have different sizes (i.e., number of individuals in a group). In practice, the search for homogeneous groups is simplified by imposing equal group sizes for all groups.

In the case where microaggregation is applied to categorical variables, the median is used to calculate the replacement value for the group.

Table 6.11: Illustrating the effect of choosing mean vs. median for microaggregation where outliers are concerned

ID	Group	Income	Microaggregation (mean)	Microaggregation (median)
1	1	2,300	2,245	2,300
2	2	2,434	3,608	2,434
3	1	2,123	2,245	2,300
4	1	2,312	2,245	2,300
5	2	6,045	3,608	2,434
6	2	2,345	3,608	2,434

In case of multiple variables that are candidates for microaggregation, one possibility is to apply univariate microaggregation to each of the variables separately. The advantage of univariate microaggregation is minimal information loss, since the changes in the variables are limited. The literature shows, however, that disclosure risk can be very high if univariate microaggregation is applied to several variables separately and no additional anonymization techniques are applied (*DMOT02*). To overcome this shortcoming, an alternative to univariate microaggregation is multivariate microaggregation.

Multivariate microaggregation is widely used in official statistics. The first step in multivariate aggregation is the creation of homogeneous groups based on several variables. Groups are formed based on multivariate distances between the individuals. Subsequently, the values of all variables for all group members are replaced with the same values. Table 6.12 illustrates this with three variables. We see that the grouping by income, expenditure and wealth leads to a different grouping, as in the case in Table 6.11, where groups were formed based only on income.

Table 6.12: Illustration of multivariate microaggregation

ID	Group	Before microaggregation			After microaggregation		
		Income	Exp	Wealth	Income	Exp	Wealth
1	1	2,300	1,714	5.3	2,285.7	1,846.3	6.3
2	1	2,434	1,947	7.4	2,285.7	1,846.3	6.3
3	1	2,123	1,878	6.3	2,285.7	1,846.3	6.3
4	2	2,312	1,950	8.0	3,567.3	2,814.0	8.3
5	2	6,045	4,569	9.2	3,567.3	2,814.0	8.3
6	2	2,345	1,923	7.8	3,567.3	2,814.0	8.3

There are several multivariate microaggregation methods that differ with respect to the algorithm used for creating groups of individuals. There is a trade-off between speed of the algorithm and within-group homogeneity, which is directly related to information loss. For large datasets, this is especially challenging. We discuss the Maximum Distance to Average Vector (MDAV) algorithm here in more detail. The MDAV algorithm was first introduced by *DoTo05* and represents a good choice with respect to the trade-off between computation time and the group homogeneity, computed by the within-group *SSE*.

The algorithm computes an average record or centroid  $C$ , which contains the average values of all included variables. We select an individual  $A$  with the largest squared Euclidean distance from  $C$ , and build a group of  $k$  records around  $A$ . The group of  $k$  records is made up of  $A$  and the  $k - 1$  records closest to  $A$  measured by the Euclidean distance. Next, we select another individual  $B$ , with the largest squared Euclidean distance from individual  $A$ . With the remaining records, we build a group of  $k$  records around  $B$ . In the same manner, we select an individual  $D$  with the largest distance from  $B$  and, with the remaining records, build a new group of  $k$  records around  $D$ . The process is repeated until we have fewer than  $2 * k$  records remaining. The MDAV algorithm creates groups of equal size with the exception of maybe one last group of remainders. The microaggregated dataset is then computed by replacing each record in the original dataset by the average values of the group to which it belongs. Equal group sizes, however, may not be ideal for data characterized by greater variability. The MDAV algorithm is illustrated in Fig. 6.6.

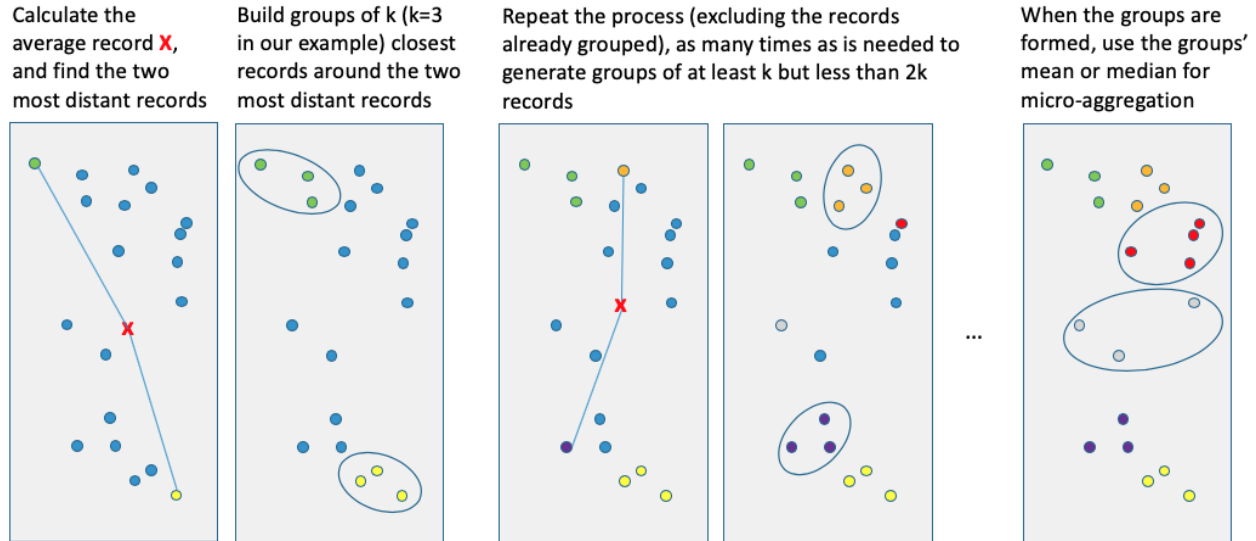


Fig. 6.6: Illustration of MDAV algorithm

It is also possible to group variables only within strata. This reduces the computation time and adds an extra layer of protection to the data, because of the greater uncertainty produced<sup>7</sup>.

Besides the method MDAV, there are few other grouping methods implemented in standard SDC software, such as *sdcMicro* (*TeMK14*). The differences are mainly the distance measure (Euclidean distance, Mahalanobis distance) or sorting based on the first principal component and whether clustering is used. The first principal component (PC) is the projection of all variables into a one-dimensional space maximizing the variance of this projection. The performance of this method depends on the share of the total variance in the data that is explained by the first PC. Using the Mahalanobis distance is computationally more intensive, but provides better results with respect to group homogeneity. It is recommended for smaller datasets (*TeMK14*).

In case of several variables to be used for microaggregation, looking first at the covariance or correlation matrix of these variables is recommended. If not all variables correlate well, but two or more sets of variables show high correlation, less information loss will occur when applying microaggregation separately to these sets of variables. In general, less information loss will occur when applying multivariate microaggregation, if the variables are highly correlated. The advantage of replacing the values with the mean of the groups rather than other replacement values has the advantage that the overall means of the variables are preserved.

### Recommended Reading Material on Microaggregation

Domingo-Ferrer, Josep, and Josep Maria Mateo-Sanz. 2002. "Practical data-oriented microaggregation for statistical disclosure control." *IEEE Transactions on Knowledge and Data Engineering* 14 (2002): 189-201.

Hansen, Stephen Lee, and Sumitra Mukherjee. 2003. "A polynomial algorithm for univariate optimal." *IEEE Transactions on Knowledge and Data Engineering* 15 (2003): 1043-1044.

Hundepool, Anco, Josep Domingo-Ferrer, Luisa Franconi, Sarah Giessing, Rainer Lenz, Jane Naylor, Eric Schulte Nordholt, Giovanni Seri, and Peter Paul de Wolf. 2006. *Handbook on Statistical Disclosure Control*. ESSNet SDC. <http://neon.vb.cbs.nl/casc/handbook.htm>

Hundepool, Anco, Josep Domingo-Ferrer, Luisa Franconi, Sarah Giessing, Eric Schulte Nordholt, Keith Spicer, and Peter Paul de Wolf. 2012. *Statistical Disclosure Control*. Chichester: John Wiley & Sons Ltd. doi:10.1002/9781118348239.

<sup>7</sup> Also the homogeneity in the groups will be generally lower, leading to larger changes, higher protection, but also more information loss, unless the strata variable correlates with the microaggregation variable.

Templ, Matthias, Bernhard Meindl, Alexander Kowarik, and Shuang Chen. 2014, August. “International Household Survey Network (IHSN).” <http://www.ihsn.org/home/software/disclosure-control-toolbox>. (accessed July 9, 2018).

### 6.3.3 Noise addition

Noise addition, or noise masking, means adding or subtracting (small) values to the original values of a variable, and is most suited to protect continuous variables (see *Bran02* for an overview). Noise addition can prevent exact matching of continuous variables. The advantages of noise addition are that the noise is typically continuous with mean zero, and exact matching with external files will not be possible. Depending on the magnitude of noise added, however, approximate interval matching might still be possible.

When using noise addition to protect data, it is important to consider the type of data, the intended use of the data and the properties of the data before and after noise addition, i.e., the distribution – particularly the mean – covariance and correlation between the perturbed and original datasets.

Depending on the data, it may also be useful to check that the perturbed values fall within a meaningful range of values. Fig. 6.8 illustrates the changes in data distribution with increasing levels of noise. For data that has outliers, it is important to note that when the perturbed data distribution is similar to the original data distribution (e.g., at low noise levels), noise addition will not protect outliers. After noise addition, these outliers can generally still be detected as outliers and hence easily be identified. An example is a single very high income in a certain region. After perturbing this income value, the value will still be recognized as the highest income in that region and can thus be used for re-identification. This is illustrated in Fig. 6.7, where 10 original observations (open circles) and the anonymized observations (red triangles) are plotted. The tenth observation is an outlier. The values of the first nine observations are sufficiently protected by adding noise: their magnitude and order has changed and exact or interval matching can be successfully prevented. The outlier is not sufficiently protected since, after noise addition, the outlier can still be easily identified. The fact that the absolute value has changed is not sufficient protection. On the other hand, at high noise levels, protection is higher even for the outliers, but the data structure is not preserved and the information loss is large, which is not an ideal situation. One way to circumvent the outlier problem is to add noise of larger magnitude to outliers than to the other values.

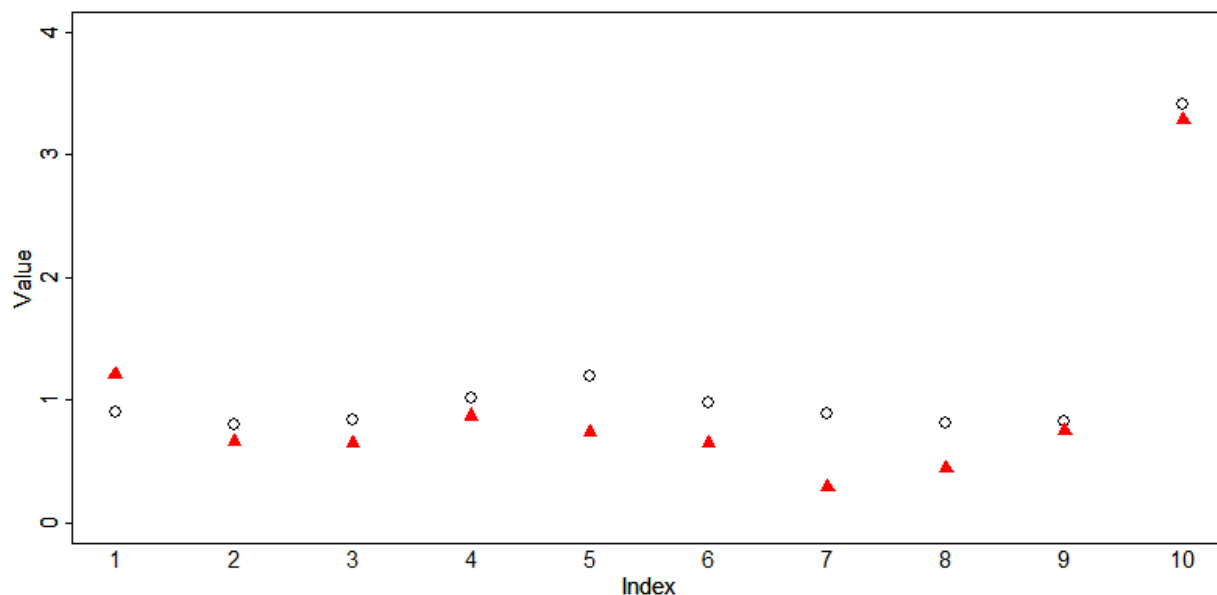


Fig. 6.7: Illustration of effect of noise addition to outliers

There are several noise addition algorithms. The simplest version of noise addition is uncorrelated additive normally

distributed noise, where  $x_j$ , the original values of variable  $j$  are replaced by

$$z_j = x_j + \varepsilon_j,$$

where  $\varepsilon_j \sim N(0, \sigma_{\varepsilon_j}^2)$  and  $\sigma_{\varepsilon_j} = \alpha * \sigma_j$  with  $\sigma_j$  the standard deviation of the original data. In this way, the mean and the covariances are preserved, but not the variances and correlation coefficient. If the level of noise added,  $\alpha$ , is disclosed to the user, many statistics can be consistently estimated from the perturbed data. The added noise is proportional to the variance of the original variable. The magnitude of the noise added is specified by the parameter  $\alpha$ , which specifies this proportion. The standard deviation of the perturbed data is  $1 + \alpha$  times the standard deviation of the original data. A decision on the magnitude of noise added should be informed by the legal situation regarding data privacy, data sensitivity and the acceptable levels of disclosure risk and information loss. In general, the level of noise is a function of the variance of the original variables, the level of protection needed and the desired value range after anonymization<sup>8</sup>. An  $\alpha$  value that is too small will lead to insufficient protection, while an  $\alpha$  value that is too high will make the data useless for data users.

Fig. 6.8 shows the frequency distribution of a numeric continuous variable and the distribution before and after noise addition with different levels of noise (0.1, 0.5, 1, 2 and 5). The first plot shows the distribution of the original values. The histograms clearly show that noise of large magnitudes (high values of alpha) lead to a distribution of the data far from the original values. The distribution of the data changes to a normal distribution when the magnitude of the noise grows respective to the variance of the data. The mean in the data is preserved, but, with an increased level of noise, the variance of the perturbed data grows. After adding noise of magnitude 5, the distribution of the original data is completely destroyed.

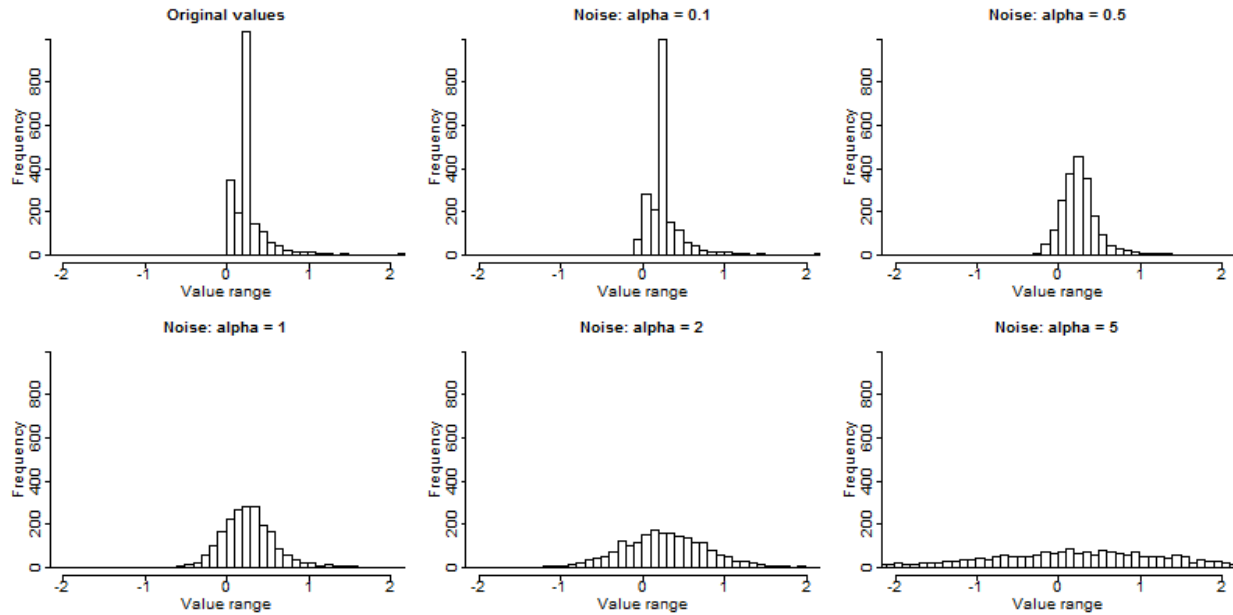


Fig. 6.8: Frequency distribution of a continuous variable before and after noise addition

Fig. 6.9 shows the value range of a variable before adding noise (no noise) and after adding several levels of noise ( $\alpha$  from 0.1 to 1.5 with 0.1 increments). In the figure, the minimum value, the 20<sup>th</sup>, 30<sup>th</sup>, 40<sup>th</sup> percentiles, the median, the 60<sup>th</sup>, 70<sup>th</sup>, 80<sup>th</sup> and 90<sup>th</sup> percentiles and the maximum value are plotted. The median (50<sup>th</sup> percentile) is indicated with the red “+” symbol. From Fig. 6.8 and Fig. 6.9, it is apparent that the range of values expands after noise addition, and the median stays roughly at the same level, as does the mean by construction. The larger the magnitude of noise added, the wider the value range. In cases where the variable should stay in a certain value range (e.g., only positive values, between 0 and 100), this can be a disadvantage of noise addition. For instance, expenditure variables typically have non-negative values, but adding noise to these variables can generate negative values, which are difficult to interpret.

<sup>8</sup> Common values for  $\alpha$  are between 0.5 and 2.

One way to get around this problem is to set any negative values to zero. This truncation of values below a certain threshold, however, will distort the distribution (mean and variance matrix) of the perturbed data. This means that the characteristics that were preserved by noise addition, such as the conservation of the mean and covariance matrix, are destroyed and the user, even with knowledge of the magnitude of the noise, can no longer use the data for consistent estimation.

Another way to avoid negative values is the application of multiplicative rather than additive noise. In that case, variables are multiplied by a random factor with expectation 1 and a positive variance. This will also lead to larger perturbations (in absolute value) of large initial values (outliers). If the variance of the noise added is small, there will be no or few negative factors and thus fewer sign changes than in case of additive noise masking. Multiplicative noise masking is not implemented in *sdcMicro*, but can be relatively easily implemented in base *R* by generating a vector of random numbers and multiplying the data with this vector. For more information on multiplicative noise masking and the properties of the data after masking, we refer to *KiWi03*.

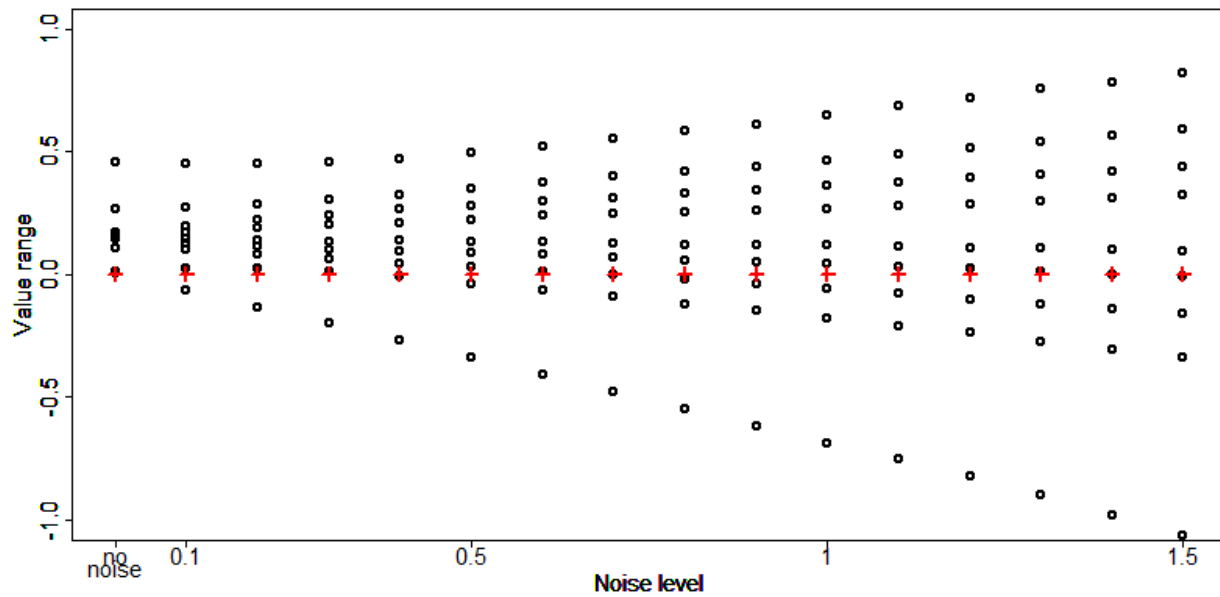


Fig. 6.9: Noise levels and the impact on the value range (percentiles)

If two or more variables are selected for noise addition, correlated noise addition is preferred to preserve the correlation structure in the data. In this case, the covariance matrix of noise  $\Sigma_{\varepsilon}$  is proportional to the covariance matrix of the original data  $\Sigma_X$  :

$$\Sigma_{\varepsilon} = \alpha \Sigma_X$$

If noise addition is applied to variables that are a ratio of an aggregate, this structure can be destroyed by noise addition. Examples are income and expenditure data with many income and expenditure categories. The categories add up to total income or total expenditures. In the original data, the aggregates match with the sum of the components. After adding noise to their components (e.g., different expenditure categories), however, their new aggregates will not necessarily match the sum of the categories anymore. One way to keep this structure is to add noise only to the aggregates and release the components as ratio of the perturbed aggregates. Subsequently, the ratios of the initial expenditure categories are used for each individual to reconstruct the perturbed values for each expenditure category.

### Recommended Reading Material on Noise Addition

Brand, Ruth. 2002. "Microdata Protection through Noise Addition." In *Inference Control in Statistical Databases - From Theory to Practice*, edited by Josep Domingo-Ferrer. Lecture Notes in Computer Science Series 2316, 97-116. Berlin Heidelberg: Springer. [http://link.springer.com/chapter/10.1007%2F3-540-47804-3\\_8](http://link.springer.com/chapter/10.1007%2F3-540-47804-3_8)

Kim, Jay J, and William W Winkler. 2003. “Multiplicative Noise for Masking Continuous Data.” *Research Report Series* (Statistical Research Division. US Bureau of the Census). <https://www.census.gov/srd/papers/pdf/rrs2003-01.pdf>

Torra, Vicenç, and Isaac Cano. 2011. “Edit Constraints on Microaggregation and Additive Noise.” In *Privacy and Security Issues in Data Mining and Machine Learning*, edited by C. Dimitrakakis, A. Gkoulalas-Divanis, A. Mitrokotsa, V. S. Verykios, Y. Saygin. Lecture Notes in Computer Science Volume 6549, 1-14. Berlin Heidelberg: Springer. <http://link.springer.com/book/10.1007/978-3-642-19896-0>

Mivule, K. 2013. “Utilizing Noise Addition for Data Privacy, An Overview.” *Proceedings of the International Conference on Information and Knowledge Engineering (IKE 2012)*, (pp.65-71).Las Vegas, USA. <http://arxiv.org/ftp/arxiv/papers/1309/1309.3958.pdf>

---

### 6.3.4 Rank swapping

Data swapping is based on interchanging values of a certain variable across records. Rank swapping is one type of data swapping, which is defined for ordinal and continuous variables. For rank swapping, the values of the variable are first ordered. The possible number of values for a variable to swap with is constrained by the values in a neighborhood around the original value in the ordered values of the dataset. The size of this neighborhood can be specified, e.g., as a percentage of the total number of observations. This also means that a value can be swapped with the same or very similar values. This is especially the case if the neighborhood is small or there are only a few different values in the variable (ordinal variable). An example is the variables “education” with only few categories: (‘none’, ‘primary’, ‘secondary’, ‘tertiary’). In these cases, rank swapping is not a suitable method.

If rank swapping is applied to several variables simultaneously, the correlation structure between the variables is preserved. Therefore, it is important to check whether the correlation structure in the data is plausible. Since rank swapping is a probabilistic method, i.e., the swapping depends on a random number generating mechanism, specifying a seed for the random number generator before using rank swapping is recommended to guarantee reproducibility of results.

Rank swapping has been found to yield good results with respect to the trade-off between information loss and data protection (*DoTo01a*). Rank swapping is not useful for variables with few different values or many missing values, since the swapping in that case will not result in altered values. Also, if the intruder knows to whom the highest or lowest value of a specific variable belongs (e.g., income), the level of this variable will be disclosed after rank swapping, because the values themselves are not altered and the original values are all disclosed. This can be solved by top and bottom coding the lowest and/or highest values.

---

#### Recommended Reading Material on Rank Swapping

Dalenius T. and Reiss S.P. 1978. Data-swapping: a technique for disclosure control (extended abstract). In Proc. ASA Section on Survey Research Methods. American Statistical Association, Washington DC, 191–194.

Domingo-Ferrer J. and Torra V. 2001. “A Quantitative Comparison of Disclosure Control Methods for Microdata.” In *Confidentiality, Disclosure and Data Access: Theory and Practical Applications for Statistical Agencies*, edited by P. Doyle, J.I. Lane, J.J.M. Theeuwes, and L. Zayatz, 111–134. Amsterdam, North-Holland.

Hundepool A., Van de Wetering A., Ramaswamy R., Franconi F., Poletini S., Capobianchi A., De Wolf P.-P., Domingo-Ferrer J., Torra V., Brand R. and Giessing S. 2007. *μ-Argus User’s Manual* version 4.1.

---

### 6.3.5 Shuffling

Shuffling as introduced by *MuSa06* is similar to swapping, but uses an underlying regression model for the variables to determine which variables are swapped. Shuffling can be used for continuous variables and is a deterministic method.



Shuffling maintains the marginal distributions in the shuffled data. Shuffling, however, requires a complete ranking of the data, which can be computationally very intensive for large datasets with several variables.

The method is explained in detail in *MuSa06*. The idea is to rank the individuals based on their original variables. Then fit a regression model with the variables to be protected as regressands and a set of variables that predict this variable well (i.e., are correlated with) as regressors. This regression model is used to generate  $n$  synthetic (predicted) values for each variable that has to be protected. These generated values are also ranked and each original value is replaced with another original value with the rank that corresponds to the rank of the generated value. This means that all original values will be in the data. Table 6.13 presents a simplified example of the shuffling method. The regressands are not specified in this example.

Table 6.13: Simplified example of the shuffling method

ID	Income (orig)	Rank (orig)	Income (pred)	Rank (pred)	Shuffled values
1	2,300	2	2,466.56	4	2,345
2	2,434	6	2,583.58	7	2,543
3	2,123	1	2,594.17	8	2,643
4	2,312	3	2,530.97	6	2,434
5	6,045	10	5,964.04	10	6,045
6	2,345	4	2,513.45	5	2,365
7	2,543	7	2,116.16	1	2,123
8	2,854	9	2,624.32	9	2,854
9	2,365	5	2,203.45	2	2,300
10	2,643	8	2,358.29	3	2,312

The suitability of shuffling depends on the predictive power of the regressors for the variables to be predicted. This can be checked with goodness-of-fit measures such as the  $R^2$  of the regression. The  $R^2$  captures only linear relations, but these are also the only relations that are captured by the linear regression model used for shuffling.

---

#### Recommended Reading Material on Shuffling

K. Muralidhar and R. Sarathy. 2006."Data shuffling - A new masking approach for numerical data," *Management Science*, 52, 658-670.

---

### 6.3.6 Comparison of PRAM, rank swapping and shuffling

PRAM, rank swapping and shuffling are all perturbative methods, i.e., they change the values for individual records and are mainly used for continuous variables. After rank swapping and shuffling, the original values are all contained in the treated dataset but might be assigned to other records. This implies that univariate tabulations are not changed. This also holds in expectation for PRAM, if a transition matrix is chosen that has the invariant property.

Choosing a method is based on the structure to be preserved in the data. In cases where the regression model fits the data well, data shuffling would work very well, as there should be sufficient (continuous) regressors available. Rank swapping works well if there are sufficient categories in the variables. PRAM is preferred if the perturbation method should be applied to only one or few variables; the advantage is the possibility of specifying restrictions on the transition matrix and applying PRAM only within strata, which can be user defined.

## 6.4 Anonymization of geospatial variables

Recently, geospatial data has become increasingly popular with researchers and wide-spread. Georeferenced data identifies the geographical location for each record with the help of a Geographical Information System (GIS), that



uses for instance GPS (Global Positioning System) coordinates or address data. The advantages of geospatial data are manifold: 1) researchers can create their own geographical areas, such as the service area of a hospital; 2) it enables researchers to measure the proximity to facilities, such as schools; 3) researchers can use the data to extract geographical patterns; and 4) it enables linking of data from different sources (see e.g., *BCRZ13*). However, geospatial data, due to the precise reference to a location, also pose a challenge to the privacy of the respondents.

One way to anonymize georeferenced data is removing the GIS variables and instead leaving in or creating other geographical variables, such as province, region. However, this approach also removes the benefits of geospatial data. Another option is the geographical displacement of areas and/or records. *BCRZ13* describe a geographical displacement procedure for a health dataset. This paper also includes the code in Python. *HuDr15* propose three different strategies for generating synthetic geocodes.

---

### Recommended Reading Material on Anonymization of Geospatial Data

C.R. Burgert, J. Colston, T. Roy and B. Zachary. 2013. “DHS Spatial Analysis Report No. 7 - Geographic Displacement Procedure and Georeferenced Data Release Policy for the Demographic and Health Surveys” (USAID). <http://dhsprogram.com/pubs/pdf/SAR7/SAR7.pdf>

J. Hu and J. Drechsler. 2015. “Generating synthetic geocoding information for public release.” <http://www.iab.de/389/section.aspx/Publikation/k150601301>

---

## 6.5 Anonymization of the quasi-identifier household size

The size of a household is an important identifier, especially for large households.<sup>9</sup> Suppression of the actual size variable, if available (e.g., number of household members), however, does not suffice to remove this information from the dataset, as a simple count of the household members for a particular household will allow this variable to be reconstructed as long as a household ID is in the data. In any case, households of a very large size or with a unique or special key (i.e., combination of values of quasi-identifiers) should be checked manually. One way to treat them is to remove these households from the dataset before release. Alternatively, the households can be split, but care should be taken to suppress or change values for these households to prevent an intruder from immediately understanding that these households have been split and reconstructing them by combining the two households with the same values.

## 6.6 Special case: census data

Census microdata are a special case because the user (and intruder) knows that all respondents are included in the dataset. Therefore, risk measures that use the sample weights and are based on uncertainty of the correctness of a match are no longer applicable. If an intruder has identified a sample unique and successfully matched, there is no doubt whether the match is correct, as it would be in the case of a sample. One approach to release census microdata is to release a stratified sample of the sample (1 – 5% of the total census).

---

**Note:** After sampling, the anonymization process has to be followed; sampling alone is not sufficient to guarantee confidentiality.

---

Several statistical offices release microdata based on census data. A few examples are:

- **The British Office for National Statistics (ONS)** released several files based on the 2011 census: 1. A micro-data teaching file for educational purposes. This file is a 1% sample of the total census with a limited set of variables. 2. Two scientific use files with 5% samples are available for registered researchers who accept

---

<sup>9</sup> Even if the dataset does not contain an explicit variable with household size, this information can be easily extracted from the data and should be taken into account. The Section [Household structure](#) shows how to create a variable “household size” based on the household IDs.

the terms and conditions of their use. 3. Two 10% samples are available in controlled research data centers for approved researchers and research goals. All these files have been anonymized prior to release.<sup>10</sup>

- **The U.S. Census Bureau** released two samples of the 2000 census: a 5% sample on the national level and a 1% sample on the state level. The national level file is more detailed, but the most detailed geographical area has at least 400,000 people. This, however, allows representation of all states from the dataset. The state-level file has less detailed variables but a more detailed geographical structure, which allows representation of cities and larger counties from the dataset (the minimum size of a geographical area is 100,000). Both files have been anonymized by using data swapping, top coding, perturbation and reducing detail by recoding.<sup>11</sup>

## References

---

<sup>10</sup> More information on census microdata at ONS is available on their website: <http://www.ons.gov.uk/ons/guide-method/census/2011/census-data/census-microdata/index.html>

<sup>11</sup> More information on the anonymization of these files is available on the website of the U.S. Census Bureau: <https://www.census.gov/population/www/cen2000/pums/index.html>

---

## Measuring Utility and Information Loss

---

SDC is a trade-off between risk of disclosure and loss of data utility and seeks to minimize the latter, while reducing the risk of disclosure to an acceptable level. Data utility in this context means the usefulness of the anonymized data for statistical analyses by end users as well as the validity of these analyses when performed on the anonymized data. Disclosure risk and its measurement are defined in the Section [Measure Risk](#) of this guide. In order to make a trade-off between minimizing disclosure risk and maximizing utility of data for end users, it is necessary to measure the utility of the data after anonymization and compare it with the utility of the original data. This section describes measures that can be used to compare the data utility before and after anonymization, or alternatively quantify the information loss. Information loss is the inverse of data utility: the larger the data utility after anonymization, the smaller the information loss.

---

**Note:** If the microdata to be anonymized is based on a sample, the data will incur a sampling error. Also other errors may be present in the data, such as nonresponse error.

---

---

**Note:** The methods discussed here only measure the information loss caused by the anonymization process relative to the original sample data and do not attempt to measure the error caused by other sources.

---

Ideally, the information loss is evaluated with respect to the needs and uses of the end users of the microdata. However, different end users of anonymized data may have very diverse uses for the released data and it might not be possible to collect an exhaustive list of these uses. Even if many uses can be identified, the characteristics in the data needed for these uses can be contradictory (e.g., one user needs a detailed geographical level whereas another is interested in a detailed age structure and does not need a detailed geographical structure). Nevertheless, as pointed out earlier, only one anonymized dataset can be released for each dataset and every type of release to avoid unintended disclosure. Releasing multiple anonymized datasets for different purposes may lead to unintended disclosure.<sup>1</sup> Therefore, it is not possible to anonymize and release a file tailored to each user's needs.

Since collecting and taking into account all data uses is often impossible, we also look at general (or generic) information loss measures besides user- and data-specific information loss measures. These measures do not take into account

---

<sup>1</sup> It is possible to release data files for different groups of users, e.g., PUF and SUF. All information in the less detailed file, however, must also be included in the more detailed file to prevent unintended disclosure. Datasets released in data enclaves can be customized for the user, since the risk that they will be combined with other version is zero.

the specific data use, but can be used as guiding measures for information loss and evaluating whether a dataset is still analytically valid after anonymization. The main idea for such measures is to compare records between the original and treated datasets and compare statistics computed from both datasets (*HDFG12*). Examples of such measures are the number of suppressions, number of changed values, changes in contingency tables and changes in mean and covariance matrices.

Many of the SDC methods discussed earlier are parametric, in the sense that their outcome depends on parameters chosen by the user. Examples are the cluster size for microaggregation (see the Section [Microaggregation](#)) or the importance vector in local suppression (see the Section [Local suppression](#)). Data utility and information loss measures are useful for choosing these parameters by comparing the impact of different parameters on the information loss. [Fig. 7.1](#) illustrates this by showing the trade-off between the disclosure risk and data utility of a hypothetical dataset. The triangle represents the original data with full utility and a certain level of disclosure risk, which is too high for disclosure. The square represents no release of microdata. Although there is no risk of disclosure, there is also no utility from the data for users since no data is released. The points in between represent the result of applying different SDC methods with different parameter specifications. We would select the SDC method corresponding to the point, which maximizes the utility, while keeping disclosure risk at an acceptable level.

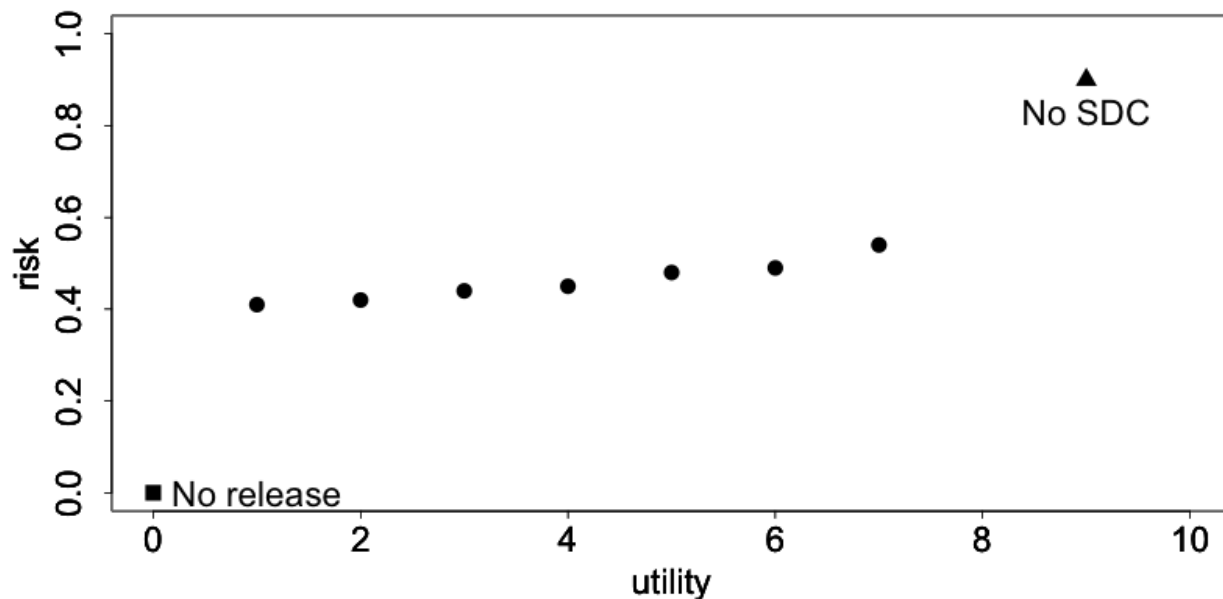


Fig. 7.1: The trade-off between risk and utility in a hypothetical dataset

In the following sections, we first propose general utility measures independent of data use, and later present an example of a specific measure useful to measure information loss with respect to specific data uses. Finally, we show how to visualize changes in the data caused by anonymization and discuss the selection of utility measures for a particular dataset.

## 7.1 General utility measures for continuous and categorical variables

General or generic measures of information loss can be divided into those comparing the actual values of the raw and anonymized data, and those comparing statistics from both datasets. All measures are a posteriori, since they measure utility after anonymization and require both the data before and after the anonymization process. General utility measures are different for categorical and continuous variables.

### 7.1.1 General utility measures for categorical variables

#### Number of missing values

An informative measure is to compare the number of missing values in the data. Missing values are often introduced after suppression and more suppressions indicate a higher degree of information loss. More generally, it is possible to count and compare the number of missing values in the original data and the treated data. This can be useful to see the proportional increase in the number of missing values. Missing values can also have other sources, such as nonresponse.

#### Number of records changed

Another useful statistic is the number of records changed per variable. The number of records changed gives a good indication of the impact of the anonymization methods on the data.

#### Comparing contingency tables

A useful way to measure information loss in categorical variables is to compare univariate tabulations and, more interestingly, contingency tables (also cross tabulations or two-way tables) between pairs of variables. To maintain the analytical validity of a dataset, the contingency tables should stay approximately the same.

*DoTo01b* propose a Contingency Table-Based Information Loss (CTBIL) measure, which quantifies the distance between the contingency tables in the original and treated data. Alternatively, visualizations of the contingency table with mosaic plots can be used to compare the impact of anonymization methods on the tabulations and contingency tables (see the Section *Mosaic plots*).

### 7.1.2 General utility measures for continuous variables

#### Statistics: mean, covariance, correlation

The statistics characterizing the dataset should not change after the anonymization. Examples of such statistics are the mean, variance, and covariance and correlation structure of the most important variables in the dataset. Other statistics characterizing the data include the principal components and the loadings. *DoTo01b* give an overview of statistics that can be considered. In order to evaluate the information loss caused by the anonymization, one should compare the appropriate statistics for continuous variables computed from the data before and after anonymization. There are several ways to evaluate the loss of utility with respect to the changes in these statistics, for instance, by comparing means and (co-)variances in the data or comparing the (multivariate) distributions of the data. Especially changes in the correlations gives valuable information on the validity of the data for regressions.

*DoTo01b* propose several measures for the discrepancy between the covariance and correlation matrices. These measures are based on the mean squared error, the mean absolute error or the mean variation of the individual cells. We refer to *DoTo01b* for a complete overview of these measures.

#### IL1s information loss measure

Alternatively, we can also compare the actual data and quantify the distance between the original dataset  $X$  and the treated dataset  $Z$ . Here  $X$  and  $Z$  contain only continuous variables. *YaWC02* introduce the distance measure IL1s, which is the sum of the absolute distances between the corresponding observations in the raw and anonymized datasets, which are standardized by the standard deviation of the variables in the original data. For the continuous variables in

the dataset, the IL1s measure is defined as

$$IL1s = \frac{1}{pn} \sum_{j=1}^p \sum_{i=1}^n \frac{|x_{ij} - z_{ij}|}{\sqrt{2}S_j},$$

where  $p$  is the number of continuous variables;  $n$  is the number of records in the dataset;  $x_{ij}$  and  $z_{ij}$ , respectively, are the values before and after anonymization for variable  $j$  and individual  $i$ ; and  $S_j$  is the standard deviation of variable  $j$  in the original data (YaWC02).

The measure is useful for comparing different methods. The smaller the value of the measure, the closer the values are to the original values and the higher the utility.

---

**Note:** This measure is related to risk measures based on distance and intervals (see the Section [Risk measures for continuous variables](#)).

---

The greater the distance between the original and anonymized values, the lower the data utility. Greater distance, however, also reduces the risk of re-identification.

### Eigenvalues

Another way to evaluate the information loss is to compare the robust eigenvalues of the data before and after anonymization. Again, the main use is to compare different methods. The greater the value, the greater the changes in the data and the information loss.

## 7.2 Utility measures based on the end user's needs

Not all needs and uses of a certain dataset can be inventoried. Nevertheless, some types of data have similar uses or important characteristics, which can be evaluated before and after anonymization. Examples of such “benchmarking indicators” (*TMKCI4*) are different for each dataset. Examples include poverty measures for income datasets and school attendance ratios. Often ideas for selecting such indicators come from the reports data users publish based on previously released microdata.

The approach is to compare the indicators calculated on the untreated data and the data after anonymization with different methods. If the differences between the indicators are not too large, the anonymized dataset can be released for use by researchers. It should be taken into account that indicators calculated on samples are estimates with a certain variance and confidence interval. Therefore, for sample data, it is more informative to compare the overlap of confidence intervals and/or to evaluate whether the point estimate calculated after anonymization is contained within the confidence interval of the original estimate. Examples of benchmark indicators and their confidence intervals are included in the case studies in these guidelines. Here we give the example of the GINI coefficient.

The GINI coefficient is a measure of statistical dispersion, which is often used to measure inequality in income. A way to measure the information loss in income data is to compare the income distribution, which can be easily done by comparing the GINI coefficients. The GINI coefficient of sample data is a random variable. Therefore, it is useful to construct a confidence interval around the coefficient to evaluate the significance of any change in the coefficient after anonymization.

## 7.3 Regression

Besides comparing covariance and correlation matrices, regressions are a useful tool to evaluate whether the structure in the data is maintained after anonymization. By comparing regressions parameters, it is also possible to compare relations between non-continuous variables (e.g., by introducing dummy variables or regression with ordinal variables).

If it is known for what purpose and in what field the data is used, common regressions can be used to compare the change in coefficients and confidence intervals.

An example of using regression to evaluate the data utility in income data is the Mincer equation. The Mincer equation explains earnings as a function of education and experience while controlling for other variables. The Mincer equation is often used to evaluate the gender pay gap and gender wage inequality by including a gender dummy. Here we show how to evaluate the impact of anonymization methods on the gender coefficient. We regress the log income on a constant, a gender dummy, years of education, years of experience, years of experience squared and other factors influencing wage.

$$\ln(\text{wage}) = \beta_0 + \beta_1 \text{gender} + \beta_2 \text{education} + \beta_3 \text{experience} + \beta_4 \text{experience}^2 + \beta_5 X$$

The parameter of interest here is  $\beta_1$ , the effect of gender on the log wage.  $X$  is a matrix with several other factors influencing wage and  $\beta$  the coefficients of these factors.

If the new estimates fall within the original confidence interval and the new and original confidence intervals are greatly overlapping, the data can be considered valid for this type of regression after anonymization. Fig. 7.2 shows the point estimates and confidence intervals for the gender coefficient in this trade-off for a sample income dataset and several SDC methods and parameters. The red dot and confidence bar (on the top) correspond to the estimates for the untreated data, whereas the other confidence bars correspond to the respective SDC methods and different parameters. The anonymization reduces the number of expected re-identifications in the data (left axis) and the point estimates and confidence intervals vary greatly for the different SDC methods. We would choose a method, which reduces the expected number of identifications, while not changing the gender coefficient and having a great overlap of the confidence interval with the confidence interval estimated from the original data.

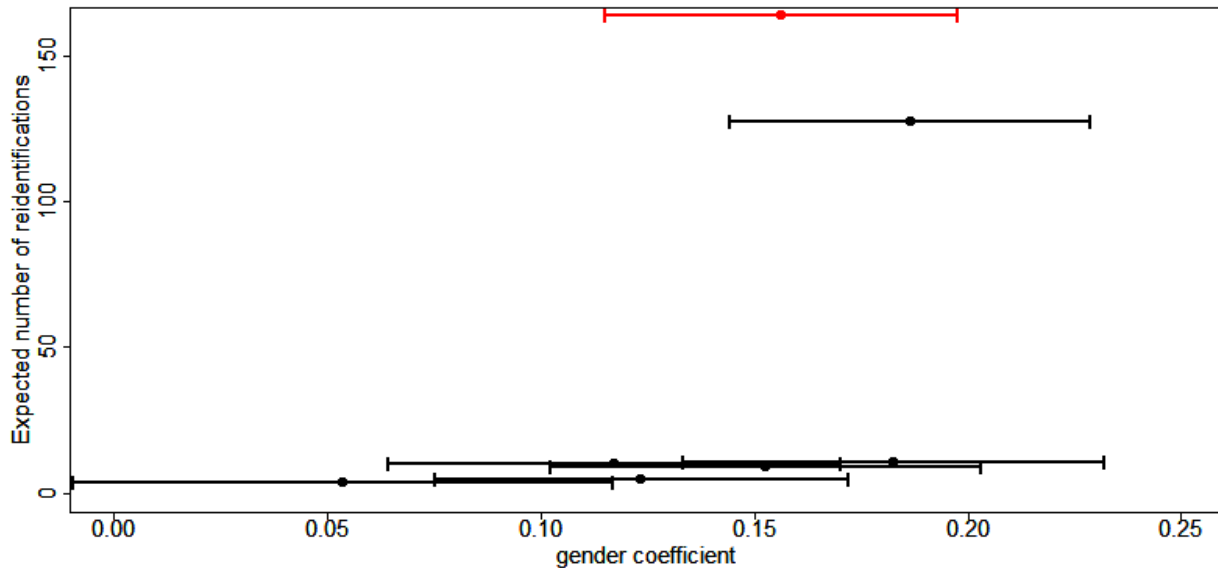


Fig. 7.2: Effect of anonymization on the point estimates and confidence interval of the gender coefficient in the Mincer equation

## 7.4 Assessing data utility with the help of data visualizations

The use of graphs and other visualization techniques is a good way to assess at a glance how much the data have changed after anonymization, and can aid the selection of appropriate anonymization techniques for the data. Visualizations can be a useful tool to assess the impact on data utility of anonymization methods and helps choose among

anonymization methods. The software package *R* provides several functions and packages that can help visualize the results of anonymization. This section provides a few examples. We present the following visualizations:

- histograms and density plots
- boxplots
- mosaic plots

To make appropriate visualizations, we need to use the raw data and the anonymized data.

### 7.4.1 Histograms and density plots

Histograms and density plots are useful for quick comparisons of variable distribution before and after anonymization. The advantage of histograms is that the results are exact. Visualization depends on the bin widths and the start point of the first bin, however. Histograms can be used for continuous and semi-continuous variables. Density plots display the kernel density of the data; therefore, the plot depends on the kernel that is chosen and whether the data fits the kernel well. Nevertheless, density plots are a good tool to illustrate the change of values and value ranges of continuous variables.

Fig. 7.3 and Fig. 7.4 show examples. The histograms and density plots give a clear indication how the values have changed: the variability of the data has increased and the shape of the distribution has changed.

---

**Note:** The vertical axes of the histograms have different scales.

---

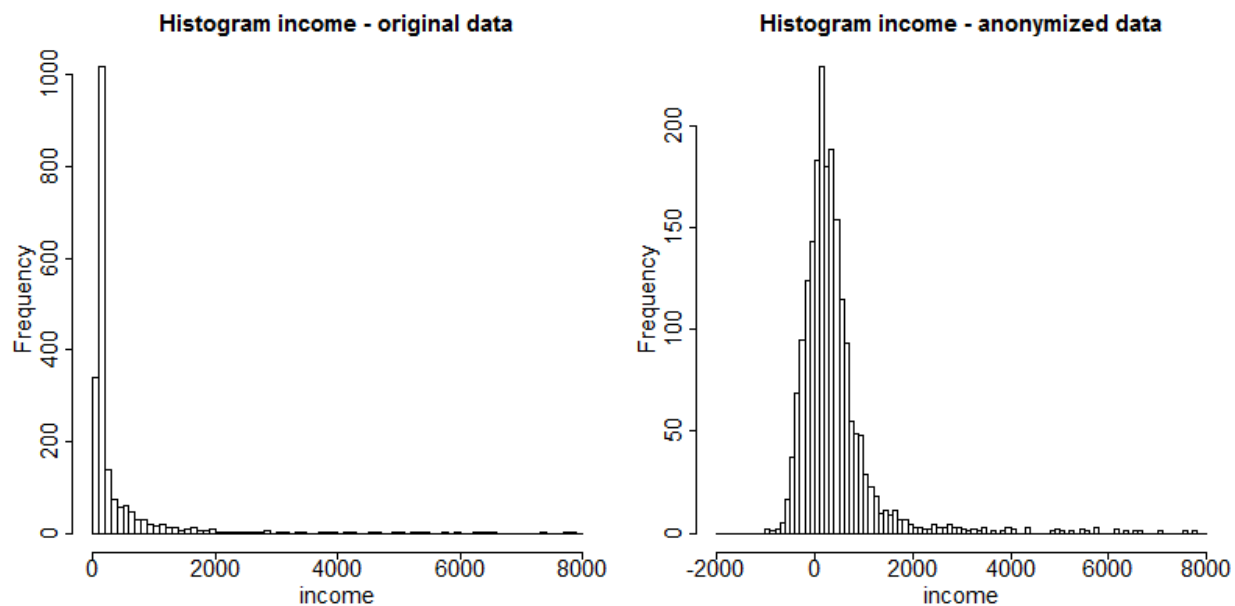


Fig. 7.3: Histograms of income before and after anonymization

### 7.4.2 Box plots

Box plots give a quick overview of the changes in the spread and outliers of continuous variables before and after anonymization. Fig. 7.5 shows an example for an expenditure variable after adding noise. The box plot shows clearly that the variability in the expenditure variable increased as a result of the anonymization methods applied.



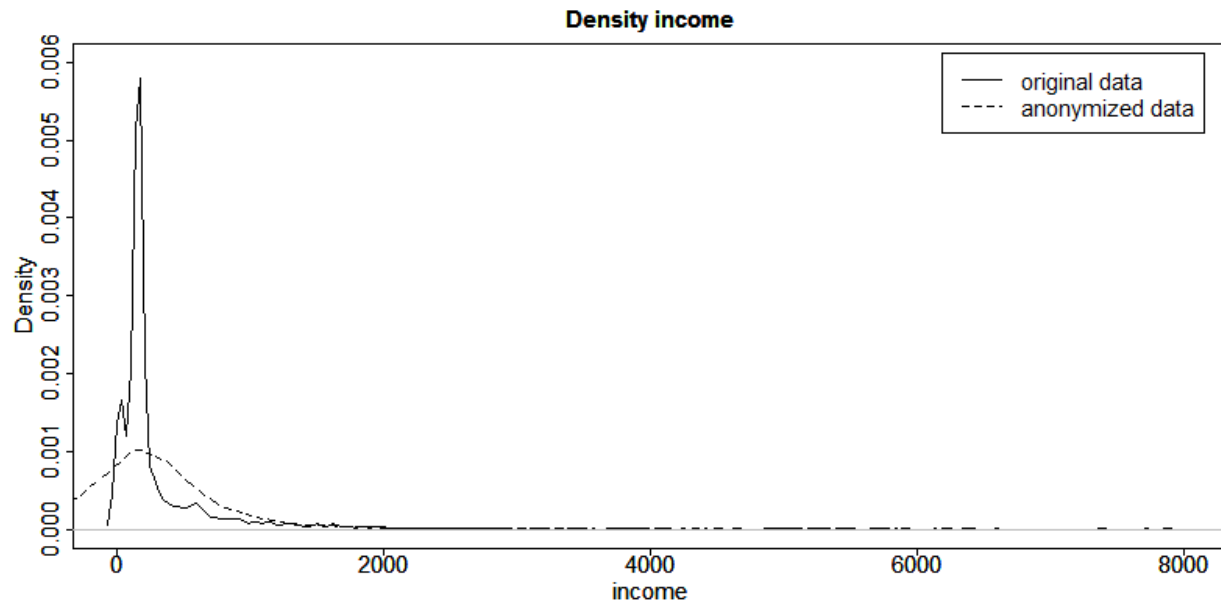


Fig. 7.4: Density plots of income before and after anonymization

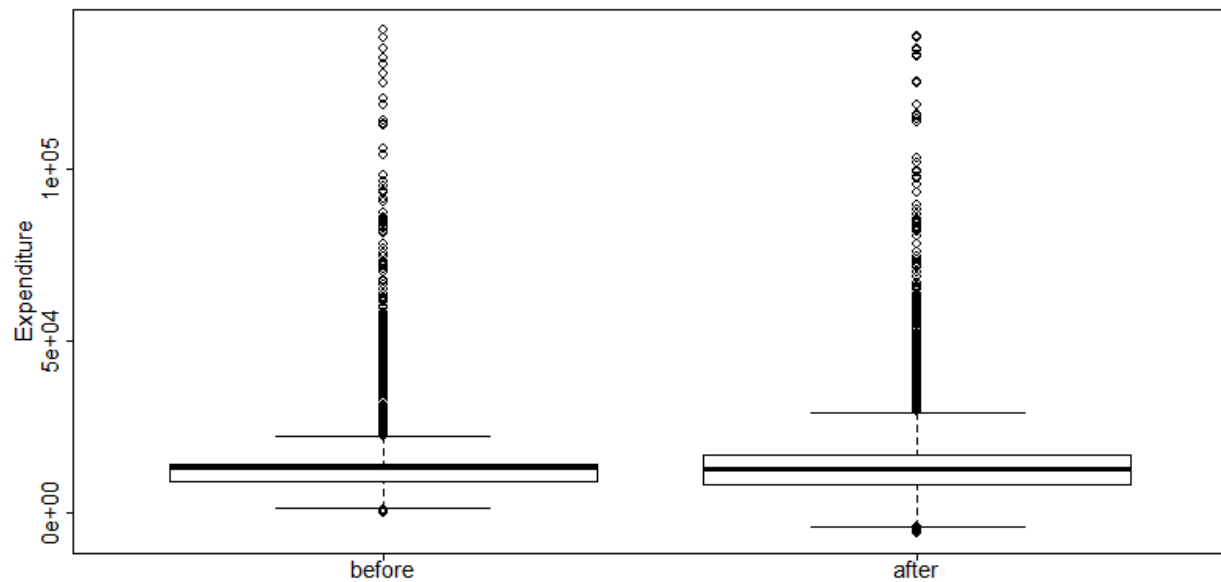


Fig. 7.5: Example of box plots of an expenditure variable before and after anonymization

7.4.3 Mosaic plots

Univariate and multivariate mosaic plots are useful for showing changes in the tabulations of categorical variables, especially when comparing several “scenarios” next to one another. A scenario here refers to the choice of anonymization methods and their parameters. With mosaic plots we can, for instance, quickly see the effect of different levels of  $k$ -anonymity or differences in the importance vectors in the local suppression algorithm (see the Section [Local suppression](#)).

We illustrate the changes in tabulations with an example of the variable “WATER” before and after applying PRAM. We can use mosaic plots to quickly see the changes for each category. Looking at the mosaic plot in [Fig. 7.6](#) we see invariant PRAM has virtually no influence on the univariate distribution.

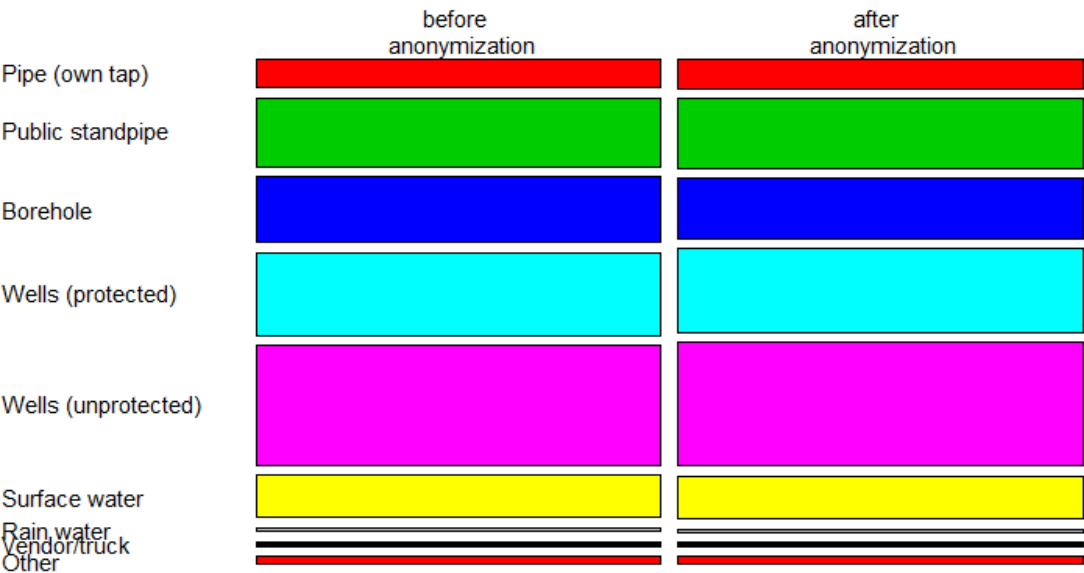


Fig. 7.6: Mosaic plot to illustrate the changes in the WATER variable

We use the variables “gender” and “relationship status” to illustrate the use of mosaic plots for the illustration of changes in univariate tabulations introduced by several sets of anonymization methods. [Table 7.1](#) provides the methods applied in each scenario. Scenario 0, the base scenario, shows the original categories of the gender and relationship status variables, while scenarios 1 to 6 show shifts in the categories after applying different anonymization techniques. [Table 6.1](#) provides a description of the anonymization methods used in each scenario. In total we visualize the impact of six different sets of anonymization methods. We can use mosaic plots to quickly see which set of methods has what impact on the gender and relationship status variables, which can be used to select the best scenario. Looking at the mosaic plots in [Fig. 7.7](#) , we see that scenarios 2, 5 and 6 give the smallest changes for the gender variable and scenarios 3 and 4 for the relationship status variable.

Table 7.1: Description of anonymization methods by scenario

Scenario	Description of anonymization methods applied
0 (base)	Original data, no treatment
1	Recode age (five-year intervals), plus local suppression (required $k = 3$ , high importance on water, toilet and literacy variables)
2	Recode age (five-year intervals), plus local suppression (required $k = 5$ , no importance vector)
3	Recode age (five-year intervals), plus local suppression (required $k = 3$ , high importance on toilet), while also recoding region, urban, education level and occupation variables
4	Recode age (five-year steps), plus local suppression (required $k = 5$ , high importance on water, toilet and literacy), while also recoding region, urban, education level and occupation variables
5	Recode age (five-year intervals), plus local suppression (required $k = 3$ , no importance vector), microaggregation (wealth index), while also recoding region, urban, education level and occupation variables
6	Recode age (five-year intervals) plus local suppression (required $k=3$ , no importance vector), PRAM literacy, while also recoding region, urban, education level and occupation variables

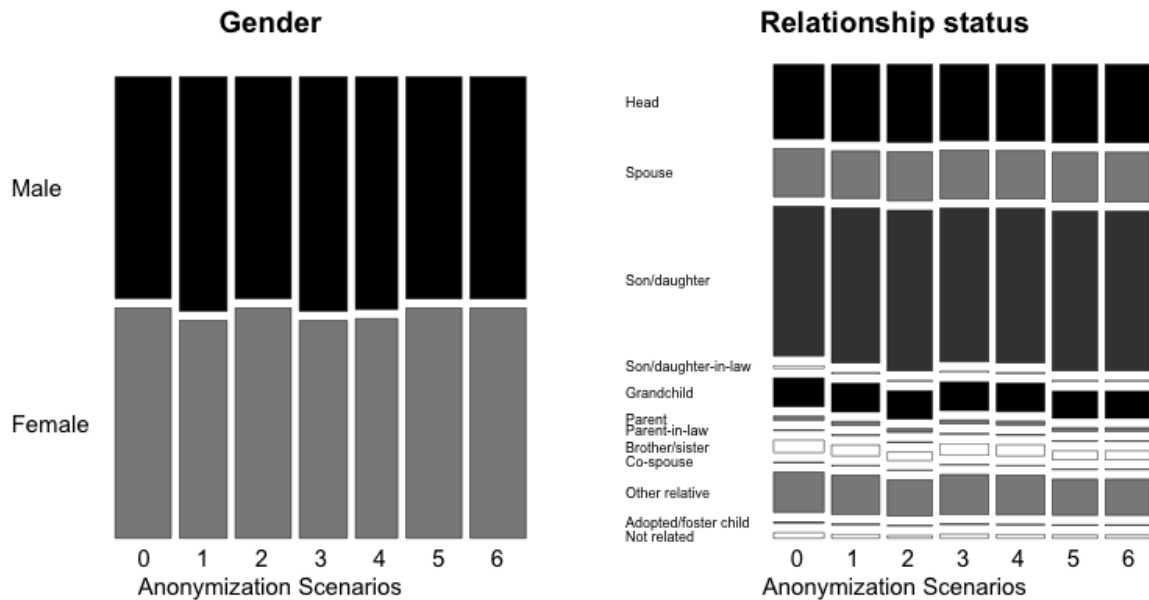


Fig. 7.7: Comparison of treated vs. untreated gender and relationship status variables with mosaic plots

As we discussed in the Section [PRAM \(Post RAndomization Method\)](#), invariant PRAM preserves the univariate distributions. Therefore, in this case it is more interesting to look at the multivariate mosaic plots. Mosaic plots are also a powerful tool to show changes in cross-tabulations/contingency tables. To compare the changes, we need to compare two different plots. [Fig. 7.8](#) and [Fig. 7.9](#) illustrate that (invariant) PRAM does not preserve the two-way tables in this case.

## 7.5 Choice of utility measure

Besides the users' requirements on the data, the utility measures should be chosen in accordance with the variable types and anonymization methods employed. The employed utility measures can be a combination of both general and user-specific measures. As discussed earlier, different utility measures should be used for continuous and categorical data.

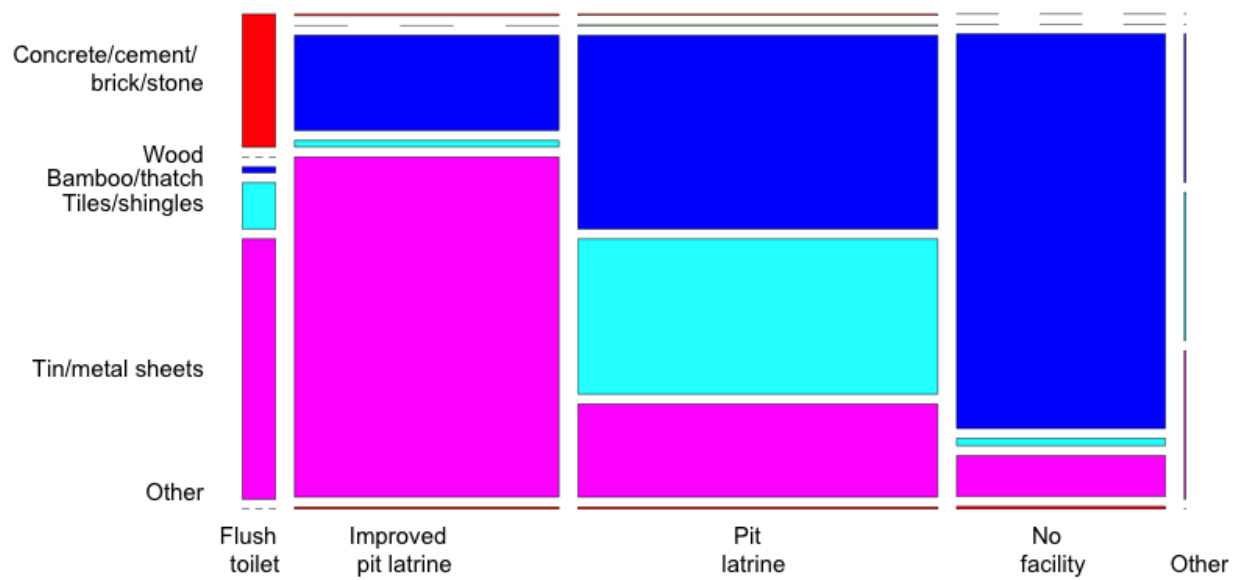


Fig. 7.8: Mosaic plot of the variables ROOF and TOILET before anonymization

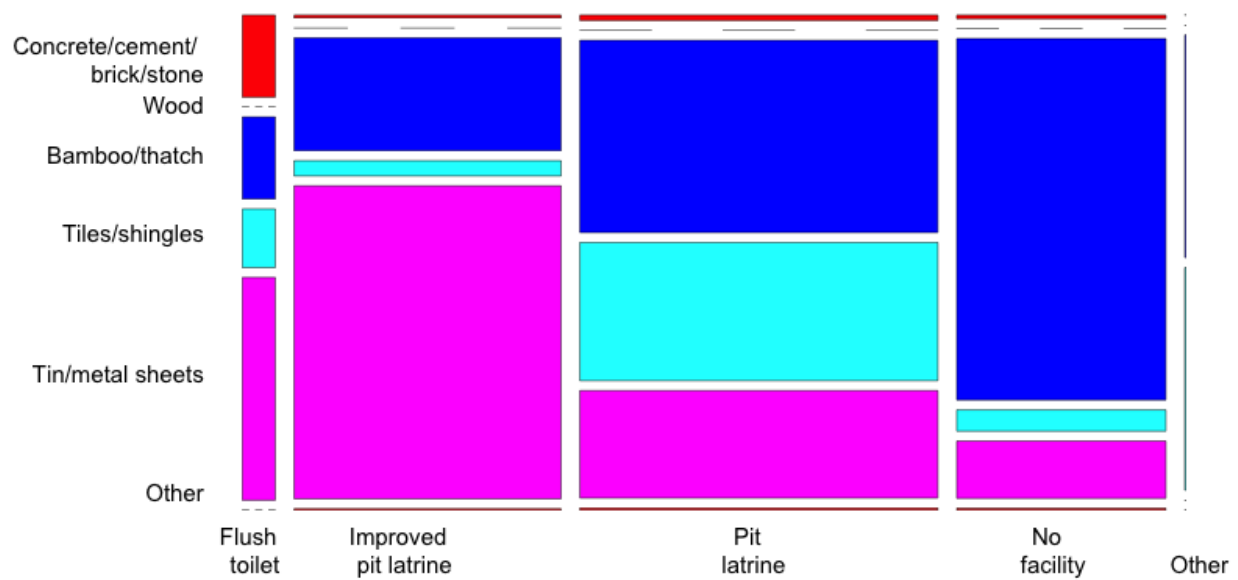


Fig. 7.9: Mosaic plot of the variables ROOF and TOILET after anonymization

Furthermore, some utility measures are not informative after certain anonymization methods have been applied. For example, after applying perturbative methods that interchange data values, comparing values directly is not useful because they will give the impression of high levels of information loss. In such cases, it is more informative to look at means, covariances and benchmarking indicators that can be computed from the data. Furthermore, it is important not only to focus on the characteristics of variables one by one, but also on the interactions between variables. This can be done by cross-tabulations and regressions. In general, when anonymizing sampled data, it is advisable to compute confidence intervals around estimates to interpret the magnitude of changes.

---

**Recommended Reading Material on Measuring Utility and Information Loss**

A.G. De Waal and L.C.R.J. Willenborg. 1999. “Information Loss through Global Recoding and Local Suppression” In Netherlands Official Statistics: Special Issue on SDC, 14, 17-10.

J. Domingo-Ferrer, J.M. Mateo-Sanz and V. Torra. 2001. “Comparing SDC Methods for Microdata on the basis of Information Loss and Disclosure Risk”. In Pre-proceedings of ETK-NTTS 2001 (vol. 2), 807-826. <http://neon.vb.cbs.nl/casc/NTTSJosep.pdf>

J. Domingo-Ferrer and V. Torra. 2001. “Disclosure Protection Methods and Information Loss for Microdata”. In P. Doyle, J.I. Lane, J.J.M. Theeuwes and L. Zayatz (eds.) *Theory and Practical Applications for Statistical Agencies*, 91-110, Amsterdam. <http://crises-deim.urv.cat/webCrises/publications/bcpi/cliatpasa01Disclosure.pdf>

---

**References**



---

### The SDC Process

---

This section presents the SDC process in a step-by-step representation and can be used as guidance for the actual SDC process. It should be noted, however, that jumping between steps and returning to previous steps is often required during the actual SDC process, as it is not necessarily a linear step-by-step process. This guidance brings together the different parts of the SDC process as discussed in the previous sections and links to these sections. The case studies in the next section follow these steps. This presentation is adapted from *HDFG12*. [Fig. 8.1](#) at the end of this section presents the entire process in a schematic way.

#### 8.1 Step 1: Need for confidentiality protection

Before starting the SDC process for a microdata set, the need for confidentiality protection has to be determined. This is closely linked to the interpretation of laws and regulations on this topic from the country in which the data originates and thus country-specific. A first step is to determine the statistical units in the dataset: if these are individuals, households or legal entities, such as companies, a need for disclosure control is likely. There are also examples of microdata for which there is no need for disclosure control. Examples could be data with climate and weather observations or data with houses as statistical units. Even if the primary statistical units are not natural or legal persons, however, the data can still contain confidential information on natural or legal persons. For example, a dataset with houses as primary statistical units can also contain information on the persons living in these houses and their income or a dataset on hospitalizations can include information about the hospitalized patients. In these cases, there is likely still a need for confidentiality protection. One option to solve this is to remove the information on the natural and legal persons in the datasets for release.

One dataset can also contain more than one type of statistical unit. The standard example here is a dataset containing both information on individuals and households. Another example is data with employees in enterprises. All types of statistical units present in the dataset have to be considered for the need of SDC. This is especially important in case the data has a hierarchical structure, such as individuals in households or employees in enterprises.

In addition, one has to evaluate whether the variables contained in the dataset are confidential or sensitive. Which variables are sensitive or confidential depends again on the applicable legislation and can differ substantially from country to country. In case the dataset includes sensitive or confidential variables, SDC is likely required. The set of sensitive variables and confidential variables together with the statistical units in the dataset determine the need for statistical disclosure control.

## 8.2 Step 2: Data preparation and exploring data characteristics

After assessing the need for statistical disclosure control, we should prepare the data and, if there are multiple, combine and consider all related data files. Then we explore the characteristics and structure in the data, which are important for the users of the data. Compiling an inventory of these characteristics is important for assessing the utility of the data after anonymization and producing an anonymized dataset, which is useful for end users.

The first step in data preparation is classifying the variables as sensitive or non-sensitive, and removing direct identifiers such as full names, passport numbers, addresses, phone numbers and GPS coordinates. In case of survey data, an inspection of the survey questionnaire is useful to classify the variables. Furthermore, it is necessary to select the variables that contain relevant information for end users and should be included in the dataset for release. At this point, it can also be useful to remove variables other than direct identifiers from the microdata set to be released. An example can be a variable with many missing values, e.g., a variable recorded only for a select group of individuals eligible for a particular survey module, and missing values for the rest. Such variables can cause a high level of disclosure risk while adding little information for end users. Examples are variables relating to education (current grade), where a missing value indicates that the individual is not currently in school, or variables relating to childbirth, where a missing value indicates that the individual has not delivered a child in the reference period. Missing values in themselves can be disclosive, especially if they indicate that the variable is not applicable. Often variables with the majority of values missing are deleted at this stage already. Other variables that might be deleted at this stage are those too sensitive to be anonymized and released or those not important to data users and that could increase the risk of disclosure.

Relationships may exist among variables in a dataset for a variety of reasons. For instance, variables can be mutually exclusive in cases where several binary variables are used for each category. An individual not in the labor force will have a missing value for the sector in which this person is employed (or more precisely not applicable). Relationships may also exist if some variables are ratios, sums or other mathematical functions of other variables. Examples are the variable household size (as a count of individuals per household), or aggregate expenditure (as a sum of all expenditure components). A certain value in one variable may also reduce the number of possible or valid values for another variable; for example, the age of an individual attending primary education or the gender of an individual having delivered a child. These relationships are important for two reasons: 1) they can be used by intruders to reconstruct anonymized values. For example, if age is suppressed but another variable indicates that they are in school, then it is still possible to infer a likely age range for that individual. Another example is if an individual is shown to be active in Sector B of the economy. Even if the labor status of this individual is suppressed, it can be inferred that this person is employed. 2) the relationships in the original data should be maintained in the anonymized dataset and inconsistencies should be avoided (e.g., SDC methods should not create 58-year-old school boys, or married 3-year-olds), or the dataset will be invalid for analysis. Another example is the case of expenditures per category, where it is important that the sum of the categories adds up to the total. One way to guarantee this is to anonymize the totals and then recalculate the sub-categories according to the original shares of the anonymized totals.

It is also useful at this stage to consolidate variables that provide the same information where possible, so as reduce the number of variables, reduce the likelihood of inconsistencies and minimize the variables an intruder can use to reconstruct the data. This is especially true if the microdata stems from an elaborate questionnaire and each variable represents one (sub-) question leading to a dataset with hundreds of variables. As an example, we take a survey with several labor force variables indicating whether a person is in the labor force, employed or unemployed, and if employed, in what sector. The data in [Table 8.1](#) illustrates this example. It is possible that each type of sector has its own binary variable. In that case, this set of variables can be summarized in two variables: one variable indicating whether a person is in labor force and another indicating the employment status, as well as the respective sector if a person is employed. These two variables contain all information contained in the previous five variables and make the anonymization process easier. If data users are used to a certain release format where including all five variables has been the norm, then it is possible to transform the variables back after the anonymization process rather than complicating the anonymization process by trying to treat more variables than is necessary. This approach also guarantees that the relationships between the variables are preserved (e.g., no individuals will be employed in several sectors).



Table 8.1: Illustration of merging variables without information loss for SDC process

Before					After	
In labor force	Employed	Sector A	Sector B	Sector C	In labor force	Employed
Yes	Yes	Missing	Yes	Missing	Yes	B
No	No	Missing	Missing	Missing	No	No
Yes	Yes	Yes	Missing	Missing	Yes	A
Yes	Yes	Missing	Yes	Missing	Yes	B
Yes	Yes	Missing	Missing	Yes	Yes	C
Yes	No	Missing	Missing	Missing	Yes	No

Besides relationships between variables, we also gather information about the survey methodology, such as strata, sampling methods, survey design and sample weights. This information is important in later stages, when estimating the disclosure risk and choosing the SDC methods. It is important to distinguish between a full census and a sample. For a full census, it is common practice to publish only a sample, as the risk of disclosure for a full sample is too high, given that we know that everyone in the country or institution is in the data (see also the Section **‘Special case: census data <anon\_methods.html#Special case: census data’**). Strata and sample weights can disclose information about the area or group to which an individual belongs (e.g., the weights can be linked with the geographical area or specific group in case of stratified sampling); this should be taken into account in the SDC process and checked before release of the dataset.

### 8.3 Step 3: Type of release

The type of release is an important factor for determining the required level of anonymization as well as the requirements end users have for the data (e.g., researchers need more detail than students for whom a teaching file might be sufficient) and should be clarified before the start of the anonymization process. Data release or dissemination by statistical agencies and data producers is often guided by the applicable law and dissemination strategies of the statistical agency, which specify the type of data that should be disseminated as well as the fashion.

Generally, there exist three types of data release methods for different target groups (the Section **Release types** provides more information on different release types):

- PUF: The data is directly available to anyone interested, e.g., on the website of the statistical agency
- SUF: The data is available to accredited researchers, who have to file their research proposals beforehand and have to sign a contract; this is also known as licensed file or microdata under contract
- Available in a controlled research data center: only on-site access to data on special computers; this is also known as data enclave

There are other data access possibilities besides these, such as teaching files or files for other specific purposes. Obviously, the required level of protection depends on the type of release: a PUF file must be protected to a greater extent than a SUF file, which in turn has to be protected more than a file which is available only in an on-site facility, since the options the intruder can use the data are limited in the latter case.

Besides the applicable legislation, the choice of the type of release depends on the type of the data and the content.

---

**Note:** Not every microdata set is suitable for release in any release type, even after SDC.

---

Some data cannot be protected sufficiently – it might always contain information that is too sensitive to be published as SUF or PUF. In such cases, the data can be released in on-site facilities, or the number of variables can be reduced by removing problematic variables.

Generally, the release of two or more anonymized datasets, e.g., tailored for different end users from the same original, is problematic because it can lead to disclosure if the two were later obtained and merged by the same user. The information contained in one dataset that is not contained in the other can lead to unintended disclosure. An exception is the simultaneous release and anonymization of a microdata set as PUF and SUF files. In this case, the PUF file is constructed from the SUF file by further anonymization. In that way, all information in the PUF file is also contained in the SUF file and the PUF file does not provide any additional information for users that have access to the SUF file.

---

**Note:** The anonymization process is an iterative process where steps can be revisited, whereas the publication of an anonymized dataset is a one-shot process.

---

Once the anonymized data is published, it is not possible to revoke and publish another dataset of the same microdata file. This would in fact mean publishing more than one anonymized file from the same microdata set, since some users might have saved the previous file.

### 8.4 Step 4: Intruder scenarios and choice of key variables

After determining the release type of the data, the possibilities of how an individual in the microdata could (realistically) be identified by an intruder under that release type should be examined. For PUF and SUF release the focus is on the use of external datasets from various source. These possibilities are described in disclosure or intruder scenarios, which specify what data an intruder could possibly have access to and how this auxiliary data can be used for identity disclosure. This leads to the specification of quasi-identifiers, which are a set of variables that are available both in the dataset to be released and in auxiliary datasets and need protection.

---

**Note:** If the number of quasi-identifiers is high, it is recommended to reduce the set of quasi-identifiers by removing some variables from the dataset for release.

---

This is especially true for PUF releases. Disclosure scenarios can also help define the required level of anonymization.

Drafting disclosure scenarios requires the support of subject matter specialists, assuming the subject specialist is not the same as the person doing the anonymization. Auxiliary datasets may contain information on the identity of the individuals and allow identity disclosure. Examples of such auxiliary data files are population registers and electoral rolls, as well as data collected by specialized firms.

---

**Note:** External datasets can come from many sources (e.g., other institutions, private companies) and it is sometimes difficult to make a full list of external data sources.

---

In addition, not all external data sources are in the public domain. Nevertheless, proprietary data can be used by the owner to re-identify individuals and should be taken into account in the SDC process, even if the exact content is not known. Also, the variables or datasets may not coincide perfectly (e.g., different definitions, more or less detailed variables, different survey period). Nevertheless, they should be considered in the SDC process.

Disclosure scenarios include both identity and inferential disclosure. The disclosure depends on the type of release, i.e., different data users have different data available and may use the data in a different way for re-identification. For example, a user in a research data center cannot match with large external datasets as (s)he is not permitted to take these into the data center. A user of a SUF is bound by an agreement specifying the use of the data and consequences if the agreement is breached (see the Section [Release types](#)). Furthermore, it should be evaluated whether, in case of a sample, possible intruders have knowledge as to which individuals are in the sample. This can be the case if it is known which schools were visited by the survey team, for example. A few examples of disclosure scenarios are (see the Section [Disclosure scenarios](#) for more information):

- **Matching:** The intruder uses auxiliary data, e.g., data on region, marital status and age from a population register, and matches them to released microdata. Individuals from the two datasets that match and are unique are successfully identified. This principle is used as an assumption in several disclosure risk measures, such as *k*-anonymity, individual and global risk, as described in the Section [Measuring Risk](#). This scenario can apply to both PUFs and SUFs.
- **Spontaneous recognition:** This scenario should be considered for SUF files, but is especially important for data available in research data centers where outliers are present in the data and data is often not strongly anonymized. The researcher might (unintentionally) recognize some individuals he knows (e.g., his colleagues, neighbors, family members, public figures, famous persons or large companies), while inspecting the data. This is especially true for rare combinations of values, such as outliers or unlikely combinations.

## 8.5 Step 5: Data key uses and selection of utility measures

In this step, we analyze the main uses of the data by the end users of the released microdata file. The data should be useful for the type of statistical analysis for which the data was collected and for which it is mostly used. The uses and requirements of data users will be different for different release types. Contacting data users directly or searching for scientific studies and papers that use similar data can be useful when collecting this information and making this assessment. Alternatively, this information can be collected from research proposals by researchers when applying for microdata access (SUF) or user groups can be set up. Furthermore, it is important to understand what level of precision the data users need and what types of categories are used. For instance, in the case of global recoding of age in years, one could recode age in groups of 10 years, e.g., 0 – 9, 10 – 19, 20 – 29, ... Many indicators relating to the labor market use categories that span the range 15 – 65, however. Therefore, constructing categories that coincide with the categories used for the indicators keeps the data much more useful while at the same time reducing the risk of disclosure in a similar way. This knowledge is important for the selection of useful utility measures, which in turn are used for selecting appropriate SDC methods.

The uses of the data depend on the release type, too. Researchers using SUF files require a higher level of detail in the data than PUF users.

---

**Note:** Anonymization will always lead to information loss and a PUF file will have reduced utility. If certain users require a high level of detail, release types other than PUF should be considered, such as SUF or release through a research data center.

---

In the case of SUFs, it is easier to find the main uses of the data since access is documented. One way to obtain information on the use of PUF files is to ask for a short description of intended use of the data before supplying the data. This is, however, useful only if microdata has been released previously.

Statistics computed from the anonymized and released microdata file should produce analytical results that agree or almost agree with previously published statistics from the original data. If, for instance, a previously published primary school enrollment rate was computed from these data and published, the released anonymized dataset should produce a very similar result to the officially published result. At the very least, the result should fall within the confidence region of the published result. It might be the case that not all published statistics can be generated from the published data. If this is the case, a choice should be made on which indicators and statistics to focus, and inform the users as to which ones have been selected and why.

As discussed in the Section [Measuring Utility and Information Loss](#), it is necessary to compute general utility measures that compare the raw and anonymized data, taking into consideration the end user's need for their analysis. In some cases the utility measures can give contradicting results, for example, a certain SDC method might lead to lower information loss for labor force figures but greater information loss for ratios relating to education. In such cases, the data uses might need to be ranked in order of importance and it should be clearly documented for the user that the prioritization of certain metrics over others means that certain metrics are no longer valid. This may be necessary, as

it is not possible to release multiple files for different users. This problem occurs especially in multi-purpose studies. For more details on utility measures, refer to the Section [Measuring Utility and Information Loss](#).

---

### Note on Steps 6 to 10

The following Steps 6 through 10 should be repeated if the data has quasi-identifiers that are on different hierarchical levels, e.g., individual and household. In that case, variables on the higher hierarchical level should be anonymized first, and then merged with the lower-level untreated variables. Subsequently, the merged dataset should be anonymized. This approach guarantees consistency in the treated data. If we neglect this procedure, the values of variables measured on the higher hierarchical level could be treated differently for observations of the same unit. For instance, the variable “region” is the same for all household members. If the value ‘rural’ would be suppressed for two members but not for the remaining three, this would lead to unintended disclosure; with the household ID the variable region would be easy to reconstruct for the two suppressed values. The Sections [Household risk](#) and [Household structure](#) provide more details on how to deal with data with household structure in *R* and *sdcMicro*.

---

## 8.6 Step 6: Assessing disclosure risk

The next step is to evaluate the disclosure risk of the raw data. Here it is important to distinguish between sample data and census data. In the case of census data, it is possible to directly calculate the risk measures when assuming that the dataset covers the entire population. If working with a sample, or a sample of the census (which is the more common case when releasing sample data), we can use the models discussed in the Section [Measuring Risk](#) to estimate the risk in the population. The main inputs for the risk measurement are the set of quasi-identifiers determined from the disclosure scenarios in Step 4 and the thresholds for risk calculations (e.g., the level of  $k$ -anonymity or the threshold for which an individual is considered at risk). If the data has a hierarchical structure (e.g., a household structure), the risk should be measured taking into account this structure as described the Section [Household risk](#).

The different risk measures described in the Section [Measuring Risk](#) each have advantages and disadvantages. Generally,  $k$ -anonymity, individual risk and global risk are used to produce an idea of the disclosure risk. These values can initially be very high but can often very easily be reduced after some simple but appropriate recoding (see [Step 8: Choice and application of SDC methods](#)). The thresholds shall be determined according to the release type. Always remember, though, that when using a sample, the risk measures based on the models presented in the literature offer a worst-case risk scenario and might therefore be an exaggeration of the real risks for some cases (see the Section [Individual risk](#)).

## 8.7 Step 7: Assessing utility measures

To quantify the information loss due to the anonymization, we first compute the utility measures selected in Step 5 using the raw data. This creates a base for comparison of results obtained when using the anonymized data – i.e., in Step 10.

---

**Note:** If the raw data is a sample, the utility measures are an estimate with a variance and therefore it is useful to construct confidence intervals in addition to the point estimates for the utility measures.

---

## 8.8 Step 8: Choice and application of SDC methods

The choice of SDC methods depends on the need for data protection (as measured by the disclosure risk), the structure of the data and the type of variables. The influence of different methods on the characteristics of the data important

for the users or the data utility should also be taken into account when selecting the SDC methods. In practice, the choice of SDC methods is partially a trial-and-error process: after applying a chosen method, disclosure risk and data utility are measured and compared to other choices of methods and parameters. The choice of methods is bound by legislation on the one hand, and a trade-off between utility and risk on the other.

The classification of methods as presented in [Table 6.1](#) gives a good overview for choosing the appropriate methods. Methods should be chosen according to the type of variable – continuous or categorical – the requirements by the users and the type of release. The anonymization of datasets with both continuous and categorical variables is discussed in the Section [Classification of variables](#).

In general for anonymization of categorical variables, it is useful to restrict the number of suppressions by first applying global recoding and/or removing variables from the microdata set. When the required number of suppressions to achieve the required level of risk is sufficiently low, the few individuals at risk can be treated by suppression. These are generally outliers. It should be noted that possibly not all variables can be released and some must be removed from the dataset (see [Step 2: Data preparation and exploring data characteristics](#)). Recoding and minimal use of suppression ensures that already published figures from the raw data can be reproduced sufficiently well from the anonymized data. If suppression is applied without sufficient recoding, the number of suppressions can be very high and the structure of the data can change significantly. This is because suppression mainly affects combinations that are rare in the data.

If the results of recoding and suppression do not achieve the required result, especially in cases where the number of select quasi-identifiers is high, an alternative is using perturbative methods. These can be used without prior recoding of variables. These methods, however, preserve data structure only partially. The preferred method depends on the requirements of the users. We refer to the Section [Anonymization Methods](#) and especially the Section [Perturbative methods](#) for a discussion of perturbative methods implemented in *sdcMicro*.

Finally, the choice of SDC methods depends on the data used since the same methods produce different results on different datasets. Therefore, the comparison of results with respect to risk and utility (Steps 9 and 10) is key to the choice made. Most methods are implemented in the *sdcMicro* package. Nevertheless, it is sometimes useful to use custom-made solutions. A few examples are presented in the Section [Anonymization Methods](#).

## 8.9 Step 9: Re-measure risk

In this step, we re-evaluate the disclosure risk with the risk measures chosen under Step 6 after applying SDC methods. Besides these risk measures, it is also important to look at individuals with high risk and/or special characteristics, combinations of values or outliers in the data. If the risk is not at an acceptable level, Steps 6 to 10 should be repeated with different methods and/or parameters.

---

**Note:** Risk measures based on frequency counts ( $k$ -anonymity, individual risk, global risk and household risk) cannot be used after applying perturbative methods since their risk estimates are not valid.

---

These methods are based on introducing uncertainty into the dataset and not on increasing the frequencies of keys in the data and will hence overestimate the risk.

## 8.10 Step 10: Re-measure utility

In this step, we re-measure the utility measures from Step 7 and compare these with the results from the raw data. Also, it is useful here to construct confidence intervals around the point estimates and compare these confidence intervals. The importance of the absolute value of a deviation can only be interpreted knowing the variance of the estimate. Besides these specific utility measures, general utility measures, as discussed in the Section [Measuring Utility and Information Loss](#), should be evaluated. This is especially important if perturbative methods have been applied. If

the data does not meet the user requirements and deviations are too large, repeat Steps 6 to 10 with different methods and/or different parameters.

---

**Note:** Anonymization will always lead to at least some information loss.

---

### 8.11 Step 11: Audit and Reporting

After anonymization, it is important to check whether all relationships in the data as identified in Step 2, such as variables that are sums of other variables or ratios, are preserved. Also, any unusual values caused by the anonymization should be detected. Examples of such anomalies are negative income or a pupil in the twentieth grade of school. This can happen after applying perturbative SDC methods. Furthermore, it is necessary to check whether previously published indicators from the raw data are reproducible from the data to be released. If this is not the case, data users might question the credibility of the anonymized dataset.

An important step in the SDC process is reporting, both internal and external. Internal reporting contains the exact description of anonymization methods used, parameters but also the risk measures before and after anonymization. This allows replication of the anonymized dataset and is important for supervisory authorities/bodies to ensure the anonymization process is sufficient to guarantee anonymity according to the applicable legislation.

External reporting informs the user that the data has been anonymized, provides information for valid analysis on the data and explains the limitations to the data as a result of the anonymization. A brief description of the methods used can be included. The release of anonymized microdata should be accompanied by the usual metadata of the survey (survey weight, strata, survey methodology) as well as information on the anonymization methods that allow researchers to do valid analysis (e.g., amount of noise added, transition matrix for PRAM).

---

**Note:** Care should be taken that this information cannot be used for re-identification (e.g., no release of random seed for PRAM).

---

The metadata must be updated to comply with the anonymized data. Variable descriptions or value labels might have changed as a result of the anonymization process. In addition, the information loss due to the anonymization process should be explained in detail to the users to make them aware of the limits to the validity of the data and their analyses.

### 8.12 Step 12: Data release

The last step in the SDC process is the actual release of the anonymized data. This step depends on the type of release chosen in Step 3. Changes to the variables made under Step 2, such as merging variables, can be undone to generate a dataset useful for users.

---

#### Recommended Reading Material on Risk Measurement

Dupriez, Olivier, and Ernie Boyko. 2010. *Dissemination of Microdata Files; Principles, Procedures and Practices*. IHSN Working Paper No. 005, International Household Survey Network (IHSN). <http://www.ihsn.org/HOME/sites/default/files/resources/IHSN-WP005.pdf>

---

#### References



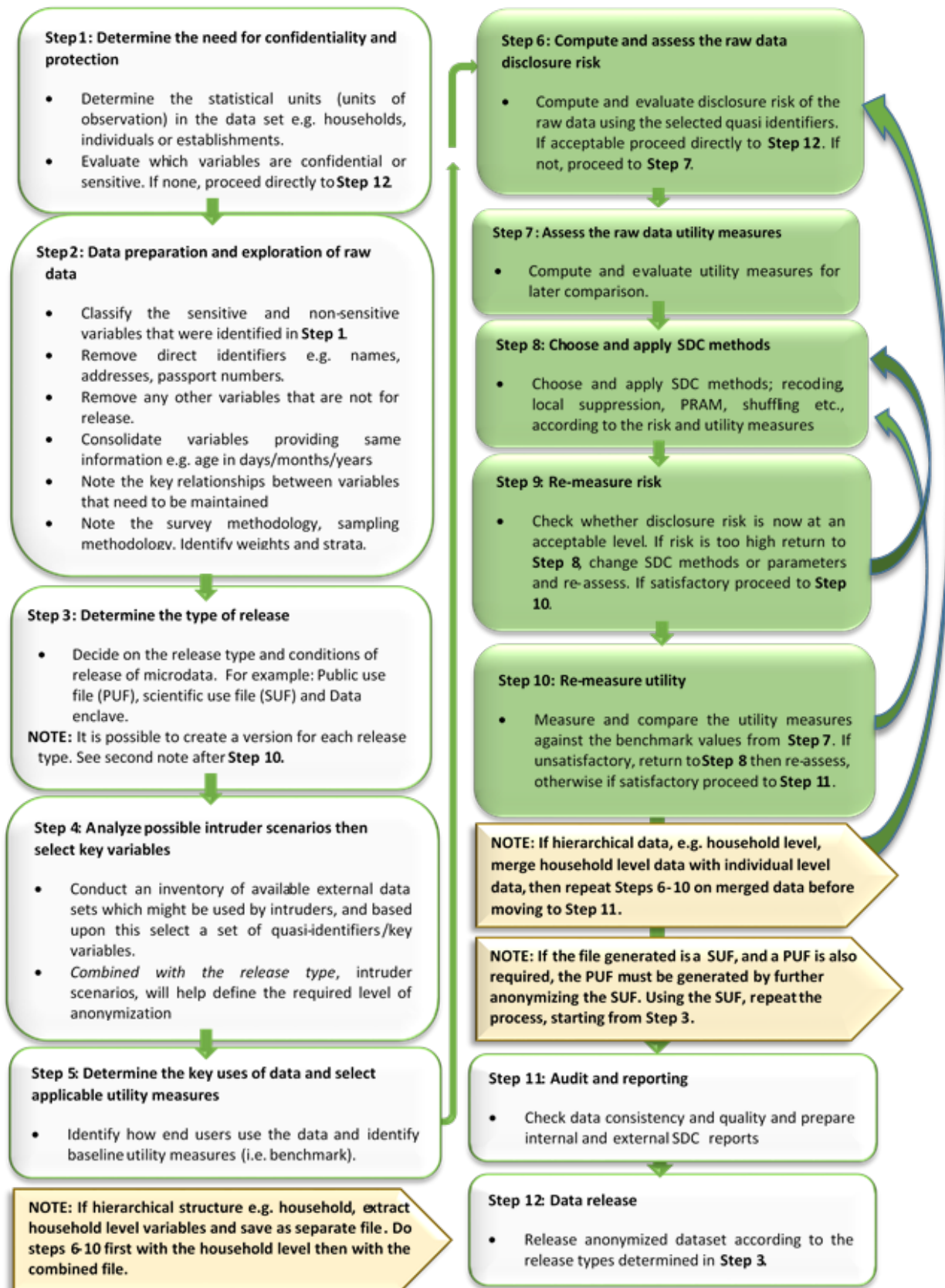


Fig. 8.1: Overview of the SDC process





---

## Case Studies (Illustrating the SDC Process)

---

In order to evaluate the use of different SDC methods on different types of survey datasets, we compared the results of the different methods applied to 75 datasets from 52 countries representing six geographic regions: Latin America and the Caribbean (LAC), Sub-Saharan Africa (AFR), South Asia (SA), Europe and Central Asia (ECA), Middle East and North Africa (MENA) as well as East Asia and the Pacific (EAP). The datasets chosen were from a mix of datasets that are already publically available, as well as data made available to the World Bank. The surveys used included, amongst others, household, demographic, and health surveys. The variables from these surveys used for the experiments were selected based on their relevance for users (e.g., for indicators, MDGs), their sensitivity, and their classification with respect to the SDC process.

The following case studies draw from knowledge gained from these experiments and try to incorporate the lessons learned. The case studies use synthetic data that mimic the structure of the survey types we used in our experiments and present the anonymization of a dataset similar to many surveys designed to measure household income and consumption, labor force participation and general demographic characteristics. The first case study creates a SUF, whereas in the second case study we take this SUF and treat it further to create a PUF.

### 9.1 Case study 1- SUF

This case study shows an example of how the anonymization process might be approached, particularly for a dataset with many continuous variables. We also show how this can be achieved using the open source and free *sdcmicro* package and *R*. A ready-to-run *R* script for this case study and the dataset are also available to reproduce the results and allow the user to adapt the code (see <http://ihsn.org/home/projects/sdc-practice>). Extracts of this code are presented in this section to illustrate several steps of the anonymization process.

---

**Note:** The choices of methods and parameters in this case study are based on this particular dataset and the results and choices might be different for other datasets.

---

The aim is to show the process, not to compare methods per se.

This example uses a dataset with a similar structure to that of a typical social survey with a focus on demographics, labor force participation and income and expenditure patterns. The dataset has been compiled using observations from

several datasets from different countries. They are considered synthetic data and as such are used only for illustrative purposes. The source datasets were already treated for disclosure control by their producers. This does not matter, as our concern is to illustrate the process only. The data from which we compiled our case study file was from surveys that contain many variables, but pay particular attention to labor force variables as well as household income and household expenditure variables. The variables in the demo dataset have already been pre-selected from the total set of variables available in the datasets. See [Appendix A](#) for the complete overview of all variables.

This case study follows the steps of the SDC process outlined in the Section [The SDC Process](#).

### 9.1.1 Step 1: Need for disclosure control

The statistical units in this dataset are individuals and households. The household structure provides a hierarchical structure in the data, which should be taken into account when measuring risk and selecting anonymization methods.

The data contains variables with demographic information, income, expenditures, education variables and variables relating to the labor status of the individual. These variables include sensitive and confidential variables. The dataset is an example of a social survey and, due to the nature of the statistical units and the variables, disclosure control is needed before release of the microdata. This is the case regardless of the legal framework, which is not specified here, as this is a hypothetical dataset.

### 9.1.2 Step 2: Data preparation and exploring data characteristics

The first step is to explore the data. To analyze the data in *R* we first have to read the data into *R*. In our case, the data is saved in a *STATA* file (.dta file). To read *STATA* files, we need to load the *R* package *foreign* (see the Section [Read functions in R](#) on importing other data formats in *R*). We also load the *sdcMicro* package and several other packages used later for the computation of the utility measures. If these packages are not yet installed, you should do so before trying to load them. The *R* code for this case study demonstrates how to do this.

Listing 9.1: Loading required packages

```
1 # Load required packages
2 library(foreign) # for read/write function for STATA files
3 library(sdcMicro) # sdcMicro package with functions for the SDC process
4 library(laeken) # for GINI
5 library(reldist) # for GINI
6 library(bootstrap) # for bootstrapping
7 library(ineq) # for Lorenz curves
```

After setting the working directory to the directory where the *STATA* file is stored, we load the data into the object called *file*. All output, unless otherwise specified, is saved in the working directory.

Listing 9.2: Loading the data

```
1 # Set working directory
2 setwd("C:/WorldBank/CaseStudy/")
3
4 # Specify file name
5 fname <- " case_1_data.dta"
6
7 # Read-in file
8 file <- read.dta(fname, convert.factors = F) # factors as numeric code
```

We check the number of variables, number of observations and variable names, as shown in [Listing 9.3](#).

Listing 9.3: Number of individuals and variables and variable names

```

1 dim(file) # Dimensions of file (observations, variables)
2 ## [1] 10574 68
3
4 colnames(file) # Variable names
5 ## [1] "REGION" "DIST" "URBRUR" "WGTHH"
6 ## [5] "WGTPOP" "IDH" "IDP" "HHSIZE"
7 ## [9] "GENDER" "REL" "MARITAL" "AGEYRS"
8 ## [13] "AGEMTH" "RELIG" "ETHNICITY" "LANGUAGE"
9 ## [17] "MORBID" "MEASLES" "MEDATT" "CHWEIGHTKG"
10 ## [21] "CHHEIGHTCM" "ATSCHOOL" "EDUCY" "EDYRS"
11 ## [25] "EDYRSCURRAT" "SCHTYP" "LITERACY" "EMPTYPI"
12 ## [29] "UNEMP1" "INDUSTRY1" "EMPCAT1" "WHOURSWEK1"
13 ## [33] "OWNHOUSE" "ROOF" "TOILET" "ELECTCON"
14 ## [37] "FUELCOOK" "WATER" "OWNAGLAND" "LANDSIZEHA"
15 ## [41] "OWNMOTORCYCLE" "CAR" "TV" "LIVESTOCK"
16 ## [45] "INCRMT" "INCWAGE" "INCBONSOCALL" "INCFARMBSN"
17 ## [49] "INCNFARMBSN" "INCRENT" "INCFIN" "INCPENSN"
18 ## [53] "INCOTHER" "INCTOTGROSSHH" "FARMEMP" "THOUSEXP"
19 ## [57] "TFOODEXP" "TALCHEXP" "TCLTHEXP" "TFURNEXP"
20 ## [61] "THLTHEXP" "TTRANSEXP" "TCOMMEXP" "TRECEXP"
21 ## [65] "TEDUEXP" "TRESTHOTEXP" "TMISCEXP" "TANHHEXP"

```

The dataset has 10,574 individuals in 2,000 households and contains 68 variables. The survey corresponds to a population of about 4.3 million individuals, which means that the sample is relatively small and the sample weights are high. This has an impact on the disclosure risk, as we will see in Steps 6a and 6b.

To get an overview of the values of the variables, we use tabulations and cross-tabulations for categorical variables and summary statistics for continuous variables. To include the number of missing values (NA or other), we use the option `useNA = "ifany"` in the `table()` function (see Listing 9.4).

In Table 9.1 the variables in the dataset are listed along with concise descriptions of the variables, the level at which they are collected (individual (IND), household (HH)), the measurement type (continuous, semi-continuous, categorical) and value ranges. Note that the dataset contains a selection of 68 variables (cf. Appendix A) of a total of 112 variables in the survey dataset. The variables have been preselected based on their relevance for data users. This allows to reduce the total numbers of variables to consider in the anonymization process and makes the process easier. The numerical values for many of the categorical variables are codes that refer to values, e.g., in the variable `URBRUR`, 1 stands for rural and 2 for urban. More information on the meanings of coded values of the categorical variables is available in the R code for this case study.

We identified the following sensitive variables in the data: `ETHNICITY`, `RELIGION`, variables related to the labor force status of the individual and the variables containing information on income and expenditures of the household. Whether variables can be identified as sensitive may vary across countries and datasets.

The case study dataset does not have any direct identifiers that, if they were present, would need to be removed at this stage. Examples of direct identifiers would be names, telephone numbers, geographical location coordinates, etc.

Listing 9.4: Tabulation of the variable 'gender' and summary statistics for the variable 'total annual expenditures' in R

```

1 # tabulation of variable GENDER (sex, categorical)
2 table(file$GENDER, useNA = "ifany")
3 ## 0 1
4 ## 5448 5126
5
6 # summary statistics for variable TANHHEXP (total annual household expenditures,

```

(continues on next page)

(continued from previous page)

```

7 # continuous)
8 summary(file$TANHHEXP)
9 ##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
10 ##      498   15550   17290   28560   29720   353200

```

Table 9.1: Overview of variables in dataset

No.	Variable name	Description	Level	Measurement	Value range
1	IDH	Household ID	HH	.	1-1000000
2	IDP	Individual ID	IND	.	1-1000000
3	REGION	Region	HH	categorical	1-10
4	DISTRICT	District	HH	categorical	10-100
5	URBRUR	Area of residence	HH	categorical	1, 2
6	WGTHH	Individual weighting coefficient	HH	weight	31-100
7	WGTPOP	Population weighting coefficient	IND	weight	45-100
8	HHSIZE	Household size	HH	semi-cont	1-10
9	GENDER	Gender	IND	categorical	0, 1
10	REL	Relationship to household head	IND	categorical	1-10
11	MARITAL	Marital status	IND	categorical	1-10
12	AGEYRS	Age in completed years	IND	semi-continuous	0-100
13	AGEMTH	Age of child in completed years	IND	semi-continuous	1-100
14	RELIG	Religion of household head	HH	categorical	1, 2
15	ETHNICITY	Ethnicity of household head	HH	categorical	all
16	LANGUAGE	Language of household head	HH	categorical	all
17	MORBID	Morbidity last x weeks	IND	categorical	0, 1
18	MEASLES	Child immunized against measles	IND	categorical	0, 1
19	MEDATT	Sought medical attention	IND	categorical	0, 1
20	CHWEIGHTKG	Weight of child (Kg)	IND	continuous	2-100
21	CHHEIGHTCM	Height of child (cms)	IND	continuous	70-150
22	ATSCHOOL	Currently enrolled in school	IND	categorical	0, 1
23	EDUCY	Highest level of education attended	IND	categorical	1-10
24	EDYEARS	Years of education	IND	semi-continuous	0-10
25	EDYRSCURRAT	Years of education for currently enrolled	IND	semi-continuous	1-10
26	SCHTYP	Type of school attending	IND	categorical	1-10
27	LITERACY	Literacy	IND	categorical	1-10
28	EMPTYPI	Type of employment	IND	categorical	1-10
29	UNEMP1	Unemployed	IND	categorical	0, 1
30	INDUSTRY1	Industry classification (1-digit)	IND	categorical	1-10
31	EMPCAT1	Employment categories	IND	categorical	11-100
32	WHOURSLASTWEEK1	Hours worked last week	IND	continuous	0-100
33	OWNHOUSE	Ownership of dwelling	HH	categorical	0, 1
34	ROOF	Main material used for roof	IND	categorical	1-10
35	TOILET	Main toilet facility	HH	categorical	1-10
36	ELECTCON	Electricity	HH	categorical	0-10
37	FUELCOOK	Main cooking fuel	HH	categorical	1-10
38	WATER	Main source of water	HH	categorical	1-10
39	OWNAGLAND	Ownership of agricultural land	HH	categorical	1-10
40	LANDSIZEHA	Land size owned by household (ha) (agric and non agric)	HH	continuous	0-100
41	OWNMOTORCYCLE	Ownership of motorcycle	HH	categorical	0, 1
42	CAR	Ownership of car	HH	categorical	0, 1
43	TV	Ownership of television	HH	categorical	0, 1

Table 9.1 – continued from previous page

No.	Variable name	Description	Level	Measurement	Value
44	LIFESTOCK	Number of large-sized livestock owned	HH	semi-continuous	0-
45	INCRMT	Income – Remittances	HH	continuous	
46	INCWAGE	Income - Wages and salaries	HH	continuous	
47	INCBONSOCAL	Income - Bonuses and social allowances derived from wage jobs	HH	continuous	
48	INCFARMBSN	Income - Gross income from household farm businesses	HH	continuous	
49	INCNFARMBSN	Income - Gross income from household nonfarm businesses	HH	continuous	
50	INCRENT	Income - Rent	HH	continuous	
51	INCFIN	Income - Financial	HH	continuous	
52	INCPENSN	Income - Pensions/social assistance	HH	continuous	
53	INCOTHER	Income - Other	HH	continuous	
54	INCTOTGROSSHH	Income - Total	HH	continuous	
55	FARMEMP				
56	TFOODEXP	Total expenditure on food	HH	continuous	
57	TALCHEXP	Total expenditure on alcoholic beverages, tobacco and narcotics	HH	continuous	
58	TCLTHEXP	Total expenditure on clothing	HH	continuous	
59	THOUSEXP	Total expenditure on housing	HH	continuous	
60	TFURNEXP	Total expenditure on furnishing	HH	continuous	
61	THLTHEXP	Total expenditure on health	HH	continuous	
62	TTRANSEXP	Total expenditure on transport	HH	continuous	
63	TCOMMEXP	Total expenditure on communication	HH	continuous	
64	TRECEXP	Total expenditure on recreation	HH	continuous	
65	TEDUEXP	Total expenditure on education	HH	continuous	
66	TRESHOTEXP	Total expenditure on restaurants and hotels	HH	continuous	
67	TMISCEXP	Total expenditure on miscellaneous spending	HH	continuous	
68	TANHHEXP	Total annual nominal household expenditures	HH	continuous	

It is always important to ensure that the relationships between variables in the data are preserved during the anonymization process and to explore and take note of these relationships before beginning the anonymization. In the final step in the anonymization process, an audit should be conducted, using these initial results, to check that these relationships are maintained in the anonymized dataset.

In our demo dataset, we identify several relationships between variables that need to be preserved during the anonymization process. The variables TANHHEXP and INCTOTGROSSHH represent the total annual nominal household expenditure and the total gross annual household income, respectively, and these variables are aggregations of existing income and expenditure components in the dataset.

The variables related to education are available only for individuals in the appropriate age groups and missing for other individuals. We make a similar observation for variables relating to children, such as height, weight and age in months. In addition, the household-level variables (cf. fourth column of [Table 9.1](#)) have the same values for all members in any particular household. The value of household size corresponds to the actual number of individuals belonging to that household in the dataset. As we proceed, we have to take care that these relationships and structures are preserved in the anonymization process.

When tabulating the variables, we notice that the variables RELIG, EMPTYP1 and LIVESTOCK have missing value codes different from the *R* standard missing value code NA. Before proceeding, we need to recode these to NA so *R* interprets them correctly. The missing value codes are resp. 99999, 99 and 9999 for these three variables. These are genuine missing value codes and not caused by the variables being not applicable to the individual. [Listing 9.5](#) shows how to make these changes.

**Note:** At the end of the anonymization process, and if desired for users, it is relatively easy to change these values back to their original missing value code.

Listing 9.5: Recoding missing value codes

```

1 # Set different NA codes to R missing value NA
2 file[, 'RELIG'][file[, 'RELIG'] == 99999] <- NA
3 file[, 'EMPTY1'][file[, 'EMPTY1'] == 99] <- NA
4 file[, 'LIVESTOCK'][file[, 'LIVESTOCK'] == 9999] <- NA

```

We also take note that the variables LANGUAGE and ETHNICITY have only missing values. Variables that contain only missing values should be removed from the dataset at this stage and excluded from the anonymization process. Removing these variables does not mean loss of data or reduction of the data utility, since these variables did not contain any information. It is, however, necessary to remove them, because keeping them can lead to errors in some of the anonymization methods in R. It is always possible to add these variables back into the dataset to be released at the end of the anonymization process. It is useful to reduce the dataset to those variables and records relevant for the anonymization process. This guarantees the best results in R and fewer errors. In Listing 9.6 we drop the variables that contain all missing values.

Listing 9.6: Dropping variables with only missing values

```

1 # Drop variables containing only missings
2 file <- file[, !names(file) %in% c('LANGUAGE', 'ETHNICITY')]

```

We assume that the data are collected in a survey that uses simple sampling of households. The data contains two weight coefficients: WGTHH and WGTPOP. The relationship between the weights is  $WGTPOP = WGTHH * HH\text{-}SIZE$ . WGTPOP is the sampling weight for the households and WGTHH is the sampling weight for the individuals to be used for disclosure risk calculations. WGTHH is used for computing individual-level indicators (such as education) and WGTPOP is used for population level indicators (such as income indicators). There are no strata variables available in the data. We will use WGTPOP for the anonymization of the household variables and WGTHH for the anonymization of the individual-level variables.

### 9.1.3 Step 3: Type of release

In this case study, we assume that data will be released as a SUF, which will be only available under license to accredited researchers with approved research proposals (see the Section [Conditions for SUFs](#) for more information of the release of a SUF). Therefore, the accepted risk level is higher and a broader set of variables can be released than would be the case when releasing a PUF. Since we do not have an overview of the requirements of all users, we restrict the utility measures to a selected number of data uses (see Step 5).

### 9.1.4 Step 4: Intruder scenarios and choice of key variables

Next, we analyze possible intruder scenarios and select quasi-identifiers or key variables based on these scenarios. Since the dataset used in this case study is a demo dataset that does not stem from an existing country (and hence we do not have information on external data sources available to possible intruders) and the original data has already been anonymized, it is not possible to define exact disclosure scenarios. Instead, we draft intruder scenarios for this demo dataset based on some hypothetical assumptions about availability of external data sources. We consider two types of disclosure scenarios: 1) matching to other publicly available datasets and 2) spontaneous recognition. The license under which the dataset will be distributed (SUF) prohibits matching to external resources. Still this can happen. However, the more important scenario is the one of spontaneous recognition. We describe both scenarios in the following two paragraphs.

For the sake of illustration, we assume that population registers are available with the demographic variables gender, age, place of residence (region, urban/rural), religion and other variables such as marital status and variables relating to education and professional status that are also present in our dataset. In addition, we assume that there is a publicly available cadastral register on land ownership. Based on this analysis of available data sources, we select the

variables REGION, URBRUR, HHSIZE, OWNAGLAND, RELIG, GENDER, REL (relationship to household head), MARITAL (marital status), AGEYRS, INDUSTRY1 and two variables relating to school attendance as categorical quasi-identifiers, the expenditure and income variables as well as LANDSIZEHA as continuous quasi-identifiers. According to our assessment, these variables might enable an intruder to re-identify an individual or household in the dataset by matching with other available datasets.

Table 9.2 gives an overview of the selected quasi-identifiers and their levels of measurement.

The decision to release the dataset as a SUF means the level of anonymization will be relatively low and consequently, the variables are more detailed and a scenario of spontaneous recognition is our main concern. Therefore, we should check for rare combinations or unusual patterns in the variables. Variables that may lead to spontaneous recognition in our sample are amongst others HHSIZE (household size), LANDSIZEHA as well as income and expenditure variables. Large households and large land ownership are easily identifiable. The same holds for extreme outliers in wealth and expenditure variables, especially when combined with other identifying variables such as region. There might be only one or a few households in a certain region with a high income, such as the local doctor. Variables that are easily observable and known by neighbors such as ROOF, TOILET, WATER, ELECTCON, FUELCOOK, OWNMOTORCYCLE, CAR, TV and LIVESTOCK may also need protection depending on what stands out in the community, since a researcher might be able to identify persons (s)he knows. This is called the nosy-neighbor scenario.

Table 9.2: List of selected quasi-identifiers

Name	Measurement
REGION (region)	Household, categorical
URBRUR (area of residence)	Household, categorical
HHSIZE (household size)	Household, categorical
OWNAGLAND (agricultural land ownership)	Household, categorical
RELIG (religion of household head)	Household, categorical
LANDSIZEHA (size of agr. and non-agr. land)	Household, continuous
TANHHEXP (total expenditures)	Household, continuous
TEXP (expenditures in category)	Household, continuous
INCTOTGROSSHH (total income)	Household, continuous
INC (income in category)	Household, continuous
GENDER (sex)	Individual, categorical
REL (relationship to household head)	Individual, categorical
MARITAL (marital status)	Individual, categorical
AGEYRS (age in completed years)	Individual, semi-continuous
EDYRSCURATT (years of education for currently enrolled)	Individual, semi-continuous
EDUCY (highest level of education completed)	Individual, categorical
ATSCHOOL (currently enrolled in school)	Individual, categorical
INDUSTRY1 (industry classification)	Individual, categorical

### 9.1.5 Step 5: Data key uses and selection of utility measures

In this case study, our aim is to create a SUF that provides sufficient information for accredited researchers. We know that the primary use of these data will be to evaluate indicators relating to income and inequality. Examples are the GINI coefficient and indicators on what share of income is spent on what type of expenditures. Furthermore, we focus on some education indicators. Table 9.3 gives an overview of the utility measures we selected. Besides these utility measures, which are specific to the data uses, we also do standard checks, such as comparing tabulations, cross-tabulations and summary statistics before and after anonymization.



Table 9.3: Overview of selected utility measures

Gini point estimates and confidence intervals for total expenditures	.
Lorenz curves for total expenditures	
Mean monthly per capita total expenditures by area of residence	
Average share of components for expenditures	
Mean monthly per capita total income by area of residence	
Average share of components for income	
Net enrollment in primary education by gender	

There are no published figures and statistics available that are calculated from this dataset because it is a demo. In general, the published figures should be re-computed based on the anonymized dataset and compared to the published figures in Step 11. Large differences would reduce the credibility of the anonymized dataset.

### 9.1.6 Hierarchical (household) structure

Our demo survey collects data on individuals in households. The household structure is important for data users and should be considered in the risk assessment. Since some variables are measured on the household level and thus have identical values for each household member, the values of the household variables should be treated in the same way for each household member (see the Section [Anonymization of the quasi-identifier household size](#)). Therefore, we first anonymize only the household variables. After this, we merge them with the individual-level variables and then anonymize the individual-level and household-level variables jointly.

Since the data has a hierarchical structure, Steps 6 through 10 are repeated twice: Steps 6a through 10a are for the household-level variables and Steps 6b through 10b for the combined dataset. In this way, we ensure that household-level variable values remain consistent across household members for each household and the household structure cannot be used to re-identify individuals. This is further explained in the Sections [Levels of risk](#) and [Randomizing order and numbering of individuals or households](#).

Before continuing to Step 6a, we select the categorical key variables, continuous key variables and any variables selected for use in PRAM routines, as well as household-level sampling weights. We extract these selected household variables and the households from the dataset and save them as *fileHH*. The choice of PRAM variables is further explained in Step 8a. [Listing 9.7](#) illustrates how these steps are done in *R* (see also the Section [Household structure](#)).

---

**Note:** In our dataset, some of the categorical variables when imported from the STATA file were not imported as factors. *sdcmicro* requires that these be converted to factors before proceeding.

---

Conversion of these variables to factors is also shown in [Listing 9.7](#).

Listing 9.7: Selecting the variables for the household-level anonymization

```

1 ### Select variables (household level)
2 # Key variables (household level)
3 selectedKeyVarsHH = c('URBRUR', 'REGION', 'HHSIZE', 'OWNHOUSE',
4                       'OWNAGLAND', 'RELIG')
5
6 # Changing variables to class factor
7 file$URBRUR <- as.factor(file$URBRUR)
8 file$REGION <- as.factor(file$REGION)
9 file$OWNHOUSE <- as.factor(file$OWNHOUSE)
10 file$OWNAGLAND <- as.factor(file$OWNAGLAND)
11 file$RELIG <- as.factor(file$RELIG)

```

(continues on next page)



(continued from previous page)

```

12
13 # Numerical variables
14 numVarsHH = c('LANDSIZEHA', 'TANHHEXP', 'TFOODEXP', 'TALCHEXP',
15              'TCLTHEXP', 'THOUSEXP', 'TFURNEXP', 'THLTHEXP',
16              'TTRANSEXP', 'TCOMMEXP', 'TRECEXP', 'TEDUEXP',
17              'TRESHOTEXP', 'TMISCEXP', 'INCTOTGROSSHH', 'INCRMT',
18              'INCWAGE', 'INCFARMBSN', 'INCNFARMBSN', 'INCRENT',
19              'INCFIN', 'INCPENSN', 'INCOTHER')
20 # PRAM variables
21 pramVarsHH = c('ROOF', 'TOILET', 'WATER', 'ELECTCON',
22              'FUELCOOK', 'OWNMOTORCYCLE', 'CAR', 'TV', 'LIVESTOCK')
23
24 # sample weight (WGTPOP) (household)
25 weightVarHH = c('WGTPOP')
26
27 # All household level variables
28 HHVars <- c('HID', selectedKeyVarsHH, pramVarsHH, numVarsHH, weightVarHH)

```

We then extract these variables from *file*, the dataframe in R that contains all variables. Every household has the same number of entries as it has members (e.g., a household of three will be repeated three times in *fileHH*). Before analyzing the household-level variables, we select only one entry per household, as illustrated in Listing 9.8. This is further explained in the Section [Household structure](#).

Listing 9.8: Taking a subset with only households

```

1 # Create subset of file with households and HH variables
2 fileHH <- file[,HHVars]
3
4 # Remove duplicated rows based on IDH, select uniques,
5 # one row per household in fileHH
6 fileHH <- fileHH[which(!duplicated(fileHH$IDH)),]
7
8 dim(fileHH)
9 ## [1] 2000 39

```

The file *fileHH* contains 2,000 households and 39 variables. We are now ready to create our *sdcMicro* object with the corresponding variables we selected in Listing 9.7. For our case study, we will create an *sdcMicro* object called *sdcHH* based on the data in *fileHH*, which we will use for steps 6a – 10a (see Listing 9.9).

---

**Note:** When the *sdcMicro* object is created, the *sdcMicro* package automatically calculates and stores the risk measures for the data.

---

This leads us to Step 6a.

Listing 9.9: Creating a *sdcMicro* object for the household variables

```

1 # Create initial SDC object for household level variables
2 sdcHH <- createSdcObj(dat = fileHH, keyVars = selectedKeyVarsHH, pramVars =
  ↪ pramVarsHH,
3                       weightVar = weightVarHH, numVars = numVarsHH)
4
5 numHH <- length(fileHH[,1]) # number of households

```

### 9.1.7 Step 6a: Assessing disclosure risk (household level)

As a first measure, we evaluate the number of households violating  $k$ -anonymity at the levels 2, 3 and 5.

Table 9.4 shows the number of violating households as well as the percentage of the total number of households. Listing 9.10 illustrates how to find these values with *sdcMicro*. The `print()` function in *sdcMicro* shows only the values for thresholds 2 and 3. Values for other thresholds can be calculated manually by summing up the frequencies smaller than the  $k$ -anonymity threshold, as shown in Listing 9.10.

Table 9.4: Number and proportion of households violating  $k$ -anonymity

k-anonymity level	Number of HH violating	Percentage of total number of HH
2	103	5.15 %
3	229	11.45 %
5	489	24.45 %

Listing 9.10: Showing number of households violating  $k$ -anonymity for levels 2, 3 and 5

```

1 # Number of observations violating k-anonymity (thresholds 2 and 3)
2 print(sdcHH)
3 ## Infos on 2/3-Anonymity:
4 ##
5 ## Number of observations violating
6 ## - 2-anonymity: 103
7 ## - 3-anonymity: 229
8 ##
9 ## Percentage of observations violating
10 ## - 2-anonymity: 5.150 %
11 ## - 3-anonymity: 11.450 %
12 -----
13
14 # Calculate sample frequencies and count number of obs. violating k(5) - anonymity
15 kAnon5 <- sum(sdcHH@risk$individual[,2] < 5)
16
17 kAnon5
18 ## [1] 489
19
20 # As percentage of total
21 kAnon5 / numHH
22 ## [1] 0.2445

```

It is often useful to view the values for the household(s) that violate  $k$ -anonymity. This might help clarify which variables cause the uniqueness of these households; this can then be used later when choosing appropriate SDC methods. Listing 9.11 shows how to assess the values of the households violating 3- and 5-anonymity. It seems that among the categorical key variables, the variable HHSIZE is responsible for many of the unique combinations and the

origin of much of the risk. Having determined this, we can flag HHSIZE as a possible first variable to treat to obtain the required risk level. In practice, with a variable like HHSIZE, this will likely involve removing large households from the dataset to be released. As explained in the Section [Anonymization of the quasi-identifier household size](#), recoding and local suppression are no valid options for the variable HHSIZE. The frequencies of household size in [Table 9.7](#) show that there are few households with more than 13 household members. This makes these households easily identifiable based on the number of household members and at high risk of re-identification, also in the context of the nosy neighbor scenario.

Listing 9.11: Showing households that violate  $k$ -anonymity

```
1 # Show values of key variable of records that violate k-anonymity
2 fileHH[sdcHH@risk$individual[,2] < 3, selectedKeyVarsHH] # for 3-anonymity
3 fileHH[sdcHH@risk$individual[,2] < 5, selectedKeyVarsHH] # for 5-anonymity
```

We also assess the disclosure risk of the categorical variables with the individual and global risk measures as described in the Sections [Individual risk](#) and [Global risk](#). In *fileHH* every entry represents a household. Therefore, we use the individual non-hierarchical risk here, where the individual refers in this case to a household. *fileHH* contains only households and has no hierarchical structure. In Step 6b, we evaluate the hierarchical risk in *file*, the dataset containing both households and individuals. The individual and global risk measures automatically take into consideration the household weights, which we defined in [Listing 9.7](#). In our file, the global risk measure calculated using the chosen key variables is 0.05%. This percentage is extremely low and corresponds to 1.03 expected re-identifications. The results are also shown in [Listing 9.12](#). This low figure can be explained by the relatively small sample size of 0.25% of the total population. Furthermore, one should keep in mind that this risk measure is based only on the categorical quasi-identifiers at the household level. [Listing 9.12](#) illustrates how to print the global risk measure.

Listing 9.12: Printing global risk measures

```
1 print(sdcHH, "risk")
2
3 ## Risk measures:
4 ##
5 ## Number of observations with higher risk than the main part of the data: 0
6 ## Expected number of re-identifications: 1.03 (0.05 %)
```

The global risk measure does not provide information about the spread of the individual risk measures. There might be a few households with relatively high risk, while the global (average) risk is low. It is therefore useful as a next step to inspect the observations with relatively high risk. The highest risk is 5.5% and only 14 households have risk larger than 1%. [Listing 9.13](#) shows how to display those households with risk over a certain threshold. Here the threshold is 0.01 (1%).

Listing 9.13: Observations with individual risk higher than 1%

```
1 # Observations with risk above certain threshold (0.01)
2 fileHH[sdcHH@risk$individual[, "risk"] > 0.01,]
```

Since the selected key variables at the household level are both categorical and numerical, the individual and global risk measures based on frequency counts do not completely reflect the disclosure risk of the entire dataset. Both categorical and continuous key variables are important for the data users, thus options like recoding the continuous variables (e.g., by creating quantiles of income and expenditure variables) to make all of them categorical will likely not satisfy the data user's needs. We therefore avoid recoding continuous variables and assess the disclosure risk of the categorical and continuous variables separately. This approach can be partly justified by the fact that any potential matching to external data sources for the continuous and categorical variables are available from different external data sources and as such will not be used simultaneously for matching.

### Continuous variables

To measure the risk of the continuous variables, we use an interval measure, which measures the number of

anonymized values that are too close to their original values. See the Section [Interval measure](#) for more information on interval-based risk measures for continuous variables. This measure is an ex-post measure, meaning that the risk can be evaluated only after anonymization and measures whether the perturbation is sufficiently large. Since it is an ex-post measure, we can evaluate it only in Step 9a after the variables have been treated. Evaluating this measure based on the original data would lead to a risk of 100%; all values would be too close to the original values since they would coincide with the original values, no matter how small the chosen intervals would be.

We also look at the distribution of LANDSIZEHA. In the variable LANDSIZEHA high values are rare and can lead to re-identification. An example is a large landowner in a specific region. To evaluate the distribution of the variable LANDSIZEHA, we look at the percentiles. Every percentile represents approximately 20 households. In addition, we look at the values of the largest 50 plots. [Listing 9.14](#) shows how to use *R* to display the quantiles and the largest landplots. [Table 9.5](#) shows the 90<sup>th</sup> – 100<sup>th</sup> percentiles and [Table 9.6](#) displays the largest 50 values for LANDSIZEHA. Based on these values, we conclude that values of LANDSIZEHA over 40 make the household very identifiable. These large households and households with large land plots need extra protection, as discussed in Step 8a.

Listing 9.14: Percentiles of LANDSIZE and listing the sizes of the largest 50 plots

```
1 # 1st - 100th percentiles of land size
2 quantile(fileHH$LANDSIZEHA, probs = (1:100)/100, na.rm= TRUE)
3
4 # Values of landsize for largest 50 plots
5 tail(sort(fileHH$LANDSIZEHA), n = 50)
```

Table 9.5: Percentiles 90-100 of the variable LANDSIZE

Percentile	90	91	92	93	94	95
Value	6.00	8.00	8.09	10.12	10.12	10.12
Percentile	96	97	98	99	100	
Value	12.14	20.23	33.83	121.41	1,214.08	

Table 9.6: 50 largest values of the variable LANDSIZE

12.14	15.00	15.37	15.78	16.19	20.00	20.23	20.23	20.23	20.23
20.23	20.23	20.23	20.23	20.23	20.23	20.23	20.23	20.23	20.23
20.23	20.23	20.50	30.35	32.38	40.47	40.47	40.47	40.47	40.47
40.47	40.47	80.93	80.93	80.93	80.93	121.41	121.41	161.88	161.88
161.88	182.11	246.86	263.05	283.29	404.69	404.69	607.04	809.39	1214.08

### 9.1.8 Step 7a: Assessing utility measures (household level)

The utility of the data does not only depend on the household level variables, but on the combination of household-level and individual-level variables. Therefore, it is not useful to evaluate all the utility measures selected in Step 5 at this stage, i.e., before anonymizing the individual level variables. We restrict the initial measurement of utility to those measures that are solely based on the household variables. In our dataset, these are the measures related to income and expenditure and their distributions. The results are presented in Step 10a, together with the results after anonymization, which allow direct comparison. If after the next anonymization step it appears that the data utility has been significantly decreased by the suppression of some household level variables, we can return to this step.

### 9.1.9 Step 8a: Choice and application of SDC methods (household variables)

This step is divided into the anonymization of the variable HHSIZE, as this is a special case, the anonymization of the other selected categorical quasi-identifiers and the anonymization of the selected continuous quasi-identifiers.

## Variable HHSIZE

The variable HHSIZE poses a problem for the anonymization of the file, since suppressing it will not anonymize this variable: a simple headcount based on the household ID would allow the reconstruction of this variable. Table 9.7 shows the absolute frequencies of HHSIZE. The number of households for each size larger than 13 is 6 or fewer and can be considered outliers with a higher risk of re-identification, as discussed in Step 6a. One way to deal with this is to remove all households of size 14 or larger from the dataset<sup>1</sup>. Removing 29 households of size 14 or larger reduces the number of 2-anonymity violations by 18, of 3-anonymity violations by 26 and of 5-anonymity violations by 29. This means that all removed households violated 5-anonymity due to the value of the variable HHSIZE and many of them 2- or 3-anonymity. In addition, the average individual risk amongst the 29 households is 0.15%, which is almost three times higher than the average individual risk of all households. The impact on the global risk measure of removing these 29 households is, however, limited, due to the relatively small number of removed households in comparison to the total number of 2,000 households. Removing the households is primarily to protect these specific households, not to reduce the global risk.

Changes, such as removing records, cannot be done in the *sdcMicro* object. Listing 9.15 illustrates the way to remove households and recreate the *sdcMicro* object.

Table 9.7: Frequencies of variable HHSIZE (household size)

HHSIZE	1	2	3	4	5	6	7	8	9	10	11	12
Frequency	152	194	238	295	276	252	214	134	84	66	34	21
HHSIZE	13	14	15	16	17	18	19	20	21	22	33	
Frequency	11	6	6	5	4	2	1	2	1	1	1	

Listing 9.15: Removing households with large (rare) household sizes

```

1 # Tabulation of variable HHSIZE
2 table(sdcHH@manipKeyVars$HHSIZE)
3
4 # Remove large households (14 or more household members) from file and fileHH
5 file <- file[!file[, 'HHSIZE'] >= 14,]
6
7 fileHHnew <- fileHH[!fileHH[, 'HHSIZE'] >= 14,]
8
9 # Create new sdcMicro object based on the file without the removed households
10 sdcHH <- createSdcObj(dat=fileHHnew, keyVars=selectedKeyVarsHH, pramVars=pramVarsHH,
11                      weightVar=weightVarHH, numVars = numVarsHH)

```

## Categorical variables

We are now ready to move on to the choice of SDC methods for the categorical variables on the household level in our dataset. As noted in our discussion of the methods, applying perturbative methods and local suppression may lead to large loss of utility. The common approach is to apply recoding to the largest possible extent as a first approach, to reach a prescribed level of risk and reduce the number of suppressions needed. Only after that should methods such as local suppression be applied. If this approach does not already achieve the desired result, we can consider perturbative methods.

Since the file is to be released as a SUF, we can keep a higher level of detail in the data. The selected categorical key variables at the household level are not suitable for recoding at this point. Due to the relatively low risk of re-identification based on the five selected categorical household level variables, it is possible in this case to use an option like local suppression to achieve our desired level of risk. Applying local suppression when initial risk is relatively low will likely only lead to suppression of few observations and thus limit the loss of utility. If, however, the data had been measured to have a relatively high risk, then applying local suppression without previous recoding would likely result in a large number of suppressions and greater information loss. Efforts such a recoding should be taken first before

<sup>1</sup> Other methods and guidance on treating datasets where household size is a quasi-identifier are discussed in the Section [Anonymization of the quasi-identifier household size](#).

suppressing in cases where risk is initially measured as high. Recoding will reduce risk with little information loss and thus the number of suppressions, if local suppression is applied as an additional step. We apply local suppression to reach 2-anonymity. The choice of the low level of two is based on the overall low re-identification risk due to the high sample weights and the release as SUF. High sample weights mean, ceteris paribus, a low level of re-identification risk. Achieving 2-anonymity is the same as removing sample uniques. This leads to 42 suppressions in the variable HHSIZE and 4 suppressions in the variable REGION. As explained earlier, suppression of the value of the variable HHSIZE does not lead to actual suppression of this information. Therefore, we redo the local suppression, but this time we tell *sdcMicro* to, if possible, not suppress HHSIZE but one of the other variables.

In *sdcMicro* it is possible to tell the algorithm which variables are important and less important for making small changes (see also the Section [Local suppression](#)). To prevent HHSIZE being suppressed, we set the importance of HHSIZE in the importance vectors to the highest (i.e., 1). [Listing 9.16](#) shows how to apply local suppression and put importance on the variable HHSIZE. The variable REGION is the type of variable that should not have any suppressions either. We also set the importance of REGION to 2 and the importance of RURURB to 3. This leads to an order of the variables to be considered for suppression by the algorithm. Instead of 42 suppressions in the variable HHSIZE, this leads one suppressed value in the variable HHSIZE, and to 6, 1, 48 and 16 suppressions respectively for the variables URBRUR, REGION, OWNAGLAND and RELIG (which we set as less important). The importance is clearly reflected in the number of suppression. The total number of suppressions is higher than without importance vector (71 vs. 46), but 2-anonymity is achieved in the dataset with fewer suppressions in the variables HHSIZE and REGION. We remove the one household with the suppressed value of HHSIZE (13) to protect this household.

**Note:** In [Listing 9.16](#) we use the `undolast()` function in *sdcMicro* to go one step back after we had first applied local suppression with no importance vector.

The `undolast()` function restores the *sdcMicro* object back to the previous state (i.e., before we applied local suppression), which allows us to rerun the same command, but this time with an importance vector set. The `undolast()` function can only be used to go one step back.

Listing 9.16: Local suppression with and without importance vector

```

1 # Local suppression
2 sdcHH <- localSuppression(sdcHH, k=2, importance = NULL) # no importance vector
3
4 print(sdcHH, "ls")
5 ## Local Suppression:
6 ##   KeyVar | Suppressions (#) | Suppressions (%)
7 ##   URBRUR |           0 |           0.000
8 ##   REGION |           4 |           0.203
9 ##   HHSIZE |          37 |           1.877
10 ##  OWNAGLAND |           0 |           0.000
11 ##    RELIG |           0 |           0.000
12
13 sdcHH <- undolast(sdcHH)
14
15 sdcHH <- localSuppression(sdcHH, k=2, importance = c(3, 2, 1, 5, 5))
16 # importance on HHSIZE (1), REGION (2) and URBRUR (3)
17
18 print(sdcHH, "ls")
19 ## Local Suppression:
20 ##   KeyVar | Suppressions (#) | Suppressions (%)
21 ##   URBRUR |           6 |           0.304
22 ##   REGION |           1 |           0.051
23 ##   HHSIZE |           1 |           0.051
24 ##  OWNAGLAND |          43 |           2.182
25 ##    RELIG |          16 |           0.812

```

The variables ROOF, TOILET, WATER, ELECTCON, FUELCOOK, OWNMOTORCYCLE, CAR, TV and LIVESTOCK are not sensitive variables and were not selected as quasi-identifiers because we assumed that there are no external data sources containing this information that could be used for matching. Values can be easily observed or be known to neighbors, however, and therefore are important, together with other variables, for the spontaneous recognition scenario and nosy neighbor scenario. To anonymize these variables, we want to introduce a low level of uncertainty in them. Therefore, we decide to use invariant PRAM for the variables ROOF, TOILET, WATER, ELECTCON, FUELCOOK, OWNMOTORCYCLE, CAR, TV and LIVESTOCK, where we treat LIVESTOCK as a semi-continuous variable due to the low number of different values. The Section [PRAM \(Post Randomization Method\)](#) provides more information on the PRAM method and its implementation in *sdcMicro*. [Listing 9.17](#) illustrates how to apply PRAM. We choose the parameter *pd*, the lower bound for the probability that a value is not changed, to be relatively high at 0.8. We can choose a high value, because the variables themselves are not sensitive and we only want to introduce a low level of changes to minimize the utility loss. Because the distribution of many of the variables chosen for PRAM depends on the REGION, we decide to use the variable REGION as a strata variable. In this way the transition matrix is computed for each region separately. Because PRAM is a probabilistic method, we set a seed for the random number generator before applying PRAM to ensure reproducibility of the results.

**Note:** In practice, it is not advisable to set a seed of 12345, but rather a longer more complex and less easy to guess sequence.

The seed should not be released, since it allows for reconstructing the original values if combined with the transition matrix. The transition matrix can be released: this allows for consistent statistical inference by correcting the statistical methods used if the researcher has knowledge about the PRAM method (at this point *sdcMicro* does not allow the retrieval of the transition matrix).

Listing 9.17: Applying PRAM

```

1 # Pram
2 set.seed(12345)
3 sdcHH <- pram(sdcHH, strata_variables = "REGION", pd = 0.8)
4
5 ## Number of changed observations:
6 ## - - - - -
7 ## ROOF != ROOF_pram : 98 (4.97%)
8 ## TOILET != TOILET_pram : 151 (7.66%)
9 ## WATER != WATER_pram : 167 (8.47%)
10 ## ELECTCON != ELECTCON_pram : 90 (4.57%)
11 ## FUELCOOK != FUELCOOK_pram : 113 (5.73%)
12 ## OWNMOTORCYCLE != OWNMOTORCYCLE_pram : 41 (2.08%)
13 ## CAR != CAR_pram : 172 (8.73%)
14 ## TV != TV_pram : 137 (6.95%)
15 ## LIVESTOCK != LIVESTOCK_pram : 149 (7.56%)

```

PRAM has changed values within the variables according to the invariant transition matrices. Since we used the invariant PRAM method (see the Section [PRAM \(Post Randomization Method\)](#)), the absolute univariate frequencies remain unchanged. This is not the case for the multivariate frequencies. In Step 10a we compare the changes in the multivariate frequencies for the PRAMmed variables.

### Continuous variables

We have selected income and expenditures variables and the variable LANDSIZEHA as numerical quasi-identifiers, as discussed in Step 4. In Step 5 we identified variables having high interest for the users of our data: many users use the data for measuring inequality and expenditure patterns.

Based on the risk evaluation in Step 6a, we decide to anonymize the variable LANDSIZEHA by top coding at the value 40 (cf. [Table 9.5](#) and [Table 9.6](#)) and round values smaller than 1 to one digit, and values larger than 1 to zero digits. Rounding the values prevents exact matching with the available cadastral register. Furthermore, we group the



values between 5 and 40 in the groups 5 – 19 and 20 – 39. After these steps, no household has a unique plot size and the number of households in the sample with the same plot size was increased to at least 7. This is shown by the tabulation of the variable LANDSIZEHA after manipulation in the last line of Listing 9.18. In addition, all outliers have been removed by top coding the values. This has reduced the risk of spontaneous recognition as discussed in Step 6. How to recode values in R is introduced in the Section [Recoding](#) and, for this particular case, shown in Listing 9.18.

Listing 9.18: Anonymizing the variable LANDSIZEHA

```

1  # Rounding values of LANDSIZEHA to 1 digit for plots smaller than 1 and
2  # to 0 digits for plots larger than 1
3  sdchH@manipNumVars$LANDSIZEHA[sdchH@manipNumVars$LANDSIZEHA <= 1 &
4      !is.na(sdchH@manipNumVars$LANDSIZEHA)] <-
5      round(sdchH@manipNumVars$LANDSIZEHA[sdchH@manipNumVars$LANDSIZEHA <= 1 &
6          !is.na(sdchH@manipNumVars
7  ↪ $LANDSIZEHA)],
8          digits = 1)
9
10 sdchH@manipNumVars$LANDSIZEHA[sdchH@manipNumVars$LANDSIZEHA > 1 &
11     !is.na(sdchH@manipNumVars$LANDSIZEHA)] <-
12     round(sdchH@manipNumVars$LANDSIZEHA[sdchH@manipNumVars$LANDSIZEHA > 1 &
13         !is.na(sdchH@manipNumVars
14 ↪ $LANDSIZEHA)],
15         digits = 0)
16
17 # Grouping values of LANDSIZEHA into intervals 5-19, 20-39
18 sdchH@manipNumVars$LANDSIZEHA[sdchH@manipNumVars$LANDSIZEHA >= 5 &
19 sdchH@manipNumVars$LANDSIZEHA < 20 & !is.na(sdchH@manipNumVars$LANDSIZEHA)] <- 13
20
21 sdchH@manipNumVars$LANDSIZEHA[sdchH@manipNumVars$LANDSIZEHA >= 20 &
22 sdchH@manipNumVars$LANDSIZEHA < 40 & !is.na(sdchH@manipNumVars$LANDSIZEHA)] <- 30
23
24 # Topcoding values of LANDSIZEHA larger than 40 (also recomputes risk after manual
25 ↪ changes)
26 sdchH <- topBotCoding(sdchH, value = 40, replacement = 40, kind = 'top', column =
27 ↪ 'LANDSIZEHA')
28
29 # Results for LANDSIZEHA
30 table(sdchH@manipNumVars$LANDSIZEHA)
31 ##    0 0.1 0.2 0.3 0.4 0.5 0.6 0.7 0.8 0.9    1    2    3    4   13   30   40
32 ## 188 109  55  30  24  65  22   7  31  16 154 258  53  60 113  18  25

```

For the expenditure and income variables we compared, **based on the actual case study data**, several methods. As mentioned earlier, the main use of the data is to compute inequality measures, such as the Gini coefficient. Recoding these variables into percentiles creates difficulties computing these measures or changes these measures to a large extent and is hence not a suitable method. Often, income and expenditure variables that are released in a SUF are anonymized by top-coding. This protects the outliers, which are the values that are the most at risk. Top-coding, however, destroys the inequality information in the data, by removing high (and low) incomes. Therefore, we decide to use noise addition. To take into account the higher risk of outliers, we add a higher level of noise to those.

Adding noise can lead to a transformation of the shape of the distribution. Depending on the magnitude of the noise (see the Section [Noise addition](#) for the definition of the magnitude of noise), the values of income can also become negative. One way to solve this would be to cut off the values below zero and set them to zero. This would, however, destroy the properties conserved by noise addition (amongst others the value of the expected mean, see also the Section [Noise addition](#)) and we chose to keep the negative values.

As mentioned before, the aggregates of income and expenditures are the sums of the components. Adding noise to each of the components might lead to violation of this condition. Therefore, one solution is to add noise to the



aggregates and remove the components. We prefer to keep the components in the data and apply noise addition to each component separately. This allows to apply a lower level of noise than when applying noise only to the aggregates. A noise level of 0.01 seems to be sufficient with extra noise of 0.05 added to the outliers. The outliers are defined by a robust Mahalanobis distance (see the Section [Noise addition](#)). After adding noise to the components, we recomputed the aggregates as the sum of the perturbed components.

**Note:** This result is only based on the actual case study dataset and is not necessarily true for other datasets.

The noise addition is shown in [Listing 9.19](#). Before applying probabilistic methods such as noise addition, we set a seed for the random number generator. This allows us to reproduce the results.

Listing 9.19: Anonymizing continuous variables

```

1  # Add noise to income and expenditure variables by category
2
3  # Anonymize components
4  compExp <- c("TFOODEXP", "TALCHEXP", "TCLTHEXP", "THOUSEXP",
5              "TFURNEXP", "THLTHEXP", "TTRANSEXP", "TCOMMEXP", "TRECEXP", "TEDUEXP",
6              "TRESHOTEXP", "TMISCEXP")
7  set.seed(123)
8
9  # Add noise to expenditure variables
10 sdchH <- addNoise(noise = 0.01, obj = sdchH, variables = compExp, method = "additive")
11
12 # Add noise to outliers
13 sdchH <- addNoise(noise = 0.05, obj = sdchH, variables = compExp, method = "outdetect")
14
15 # Sum over expenditure categories to obtain consistent totals
16 sdchH@manipNumVars[, 'TANHHEXP'] <- rowSums(sdchH@manipNumVars[, compExp])
17 compInc <- c('INCRMT', 'INCWAGE', 'INCFARMBN', 'INCNFARMBN',
18             'INCRENT', 'INCFIN', 'INCPENSN', 'INCOTHER')
19
20 # Add noise to income variables
21 sdchH <- addNoise(noise = 0.01, obj = sdchH, variables = compInc, method = "additive")
22
23 # Add noise to outliers
24 sdchH <- addNoise(noise = 0.05, obj = sdchH, variables = compInc, method = "outdetect")
25
26 # Sum over income categories to obtain consistent totals
27 sdchH@manipNumVars[, 'INCTOTGROSSH'] <- rowSums(sdchH@manipNumVars[, compInc])
28
29 # recalculate risks after manually changing values in sdcmicro object
30 calcRisks(sdchH)

```

### 9.1.10 Step 9a: Re-measure risk

For the categorical variables, we conclude that we have achieved 2-anonymity in the data with local suppression. Only 104 households, or about 5% of the total number, violate 3-anonymity. [Table 9.8](#) gives an overview of these risk measures. The global risk is reduced to 0.02% (expected number of re-identifications 0.36), which is extremely low. Therefore, we conclude that based on the categorical variables, the data has been sufficiently anonymized. No household has a risk of re-identification higher than 0.01 (1%). By removing households with rare values (or outliers) of the variable HHSIZE, we have reduced the risk of spontaneous recognition of these households. This reasoning can also be applied to the result of the risk of recoding the variable LANDSIZEHA and PRAMming the variables identified to be important in the nosy neighbor scenario. An intruder cannot know with certainty whether a household

that he recognizes in the data is the correct household, due to the noise.

Table 9.8: Number and proportion of households violating k-anonymity after anonymization

k-anonymity	Number HH violating	Percentage
2	0	0 %
3	104	5.28 %
5	374	18.70 %

These measures refer only to the categorical variables. To evaluate the risk of the continuous variables we could use an interval measure or closest neighbor algorithm. These risk measures are discussed in the Section [Risk measures for continuous variables](#). We chose to use an interval measure, since exact value matching is not our largest concern based on the assumed scenarios and external data sources. Instead, datasets with similar values but not the exact same values could be used for matching. Here the main concern is that the values are sufficiently far from the original values, which is measured with an interval measure.

[Listing 9.20](#) shows how to evaluate the interval measure for each of the expenditure variables, which are contained in the vector *compExp*<sup>2</sup>. The different values of the parameter *k* in the function *dRisk()* define the size of the interval around the original value, as explained in the Section [Interval measure](#). The larger *k*, the larger the intervals, the higher the probability that a perturbed value is in the interval around the original value and the higher the risk measure. The result is satisfactory with relatively small intervals (*k* = 0.01), but not when increasing the size of the intervals. In our case, *k* = 0.01 is sufficiently large, since we are looking at the components, not the aggregates. We have to pay special attention to the outliers. Here the value 0.01 for *k* is too small to assume that they are protected when outside this small interval. It would be necessary to check outliers and their perturbed values and there might be a need for a higher level of perturbation for outliers. We conclude that, from a risk perspective and based on the interval measure, the chosen levels of noise are acceptable. In the next step, we will look at the impact on the data utility of the noise addition.

Listing 9.20: Measuring risk of re-identification of continuous variables

```

1 dRisk(sdcHH@origData[,compExp], xm = sdcHH@manipNumVars[,compExp], k = 0.01)
2 [1] 0.0608828
3
4 dRisk(sdcHH@origData[,compExp], xm = sdcHH@manipNumVars[,compExp], k = 0.02)
5 [1] 0.9025875
6
7 dRisk(sdcHH@origData[,compExp], xm = sdcHH@manipNumVars[,compExp], k = 0.05)
8 [1] 1

```

## 9.1.11 Step 10a: Re-measure utility

None of the variables has been recoded and the original level of detail in the data is kept, except for the variable LANDSIZEHA. As described in Step 8a, local suppression has only removed a few values in the other variables, which has not greatly reduced the validity of the data.

The univariate frequency distributions of the variables ROOF, TOILET, WATER, ELECTCON, FUELCOOK, OWN-MOTORCYCLE, CAR, TV and LIVESTOCK did not, by definition of the invariant PRAM method (see the Section [PRAM \(Post Randomization Method\)](#)), change to a large extent. The tabulations are presented in [Table 9.9](#) (the values 1 – 9 and NA in the first row are the values of the variables and .m after the variable name refers to the values after anonymization).

**Note:** Although the frequencies are almost the same, this does not mean that the values of particular households did

<sup>2</sup> For illustrative purposes, we only show this evaluation for the expenditure variables. It can be easily copied for the income variables. The results are similar.

not change.

Values have been swapped between households. This becomes apparent when looking at the multivariate frequencies of the WATER with the variable URBRUR in Table 9.10. The multivariate frequencies of the PRAMmed with the variable URBRUR could be of interest for users, but these are not preserved. Since we applied PRAM within the regions, the multivariate frequencies of the PRAMmed variables with REGION are preserved.

Table 9.9: Univariate frequencies of the PRAMmed variable before and after anonymization

.	0	1	2	3	4	5	6	7	8	9	NA
ROOF		27	1	914	307	711				10	1
ROOF.m		25	1	907	319	712				6	1
TOILET		76	594	817	481					3	
TOILET.m		71	597	816	483					4	
WATER		128	323	304	383	562	197	18	21	35	
WATER.m		134	319	308	378	573	188	16	21	34	
ELECTCON	768	216	8	2							977
ELECTCON.m	761	218	8	3							981
FUELCOOK		1289	21	376	55	36				139	55
FUELCOOK.m		1284	22	383	50	39				143	50
OWNMOTORCYCLE	1883	86									2
OWNMOTORCYCLE.m	1882	86									2
CAR	963	31									977
CAR.m	966	25									
TV	1216	264									491
TV.m	1203	272									496

Table 9.10: Multivariate frequencies of the variables WATER with RU-RURB before and after anonymization

.	1	2	3	4	5	6	7	8	9
WATER/URB	11	49	270	306	432	183	12	15	21
WATER/RUR	114	274	32	76	130	14	6	6	14
WATER/URB.m	79	220	203	229	402	125	10	12	19
WATER/RUR.m	54	98	105	147	169	63	6	9	15

For conciseness, we restrict ourselves to the analysis of the expenditure variables. The analysis of the income variables can be done in the same way and leads to similar results.

We look at the effect of anonymization on some indicators as discussed in Step 5. Table 9.11 presents the point estimates and bootstrapped confidence interval of the GINI coefficient<sup>3</sup> for the sum of the expenditure components. The calculation of the GINI coefficient and the confidence interval are based on the positive expenditure values. We observe very small changes in the Gini coefficient, that are statistically negligible. We use a visualization to illustrate the impact on utility of the anonymization. Visualizations are discussed in the Section [Assessing data utility with the help of data visualizations \(in R\)](#) and the specific R code for this case study is available in the R script. The change in the inequality measures is illustrated in Fig. 9.1, which shows the Lorenz curves based on the positive expenditure values before and after anonymization.

<sup>3</sup> To compute the GINI coefficient, bootstrap to construct the confidence intervals and plot the Lorenz curve we used the R packages *laeken*, *reldist*, *bootstrap* and *ineq*.

Table 9.11: GINI point estimates and bootstrapped confidence intervals for sum of expenditure components

.	before	after
Point estimate	0.510	0.508
Left bound of CI	0.476	0.476
Right bound of CI	0.539	0.538

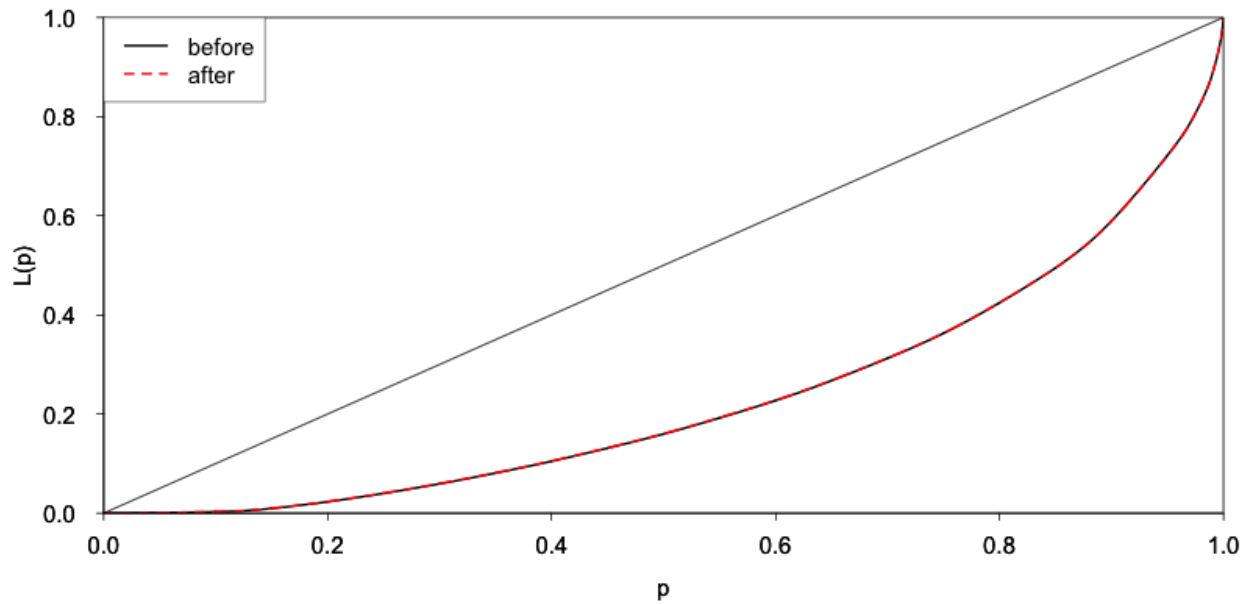


Fig. 9.1: Lorenz curve based on positive total expenditures values

We compare the mean monthly expenditures (MME) and mean monthly income (MMI) for rural, urban and total population. The results are shown in Table 9.12. We observe that the chosen levels of noise add only small distortions to the MME and slightly larger changes to the MMI.

Table 9.12: Mean monthly expenditure and mean monthly income per capita by rural/urban

.	before	after
MME rural	400.5	398.5
MME urban	457.3	459.9
MME total	412.6	412.6
MMI rural	397.1	402.2
MMI urban	747.6	767.8
MMI total	472.1	478.5

Table 9.13 shows the share of each of the components of the expenditure variables before and after anonymization.

Table 9.13: Shares of expenditures components

.	TFOODEXP	TALCHEXP	TCLTHEXP	THOUSEXP	TFURNEXP	THLTHEXP
before	0.58	0.01	0.03	0.09	0.02	0.03
after	0.59	0.01	0.03	0.09	0.02	0.03
.	TTRANSEXP	TCOMMEXP	TRECEXP	TEDUEXP	TRESHOTEXP	TMISCEXP
before	0.04	0.02	0.00	0.08	0.03	0.05
after	0.04	0.02	0.00	0.08	0.03	0.05

Anonymization for the creation of a SUF will inevitably lead to some degree of utility loss. It is important to describe this loss in the external report, so that users are aware of the changes in the data. This is described in Step 11 and presented in [Appendix C](#). [Appendix C](#) also shows summary statistics and tabulations of the household level variables before and after anonymization.

### Merging the household- and individual-level variables

The next step is to merge the treated household variables with the untreated individual variables for the anonymization of the individual level variables. [Listing 9.21](#) shows the steps to merge these files. This also includes the selection of variables used in the anonymization of the individual-level variables. We create the *sdcMicro* object for the anonymization of the individual variables in the same way as for the household variable in [Listing 9.9](#). Subsequently, we repeat Steps 6-10 for the individual-level variables.

Listing 9.21: Merging the files with household and individual-level variables and creating an *sdcMicro* object for the anonymization of the individual-level variables

```

1  ### Select variables (individual level)
2  # Key variables (individual level)
3  selectedKeyVarsIND = c('GENDER', 'REL', 'MARITAL', 'AGEYRS',
4                        'EDUCY', 'ATSCHOOL', 'INDUSTRY1') # list of selected key_
5  ↪ variables
6
7  # Sample weight (WGTHH, individual weight)
8  selectedWeightVarIND = c('WGTHH')
9
10 # Household ID
11 selectedHouseholdID = c('IDH')
12
13 # No strata
14
15 # Recombining anonymized HH datasets and individual level variables
16 indVars <- c("IDH", "IDP", selectedKeyVarsIND, "WGTHH") # HID and all non HH variables
17 fileInd <- file[indVars] # subset of file without HHVars
18
19 HHmanip <- extractManipData(sdcHH) # manipulated variables HH
20 HHmanip <- HHmanip[HHmanip[, 'IDH'] != 1782,]
21
22 fileCombined <- merge(HHmanip, fileInd, by.x= c('IDH'))
23
24 fileCombined <- fileCombined[order(fileCombined[, 'IDH'],
25 fileCombined[, 'IDP']),]
26
27 dim(fileCombined)
28
29 # SDC objects with all variables and treated HH vars for
30 # anonymization of individual level variables
31 sdcCombined <- createSdcObj(dat = fileCombined, keyVars = selectedKeyVarsIND,

```

(continues on next page)

(continued from previous page)

```

31 weightVar = selectedWeightVarIND, hhId = selectedHouseholdID
    ↪selectedHouseholdID)

```

### 9.1.12 Step 6b: Assessing disclosure risk (individual level)

All key variables at the individual level are categorical. Therefore, we can use  $k$ -anonymity and the individual and global risk measures (see the Sections [Individual risk](#) and [Global risk](#)). The hierarchical risk is now of interest, given the household structure in the dataset *fileCombined*, which includes both household- and individual-level variables. The number of individuals (absolute and relative) that violate  $k$ -anonymity at the levels 2, 3 and 5 are shown in [Table 9.14](#).

**Note:**  $k$ -anonymity does not consider the household structure and therefore underestimates the risk. Therefore, we are more interested in the individual and global hierarchical risk measures.

Table 9.14:  $k$ -anonymity violations

$k$ -anonymity	Number HH violating	Percentage
2	998	9.91%
3	1,384	13.75%
5	2,194	21.79%

The global risk measures can be found using  $R$  as illustrated in [Listing 9.22](#). The global risk is 0.24%, which corresponds to 24 expected re-identifications. Accounting for the hierarchical structure, this rises to 1.26%, or 127 expected re-identifications. The global risk measures are low compared to the number of  $k$ -anonymity violators due to the low sampling weights. The high number of  $k$ -anonymity violators is mainly due to the very detailed age variable. The risk measures are based only on the individual level variables, since we assume that the individual and household level variables are be used simultaneously by an intruder. If we would consider an intruder scenario where these variables are used simultaneously by an intruder to re-identify individuals, the household level variables should also be taken into account here. This would results in a high number of key variables.

Listing 9.22: Global risk of the individual-level variables

```

1 print(sdcCombined, 'risk')
2 ## Risk measures:
3 ##
4 ## Number of observations with higher risk than the main part of the data: 0
5 ## Expected number of re-identifications: 23.98 (0.24 %)
6 ##
7 ## Information on hierarchical risk:
8 ## Expected number of re-identifications: 127.12 (1.26 %)

```

### 9.1.13 Step 7b: Assessing utility (individual level)

We evaluate the utility measures selected in Step 5 besides some general utility measures. The values computed from the raw data are presented in step 10b to allow for direct comparison with the values computed from the anonymized data.

### 9.1.14 Step 8b: Choice and application of SDC methods (individual level)

We use the same approach for the anonymization of the individual-level categorical key variables as for the household level categorical variables described earlier: first use global recoding to limit the necessary number of suppressions, then apply local suppressions and finally, if necessary, use of perturbative methods.

The variable AGEYRS (i.e., age in years) has many different values (age in months for children 0 – 1 years and age in years for individuals over 1 year). This level of detail leads to a high level of re-identification risk, given external datasets with exact age as well as knowledge of the exact age of close relatives. We have to reduce the level of detail in the age variables by recoding the age values (see the Section [Recoding](#)). First, we recode the values from 15 to 65 in ten-year intervals. Since some indicators related to education are computed from the survey dataset, our first approach is not to recode the age range 0 – 15 years. For children under the age of 1 year, we reduce the level of detail and recode these to 0 years. These recodes are shown in [Listing 9.23](#). We also top-code age at the age of 65 years. This protects individuals with high (rare) age values.

Listing 9.23: Recoding age in 10-year intervals in the range 15 – 65 and top code age over 65 years

```

1  # Recoding age and top coding age (top code 65), below that 10 year age
2  # groups, children aged under 1 are recoded 0 (previously in months)
3
4  sdcCombined@manipKeyVars$AGEYRS[sdcCombined@manipKeyVars$AGEYRS >= 0 &
5  sdcCombined@manipKeyVars$AGEYRS < 1] <- 0
6
7  sdcCombined@manipKeyVars$AGEYRS[sdcCombined@manipKeyVars$AGEYRS >= 15 &
8  sdcCombined@manipKeyVars$AGEYRS < 25] <- 20
9
10 ...
11
12 sdcCombined@manipKeyVars$AGEYRS[sdcCombined@manipKeyVars$AGEYRS >= 55 &
13 sdcCombined@manipKeyVars$AGEYRS < 66] <- 60
14
15 # topBotCoding also recalculates risk based on manual recoding above
16 sdcCombined <- **topBotCoding(obj = sdcCombined, value = 65,
17 replacement = 65, kind = 'top', column = 'AGEYRS')
18
19 table(sdcCombined@manipKeyVars$AGEYRS) # check results
20 ##      0      1      2      3      4      5      6      7      8      9     10     11     12     13     14
21 ##  311   367   340   332   260   334   344   297   344   281   336   297   326   299   263
22 ##    20    30    40    50    60    65
23 ## 1847 1220  889  554  314  325

```

These recodes already reduce the risk to 531 individuals violating 3-anonymity. We could recode the values of age in the lower range according to the age categories users require (e.g., 8 – 11 for education). There are many different categories for different indicators, however, including education indicators. This would reduce the utility of the data for some users. Therefore, we decide to look first at the number of suppressions needed in local suppression after this limited recoding. If the number of suppressions is too high, we can go back and recode age in the range 1 – 14 years.

In [Listing 9.24](#) we demonstrate how one might experiment with local suppression to find the best option. We use local suppression to achieve 3-anonymity (see the Section [Local suppression](#)). On the first attempt, we do not specify any importance vector; this leads to many suppressions in the variable AGEYRS (see [Table 9.15](#) below, first row), however. This is undesirable from a utility point of view. Therefore, we decide to specify an importance vector to prevent suppressions in the variable AGEYRS. Suppressing the variable GENDER is also undesirable from the utility point of view. The variable GENDER is a type of variable that should not have suppressions. We set GENDER as variable with the second highest importance. After specifying the importance vector to prevent suppressions of the age variable, there are no age suppressions (see [Table 9.15](#), second row). The total number of suppressions in the other variables increased, however, from 253 to 323 because of the importance vector. This is to be expected because

the algorithm without the importance vector minimizes the total number of suppressions by first suppressing values in variables with many categories – in this case, age and gender. Specifying an importance vector prevents reaching this optimality and hence leads to a higher total number of suppressions. There is a trade-off between which variables are suppressed and the total number of suppressions. After specifying an importance vector, the variable REL has many suppressions (see Table 9.15, second row). We choose this second option.

Listing 9.24: Experimenting with different options in local suppression

```

1  # Copy of sdcMicro object to later undo steps
2  sdcCopy <- sdcCombined
3
4  # Importance vectors for local suppression (depending on utility measures)
5  impVec1 <- NULL # for optimal suppression
6  impVec2 <- rep(length(selectedKeyVarsIND), length(selectedKeyVarsIND))
7  impVec2[match('AGEYRS', selectedKeyVarsIND)] <- 1 # AGEYRS
8  impVec2[match('GENDER', selectedKeyVarsIND)] <- 2 # GENDER
9
10 # Local suppression without importance vector
11 sdcCombined <- localSuppression(sdcCombined, k = 2, importance = impVec1)
12
13 # Number of suppressions per variable
14 print(sdcCombined, "ls")
15
16 ## Local Suppression:
17 ##      KeyVar | Suppressions (#) | Suppressions (%)
18 ##      GENDER |           0 |           0.000
19 ##      REL    |          34 |           0.338
20 ##      MARITAL |           0 |           0.000
21 ##      AGEYRS |         195 |           1.937
22 ##      EDUCY  |           0 |           0.000
23 ##      EDYRSCURRAT |         3 |           0.030
24 ##      ATSCHOOL |           0 |           0.000
25 ##      INDUSTRY1 |         21 |           0.209
26
27 # Number of suppressions per variable for each value of AGEYRS
28 table(sdcCopy@manipKeyVars$AGEYRS) - table(sdcCombined@manipKeyVars$AGEYRS)
29
30 ##  0  1  2  3  4  5  6  7  8  9 10 11 12 13 14 20 30 40 50 60 65
31 ##  0  0  0  0  0  0  2  0  2  1  0  1  4  1  5 25 53 37 36 15 13
32
33 # Undo local suppression
34 sdcCombined <- undolast(sdcCombined)
35
36 # Local suppression with importance vector on AGEYRS and GENDER
37 sdcCombined <- localSuppression(sdcCombined, k = 2, importance = impVec2)
38
39 # Number of suppressions per variable
40 print(sdcCombined, "ls")
41 ## Local Suppression:
42 ##      KeyVar | Suppressions (#) | Suppressions (%)
43 ##      GENDER |           0 |           0.000
44 ##      REL    |         323 |           3.208
45 ##      MARITAL |           0 |           0.000
46 ##      AGEYRS |           0 |           0.000
47 ##      EDUCY  |           0 |           0.000
48 ##      EDYRSCURRAT |           0 |           0.000
49 ##      ATSCHOOL |           0 |           0.000
50 ##      INDUSTRY1 |           0 |           0.000

```

(continues on next page)



(continued from previous page)

```

51 # Number of suppressions for each value of the variable AGEYRS
52 table(sdcCopy@manipKeyVars$AGEYRS) - table(sdcCombined@manipKeyVars$AGEYRS)
53 ## 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 20 30 40 50 60 65
54 ## 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0

```

Table 9.15: Number of suppressions by variable for different variations of local suppression

Local suppression options	GEN-DER	REL	MARI-TAL	AGEYRS	EDUCY	EDYRS	SCU-RATT	AT	SCHOOLINDUS-TRY1
k = 2, no imp	0	34	0	195	0	3	0	21	
k = 2, imp on AGEYRS	0	323	0	0	0	0	0	0	0

### 9.1.15 Step 9b: Re-measure risk (individual level)

We re-evaluate the risk measures selected in Step 6b. Table 9.16 shows that local suppression, not surprisingly, has reduced the number of individuals violating 2-anonymity to 0.

Table 9.16: k-anonymity violations

k-anonymity	Number HH violating	Percentage
2	0	0.00 %
3	197	1.96 %
5	518	5.15 %

The hierarchical global risk was reduced to 0.11%, which corresponds to 11.3 expected re-identifications. The highest individual hierarchical re-identification risk is 1.21%. These risk levels would seem acceptable for a SUF.

### 9.1.16 Step 10b: Re-measure utility (individual level)

We selected two utility measures for the individual variables: primary and secondary education enrollment, both also by gender. These two measures are sensitive to changes in the variables gender (GENDER), age (AGEYRS) and education (EDUCY and EDYRSATCURR), and therefore give a good overview of the impact of the anonymization. As shown in Table 9.17 the anonymization did not change the results. The results of the tabulations in Appendix C confirm these results.

Table 9.17: Net enrollment in primary and secondary education by gender

.	Primary education			Secondary education		
Before	72.6%	74.2%	70.9%	42.0%	44.8%	39.1%
After	72.6%	74.2%	70.9%	42.0%	44.8%	39.1%

### 9.1.17 Step 11: Audit and reporting

In the audit step, we check whether the data allow for reproduction of published figures from the original dataset and relationships between variables and other data characteristics are preserved in the anonymization process. In short, we

check whether the dataset is valid for analytical purposes. There are no figures available that were published from the dataset and need to be reproducible from the anonymized data.

In Step 2, we explored the data characteristics and relationships between variables. These data characteristics and relationships have been mainly preserved, since we took them into account when choosing the appropriate anonymization methods. The variables `TANHHEXP` and `INCTOTGROSSHH` are the sums of the individual components, because we added noise to the components and reconstructed the aggregates by summing over the components. Initially, the income variables were all positive. This characteristic has been violated, as a result of noise addition. Since values of the variable `AGEYRS` were not perturbed, but only recoded and suppressed, we did not introduce unlikely combinations, such as a 60-year-old individual enrolled in primary education. Also, by separating the anonymization process into two parts, one for household-level variables and one for individual-level variables, the values of variables measured at the household level agree for all members of each household.

Furthermore, we drafted two reports, internal and external, on the anonymization of the case study dataset. The internal report includes the methods used, the risk before and after anonymization as well as the reasons for the selected methods and their parameters. The external report focuses on the changes in the data and the loss in utility. Focus here should be on the number of suppressions as well as the perturbative methods (PRAM). This is described in the previous steps.

---

**Note:** When creating a SUF, it is inevitable that there will be a loss of information and it is very important for the users to be aware of these changes and release them in a report that accompanies the data.

---

[Appendix C](#) provides examples of an internal and external report of the anonymization process of this dataset. Depending on the users and readers of the reports, the content may differ. The code to this case study shows how to obtain the information for the reports. Some measures are also available in the standard reports generated with the `report()` function. This is shown in [Listing 9.25](#). The `report()` function will only use the data available in the `sdcMicro` object, which does not contain all households for `sdcHH`.

Listing 9.25: Using the `report()` function for internal and external reports

```
1 # Create reports with sdcMicro report() function
2 report(sdcHH, internal = F) # external (brief) report
3 report(sdcHH, internal = T) # internal (extended) report
4
5 # Create reports with sdcMicro report() function
6 report(sdcCombined, internal = F) # external (brief) report
7 report(sdcCombined, internal = T) # internal (extended) report
```

### 9.1.18 Step 12: Data release

The final step is the release of the anonymized dataset together with the external report. [Listing 9.26](#) shows how to collect the data from the `sdcMicro` object with the `extractManipData()` function. Before releasing the file, we add an individual ID to the file (line number in household). We export the anonymized dataset in as *STATA* file. The [Section Read functions in R](#) presents functions for exporting files in other data formats.

Listing 9.26: Exporting the anonymized dataset

```
1 # Anonymized dataset
2 # Household variables and individual variables
3 # extracts all variables, not just the manipulated ones
4 dataAnon <- extractManipData(sdcCombined, ignoreKeyVars = F, ignorePramVars = F,
5                             ignoreNumVars = F, ignoreStrataVar = F)
6
```

(continues on next page)

(continued from previous page)

```

7 # Create STATA file
8 write.dta(dataframe = dataAnon, file= 'Case1DataAnon.dta', convert.dates=TRUE)

```

## 9.2 Case study 2 - PUF

This case study is a continuation of case study 1 in the Section *Case study 1- SUF* . Case study 1 produces a SUF file. In this case study we use this SUF file to produce a PUF file of the same dataset, which can be freely distributed. The structure of the SUF and PUF releases will be the same. However, the PUF will contain fewer variables and less (detailed) information than the SUF. We refer to the Section *Case study 1- SUF* for a description of the dataset.

---

**Note:** It is also possible to directly produce a PUF from a dataset without first creating a SUF.

---

As in case study 1, we show how the creation of a PUF can be achieved using the open source and free *sdcMicro* package and *R*. A ready-to-run *R* script for this case study and the dataset are also available to reproduce the results and allow the user to adapt the code (see <http://ihnsn.org/home/projects/sdc-practice>). Extracts of this code are presented in this section to illustrate several steps of the anonymization process.

---

**Note:** The choices of methods and parameters in this case study are based on this particular dataset and the results and choices might be different for other datasets.

---

This case study follows the steps of the SDC process outlined in *The SDC Process*.

### 9.2.1 Step 1: Need for disclosure control

The same reasoning as in case study 1 applies: the SUF dataset produced in case study 1 contains data on individuals and households and some variables are confidential and/or sensitive. The decisions made in case study 1 are based on the disclosure scenarios for a SUF release. The anonymization applied for the SUF does not provide sufficient protection for the release as PUF and the SUF file cannot be released as PUF without further treatment. Therefore, we have to repeat the SDC process with a different set of disclosure scenarios based on the characteristics of a PUF release (see Step 4). This leads to different risk measures, lower accepted risk levels and different SDC methods.

### 9.2.2 Step 2: Data preparation and exploring data characteristics

In order to guarantee consistency between the released PUF and SUF files, which is required to prevent intruders from using the datasets together (SUF users have also access to the PUF file), we have to use the anonymized SUF file to create the PUF file (see also the Section *Step 3: Type of release*). In this way all information in the PUF file is also contained in the SUF, and the PUF does not provide additional information to an intruder with access to the SUF. We load the required packages to read the data (*foreign* package for *STATA* files) and load the SUF dataset into “file” as illustrated in *Listing 9.27*. We also load the original data file (raw data) as “fileOrig”. We need the raw data to undo perturbative methods used in case study 1 (see Step 8) and to compare data utility measures (see Step 5). To evaluate the utility loss in the PUF, we have to compare the information in the anonymized PUF file with the information in the raw data. For an overview of the data characteristics and a description of the variables in both files, we refer to Step 2 of *Case study 1- SUF* .

Listing 9.27: Loading required packages and datasets

```

1 # Load required packages
2 library(foreign) # for read/write function for STATA
3 library(sdcMicro) # sdcMicro package
4
5 # Set working directory - set to the path on your machine
6 setwd("/Users/CaseStudy2")
7
8 # Specify file name of SUF file from case study 1
9 fname <- "CaseDataAnon.dta"
10
11 # Specify file name of original dataset (raw data)
12 fnameOrig <- "CaseA.dta"
13
14 # Read-in files
15 file <- read.dta(fname, convert.factors = TRUE) # SUF file case study 1
16 fileOrig <- read.dta(fnameOrig, convert.factors = TRUE) # original data

```

We check the number of variables and number of observations of both files and the variable names of the SUF file, as shown in Listing 9.28. The PUF file has fewer records and fewer variables than the original data file, since we removed large households and several variables to generate the SUF file.

Listing 9.28: Number of individuals and variables and variable names

```

1 # Dimensions of file (observations, variables)
2 dim(file)
3 ## [1] 10068    49
4
5 dim(fileOrig)
6 ## [1] 10574    68
7
8 colnames(file) # Variable names
9 ## [1] "IDH" "URBRUR" "REGION" "HHSIZE"
10 ## [5] "OWNAGLAND" "RELIG" "ROOF" "TOILET"
11 ## [9] "WATER" "ELECTCON" "FUELCOOK" "OWNMOTORCYCLE"
12 ## [13] "CAR" "TV" "LIVESTOCK" "LANDSIZEHA"
13 ## [17] "TANHHEXP" "TFOODEXP" "TALCHEXP" "TCLTHEXP"
14 ## [21] "THOUSEXP" "TFURNEXP" "THLTHEXP" "TTRANSEXP"
15 ## [25] "TCOMMEXP" "TRECEXP" "TEDUEXP" "TRESTHOTEXP"
16 ## [29] "TMISCEXP" "INCTOTGROSSHH" "INCRMT" "INCWAGE"
17 ## [33] "INCFARMBSN" "INCNFARMBSN" "INCRENT" "INCFIN"
18 ## [37] "INCPENSN" "INCOTHER" "WGTPOP" "IDP"
19 ## [41] "GENDER" "REL" "MARITAL" "AGEYRS"
20 ## [45] "EDUCY" "EDYRSCURRAT" "ATSCHOOL" "INDUSTRY1"
21 ## [49] "WGTHH"

```

To get an overview of the values of the variables, we use tabulations and cross-tabulations for categorical variables and summary statistics for continuous variables. To include the number of missing values ('NA' or other), we use the option `useNA = "ifany"` in the `table()` function. For some variables, these tabulations differ from the tabulations of the raw data, due to the anonymization of the SUF file.

In Table 9.18 the variables in the dataset "file" are listed along with concise descriptions of the variables, the level at which they are collected (individual level (IND) or household level (HH)), the measurement type (continuous, semi-continuous or categorical) and value ranges. Note that the dataset contains a selection of 49 variables of the 68 variable selected for the SUF release. The variables have been preselected based on their relevance for data users and some variables were removed while creating a SUF file. The numerical values for many of the categorical variables are

codes that refer to values, e.g., in the variable “URBRUR”, 1 stands for ‘rural’ and 2 for ‘urban’. More information on the meanings of coded values of the categorical variables is available in the *R* code for this case study.

Any data cleaning, such as recoding missing value codes and removing empty variables, was already done in case study 1. The same holds for removing any direct identifiers. Direct identifiers are not released in the SUF file.

We identified the following sensitive variables in the dataset: variables related to schooling and labor force status as well as income and expenditure related variables. These variables need protection. Whether a variable is considered sensitive may depend on the release type, country and the dataset itself.

Table 9.18: Overview of the variables in the dataset

No.	Variable name	Description	Level	Measurement	Values
1	IDH	Household ID	HH	.	1-2,000
2	IDP	Individual ID	IND	.	1-13
3	REGION	Region	HH	categorical	1-6
4	URBRUR	Area of residence	HH	categorical	1, 2
5	WGTHH	Individual weighting coefficient	HH	weight	31.2-84
6	WGTPOP	Population weighting coefficient	IND	weight	45.8-93
7	HHSIZE	Household size	HH	semi-continuous	1-33
8	GENDER	Gender	IND	categorical	0, 1
9	REL	Relationship to household head	IND	categorical	1-9
10	MARITAL	Marital status	IND	categorical	1-6
11	AGEYRS	Age in completed years	IND	semi-continuous	0-65
12	RELIG	Religion of household head	HH	categorical	1, 5-7, 9
13	ATSCHOOL	Currently enrolled in school	IND	categorical	0, 1
14	EDUCY	Highest level of education attended	IND	categorical	1-6
15	EDYRSCUR AT	Years of education for currently enrolled	IND	semi-continuous	1-18
16	INDUSTRY	Industry classification (1-digit)	IND	categorical	1-10
17	ROOF	Main material used for roof	IND	categorical	1-5, 9
18	TOILET	Main toilet facility	HH	categorical	1-4, 9
19	ELECTCON	Electricity	HH	categorical	0-3
20	FUELCOOK	Main cooking fuel	HH	categorical	1-5, 9
21	WATER	Main source of water	HH	categorical	1-9
22	OWNAGLAN	Ownership of agricultural land	HH	categorical	1-3
23	LANDSIZE	Land size owned by household (ha) (agric and non agric)	HH	continuous	0-40
24	OWNMOTORCYCLE	Ownership of motorcycle	HH	categorical	0, 1
25	CAR	Ownership of car	HH	categorical	0, 1
26	TV	Ownership of television	HH	categorical	0, 1
27	LIVESTOC	Number of large-sized livestock owned	HH	semi-continuous	0-25
28	INCRMT	Income – Remittances	HH	continuous	
29	INCWAGE	Income - Wages and salaries	HH	continuous	
30	INCFARMBSN	Income - Gross income from household farm businesses	HH	continuous	
31	INCNFARMBSN	Income - Gross income from household nonfarm businesses	HH	continuous	
32	INCRENT	Income - Rent	HH	continuous	
33	INCFIN	Income - Financial	HH	continuous	
34	INCPENSN	Income - Pensions/social assistance	HH	continuous	
35	INCOTHER	Income - Other	HH	continuous	
36	INCTOTGROSHH	Income - Total	HH	continuous	
37	FARMEMP				
38	TFOODEXP	Total expenditure on food	HH	continuous	
39	TALCHEXP	Total expenditure on alcoholic beverages, tobacco and narcotics	HH	continuous	
40	TCLTHEXP	Total expenditure on clothing	HH	continuous	

Continued on next

Table 9.18 – continued from previous page

No.	Variable name	Description	Level	Measurement	Values
41	THOUSEXP	Total expenditure on housing	HH	continuous	
42	TFURNEXP	Total expenditure on furnishing	HH	continuous	
43	THLTHEXP	Total expenditure on health	HH	continuous	
43	TTRANSEXP	Total expenditure on transport	HH	continuous	
44	TCOMMEXP	Total expenditure on communication	HH	continuous	
45	TRECEXP	Total expenditure on recreation	HH	continuous	
46	TEDUEXP	Total expenditure on education	HH	continuous	
47	TRESHOTEXP	Total expenditure on restaurants and hotels	HH	continuous	
48	TMISCEXP	Total expenditure on miscellaneous spending	HH	continuous	
49	TANHHEXP	Total annual nominal household expenditures	HH	continuous	

It is always important to ensure that the relationships between variables in the data are preserved during the anonymization process and to explore and take note of these relationships before beginning the anonymization. At the end of the anonymization process before the release of the data, an audit should be conducted, using these initial results, to check that these relationships are maintained in the anonymized dataset (see Step 11).

In our dataset, we identify several relationships between variables that need to be preserved during the anonymization process. The variables “TANHHEXP” and “INCTOTGROSSHH” represent the total annual nominal household expenditure and the total gross annual household income, respectively, and these variables are aggregations of existing income and expenditure components in the dataset.

The variables related to education are available only for individuals in the appropriate age groups and missing for other individuals. In addition, the household-level variables (cf. fourth column of Table 9.18) have the same values for all members in any particular household. The value of household size corresponds to the actual number of individuals belonging to that household in the dataset. As we proceed, we have to take care that these relationships and structures are preserved in the anonymization process.

We assume that the data are collected in a survey that uses simple sampling of households. The data contains two weight coefficients: “WGTHH” and “WGTPOP”. The relationship between the weights is  $WGTPOP = WGTHH * HHSIZE$ . “WGTPOP” is the sampling weight for the households and “WGTHH” is the sampling weight for the individuals to be used for disclosure risk calculations. “WGTHH” is used for computing individual-level indicators (such as education) and “WGTPOP” is used for population level indicators (such as income indicators). There are no strata variables available in the data. We will use “WGTPOP” for the anonymization of the household variables and “WGTHH” for the anonymization of the individual-level variables.

### 9.2.3 Step 3: Type of release

In this case study, we assume that the file will be released as a PUF, which will be freely available to anyone interested in the data (see the Section [Conditions for PUFs](#) for the conditions and more information on the release of PUFs). The PUF release is intended for users with lower information requirements (e.g., students) and researchers interested in the structure of the data and willing to do preliminary research. The PUF file can give an idea to the researcher whether it is worthwhile for their research to apply for access to the SUF file. Researchers willing to do more in-depth research will most likely apply for SUF access. Generally, users of a PUF file are not restricted by an agreement that prevents them from using the data to re-identify individuals and hence the accepted risk level is much lower than in the case of the SUF and the set of released variables is limited.

### 9.2.4 Step 4: Intruder scenarios and choice of key variables

Next, based on the release type, we reformulate the intruder scenarios for the PUF release. This leads to the selection of a set of quasi-identifiers. Since this case study is based on a demo dataset, we do not have a real context and we cannot define exact disclosure scenarios. Therefore, we make hypothetical assumptions on possible disclosure

scenarios. We consider two types of disclosure scenarios: 1) matching with other publically available datasets and 2) spontaneous recognition. Since the dataset will be distributed as PUF, there are de facto no restrictions on the use of the dataset by intruders.

For the sake of illustration, we assume that population registers are available with the demographic variables gender, age, place of residence (region, urban/rural), religion and other variables such as marital status and variables relating to education and professional status that are also present in our dataset. In addition, we assume that there is a publically available cadastral register on land ownership. Based on this analysis of available data sources, we have selected in case study 1 the variables “REGION”, “URBRUR”, “HHSIZE”, “OWNAGLAND”, “RELIG”, “GENDER”, “REL” (relationship to household head), “MARITAL” (marital status), “AGEYRS”, “INDUSTRY1” and two variables relating to school attendance as categorical quasi-identifiers, the expenditure and income variables as well as LANDSIZEHA as continuous quasi-identifiers. According to our assessment, these variables might enable an intruder to re-identify an individual or household in the dataset by matching with other available datasets. The key variables for PUF release generally coincide with the key variables for the SUF release. Possibly, more variables could be added, since the user has more possibilities to match the data extensively and is not bound by any contract, as is in the case of the SUF file. Equally, some key variables in the SUF file may not be released in the PUF file and, as a consequence, these variables are removed from the list of key variables.

Upon further consideration, this initial set of identifying variables is too large for a PUF release, as the number of possible combinations (keys) is very high and hence many respondents could be identified based on these variables. Therefore, we decide to limit the set of key variables, by excluding variables from the dataset for PUF release. The choice of variables to be removed is led by the needs of the intended PUF users. Assuming the typical users are mainly interested in aggregate income and expenditure data, we can therefore remove from the initial set of key variables “OWNAGLAND”, “RELIG” and “LANDSIZEHA” at the household level and “EDYRSCURRAT” and “ATSCHOOL” at the individual level.

---

**Note:** These variables will not be released in the PUF file.

---

We also remove the income and expenditure components from the list of key variables, since we reduce their information content by building proportions (see Step 8a). The list of the remaining key variables is presented in [Table 9.19](#).

Table 9.19: Overview of selected key variables for PUF file

Variable name	Variable description	Measurement level
REGION	region	Household, categorical
URBRUR	area of residence	Household, categorical
HHSIZE	household size	Household, categorical
TANHHEXP	total expenditure	Household, continuous
INCTOTGROSSHH	total income	Household, continuous
GENDER	gender	Individual, categorical
REL	relationship to household head	Individual, categorical
MARITAL	marital status	Individual, categorical
AGEYRS	age in completed years	Individual, semi-continuous/categorical
EDUCY	highest level of education completed	Individual, categorical
INDUSTRY1	industry classification	Individual, categorical

The decision to release the dataset as a PUF means the level of anonymization will be relatively high and consequently, the variables are less detailed (e.g., after recoding) and a scenario of spontaneous recognition is less likely. Nevertheless, we should check for rare combinations or unusual patterns in the variables. Variables that may lead to spontaneous recognition in our sample are amongst others “HHSIZE” (household size) as well as “INCTOTGROSSHH” (aggregate income) and “TANHHEXP” (total expenditure). Large households and households with high income are easily identifiable, especially when combined with other identifying variables such as a geographical identifier (“REGION”). There might be only one or a few households/individuals in a certain region with a high income, such as the



local doctor. Variables that are easily observable and known by neighbors such as “ROOF”, “TOILET”, “WATER”, “ELECTCON”, “FUELCOOK”, “OWNMOTORCYCLE”, “CAR”, “TV” and “LIVESTOCK” may also need protection depending on what stands out in the community, since a user might be able to identify persons (s)he knows. This is called the nosy-neighbor scenario.

## 9.2.5 Step 5: Data key uses and selection of utility measures

A PUF file contains less information and the file is generally used by students as a teaching file, by researchers to get an idea about the data structure, and for simple analyses. The users have generally lower requirements than for a SUF file and the results of analysis may be less precise. The researcher interested in a more detailed dataset, would have to apply for access to the SUF file. Therefore, we select more aggregate utility measures for the PUF file that reflect the intended use of a PUF file. Data intensive measures, such as the Gini coefficient, cannot be computed from the PUF file. Besides the standard utility measures, such as tabulations, we evaluate the decile dispersion ratio and a regression with the income deciles as regressand.

To measure the information loss, we should compare the initial data file before any anonymization (including the anonymization for the SUF) with the file after the anonymization for the PUF. Comparing the files directly before and after the PUF anonymization would underestimate the information loss, as this would omit the information loss due to SUF anonymization. Therefore, in Step 2, we also loaded the raw dataset.

### Hierarchical (household) structure

As noted in case study 1, the data has a household structure. For the SUF release, we protected large households by removing these from the dataset. Since some variables are measured on the household level and thus have identical values for each household member, the values of the household variables should be treated in the same way for each household member (see the Section [Anonymization of the quasi-identifier household size](#)). Therefore, we first anonymize only the household variables. After this, we merge them with the individual-level variables and then anonymize the individual-level and household-level variables jointly.

Since the data has a hierarchical structure, Steps 6 through 10 are repeated twice: Steps 6a through 10a are for the household-level variables and Steps 6b through 10b for the combined dataset. In this way, we ensure that household-level variable values remain consistent across household members for each household and the household structure cannot be used to re-identify individuals. This is further explained in the Sections [Levels of risk](#) and [Household structure](#).

Before continuing to Step 6a, we select the categorical key variables, continuous key variables and any variables selected for use in PRAM routines, as well as household-level sampling weights in *R*. We also collect the variable names of the variables that will not be released. The PRAM variables are variables select for the PRAM routine, which we discuss further in Step 8a. We extract these selected household variables from the SUF dataset and save them as “fileHH”. The choice of PRAM variables is further explained in Step 8a. [Listing 9.29](#) illustrates how these steps are done in *R* (see also the Section [Household structure](#)).

Listing 9.29: Selecting the variables for the household-level anonymization

```

1 # Categorical key variables at household level
2 selectedKeyVarsHH <- c('URBRUR', 'REGION', 'HHSIZE')
3
4 # Continuous key variables
5 numVarsHH <- c('TANHHEXP', 'INCTOTGROSSHH')
6
7 # PRAM variables
8 pramVarsHH <- c('ROOF', 'TOILET', 'WATER', 'ELECTCON',
9               'FUELCOOK', 'OWNMOTORCYCLE', 'CAR', 'TV', 'LIVESTOCK')
10
11 # Household weight

```

(continues on next page)



(continued from previous page)

```

12 weightVarHH <- c('WGTPOP')
13
14 # Variables not suitable for release in PUF (HH level)
15 varsNotToBeReleasedHH <- c("OWNAGLAND", "RELIG", "LANDSIZEHA")
16
17 # Vector with names of all HH level variables
18 HHVars <- c('IDH', selectedKeyVarsHH, pramVarsHH, numVarsHH, weightVarHH)
19
20 # Create subset of file with only HH level variables
21 fileHH <- file[,HHVars]

```

Every household has the same number of entries as it has members (e.g., a household of three will be repeated three times in “fileHH”). Before analyzing the household-level variables, we select only one entry per household, as illustrated in [Listing 9.30](#). This is further explained in the [Section Household structure](#). In the same way we extract “fileOrigHH” from “fileOrig”. “fileOrigHH” contains all variables from the raw data, but contains every household only once. We need “fileOrigHH” in Steps 8a and 10a for undoing some perturbative methods used in the SUF file and computing utility measures from the raw data respectively.

Listing 9.30: Taking a subset with only households

```

1 # Remove duplicated rows based on IDH, one row per household in fileHH
2 fileHH <- fileHH[which(!duplicated(fileHH$IDH)),] # SUF file
3 fileOrigHH <- fileOrig[which(!duplicated(fileOrig$IDH)),] # original dataset
4
5 # Dimensions of fileHH
6 dim(fileHH)
7 ## [1] 1970 16
8
9 dim(fileOrigHH)
10 ## [1] 2000 68

```

The file “fileHH” contains 1,970 households and 16 variables. We are now ready to create our *sdcMicro* object with the corresponding variables we selected in [Listing 9.28](#). For our case study, we will create an *sdcMicro* object called “sdcHH” based on the data in “fileHH”, which we will use for steps 6a – 10a (see [Listing 9.34](#)).

---

**Note:** When the *sdcMicro* object is created, the *sdcMicro* package automatically calculates and stores the risk measures for the data.

---

This leads us to Step 6a.

Listing 9.31: Creating a *sdcMicro* object for the household variables

```

1 # Create initial sdcMicro object for household level variables
2 sdcHH <- createSdcObj(dat = fileHH, keyVars = selectedKeyVarsHH,
3                      pramVars = pramVarsHH, weightVar = weightVarHH, numVars =
4                      ↪ numVarsHH)
5 numHH <- length(fileHH[,1]) # number of households

```

## 9.2.6 Step 6a: Assessing disclosure risk (household level)

Based on the key variables selected in the disclosure scenarios, we can evaluate the risk at the household level. The PUF risk measures show a lower risk level than in the SUF file after anonymization in case study 1. The reason is that

the set of key variables is smaller, since some variables will not be released in the PUF file. Removing (key) variables reduces the risk, and it is one of the most straightforward SDC methods.

As a first measure, we evaluate the number of households violating  $k$ -anonymity at the levels 2, 3 and 5. Table 9.20 shows the number of violating households as well as the percentage of the total number of households. Listing 9.32 illustrates how to find these values with *sdcMicro*. The `print()` function in *sdcMicro* shows only the values for thresholds 2 and 3. Values for other thresholds can be calculated manually by summing up the frequencies smaller than the  $k$ -anonymity threshold, as shown in Listing 9.32. The number of violators is already at a low level, due to the prior anonymization of the SUF file and the reduced set of key variables.

Table 9.20: Number and proportion of households violating  $k$ -anonymity

k-anonymity level	Number of HH violating	Percentage of total number of HH
2	0	0.0%
3	18	0.9%
5	92	4.7%

Listing 9.32: Showing number of households violating  $k$ -anonymity for levels 2, 3 and 5

```

1 # Number of observations violating k-anonymity (thresholds 2 and 3)
2 print(sdcHH)
3 ## Infos on 2/3-Anonymity:
4 ##
5 ## Number of observations violating
6 ##   - 2-anonymity: 0
7 ##   - 3-anonymity: 18
8 ##
9 ## Percentage of observations violating
10 ##   - 2-anonymity: 0.000 %
11 ##   - 3-anonymity: 0.914 %
12 -----
13
14 # Calculate sample frequencies and count number of obs. violating k(5) - anonymity
15 kAnon5 <- sum(sdcHH@risk$individual[,2] < 5)
16 kAnon5
17 ## [1] 92
18
19 # As percentage of total
20 kAnon5 / numHH
21 ## [1] 0.04670051

```

It is often useful to view the records of the household(s) that violate  $k$ -anonymity. This might help to find which variables cause the uniqueness of these households; this can then be used later when choosing appropriate SDC methods. Listing 9.32 shows how to access the values of the households violating 3 and 5-anonymity. Not surprisingly, the variable “HHSIZE” is responsible for many of the unique combinations and the origin of much of the risk. This is even the case after removing large households for the SUF release.

Listing 9.33: Showing records of households that violate  $k$ -anonymity

```

1 # Show values of key variable of records that violate k-anonymity
2 fileHH[sdcHH@risk$individual[,2] < 3, selectedKeyVarsHH] # for 3-anonymity
3 fileHH[sdcHH@risk$individual[,2] < 5, selectedKeyVarsHH] # for 5-anonymity

```

We also assess the disclosure risk of the categorical variables with the individual and global risk measures as described in the Sections [Individual risk](#) and [Global risk](#). In “fileHH” every entry represents a household. Therefore, we use the individual non-hierarchical risk here, where the individual refers in this case to a household. “fileHH” is a subset of the

complete dataset and contains only households and has, contrary to the complete dataset, no hierarchical structure. In Step 6b, we evaluate the hierarchical risk in the dataset “file”, the dataset containing both households and individuals. The individual and global risk measures automatically take into consideration the household weights, which we defined in Listing 9.29. In our file, the global risk measure calculated using the chosen key variables is lower than 0.01% (the smallest reported value is 0.01%, in fact the global risk is 0.0000642 %). This percentage is extremely low and corresponds to 0.13 expected re-identifications. The results are also shown in Listing 9.34. This low figure can be explained by the relatively small sample size of 0.25% of the total population (see case study 1). Furthermore, one should keep in mind that this risk measure is based only on the categorical quasi-identifiers at the household level.

Listing 9.34: Printing global risk measures

```
1 print(sdcHH, "risk")
2 ## Risk measures:
3 ##
4 ## Number of observations with higher risk than the main part of the data: 0
5 ## Expected number of re-identifications: 0.13 (0.01 %)
```

The global risk measure does not provide information about the spread of the individual risk measures. There might be a few households with relatively high risk, while the global (average) risk is low. Therefore we check the highest individual risk as shown in Listing 9.35. The individual risk of the household with the highest risk is 0.1 %, which is still very low.

Listing 9.35: Determining the highest individual risk

```
1 # Highest individual risk
2 max(sdcHH@risk$individual[, "risk"])
3 ## [1] 0.001011633
```

Since the selected key variables at the household level are both categorical and numerical, the individual and global risk measures based on frequency counts do not completely reflect the disclosure risk of the entire dataset. When generating the SUF file, we concluded that recoding of continuous variables to make them all categorical would likely not satisfy the needs of the SUF users. For the PUF file it is acceptable to recode continuous variables, such as income and expenditures since PUF content is typically less detailed. In Step 8a we will recode these variables into deciles and convert them into categorical variables. Therefore, we exclude these variables from the risk calculations now. We take these variables into account while remeasuring risk after anonymization.

### 9.2.7 Step 7a: Assessing utility measures (household level)

The utility of the data does not only depend on the household level variables, but on the combination of household-level and individual-level variables. Therefore, it is not useful to evaluate all the utility measures selected in Step 5 at this stage, i.e., before anonymizing the individual level variables. We restrict the initial measurement of utility to those measures that are solely based on the household variables. In our dataset, these are the measures related to income and expenditure and their distributions. The results are presented in Step 10a, together with the results after anonymization, which allow direct comparison. If after the next anonymization step it appears that the data utility has been significantly decreased by the suppression of some household level variables, we can return to this step.

---

**Note:** To analyze the utility loss, the utility measures before anonymization have to be calculated from the raw data and not from the anonymized SUF file.

---

Not all measures from case study 1 can be computed from the PUF file, since the information content is lower. The set of utility measures we use to evaluate the information loss in the PUF file consists of measures that need less detailed variables. This reflects the lower requirements a PUF user has on the dataset.

### 9.2.8 Step 8a: Choice and application of SDC methods (household level)

This step is divided into the anonymization of the categorical key variables and the continuous key variables, since different methods are used for both sets of variables. As already discussed in Step 4, we do not release all variables in the PUF file. At the household level “RELIG” (religion of household head), “OWNAGLAND” (land ownership) and “LANDSIZEHA” (plot size in ha) are not released in addition to the variables removed for the SUF release. For the sake of illustration, we assume that the variable “RELIG” is too sensitive and the variables “OWNAGLAND” and “LANDSIZEHA” are too identifying.

#### Categorical variables

We are now ready to move on to the choice of SDC methods for the categorical variables on the household level in our dataset. In the SUF file we already recoded some of the key variables and used local suppression. We only have three categorical key variables at the household level; “URBRUR”, “REGION” and “HHSIZE”. The selected categorical key variables at the household level are not suitable for recoding at this point, since the values cannot be grouped further. “URBRUR” has only two distinct categories and “REGION” has only six non-combinable regions. As noted before, the variable “HHSIZE” can be reconstructed by a headcount per household. Therefore, recoding of this variable alone does not lead to disclosure control.

Due to the relatively low risk of re-identification based on the three selected categorical household level variables, it is possible in this case to use an option like local suppression to achieve our desired level of risk. Applying local suppression when initial risk is relatively low will likely only lead to suppression of few observations and thus limit the loss of utility. If, however, the data had been measured to have a relatively high risk, then applying local suppression without previous recoding would likely result in a large number of suppressions and greater information loss. Efforts such a recoding should be taken first before suppressing values in cases where risk is initially measured as high. Recoding will reduce risk with little information loss and thus the number of suppressions, if local suppression is applied as an additional step.

We apply local suppression to reach 5-anonymity. The chosen level of five is higher than in the SUF release and is based on the release type as PUF. This leads to a total of 39 suppressions, all in the variable “HHSIZE”. As explained earlier, suppression of the value of the variable “HHSIZE” does not lead to actual suppression of this information. Therefore, we redo the local suppression, but this time we tell *sdcMicro* to, if possible, not suppress “HHSIZE” but one of the other variables. Alternatively, we could remove households with suppressed values of the variable “HHSIZE”, remove large households or split households.

In *sdcMicro* it is possible to tell the algorithm which variables are important and less important for making small changes (see also the Section [Local suppression](#)). To prevent values of the variable “HHSIZE” being suppressed, we set the importance of “HHSIZE” in the importance vectors to the highest (i.e., 1). We try two different importance vectors: the first where “REGION” is more important than “URBRUR” and the second with the importance of “REGION” and “URBRUR” swapped. [Listing 9.36](#) shows how to apply local suppression and put importance on the variable “HHSIZE”.

---

**Note:** In [Listing 9.36](#) we use the `undolast()` function in *sdcMicro* to go one step back after we had first applied local suppression with no importance vector.

---

The `undolast()` function restores the *sdcMicro* object back to the previous state (i.e., before we applied local suppression), which allows us to rerun the same command, but this time with an importance vector set. The `undolast()` function can only be used to go one step back.

The suppression patterns of the three different options are shown in [Table 9.21](#). The importance is clearly reflected in the number of suppressions per variable. The total number of suppressions is with an importance vector higher than without an importance vector (44/73 vs. 39), but 5-anonymity is achieved in the dataset with no suppressions in the variable “HHSIZE”. This means that we do not have to remove or split households. The variable “REGION” is the type of variable that should not have any suppressions either. From that perspective we chose the third option. This leads to more suppressions, but no suppressions in “HHSIZE” and as few as possible in “REGION”.

Table 9.21: Number of suppressions by variable after local suppression with and without importance vector

Key variable	Number of suppressions and proportion of total		
.	No importance vector	Importance HHSIZE, URBUR, REGION	Importance HHSIZE, REGION, URBUR
URBUR	0 (0.0 %)	2 (0.1 %)	61 (3.1 %)
REGION	0 (0.0 %)	42 (2.1 %)	12 (0.6 %)
HHSIZE	39 (2.0 %)	0 (0.0 %)	0 (0.0 %)

Listing 9.36: Local suppression with and without importance vector

```

1  # Local suppression to achieve 5-anonymity
2  sdchH <- localSuppression(sdchH, k = 5, importance = NULL) # no importance vector
3  print(sdchH, "ls")
4  ## Local Suppression:
5  ##   KeyVar | Suppressions (#) | Suppressions (%)
6  ##   URBUR |           0 |           0.000
7  ##   REGION |           0 |           0.000
8  ##   HHSIZE |          39 |           1.980
9  ## -----
10
11 # Undo suppressions to see the effect of an importance vector
12 sdchH <- undolast(sdchH)
13
14 # Redo local suppression minimizing the number of suppressions in HHSIZE
15 sdchH <- localSuppression(sdchH, k = 5, importance = c(2, 3, 1))
16
17 print(sdchH, "ls")
18 ## Local Suppression:
19 ##   KeyVar | Suppressions (#) | Suppressions (%)
20 ##   URBUR |           2 |           0.102
21 ##   REGION |          42 |           2.132
22 ##   HHSIZE |           0 |           0.000
23 ## -----
24
25 # Undo suppressions to see the effect of a different importance vector
26 sdchH <- undolast(sdchH)
27
28 # Redo local suppression minimizing the number of suppressions in HHSIZE
29 sdchH <- localSuppression*(sdchH, k = 5, importance = c(3, 2, 1))
30
31 print(sdchH, "ls")
32 ## Local Suppression:
33 ##   KeyVar | Suppressions (#) | Suppressions (%)
34 ##   URBUR |          61 |           3.096
35 ##   REGION |          12 |           0.609
36 ##   HHSIZE |           0 |           0.000
37 ## -----

```

In case study 1 we applied invariant PRAM to the variables “ROOF”, “TOILET”, “WATER”, “ELECTCON”, “FUELCOOK”, “OWNMOTORCYCLE”, “CAR”, “TV” and “LIVESTOCK”, since these variables are not sensitive and were not selected as quasi-identifiers because we assumed that there are no external data sources containing this information that could be used for matching. Values can be easily observed or be known to neighbors, however, and therefore are important, together with other variables, for the nosy neighbor scenario. For the PUF release we would like to level of uncertainty by increasing the number of changes. Therefore, we redo PRAM with a different transi-

tion matrix. As discussed in the Section [PRAM \(Post Randomization Method\)](#), the invariant PRAM method has the property that the univariate distributions do not change. To maintain this property, we reapply PRAM to the raw data, rather than to the already PRAMmed variables in the SUF file.

[Listing 9.37](#) illustrates how to apply PRAM. We use the original values to apply PRAM and replace the values in the *sdcMicro* object with these values. We choose the parameter ‘pd’, the lower bound for the probability that a value is not changed, to be relatively low at 0.6. This is a lower value than the 0.8 used in the SUF file and will lead to a higher number of changes (cf. [Listing 9.17](#)). This is acceptable for a PUF file and introduces more uncertainty as required for a PUF release. [Listing 9.37](#) also shows the number of changed records per variables. Because PRAM is a probabilistic method, we set a seed for the random number generator before applying PRAM to ensure reproducibility of the results.

---

**Note:** In some cases the choice of the seed matters. The choice of seed changes the results.

---

The seed should not be released, since it allows for reconstructing the original values if combined with the transition matrix. The transition matrix can be released: this allows for consistent statistical inference by correcting the statistical methods used if the researcher has knowledge about the PRAM method (at this point *sdcMicro* does not allow to retrieve the transition matrix).

Listing 9.37: Applying PRAM

```
1 # PRAM
2 set.seed(10987)
3
4 # Replace PRAM variables in sdcMicro object sdcHH with the original raw values
5 sdcHH@origData[,pramVarsHH] <- fileHH[match(fileHH$IDH, fileOrigHH$IDH), pramVarsHH]
6 sdcHH@manipPramVars <- fileHH[match(fileHH$IDH, fileOrigHH$IDH), pramVarsHH]
7
8 sdcHH <- pram(obj = sdcHH, pd = 0.6)
9 ## Number of changed observations:
10 ## - - - - -
11 ## ROOF != ROOF_pram : 305 (15.48%)
12 ## TOILET != TOILET_pram : 260 (13.2%)
13 ## WATER != WATER_pram : 293 (14.87%)
14 ## ELECTCON != ELECTCON_pram : 210 (10.66%)
15 ## FUELCOOK != FUELCOOK_pram : 315 (15.99%)
16 ## OWNMOTORCYCLE != OWNMOTORCYCLE_pram : 95 (4.82%)
17 ## CAR != CAR_pram : 255 (12.94%)
18 ## TV != TV_pram : 275 (13.96%)
19 ## LIVESTOCK != LIVESTOCK_pram : 109 (5.53%)
```

PRAM has changed values within the variables according to the invariant transition matrices. Since we used the invariant PRAM method (see the Section [PRAM \(Post Randomization Method\)](#)), the absolute univariate frequencies remain approximately unchanged. This is not the case for the multivariate frequencies. In Step 10a we compare the multivariate frequencies before and after anonymization for the PRAMmed variables.

### Continuous variables

We have selected the variables “INCTOTGROSSHH” (total income) and “TANHHEXP” (total expenditure) as numerical quasi-identifiers, as discussed in Step 4. In Step 5 we identified variables having high interest for the users of our data: many users use the data for measuring inequality and expenditure patterns. The noise addition in the SUF file does not protect these variables sufficiently, especially, because outliers are not protected. Therefore, we decide to recode total income and total expenditure into deciles.

As with PRAM, we want to compute the deciles from the raw data rather than from the perturbed values in the SUF file. We compute the deciles directly from the raw data and overwrite these values in the *sdcMicro* object. Subsequently, we compute the mean of each decile from the raw data and replace the values for total income and total expenditures

with the mean of the respective decile. In this way the mean of both variables does not change. This approach can be interpreted as univariate microaggregation with very large groups (group size  $n/10$ ) with the mean as replacement value (see the Section [Microaggregation](#)).

The information in the income and expenditure variables by component is too sensitive to release as PUF, and, summing those variables would allow an intruder to reconstruct the totals. PUF users might however be interested in the shares. Therefore, we decide to keep the income and expenditure components as proportions of the raw totals, rounded to two digits. The anonymization of the income and expenditure variables is shown in [Listing 9.38](#).

Listing 9.38: Anonymization of income and expenditure variables

```

1  # Create bands (deciles) for income and expenditure variables
2  (aggregates) based on the original data
3  decExp <- as.numeric(cut(fileOrigHH[match(fileHH$IDH, fileOrigHH$IDH), "TANHHEXP"],
4                        quantile(fileOrigHH[match(fileHH$IDH, fileOrigHH$IDH),
5                                "TANHHEXP"],
6                                (0:10)/10, na.rm = T),
7                                include.lowest = TRUE, labels = c(1, 2, 3, 4, 5, 6, 7, 8, 9,
8                                ↪10)))
9  decInc <- as.numeric(cut(fileOrigHH[match(fileHH$IDH, fileOrigHH$IDH), "INCTOTGROSSHH
10 ↪"],
11                        quantile(fileOrigHH[match(fileHH$IDH, fileOrigHH$IDH),
12                                "INCTOTGROSSHH"],
13                                (0:10)/10, na.rm = T),
14                                include.lowest = TRUE, labels = c(1, 2, 3, 4, 5, 6, 7, 8,
15                                ↪9, 10)))
16
17  # Mean values of deciles
18  decExpMean <- round(sapply(split(fileOrigHH[match(fileHH$IDH, fileOrigHH$IDH),
19 ↪"TANHHEXP"],
20                                decExp), mean))
21  decIncMean <- round(sapply(split(fileOrigHH[match(fileHH$IDH, fileOrigHH$IDH),
22 ↪"INCTOTGROSSHH"],
23                                decInc), mean))
24
25  # Replace with mean value of decile
26  sdCHH@manipNumVars$TANHHEXP <- decExpMean[match(decExp, names(decExpMean))]
27  sdCHH@manipNumVars$INCTOTGROSSHH <- decIncMean[match(decInc, names(decIncMean))]
28
29  # Recalculate risks after manually changing values in sdMicro object calcRisks(sdCHH)
30  calcRisks(sdCHH)
31
32  # Extract data from sdCHH
33  HHmanip <- extractManipData(sdCHH) # manipulated variables HH
34
35  # Keep components of expenditure and income as share of total,
36  # use original data since previous data was perturbed
37  compExp <- c('FOODEXP', 'TALCHEXP', 'TCLTHEXP', 'THOUSEXP',
38              'TFURNEXP', 'THLTHEXP', 'TTRANSEXP', 'TCOMMEXP',
39              'TRECEXP', 'TEDUEXP', 'TRESTHOTEXP', 'TMISCEXP')
40  compInc <- c('INCRMT', 'INCWAGE', 'INCFARMBSN', 'INCNFARMBSN',
41              'INCRENT', 'INCFIN', 'INCPENSN', 'INCOTHER')
42  HHmanip <- cbind(HHmanip, round(fileOrigHH[match(fileHH$IDH, fileOrigHH$IDH),
43 ↪compExp] /
44                                fileOrigHH[match(fileHH$IDH, fileOrigHH$IDH),
45                                "TANHHEXP"], 2))
46  HHmanip <- cbind(HHmanip, round(fileOrigHH[match(fileHH$IDH, fileOrigHH$IDH),
47 ↪compInc] /

```

(continues on next page)



(continued from previous page)

```

39         fileOrigHH[match(fileHH$IDH, fileOrigHH$IDH),
40                        "INCTOTGROSSHH"], 2))

```

### 9.2.9 Step 9a: Re-measure risk (household level)

For the categorical variables, we conclude that we have achieved 5-anonymity in the data with local suppression. 5-anonymity also implies 2- and 3-anonymity. The global risk stayed close to zero (as the expected number of re-identifications), which is very low. Therefore, we conclude that based on the categorical variables, the data has been sufficiently anonymized. One should keep in mind that the anonymization methods applied are complementing the ones used for the SUF.

**Note:** The methods selected methods in this case study alone would not be sufficient to protect the data set for a PUF release.

We have reduced the risk of spontaneous recognition of households, by removing the variable “LANDSIZEHA” and PRAMming the variables identified to be important in the nosy neighbor scenario. An intruder cannot know with certainty whether a household that (s)he recognizes in the data is the correct household, due to the noise in these variables.

These measures refer only to the categorical variables. To evaluate the risk of the continuous variables we could use an interval measure or closest neighbor algorithm. These risk measures are discussed in the Section [Risk measures for continuous variables](#). We chose to use an interval measure, since exact value matching is not our largest concern based on the assumed scenarios and external data sources. Instead, datasets with similar values but not the exact same values could be used for matching. Here the main concern is that the values are sufficiently far from the original values, which is measured with an interval measure.

[Listing 9.39](#) shows how to evaluate the interval measure for the variables “INCTOTGROSSHH” and “TANHHEXP” (total income and expenditure). The different values of the parameter  $k$  in the function `dRisk()` define the size of the interval around the original value as a function of the standard deviation, as explained in the Section [Interval measure](#). The larger  $k$ , the larger the intervals, the higher the probability that a perturbed value is in the interval around the original value and the higher the risk measure. The results are satisfactory, especially when keeping in mind that there are only 10 distinct values in the dataset (the means of each of the deciles). All outliers have been recoded. Looking at the proportions of the components, we do not detect any outliers (households with an unusual high or low spending pattern in one component).

Listing 9.39: Measuring risk of re-identification of continuous variables

```

1  # Risk evaluation continuous variables
2  dRisk(sdcHH@origData[,c("TANHHEXP", "INCTOTGROSSHH")],
3        xm = sdcHH@manipNumVars[,c("TANHHEXP", "INCTOTGROSSHH")], k = 0.01)
4  ## [1] 0.4619289
5
6  dRisk(sdcHH@origData[,c("TANHHEXP", "INCTOTGROSSHH")],
7        xm = sdcHH@manipNumVars[,c("TANHHEXP", "INCTOTGROSSHH")], k = 0.02)
8  ## [1] 0.642132
9
10 dRisk(sdcHH@origData[,c("TANHHEXP", "INCTOTGROSSHH")],
11        xm = sdcHH@manipNumVars[,c("TANHHEXP", "INCTOTGROSSHH")], k = 0.05)
12 ## [1] 0.8258883

```



### 9.2.10 Step 10a Re-measure utility (household level)

The utility in the data has decreased compared to the raw data, mainly because variables were completely removed. Many of the utility measures used in case study 1 are not applicable to the PUF file. However, by replacing the deciles with their means, we can still use the income and expenditure variables for arithmetic operations. Also the shares of the income and expenditure components can still be used, since they are based on the raw data.

We select two additional utility measures: the decile dispersion ratio and the share of total consumption by the poorest decile. The decile dispersion ratio is the ratio of the average income of the top decile and the average income of the bottom decile. [Listing 9.40](#) shows how to compute these from the raw data and the household variables after anonymization. [Table 9.22](#) presents the estimated values. The differences are small and mainly due to the removed households.

Table 9.22: Comparison of utility measures

.	Raw data	PUF file
Decile dispersion ratio	24.12	23.54
Share of consumption by the poorest decile	0.0034	0.0035

Listing 9.40: Computation of decile dispersion ratio and share of total consumption by the poorest decile

```

1 # Decile dispersion ratio
2 # raw data
3 mean(tail(sort(fileOrigHH$INCTOTGROSSHH), n = 200)) /
4   mean(head(sort(fileOrigHH$INCTOTGROSSHH), n = 200))
5 ## [1] 24.12152
6
7 mean(tail(sort(HHmanip$INCTOTGROSSHH), n = 197)) /
8   mean(head(sort(HHmanip$INCTOTGROSSHH), n = 197))
9 ## [1] 23.54179
10
11 # Share of total consumption by the poorest decile households
12 sum(head(sort(fileOrigHH$TANHHEXP), n = 200)) / sum(fileOrigHH$TANHHEXP)
13 ## [1] 0.003411664
14
15 sum(head(sort(HHmanip$TANHHEXP), n = 197)) / sum(HHmanip$TANHHEXP)
16 ## [1] 0.003530457

```

#### Merging the household- and individual-level variables

The next step is to merge the treated household variables with the untreated individual variables for the anonymization of the individual level variables. [Listing 9.41](#) shows the steps to merge these files. This also includes the selection of variables used in the anonymization of the individual-level variables. We create the *sdcMicro* object for the anonymization of the individual variables in the same way as for the household variable in [Listing 9.31](#). Generally, at this stage, the household level and individual level variables should be combined and quasi-identifiers at both levels be used (see the Section [Levels of risk](#)). Unfortunately, in our dataset, this leads to long computation times. Therefore, we create two *sdcMicro* objects, one with all key variables (“sdcCombinedAll”) and one with only the individual level key variables (“sdcCombined”). In Step 6b we compare the risk measures for both cases and in Step 8b we discuss alternative approaches to keeping the complete set of variables. We now repeat Steps 6-10 for the individual-level variables.

Listing 9.41: Merging the files with household and individual-level variables and creating an *sdcMicro* object for the anonymization of the individual-level variables

```

1  ### Select variables (individual level)
2  selectedKeyVarsIND = c('GENDER', 'REL', 'MARITAL',
3                        'AGEYRS', 'EDUCY', 'INDUSTRY1') # list of selected key_
4  ↪variables
5
6  # sample weight (WGTHH, individual weight)
7  selectedWeightVarIND = c('WGTHH')
8
9  # Household ID
10 selectedHouseholdID = c('IDH')
11
12 # Variables not suitable for release in PUF (IND level)
13 varsNotToBeReleasedIND <- c("ATSCHOOL", "EDYRSCURRAT")
14
15 # All individual level variables
16 INDVars <- c(selectedKeyVarsIND)
17
18 # Recombining anonymized HH data sets and individual level variables
19 indVars <- c("IDH", "IDP", selectedKeyVarsIND, "WGTHH") # HID and all non HH vars
20 fileInd <- file[indVars] # subset of file without HHVars
21 fileCombined <- merge(HHmanip, fileInd, by.x = c('IDH'))
22 fileCombined <- fileCombined[order(fileCombined[, 'IDH'], fileCombined[, 'IDP']),]
23
24 dim(fileCombined)
25 ## [1] 10068    44
26
27 # SDC objects with only IND level variables
28 sdcCombined <- createSdcObj(dat = fileCombined, keyVars = c(selectedKeyVarsIND),
29                           weightVar = selectedWeightVarIND, hhId = _
30 ↪selectedHouseholdID)
31
32 # SDC objects with both HH and IND level variables
33 sdcCombinedAll <- createSdcObj(dat = fileCombined,
34                               keyVars = c(selectedKeyVarsIND, selectedKeyVarsHH ),
35                               weightVar = selectedWeightVarIND, hhId = _
36 ↪selectedHouseholdID)

```

### 9.2.11 Step 6b: Assessing disclosure risk (individual level)

As first measure, we evaluate the number of records violating  $k$ -anonymity at the levels 2, 3 and 5. Table 9.23 shows the number of violating individuals as well as the percentage of the total number of households. The second and third column refer to “sdcCombined” and the fourth and fifth column to “sdcCombinedAll”. We see that combining the individual level and household level variables leads to a large increase in the number of  $k$ -anonymity violators. The choice not to include the household level variables is pragmatically driven by the computation time and can be justified by the different type of variables on the household level and individual level. One could assume that these variables are not available in the same dataset and can therefore not simultaneously be used by an intruder to re-identify individuals.

Table 9.23: Number of records violating  $k$ -anonymity

.	sdcCombined		sdcCombinedAll	
k-anonymity	Number of records violating	Percentage of total records	Number of records violating	Percentage of total records
2	0	0.0 %	4,048	40.2 %
3	167	1.7 %	6,107	60.7 %
5	463	4.6 %	8,292	82.4 %

The global hierarchical risk measure is 0.095%, which corresponds to approximately 10 expected re-identifications. We use here the hierarchical risk measure, since the re-identification of a single household member would lead to the re-identification of the other members of the same household too. This number is low compared to the number of  $k$ -anonymity violations, due to the high sample weights, which protect the data already to a large extent. Only 24 observations have an individual hierarchical risk higher than 1%, with a maximum of 1.17%. This is mainly because of the lower sample weights of these records. Listing 9.42 shows how to retrieve these measures in R.

Listing 9.42: Risk measures before anonymization

```

1 numIND <- length(fileCombined[,1]) # number of households
2
3 # Number of observations violating k-anonymity
4 print(sdcCombined)
5 ## Infos on 2/3-Anonymity:
6 ##
7 ## Number of observations violating
8 ## - 2-anonymity: 0
9 ## - 3-anonymity: 167
10 ##
11 ## Percentage of observations violating
12 ## - 2-anonymity: 0.000 %
13 ## - 3-anonymity: 1.659 %
14 ## -----
15
16 # Calculate sample frequencies and count number of obs. violating k(3,5) - anonymity
17 kAnon5 <- sum(sdcCombined@risk$individual[,2] 5)
18 kAnon5
19 ## [1] 463
20
21 # As percentage of total
22 kAnon5 / numIND
23 ## [1] 0.04598729
24
25 # Global risk on individual level
26 print(sdcCombined, 'risk')
27 ## Risk measures:
28 ##
29 ## Number of observations with higher risk than the main part of the data: 0
30 ## Expected number of re-identifications: 1.69 (0.02 %)
31 ##
32 ## Information on hierarchical risk:
33 ## Expected number of re-identifications: 9.57 (0.10 %)
34 ## -----
35
36 # Number of observation with relatively high risk
37 dim(fileCombined[sdcCombined@risk$individual[, "hier_risk"] > 0.01,])
38 ## [1] 24 44

```

(continues on next page)

(continued from previous page)

```
39  
40 # Highest individual risk  
41 max(sdcCombined@risk$individual[, "hier_risk"])  
42 ## [1] 0.01169091
```

### 9.2.12 Step 7b: Assessing utility measures (individual level)

We evaluate the utility measures as discussed in Step 5 based on the raw data (before applying any anonymization measures). The results are presented in Step 10b together with the values after anonymization to allow for direct comparison.

### 9.2.13 Step 8b: Choice and application of SDC methods (individual level)

In this step, we discuss four different techniques used for anonymization: 1) removing variables from the dataset to be released, 2) recoding of categorical variables to reduce the level of detail, 3) local suppression to achieve the required level of  $k$ -anonymity, 4) randomization of the order of the records in the file. Finally, we discuss some alternative options for treating the household structure in the dataset.

#### Removing variables

Additional to the variables removed from the dataset for the SUF release (see case study 1), we further reduce the number of variables in the dataset to be released. This is normal practice for PUF releases. Sensitive or identifying variables are removed, which allows to release other variables at a more detailed level. In a PUF release, the set of key variables should be limited.

In our case, we decide to remove at the individual level the variables “EDYRSCURRAT”, as this variable is too identifying (identifies whether there are school-going children in the household). We keep the variable “EDUCY” (highest level of education attended) for information on education.

---

**Note:** As an alternative to removing the variables from the dataset, one could also set all values to missing. This would allow the user to see the structure and variables contained in the SUF file.

---

#### Recoding

As noted before, PUF users require a lower level of information and therefore we can recode the key variables even further to reduce the disclosure risk. The recoding of variables in case study 1 is not sufficient for a PUF release. Therefore, we recode most of the categorical key variables from [Table 9.19](#) to reduce the risk and number of necessary suppressions by local suppression. [Table 9.24](#) gives an overview of the recodes made. All new categories are formed with the needs of the data user in mind. [Listing 9.43](#) shows how to do this in *R* and also shows value labels and the univariate tabulations of these variables before and after recoding.

Table 9.24: Overview of recodes of categorical variables at individual level

Variable	Recoding
REL (relation to household head)	recode 'Father/Mother', 'Grandchild', 'Son/Daughter in law', 'Other relative' to 'Other relative' and recode 'Domestic help' and 'Non-relative' to 'Other'
MARITAL (marital status)	recode 'Married monogamous', 'Married polygamous', 'Common law, union coutumiere, union libre, living together' to 'Married/living together' and 'Divorced/Separated' and 'Widowed' to 'Divorced/Separated/Widowed'
AGEYRS (age in completed years)	recode values under 15 to 7 (other values have been recoded for SUF)
EDUCY (highest level of education completed)	recode 'Completed lower secondary (or post-primary vocational education) but less than completed upper secondary', 'Completed upper secondary (or extended vocational/technical education)', 'Post secondary technical' and 'University and higher' to 'Completed lower secondary or higher'
INDUSTRY1	recode to 'primary', 'secondary' and 'tertiary'

Listing 9.43: Recoding the categorical and continuous variables

```

1  # Recode REL (relation to household head)
2  table(sdcCombined@manipKeyVars$REL, useNA = "ifany")
3  ##
4  ##      1      2      3      4      5      6      7      8      9 <NA>
5  ## 1698 1319 4933   52   765   54   817   40   63   327
6
7  # 1 - Head, 2 - Spouse, 3 - Child, 4 - Father/Mother, 5 - Grandchild, 6 - Son/
8  ↪ Daughter in law
9  # 7 - Other relative, 8 - Domestic help, 9 - Non-relative
10 sdcCombined <- groupVars(sdcCombined, var = "REL", before = c("4", "5", "6", "7"),
11                          after = c("7", "7", "7", "7")) # other relative
12 sdcCombined <- groupVars(sdcCombined, var = "REL", before = c("8", "9"),
13                          after = c("9", "9")) # other
14
15 table(sdcCombined@manipKeyVars$REL, useNA = "ifany")
16 ##
17 ##      1      2      3      7      9 <NA>
18 ## 1698 1319 4933 1688  103   327
19
20 # Recode MARITAL (marital status)
21 table(sdcCombined@manipKeyVars$MARITAL, useNA = "ifany")
22 ##
23 ##      1      2      3      4      5      6 <NA>
24 ## 3542 2141  415   295   330  329 3016
25
26 # 1 - Never married, 2 - Married monogamous, 3 - Married polygamous,
27 # 4 - Common law, union coutumiere, union libre, living together, 5 - Divorced/
28 ↪ Separated,
29 # 6 - Widowed
30 sdcCombined <- groupVars(sdcCombined, var = "MARITAL", before = c("2", "3", "4"),
31                          after = c("2", "2", "2")) # married/living together
32 sdcCombined <- groupVars(sdcCombined, var = "MARITAL", before = c("5", "6"),
33                          after = c("9", "9")) # divorced/seperated/widowed*

```

(continues on next page)

(continued from previous page)

```

34 table(sdcCombined@manipKeyVars$MARITAL, useNA = "ifany")
35 ##
36 ##      1      2      9 <NA>
37 ## 3542 2851  659 3016
38
39 # Recode AGEYRS (0-15 years)
40 table(sdcCombined@manipKeyVars$AGEYRS, useNA = "ifany")
41 ##
42 ##      0      1      2      3      4      5      6      7      8      9     10     11     12     13     14
43 ##  311   367   340   332   260   334   344   297   344   281   336   297   326   299   263
44 ##    20    30    40    50    60    65 <NA>
45 ## 1847 1220  889  554  314  325  188
46
47 sdcCombined <- groupVars(sdcCombined, var = "AGEYRS",
48                          before = c("0", "1", "2", "3", "4", "5", "6", "7", "8", "9",
49                                     "10", "11", "12", "13", "14"),
50                          after = rep("7", 15))
51
52 table(sdcCombined@manipKeyVars$AGEYRS, useNA = "ifany")
53 ##
54 ##      7     20     30     40     50     60     65 <NA>
55 ## 4731 1847 1220  889  554  314  325  188
56
57 sdcCombined <- calcRisks(sdcCombined)
58 # Recode EDUCY (highest level of educ compl)
59 table(sdcCombined@manipKeyVars$EDUCY, useNA = "ifany")
60 ##
61 ##      0      1      2      3      4      5      6 <NA>
62 ## 1582 4755 1062  330  139   46  104 2050
63
64 # 0 - No education, 1 - Pre-school/ Primary not completed,
65 # 2 - Completed primary, but less than completed lower secondary
66 # 3 - Completed lower secondary (or post-primary vocational education)
67 #   but less than completed upper secondary
68 # 4 - Completed upper secondary (or extended vocational/technical education),
69 # 5 - Post secondary technical
70 # 6 - University and higher
71 sdcCombined <- groupVars(sdcCombined, var = "EDUCY", before = c("3", "4", "5", "6"),
72                          after = c("3", "3", "3", "3")) # completed lower secondary_
73 ↪ or higher
74 table(sdcCombined@manipKeyVars$EDUCY, useNA = "ifany")
75 ##
76 ##      0      1      2      3 <NA>
77 ## 1582 4755 1062  619 2050
78
79 # Recode INDUSTRY1 ()
80 table(sdcCombined@manipKeyVars$INDUSTRY1, useNA = "ifany")
81 ##
82 ##      1      2      3      4      5      6      7      8      9     10 <NA>
83 ## 5300   16   153    2    93  484   95   17    70   292 3546
84
85 # 1 - Agriculture and Fishing, 2 - Mining, 3 - Manufacturing, 4 - Electricity and_
86 ↪ Utilities
87 # 5 - Construction, 6 - Commerce, 7 - Transportation, Storage and Communication, 8 -_
88 ↪ Financial, Insurance and Real Estate
89 # 9 - Services: Public Administration, 10 - Other Services, 11 - Unspecified
90 sdcCombined <- groupVars(sdcCombined, var = "INDUSTRY1", before = c("1", "2"),

```

(continues on next page)

(continued from previous page)

```

88         after = c("1", "1")) # primary
89 sdcCombined <- groupVars(sdcCombined, var = "INDUSTRY1", before = c("3", "4", "5"),
90         after = c("2", "2", "2")) # secondary
91 sdcCombined <- groupVars(sdcCombined, var = "INDUSTRY1", before = c("6", "7", "8", "9
92   ↪", "10"),
93         after = c("3", "3", "3", "3", "3")) # tertiary
94 table(sdcCombined@manipKeyVars$INDUSTRY1, useNA = "ifany")
95 ##
96 ##      1      2      3 <NA>
97 ## 5316  248  958 3546

```

### Local suppression

The recoding has reduced the risk already considerably. We use local suppression to achieve the required level of  $k$ -anonymity. Generally, the required level of  $k$ -anonymity for PUF files is 3 or 5. In this case study, we require 5-anonymity. Listing 9.44 shows the suppression pattern without specifying an importance vector. All suppressions are made in the variable “AGEYRS”. This is the variable with the highest number of different values, and hence considered first by the algorithm. We try different suppression patterns by specifying importance vectors, but we decide that the pattern without importance vector yields the best result. This is also the result with the lowest total number of suppressions. Less than 1 percent suppression in the age variable is acceptable. We could reduce this number by further recoding the variable “AGEYRS”.

Listing 9.44: Local suppression to reach 5-anonymity

```

1 # Local suppression without importance vector
2 sdcCombined <- localSuppression(sdcCombined, k = 5, importance = NULL)
3
4 # Number of suppressions per variable
5 print(sdcCombined, "ls")
6 ## Local Suppression:
7 ##      KeyVar | Suppressions (#) | Suppressions (%)
8 ##      GENDER |                0 |             0.000
9 ##      REL    |                0 |             0.000
10 ##     MARITAL |                0 |             0.000
11 ##     AGEYRS  |               91 |             0.904
12 ##     EDUCY   |                0 |             0.000
13 ##  INDUSTRY1 |                0 |             0.000

```

### Randomization of order of records

The records in the dataset are ordered by region and household ID. There is a certain geographical order of the households within the regions, due to the way the households IDs were assigned. Intruders could reconstruct suppressed values by using this structure. To prevent this, we randomly reorder the records within the regions. Listing 9.45 shows how to do this in R. We first count the number of records per region

---

**Note:** Some records have their region value suppressed, so we include the count of NAs.

---

Subsequently, we draw randomly household IDs, in such way that the regional division is respected. Finally, we sort the file by the new, randomized, individual ID (“IDP”). Households with suppressed values for “REGION” will be last in the reordered file. Before randomizing the order, we extract the data from the *sdcMicro* object “sdcCombined” as shown in Listing 9.45.

Listing 9.45: Randomizing the order of records within regions

```

1  # Randomize order of households dataAnon and recode IDH to random
2  number (sort file by region)
3  set.seed(97254)
4  # Sort by region
5  dataAnon <- dataAnon[order(dataAnon$REGION),]
6
7  # Number of households per region
8  hhperregion <- table(dataAnon[match(unique(dataAnon$IDH), dataAnon$IDH), "REGION"],
9                      useNA = "ifany")
10
11 # Randomized IDH (household ID)
12 randomHHid <- c(sample(1:hhperregion[1], hhperregion[1]),
13                unlist(lapply(1:(length(hhperregion)-1),
14                             function(i){sample((sum(hhperregion[1:i]) + 1):
15                                     ↪sum(hhperregion[1:(i+1)]), hhperregion[(i+1)]))})))
16
17 dataAnon$IDH <- rep(randomHHid, table(dataAnon$IDH)[match(unique(dataAnon$IDH),
18                                                             as.numeric(names(table(dataAnon$IDH))))])
19
20 # Sort by IDH (and region)
21 dataAnon <- dataAnon[order(dataAnon$IDH),]

```

### Alternative options for dealing with household structure

In Step 6b we compared the disclosure risk for two cases: one with only individual level key variables and another with individual level and household level key variables combined. We decided to use only the individual level key variables to reduce the computation time and justified this choice by arguing that intruders cannot use household and individual level variables simultaneously. This might not always be the case. Therefore we explore other options to reduce the risk when taking both individual level and household level variables into account. We present two options: removing the household structure; and using options in the local suppression algorithm.

#### *Removing household structure*

We consider the risk emanating from the household structure in the dataset to be very high. We can remove the hierarchical household structure completely and treat all variables at the individual level. This entails, besides removing the household id (“IDH”), also treating variables that could be used for reconstructing households. These are, for instance, “REL” (relation to household head), “HHSIZE” (household size), and any of the household level variables, such as income and expenditure. However, not all household level variables need to be treated. For example, “REGION” is a household level variable, but the probability that this variable leads to the reconstruction of a household is low. Also, we need to reorder the records in the file, as they are sorted by households. Note that by removing the household structure, we interpret all variables as individual level variables for measuring disclosure risk. This leads to a lower need for recoding and suppression, since the hierarchical risk disappears. The reason why we did not opt for this approach is the loss of utility for the user. The household structure is an important feature of the data, and should be kept in the PUF file.

#### *Using different options for local suppression*

The long running time is mainly due to the local suppression algorithm. In the Section [Local suppression](#) we discuss options to reduce the running time of the local suppression in case of many key variables. The all- $m$  approach reduces the running time by first considering subsets of the complete set of key variables. This reduces the complexity of the problem and leads to lower computation times. However, the total number of suppressions made is likely to be higher. Also, if not explicitly specified, it is not guaranteed that the required level for  $k$ -anonymity is automatically achieved on the complete set of key variables. It is therefore important to check the results.



### 9.2.14 Step 9b: Re-measure risk

We re-evaluate the risk measures selected in Step 6. Table 9.25 shows that local suppression, not surprisingly, has reduced the number of individuals violating 5-anonymity to 0. The global hierarchical risk was reduced to 0.02%, which corresponds to approximately 2 correct re-identifications. The highest individual hierarchical re-identification risk is 0.2%. These risk levels are acceptable for a PUF release. Furthermore, the recoding has removed any unusual combinations in the data.

---

**Note:** The risk may be underestimated by excluding the household level variables.

---

Table 9.25: k-anonymity violations

k-anonymity	Number of records violating	Percentage
2	0	0.0 %
3	0	0.0 %
5	0	0.0 %

### 9.2.15 Step 10b: Re-measure utility

We compare (cross-)tabulations before and after anonymization, which are illustrated in the *R* code to this case study. We note that due to the recoding in Step 8b, the detail in the variables is reduced. This reduces the number of necessary suppressions and is acceptable for a public use file.

### 9.2.16 Step 11: Audit and reporting

In the audit step, we check whether the data allow for reproduction of published figures from the original dataset and relationships between variables and other data characteristics are preserved in the anonymization process. In short, we check whether the dataset is valid for analytical purposes. There are no figures available that were published from the dataset and need to be reproducible from the anonymized data.

In Step 2, we explored the data characteristics and relationships between variables. These data characteristics and relationships have been mainly preserved, since we took them into account when choosing the appropriate anonymization methods. Since values of the variable “AGEYRS” were not perturbed, but only recoded and suppressed, we did not introduce unlikely combinations, such as a 60-year-old individual enrolled in primary education. Also, by separating the anonymization process into two parts, one for household-level variables and one for individual-level variables, the values of variables measured at the household level agree for all members of each household.

Furthermore, we drafted two reports, internal and external, on the anonymization of the case study dataset. The internal report includes the methods used, the risk before and after anonymization as well as the reasons for the selected methods and their parameters. The external report focuses on the changes in the data and the loss in utility. Focus here should be on the number of suppressions as well as the perturbative methods (PRAM). This is described in the previous steps.

---

**Note:** When creating a PUF, it is inevitable that there will be a loss of information and it is very important for the users to be aware of these changes and release them in a report that accompanies the data.

---

Appendix C provides examples of an internal and external report of the anonymization process of this dataset. Depending on the users and readers of the reports, the content may differ.

---

**Note:** The `report()` function in `sdcmicro` is at this point not useful, since this will only report on the SDC measures in the second case study.

---

However, the report should contain the entire process, including the measures applied in case study 1.

### 9.2.17 Step 12: Data release

The final step is the release of the anonymized dataset together with the external report. [Listing 9.46](#) shows how to export the anonymized dataset as *STATA* file. The Section [Read functions in R](#) presents functions for exporting files in other data formats.

Listing 9.46: Exporting the anonymized PUF file

```
1 # Create STATA file
2 write.dta(dataframe = dataAnon, file= 'Case2DataAnon.dta', convert.dates=TRUE)
```

# CHAPTER 10

## Appendices

### 10.1 Appendix A: Overview of Case Study Variables

.	Variable	Description
1	REGION	Region
2	DIST	District
3	URBRUR	Area of residence
4	WGTHH	Individual weighting coefficient (country-specific weighting co-efficient to derive individual-level indi
5	WGTPOP	Population weighting coefficient (weighting co-efficient to derive population-level indicators)
6	IDH	Household unique identification
7	IDP	Individual identification
8	HHSIZE	Household members
9	GENDER	Sex
10	REL	Relationship to household head
11	MARITAL	Marital status
12	AGEYRS	Age in completed years
13	AGEMTH	Age of child in completed months
14	RELIG	Religion of household head
15	ETHNICITY	Ethnicity
16	LANGUAGE	Language
17	MORBID	Morbidity last x weeks
18	MEASLES	Child immunized against Measles
19	MEDATT	Sought medical attention
20	CHWEIGHTKG	Weight of the child (Kg)
21	CHHEIGHTCM	Height of the child (cms)
22	ATSCHOOL	Current school enrolment
23	EDUCY	Highest level of education completed
24	EDYRS	Years of education
25	EDYRSCURRAT	Years of education for currently enrolled
26	SCHTYP	Type of school attending

Continued

Table 10.1 – continued from previous page

.	Variable	Description
27	LITERACY	Literacy status
28	EMPTYPI	Type of employment, Primary job
29	UNEMP1	Unemployed
30	INDUSTRY1	1 digit industry classification, Primary job
31	EMPCAT1	Employment categories, Primary job
32	WHOURSWEK1	Hours worked last week, Primary job
33	OWNHOUSE	Ownership of dwelling unit
34	ROOF	Main material used for roof
35	TOILET	Main toilet facility
36	ELECTCON	Connection of electricity in dwelling
37	FUELCOOK	Main cooking fuel
38	WATER	Main source of water
39	OWNAGLAND	Ownership of agricultural land
40	LANDSIZEHA	Land size owned by household (ha)
41	OWNMOTORCYCLE	Ownership of motorcycle
42	CAR	Ownership of car
43	TV	Ownership of television
44	LIVESTOCK	Number of large-sized livestock owned
45	INCRMT	Total amount of remittances received from remittance sending members
46	INCWAGE	Wage and salaries (annual)
47	INCBONSOCALL	Bonus and social allowance from wage job (annual)
48	INCFARMBSN	Gross income from household farm businesses (annual)
49	INCNFARMBSN	Gross income from household non-farm businesses (annual)
50	INCRENT	Rental income (annual)
51	INCFIN	Financial income from savings, loans, tax refunds, maturity payments on insurance
52	INCPENSIN	Pension and other social assistance (annual)
53	INCOTHER	Other income(annual)
54	INCTOTGROSSHH	Total gross household income (annual)
55	FARMEMP	Farm employment
56	THOUSEXP	Total expenditure on housing
57	TFOODEXP	Total food expenditure
58	TALCHEXP	Total alcohol expenditure
59	TCLTHEXP	Total expenditure on clothing and footwear
60	TFURNEXP	Total expenditure on furnishing
61	THLTHEXP	Total expenditure on health
62	TTRANSEXP	Total expenditure on transport
63	TCOMMEXP	Total expenditure on communications
64	TRECEXP	Total expenditure on recreation
65	TEDUEXP	Total expenditure on education
66	TRESTHOTEXP	Total expenditure on restaurants and hotel
67	TMISCEXP	Total miscellaneous expenditure
68	TANHHEXP	Total annual nominal household expenditures

## 10.2 Appendix B: Example of Blanket Agreement for SUF

Agreement between [providing agency] and [receiving agency] regarding the deposit and use of microdata

A. This agreement relates to the following microdatasets:

1. \_\_\_\_\_

2. \_\_\_\_\_
3. \_\_\_\_\_
4. \_\_\_\_\_
5. \_\_\_\_\_

B. Terms of the agreement:

As the owner of the copyright in the materials listed in section A, or as duly authorized by the owner of the copyright in the materials, the representative of [providing agency] grants the [receiving agency] permission for the datasets listed in section A to be used by [receiving agency] employees, subject to the following conditions:

1. Microdata (including subsets of the datasets) and copyrighted materials provided by the [providing agency] will not be redistributed or sold to other individuals, institutions or organisations without the [providing agency]'s written agreement. Non-copyrighted materials which do not contain microdata (such as survey questionnaires, manuals, codebooks, or data dictionaries) may be distributed without further authorization. The ownership of all materials provided by the [providing agency] remains with the [providing agency].
2. Data will be used for statistical and scientific research purposes only. They will be employed solely for reporting aggregated information, including modeling, and not for investigating specific individuals or organisations.
3. No attempt will be made to re-identify respondents, and there will be no use of the identity of any person or establishment discovered inadvertently. Any such discovery will be reported immediately to the [providing agency].
4. No attempt will be made to produce links between datasets provided by the [providing agency] or between [providing agency] data and other datasets that could identify individuals or organisations.
5. Any books, articles, conference papers, theses, dissertations, reports or other publications employing data obtained from the [providing agency] will cite the source, in line with the citation requirement provided with the dataset.
6. An electronic copy of all publications based on the requested data will be sent to the [providing agency].
7. The [providing agency] and the relevant funding agencies bear no responsibility the data's use or for interpretation or inferences based upon it.
8. An electronic copy of all publications based on the requested data will be sent to the [providing agency].
9. Data will be stored in a secure environment, with adequate access restrictions. The [providing agency] may at any time request information on the storage and dissemination facilities in place.
10. The [recipient agency] will provide an annual report on uses and users of the listed microdatasets to the [providing agency], with information on the number of researchers having accessed each dataset, and on the output of this research.
11. This access is granted for a period of [provide information on this period, or state that the agreement is open ended].

C. Communications:

The [receiving organisation] will appoint a contact person who

will act as unique focal person for this agreement. Should the focal person be replaced, the [recipient agency] will immediately communicate the name and coordinates of the new contact person to the [providing agency]. Communications for administrative and procedural purposes may be made by email, fax or letter as follows:

Communications made by [providing agency] to [recipient agency] will be directed to:

Name of contact person:

Title of contact person:

Address of the recipient agency:

Email:

Tel:

Fax:

Communications made by [recipient agency] to [depositor agency]  
will be directed to:

Name of contact person:

Title of contact person:

Address of the recipient agency:

Email:

Tel:

Fax:

### D. Signatories

The following signatories have read and agree with the Agreement as presented above:

#### Representative of the [providing agency]

Name \_\_\_\_\_

Signature \_\_\_\_\_ Date \_\_\_\_\_

#### Representative of the [recipient agency]

Name \_\_\_\_\_

Signature \_\_\_\_\_ Date \_\_\_\_\_

Source: *DuBo10*

## 10.3 Appendix C: Internal and External Reports for Case Studies

This appendix provides example of internal and external reports on the anonymization process for the case studies in Section 9.1. The internal report consists of two parts: the first is for the anonymization of the household-level variables and the second is for the anonymization of the individual-level variables.

## 10.4 Case study 1 - Internal report

**SDC report** (adapted from the report function in *sdcMicro*)

The dataset consists of 10,574 observations (i.e., 10,574 individuals in 2,000 households).

Household-level variables

Anonymization methods applied to household-level variables:

- Removing households of size larger than 13 (29 households)
- Local suppression to achieve 2-anonymity, with importance vector to prevent suppressing values of the variables HHSIZE, REGION and URBRUR

- Recoding the variable LANDSIZEHA: rounding to one digit for values smaller than 1, rounding to zero digits for other values, grouping values 5-19 and 20-40, topcoding at 40
- PRAMming the variables ROOF, TOILET, WATER, ELECTCON, FUELCOOK, OWNMOTORCYCLE, CAR, TV and LIVESTOCK
- Noise addition (level 0.01 and 0.05 for outliers) to the income and expenditure components, replacing aggregates by sum of perturbed components

Selected (key) variables:

- Modifications on categorical key variables: **TRUE**
- Modifications on continuous key variables: **TRUE**
- Modifications using PRAM: **TRUE**
- Local suppressions: **TRUE**

Disclosure risk (household-level variables):

Frequency analysis for categorical key variables:

Number of observations violating

2-Anonymity: 0 (unmodified data: 103)

3-Anonymity: 104 (unmodified data: 229)

5-Anonymity: 374 (unmodified data: 489)

Percentage of observations violating

2-Anonymity: 0% (unmodified data: 5.15%)

3-Anonymity: 5.28% (unmodified data: 11.45%)

5-Anonymity: 18.7% (unmodified data: 24.45%)

Disclosure risk categorical variables:

Expected Percentage of Re-identifications: 0.05161614% (~ 1.0 observations)

(unmodified data: 0.001820465% (~ 0.36 observations))

10 combinations of categories with highest risk:

	URBRUR	REGION	HHSIZE	OWNAGLAND	RELIG	fk	Fk
1	2	6	2	3	7	1	372.37
2	1	5	1	1	6	1	226.35
3	2	5	2	3	6	1	430.21
4	2	2	1	1	NA	1	173.05
5	2	6	1	1	5	1	80.05
6	1	6	1	3	5	1	343.27
7	2	5	1	2	NA	1	140.60
8	2	6	1	3	7	1	230.29
9	2	5	12	1	9	1	475.01
10	2	6	3	1	1	1	338.57

Disclosure risk continuous scaled variables:

Distance-based Disclosure Risk for Continuous Key Variables:

Disclosure Risk is between 0% and 100% in the modified data. In the original data, the risk is approximately 100%.

Data Utility (household-level variables):

**29 households have been removed due to their household sizes**

**Frequencies categorical key variables**

**URBRUR**

categories1	1	2	NA
orig	1316	684	0
categories2	1	2	NA
recoded	1299	666	6

**REGION**

categories1	1	2	3	4	5	6	NA
orig	324	334	371	375	260	336	0
categories2	1	2	3	4	5	6	NA
recoded	315	328	370	370	257	330	1

**HHSIZE**

categories1	1	2	3	4	5	6	7	8	9	10	11	12
orig	152	194	238	295	276	252	214	134	84	66	34	21
categories1	13	14	15	16	17	18	19	20	21	22	33	
orig	11	6	6	5	4	2	1	2	1	1	1	
categories2	1	2	3	4	5	6	7	8	9	10	11	12
recoded	152	194	238	295	276	252	214	134	84	66	34	21
categories2	13											
recoded	10											

**OWNAGLAND**

categories1	1	2	3	NA
orig	763	500	332	405
categories2	1	2	3	NA
recoded	735	482	310	444

**RELIG**

categories1	1	5	6	7	9	NA
orig	179	383	267	7	154	1010
categories2	1	5	6	7	9	NA
recoded	175	380	260	5	148	1003

Local suppressions

Number of local suppressions:

	URBRUR	REGION	HHSIZE	OWNAGLAND	RELIG
absolute	6	1	1	48	16
relative (in percent)	0.304%	0.051%	0.051%	2.435%	0.812%



Data utility of continuous scaled key variables:

#### Univariate summary:

	Min.	1st Qu	Median	Mean	3rd Qu	Max.
TANHHEX P	0	0,2	1	6,689	2,421	1214
TANHHEX P.m	0	0,2	1	3,427	2	40
TFOODEX P	498	15170	17090	24340	23260	353200
TFOODEX P.m	127,1	15100	17060	23410	22110	275300
TALCHEX P	0	8438	11890	12920	13070	127900
TALCHEX P.m	-209,7	8377	11880	12570	13030	124800
TCLTHEX P	0	0	0	401,7	0	85280
TCLTHEX P.m	-77,53	-13,59	6,42	404,7	30,69	85280
THOUSEX P	0	121	131	733,8	672,8	28400
THOUSEX P.m	-54,65	111,4	138,8	706,1	618,9	28410
TFURNEX P	0	1211	1340	2233	1970	197500
TFURNEX P.m	-39,54	1198	1340	2066	1933	49230
THLTHEX P	0	153,8	167	479,8	302	17780
THLTHEX P.m	-18,79	146,8	168,6	453,1	295,2	15720
TTRANSE XP	0	1	634	961	687	49650
TTANSEX P.m	-80,58	26,66	627,1	917,2	692,4	49640
TCOMMEX P	0	146	241	1158	434	91920
TCOMMEX P.m	-115,2	139,1	238,3	1104	403,2	91920
TRECEXP	0	3	95	577,2	107	34000
TRECEXP .m	-61,27	21,35	92,28	555,4	128,8	33960
TEDUEXP	0	0	0	123,7	0	15880
TEDUEXP .m	-29,23	-5,06	1,213	121,8	9,748	15860
TRESHOT EXP	0	154	722	2730	784	240300
TRESHOT EXP.m	-396,1	190,5	671,6	2568	872	240400
TMISCEX P	0	0	467	875,1	528	63700
TMISCEX P.m	-93,39	0,7588	442,7	860,7	531,9	63680
INCTOTG ROSSHH	0	444	1041	1148	1126	67420
INCTOTG ROSSHH. m	-24,92	446	1041	1087	1124	14940
INCRMT	5000	12400	13390	30840	24200	683900
INCRMT. m	4069	9071	17000	33040	36680	570000
INCWAGE	0	0	0	1276	0	300000
INCWAGE .m	-295,1	-46,95	20,93	1261	114,4	300100
INCFARM BSN	0	9262	12950	23460	14570	683900
INCFARM BSN.m	-1466	9336	12980	23420	14750	684000
INCNFAR MBSN	0	0	0	3809	3900	165400
INCNFAR MBSN.m	-232,4	-10,69	142,6	3415	3846	160100
INCRENT	0	0	827,5	9166	7307	400000
INCRENT .m	-757,4	43,89	783,7	8637	7267	394800
INCFIN	0	0	0	1783	0	120000
INCFIN. m	-248,5	-56,57	11,54	1608	90,27	120000
INCPENS N	0	0	0	74,58	0	14400
INCPENS N.m	-20,2	-4,591	0,1964	76,62	5,796	14380
INCOTHE R	0	0	0	331,3	0	60000
INCOTHE R.m	-123,3	-24,78	-0,0261 7	331,1	26,75	60050
LANDSIZ EHA	0	0	0	549,1	0	82300
LANDSIZ EHA.m	-126,2	-21,91	3,4	486,7	30,88	79670

### Information loss:

Criteria IL1: 0.01219892

### Individual-level variables

- Modifications on categorical key variables: **TRUE**
- Modifications on continuous key variables: **FALSE**
- Modifications using PRAM: **FALSE**
- Local suppressions: **TRUE**

Disclosure risk (individual-level variables):

Anonymization methods applied to individual-level variables:

- Recoding AGEYRS from months to years for age under 1, and to ten-year intervals for age values between 15 and 65, topcoding age at 65
- Local suppression to achieve 2-anonymity

Frequency analysis for categorical key variables:

Number of observations violating

2-Anonymity: 0 (unmodified data: 998)

3-Anonymity: 0 (unmodified data: 1384)

5-Anonymity: 935 (unmodified data: 2194)

Percentage of observations violating

2-Anonymity: 0% (unmodified data: 9.91%)

3-Anonymity: 0% (unmodified data: 13.75%)

5-Anonymity: 6.23% (unmodified data: 21.79%)

Disclosure risk categorical variables:

Expected Percentage of Reidentifications: 0.02% (~ 2.66 observations)

(unmodified data: 0.24% (~23.98 observations))

Expected Percentage of Reidentifications (hierarchical risk): 0.1% (~ 15.34 observations)

(unmodified data: 1.26 % (~ 127.12 observations))

10 combinations of categories with highest risk:

	GEN- DER	REL	MARI- TAL	AGEYRS	EDUCY	EDYR- SATCURRAT	ATSCHOOL	INDUS- TRY1	fk	Fk
1	1	1	3	38	6	NA	0	9	1	73. 31
2	1	1	3	20	1	NA	0	6	1	69. 53
3	1	1	2	39	2	NA	0	5	1	54. 63
4	1	1	1	36	6	NA	0	9	1	73. 31
5	1	1	3	42	2	NA	0	1	1	39. 58
6	0	1	6	74	1	NA	0	1	1	58. 12
7	0	1	6	34	2	NA	0	1	1	57. 40
8	1	1	1	26	4	NA	0	5	1	66. 21
9	1	1	4	35	1	NA	0	10	1	57. 13
10	1	6	1	12	1	NA	0	5	1	57. 13

Data utility (individual-level variables):

#### Frequencies categorical key variables

##### GENDER

categories1	0	1	NA
orig	5197	4871	0
categories2	0	1	NA
recoded	5197	4871	0

##### REL

categories1	1	2	3	4	5	6	7	8	9	NA
orig	1970	1319	4933	57	765	89	817	51	63	4
categories2	1	2	3	4	5	6	7	8	9	NA
recoded	1698	1319	4933	52	765	54	817	40	63	327

##### MARITAL

categories1	1	2	3	4	5	6	NA
orig	3542	2141	415	295	330	329	3016
categories2	1	2	3	4	5	6	NA
recoded	3542	2141	415	295	330	329	3016

##### AGEYRS

categories1	0	1/12	2/12	3/12	4/12	5/12	6/12	7/12	8/12	9/12
orig	178	8	1	14	15	19	17	21	18	7
categories1	10/12	11/12	1	2	3	4	5	6	7	8
orig	5	8	367	340	332	260	334	344	297	344
categories1	9	10	11	12	13	14	15	16	17	18
orig	281	336	297	326	299	263	243	231	196	224
categories1	19	20	21	22	23	24	25	26	27	28
orig	202	182	136	146	150	137	128	139	117	152
categories1	29	30	31	32	33	34	35	36	37	38
orig	111	143	96	123	104	107	148	91	109	87
categories1	39	40	41	42	43	44	45	46	47	48
orig	89	93	58	78	72	64	84	74	48	60
categories1	49	50	51	52	53	54	55	56	57	58
orig	58	66	50	55	29	30	34	38	33	44
categories1	59	60	61	62	63	64	65	66	67	68
orig	35	36	25	33	21	15	30	18	13	29
categories1	69	70	71	72	73	74	75	76	77	78
orig	26	36	17	16	12	3	16	10	8	18
categories1	79	80	81	82	83	84	85	86	87	88
orig	11	13	5	2	7	7	7	3	2	2
categories1	89	90	91	92	93	95	NA			
orig	4	4	3	1	1	1	188			
categories2	0	1	2	3	4	5	6	7	8	9
recoded	311	367	340	332	260	334	344	297	344	281
categories2	10	11	12	13	14	20	30	40	50	60
recoded	336	297	326	299	263	1847	1220	889	554	314
categories2	65	NA								
recoded	325	188								

## EDUCY

categories1	0	1	2	3	4	5	6	NA
orig	1582	4755	1062	330	139	46	104	2050
categories2	0	1	2	3	4	5	6	NA
recoded	1582	4755	1062	330	139	46	104	2050

## EDYRSATCURR

categories1	0	1	2	3	4	5	6	7	8	9
orig	177	482	445	446	354	352	289	266	132	127
categories1	10	11	12	13	15	16	18	NA		
orig	143	58	46	27	18	10	54	6642		
categories2	0	1	2	3	4	5	6	7	8	9
recode	177	482	445	446	354	352	289	266	132	127
categories2	10	11	12	13	15	16	18	NA		
recode	143	58	46	27	18	10	54	6642		

## ATSCHOOL

categories1	0	1	NA
orig	4696	3427	1945
categories2	0	1	NA
recoded	4696	3427	1945

**INDUSTRY1**

categories1	1	2	3	4	5	6	7	8	9	10	NA
orig	5300	16	153	2	93	484	95	17	70	292	3546
categories2	1	2	3	4	5	6	7	8	9	10	NA
recoded	5300	16	153	2	93	484	95	17	70	292	3546

Local suppressions

Number of local suppressions:

	GENDER	REL	MARITAL	AGEYRS	EDUCY
absolute	0	323	0	0	0
relative (in percent)	0	3.21%	0	0	0

	EDYRSATCURR	ATSCHOOL	INDUSTRY1		
absolute	0	0	0		
relative (in percent)	0	0	0		

## 10.5 Case study 1 - External report

This case study microdata set has been treated to protect confidentiality. Several methods have been applied to protect the confidentiality: removing variables from the original dataset, removing records from the dataset, reducing detail in variables by recoding and top-coding, removing particular values of individuals at risk (local suppression) and perturbing values of certain variables.

### *Removing variables*

The released microdata set has only a selected number of variables contained in the initial survey. Not all variables could be released in this SUF without breaching confidentiality rules.

### *Removing records*

To protect confidentiality, records of households larger than 13 were removed. Thirty households out of a total of 2,000 households in the dataset were removed.

### *Reducing detail in variables by recoding and top-coding*

The variable LANDSIZEHA was rounded to one digit for values smaller than 1, rounded to zero digits for other values, grouped for values 5-19 and 20-40 and topcoded at 40. The variable AGEYRS was recoded to ten-year age intervals for values in the age range 15 to 65.

### *Local suppression*

Values of certain variables for particular households and individuals were deleted. In total, six values of the variable URBRUR, one of the REGION variable, 48 for the OWNAGLAND variable, 16 for the RELIG variable and 323 values of the variable REL were deleted.

### *Perturbing values*

Uncertainty was introduced in the variables ROOF, TOILET, WATER, ELECTCON, FUELCOOK, OWNMOTORCYCLE, CAR, TV and LIVESTOCK by using the PRAM method. This method changes a certain percentage of values of variables within each variable. Here invariant PRAM was used, which guarantees that the univariate tabulations stay unchanged. Multivariate tabulations may be changed. Unfortunately, the transition matrix cannot be published.

The income and expenditure variables were perturbed by adding noise (adding small random values to the original values). The noise added was 0.01 times the standard deviation in the original data and 0.05 for outliers. Noise was added to the components and the aggregates were recomputed to guarantee that the proportions of the different components did not change.

## 10.6 Case study 2 - Internal report

**SDC report** (adapted from the report function in *sdcMicro*)

This report describes the anonymization measures for the PUF release additional to those already taken in the first case study. Therefore, this report should be read in conjunction with the internal report for case study 1. The original dataset consists of 10,574 observations (i.e., 10,574 individuals in 2,000 households). The dataset used for the anonymization of the PUF file is the anonymized SUF file from case study 1. This dataset consists of 10,068 observations in 1,970 households. The difference is due to the removal of large households and sensitive or identifying variables in the first case study.

Household-level variables

Anonymization methods applied to household-level variables:

- For SUF release (see case study 1):
  - Removing households of size larger than 13 (29 households)
  - Local suppression to achieve 2-anonymity, with importance vector to prevent suppressing values of the variables HHSIZE, REGION and URBRUR
- For PUF release:
  - Remove variables OWNLANDAG, RELIG and LANDSIZEHA
  - Local suppression to achieve 5-anonymity, with importance vector to prevent suppressing values of the variables HHSIZE and REGION
  - PRAMming the variables ROOF, TOILET, WATER, ELECTCON, FUELCOOK, OWNMOTORCYCLE, CAR, TV and LIVESTOCK
  - Create deciles for aggregate income and expenditure (TANNEXP and INCTOTGROSSHH) and replace the actual values with the mean of the corresponding decile. Replace income and expenditure components with the proportion of original totals.

Selected (key) variables:

<i>categorical</i>	URBRUR	REGION	HHSIZE
<i>continuous</i>	TANHHEXP	INCTOTGROSSHH	
<i>weight</i>	WGTPOP		
<i>hhID</i>	not defined		
<i>strata</i>	not defined		

- Modifications on categorical key variables: **TRUE**
- Modifications on continuous key variables: **TRUE**
- Modifications using PRAM: **TRUE**

- Local suppressions: **TRUE**

Disclosure risk (household-level variables):

Frequency analysis for categorical key variables:

Number of observations violating

2-Anonymity: 0 (PUF file: 0, unmodified data: 103)

3-Anonymity: 0 (PUF file: 18, unmodified data: 229)

5-Anonymity: 0 (PUF file: 92, unmodified data: 489)

Percentage of observations violating

2-Anonymity: 0.00% (PUF file: 0.00%, unmodified data: 5.15%)

3-Anonymity: 0.00% (PUF file: 0.91%, unmodified data: 11.45%)

5-Anonymity: 0.00% (PUF file: 4.67%, unmodified data: 24.45%)

Disclosure risk categorical variables:

Expected Percentage of Re-identifications: 0.0000526% (~ 0.10 observations),

PUF file: 0.0000642% (~ 0.13 observations), unmodified data: 0.001820465% (~ 0.36 observations)

11 combinations of categories with highest risk in PUF file:

	URBRUR	REGION	HHSIZE	fk	Fk
1	2	4	1	7	1152.084
2	2	4	1	7	1152.084
3	2	2	9	2	2356.926
4	2	4	1	7	1152.084
5	2	4	1	7	1152.084
6	2	4	1	7	1152.084
7	2	5	12	2	2978.454
8	2	4	1	7	1152.084
9	2	4	1	7	1152.084
10	2	5	12	2	2978.454
11	2	2	9	2	2356.926

Disclosure risk continuous scaled variables:

Distance-based Disclosure Risk for Continuous Key Variables:

Disclosure Risk is between 0% and 100% in the modified data. In the original data, the risk is approximately 100%.

Data Utility (household-level variables):

**Frequencies categorical key variables**

**URBRUR**

categories1	1	2	NA
orig	1316	684	0
categories2	1	2	NA
recoded	1280	623	67

**REGION**

categories1	1	2	3	4	5	6	NA
orig	324	334	371	375	260	336	0
categories2	1	2	3	4	5	6	NA
recoded	311	325	369	370	253	329	13

## HHSIZE

categories1	1	2	3	4	5	6	7	8	9	10	11	12
orig	152	194	238	295	276	252	214	134	84	66	34	21
categories1	13	14	15	16	17	18	19	20	21	22	33	
orig	11	6	6	5	4	2	1	2	1	1	1	
categories2	1	2	3	4	5	6	7	8	9	10	11	12
recoded	152	194	238	295	276	252	214	134	84	66	34	21
categories2	13											
recoded	10											

Local suppressions

Number of local suppressions:

	URBRUR	REGION	HHSIZE
absolute	61	125	0
relative (in percent)	3.096%	0.609%	0.000%

Data utility of continuous scaled key variables:

### Univariate summary:

	Min.	1st Qu	Median	Mean	3rd Qu	Max.
TANHHEX P	498	15,170	17,090	24,340	23,260	353,230
TANHHEX P.m	827	14,700	17,060	23,420	22,750	83,963
INCTOTG ROSSHH	5,000	12,400	13,390	30,840	24,200	683,900
INCTOTG ROSSHH. m	6353	12,390	13,400	30,250	24,240	149,561

### Information loss:

Criteria IL1: 0.2422625

Disclosure risk (individual-level variables):

Anonymization methods applied to individual-level variables:

- For SUF release (see case study 1):
  - Recoding AGEYRS from months to years for age under 1, and to ten-year intervals for age values between 15 and 65, topcoding age at 65
  - Local suppression to achieve 2-anonymity
- For PUF release:
  - Remove variable EDYRSCURRAT
  - Recode REL to ‘Head’, ‘Spouse’, ‘Child’, ‘Other relative’, ‘Other’
  - Recode MARITAL to ‘Never married’, ‘Married/Living together’, ‘Divorced/Separated/Widowed’



- Recode AGEYRS for values under 15 to 7
- Recode EDUCY to ‘No education’, ‘Pre-school/ Primary not completed’, ‘Completed lower secondary or higher’
- Recode INDUSTRY1 to ‘Primary sector’, ‘Secondary sector’, ‘Tertiary sector’

Frequency analysis for categorical key variables:

Number of observations violating

2-Anonymity: 0 (PUF file: 0, unmodified data: 998)

3-Anonymity: 0 (PUF file: 167, unmodified data: 1384)

5-Anonymity: 0 (PUF file: 463, unmodified data: 2194)

Percentage of observations violating

2-Anonymity: 0.00% (PUF file: 0.00%, unmodified data: 9.91%)

3-Anonymity: 0.00% (PUF file: 1.66%, unmodified data: 13.75%)

5-Anonymity: 0.00% (PUF file: 4.60%, unmodified data: 21.79%)

Disclosure risk categorical variables:

Expected Percentage of Re-identifications: 0.00% (~0.41 observations)

(PUF file: 0.02 % (~ 1.69 observations), unmodified data: 0.24% (~23.98 observations))

Expected Percentage of Re-identifications (hierarchical risk): 0.02% (~2.29 observations)

(PUF file: 0.10 % (~ 9.57 observations), unmodified data: 1.26 % (~ 127.12 observations))

10 combinations of categories with highest risk:

	GEDNER	REL	MARIT AL	AGEYR S	EDUCY	INDUS TRY1	fk	Fk
1	1	1	2	50	1	7	2	324.9 275
2	0	1	3	40	3	6	2	330.0 521
3	0	1	6	60	0	3	2	350.5 000
4	0	1	3	40	3	6	2	330.0 521
5	1	1	2	30	4	5	2	253.7 431
6	1	1	2	50	1	7	2	324.9 275
7	0	1	6	50	1	6	2	255.6 142
8	1	1	4	40	1	10	2	175.0 797
9	1	1	4	40	1	10	2	175.0 797
10	1	1	3	30	1	6	2	323.4 879

Data utility (individual-level variables):

**Frequencies categorical key variables**

**GENDER**

categories1	0	1	NA
orig	5,197	4,871	0
categories2	0	1	NA
recoded	5,197	4,871	0

**REL**

categories1	1	2	3	4	5	6	7	8	9	NA
orig	1,970	1,319	4,933	57	765	89	817	51	63	4
categories2	1	2	3	7	9	NA				
recoded	1,698	1,319	4,933	1,688	103	327				

**MARITAL**

categories1	1	2	3	4	5	6	NA
orig	3,542	2,141	415	295	330	329	3,016
categories2	1	2	9	NA			
recoded	3,542	2,851	659	3,016			

**AGEYRS**

categories1	0	1/12	2/12	3/12	4/12	5/12	6/12	7/12	8/12	9/12
orig	178	8	1	14	15	19	17	21	18	7
categories1	10/12	11/12	1	2	3	4	5	6	7	8
orig	5	8	367	340	332	260	334	344	297	344
categories1	9	10	11	12	13	14	15	16	17	18
orig	281	336	297	326	299	263	243	231	196	224
categories1	19	20	21	22	23	24	25	26	27	28
orig	202	182	136	146	150	137	128	139	117	152
categories1	29	30	31	32	33	34	35	36	37	38
orig	111	143	96	123	104	107	148	91	109	87
categories1	39	40	41	42	43	44	45	46	47	48
orig	89	93	58	78	72	64	84	74	48	60
categories1	49	50	51	52	53	54	55	56	57	58
orig	58	66	50	55	29	30	34	38	33	44
categories1	59	60	61	62	63	64	65	66	67	68
orig	35	36	25	33	21	15	30	18	13	29
categories1	69	70	71	72	73	74	75	76	77	78
orig	26	36	17	16	12	3	16	10	8	18
categories1	79	80	81	82	83	84	85	86	87	88
orig	11	13	5	2	7	7	7	3	2	2
categories1	89	90	91	92	93	95	NA			
orig	4	4	3	1	1	1	188			
categories2	7	20	30	40	50	60	65	NA		
recoded	4,731	1,847	1,220	889	554	314	325	188		

**EDUCY**

categories1	0	1	2	3	4	5	6	NA
orig	1582	4755	1062	330	139	46	104	2050
categories2	0	1	2	3	4	5	6	NA
recoded	1,582	4,755	1,062	330	139	46	104	2,050

**INDUSTRY1**

categories1	1	2	3	4	5	6	7	8	9	10	NA
orig	5,300	16	153	2	93	484	95	17	70	292	3,546
categories2	1	2	3	NA							
recoded	5,316	248	958	3,546							

Local suppressions

Number of local suppressions:

	GENDER	REL	MARITAL	AGEYRS	EDUCY	INDUSTRY1
absolute	0	0	0	91	0	0
relative (in percent )	0.00%	0.00%	0.00%	0.90%	0.00%	0.00%

## 10.7 Case study 2- External report

This case study microdata set has been treated to protect confidentiality. Several methods have been applied to protect the confidentiality: removing variables from the original dataset, removing records from the dataset, reducing detail in variables by recoding and top-coding, removing particular values of individuals at risk (local suppression) and perturbing values of certain variables.

### *Removing variables*

The released microdata set has only a selected number of variables contained in the initial survey. Not all variables could be released in this PUF without breaching confidentiality rules.

### *Removing records*

To protect confidentiality, records of households larger than 13 were removed. Twenty-nine households out of a total of 2,000 households in the dataset were removed.

### *Reducing detail in variables by recoding and top-coding*

The variable AGEYRS was recoded to ten-year age intervals for values in the age range 15 to 65 and bottom- and top-coded at 15 and 65. The variables REL, MARITAL, EDUCY and INDUSTRY1 were recoded to less detailed categories. The total income and expenditure variables were recoded to the mean of the corresponding deciles and the income and expenditure components to the proportion of the totals.

### *Local suppression*

Values of certain variables for particular households and individuals were deleted. In total, 67 values of the variable URBRUR, 126 of the REGION variable, 91 for the AGEYRS variable and 323 values of the variable REL were deleted.

### *Perturbing values*

Uncertainty was introduced in the variables ROOF, TOILET, WATER, ELECTCON, FUELCOOK, OWNMOTORCYCLE, CAR, TV and LIVESTOCK by using the PRAM method. This method changes a certain percentage of values of variables within each variable. Here invariant PRAM was used, which guarantees that the univariate tabulations stay unchanged. Multivariate tabulations may be changed. Unfortunately, the transition matrix cannot be published.

## 10.8 Appendix D: Execution Times for Multiple Scenarios Tested using Selected Sample Data

Datasets	1	2	3	4	5	6
Total Observations	40,449	102,429	20,261	34,785	23,022	26,411
Quasi-identifiers	Number of categories*					
RURAL/URBAN	2	2	2	2	2	2
REGION	8	2	10	5	4	10
HOUSEHOLD SIZE	25	26	18	17	20	18
OWNHOUSE	2	NA	2	2	3	2
SEX	2	2	2	2	2	2
RELATION TO HEAD OF HOUSEHOLD	8	7	8	7	9	7
MARITAL STATUS	5	6	6	6	6	6
AGE IN COMPLETED YEARS	99	98	97	98	109	99
ETHNICITY	NA	NA	NA	NA	10	20
RELIGION	NA	NA	NA	4	10	11
HIGHEST LEVEL OF EDUC COMPLETED	7	7	9	6	8	8
YEARS OF EDUCATION - ENROLLED	16	17	17	18	14	NA
EMPLOYMENT TYPE	7	7	6	5	6	7
LITERACY	2	NA	2	3	2	2
INDUSTRY CODE	10	464	10	11	10	10
CURRENTLY ENROLLED AT SCHOOL	2	2	2	2	2	2
Processing Platform	Server: 16 Processors	Desktop: 2 Processors	Desktop : 2 Processors	Desktop: 2 Processors	Server: 16 Processors	Laptop: 4 Processors
Total execution time with parallelization (hours)	6.11	10.75	23.70	67.58	77.14	282.58
Scenarios						
base	0	0	0	0	0	0
1	0.35	0.72	0.15	0.31	0.36	0.30
2	0.37	0.75	0.16	0.32	0.37	0.31
3	0.38	0.78	0.16	0.33	0.37	0.31
4	0.39	0.81	0.16	0.33	0.38	0.32
5	4.14	6.46	2.27	6.07	34.85	163.59
6	1.34	1.37	1.02	2.36	22.58	131.91
7	0.82	0.51	0.96	2.34	16.28	70.63
8	0.88	0.63	0.82	2.17	19.88	130.50
9	2.99	2.75	1.77	5.01	43.21	177.47
10	1.51	0.80	1.23	4.06	17.99	84.16
11	1.52	0.61	1.11	4.19	24.39	114.03
12	4.03	6.45	2.25	5.94	31.38	163.42
13	1.36	1.34	0.99	2.33	23.47	129.24
14	0.78	0.54	0.96	2.31	15.83	74.24
15	0.90	0.67	0.83	2.13	21.59	123.70
16	2.82	2.84	1.74	5.00	40.50	179.26
17	1.61	0.79	1.21	4.04	19.00	82.80
18	1.48	0.62	1.11	4.31	24.11	118.77
19	4.06	6.51	2.21	6.07	31.79	159.04
20	1.25	1.43	0.96	2.35	21.36	130.83
21	0.78	0.54	0.94	2.36	16.44	69.55
22	0.91	0.68	0.80	2.18	22.74	131.49
23	2.82	2.86	1.73	4.98	42.32	168.45
24	1.49	0.81	1.20	4.04	19.37	84.13
25	1.40	0.64	1.08	4.26	24.54	113.05
26	3.94	6.36	2.19	5.99	31.00	145.07
27	1.36	1.43	0.96	2.31	23.18	107.94
28	0.79	0.54	0.93	2.30	16.96	73.40
29	0.90	0.70	0.81	2.13	22.25	118.46
30	2.86	2.80	1.79	5.05	40.63	137.30
31	1.45	0.80	1.24	4.11	19.01	80.12
32	1.60	0.65	1.14	4.39	23.87	95.11
33	3.92	6.40	2.31	6.13	29.99	110.22
34	1.27	1.43	1.02	2.35	22.61	97.67
35	0.80	0.54	0.99	2.32	16.32	72.85
36	0.89	0.70	0.85	2.14	19.68	101.54
37	2.57	2.82	1.92	4.99	37.70	104.72
38	1.34	0.80	1.35	4.03	18.64	70.92
39	1.42	0.65	1.17	4.17	23.10	70.46
Approximated execution time if sequentially processed (hours)	65.48	69.52	46.46	132.20	880.06	3987.25

\*NA indicates that the variable is not available in that dataset and, hence, could not be selected as quasi-identifier.

Description of Scenarios				
Scenario	Recode	Local suppression	Expenditure variables	PRAM
base	-	-	-	-
1	-	k=3	-	PRAM on roof, water, toilet
2	-	k=3	add noise to expenditure variables	PRAM on roof, water, toilet
3	-	k=3	shuffle expenditure variables	PRAM on roof, water, toilet
4	topcode age	k=3	add noise to aggregate expenditures	-
5	recode age	k=3	-	-
6	recode age	k=3	-	-
7	recode age	k=3, importance age	-	-
8	recode age	k=3, importance ethnicity and religion	-	-
9	recode age	k=5	-	-
10	recode age	k=5, importance age	-	-
11	recode age	k=5, importance ethnicity and religion	-	-
12	topcode age	k=3	-	PRAM on roof, water, toilet
13	topcode age	k=3	-	PRAM on roof, water, toilet
14	recode age	k=3, importance age	-	PRAM on roof, water, toilet
15	recode age	k=3, importance ethnicity and religion	-	PRAM on roof, water, toilet
16	recode age	k=5	-	PRAM on roof, water, toilet
17	recode age	k=5, importance age	-	PRAM on roof, water, toilet
18	recode age	k=5, importance ethnicity and religion	-	PRAM on roof, water, toilet
19	recode age	k=3	add noise to expenditure variables	PRAM on roof, water, toilet
20	recode age	k=3	add noise to expenditure variables	PRAM on roof, water, toilet
21	recode age	k=3, importance age	add noise to expenditure variables	PRAM on roof, water, toilet
22	recode age	k=3, importance ethnicity and religion	add noise to expenditure variables	PRAM on roof, water, toilet
23	recode age	k=5	add noise to expenditure variables	PRAM on roof, water, toilet
24	recode age	k=5, importance age	add noise to expenditure variables	PRAM on roof, water, toilet
25	recode age	k=5, importance ethnicity and religion	add noise to expenditure variables	PRAM on roof, water, toilet
26	topcode age	k=3	shuffle expenditure variables	PRAM on roof, water, toilet
27	recode age	k=3	shuffle expenditure variables	PRAM on roof, water, toilet
28	recode age	k=3, importance age	shuffle expenditure variables	PRAM on roof, water, toilet
29	recode age	k=3, importance ethnicity and religion	shuffle expenditure variables	PRAM on roof, water, toilet
30	recode age	k=5	shuffle expenditure variables	PRAM on roof, water, toilet
31	recode age	k=5, importance age	shuffle expenditure variables	PRAM on roof, water, toilet
32	recode age	k=5, importance ethnicity and religion	shuffle expenditure variables	PRAM on roof, water, toilet
33	topcode age	k=3	add noise to aggregate expenditures	PRAM on roof, water, toilet
34	recode age	k=3	add noise to aggregate expenditures	PRAM on roof, water, toilet
35	recode age	k=3, importance age	add noise to aggregate expenditures	PRAM on roof, water, toilet
36	recode age	k=3, importance ethnicity and religion	add noise to aggregate expenditures	PRAM on roof, water, toilet
37	recode age	k=5	add noise to aggregate expenditures	PRAM on roof, water, toilet
38	recode age	k=5, importance age	add noise to aggregate expenditures	PRAM on roof, water, toilet
39	recode age	k=5, importance ethnicity and religion	add noise to aggregate expenditures	PRAM on roof, water, toilet

Fig. 10.1: Description of anonymization scenarios

## **References**

## CHAPTER 11

---

### Acknowledgements

---

Authors of this guide:

- Thijs Benschop, Consultant The World Bank
- Cathrine Machingauta, The World Bank
- Matthew Welch, The World Bank

Correspondence: Matthew Welch, [mwelch@worldbank.org](mailto:mwelch@worldbank.org)

Acknowledgments: The authors thank Olivier Dupriez (The World Bank) for his technical comments and inputs throughout the process.

The production of this guide was made possible through a World Bank Knowledge for Change II Grant: KCP II - A microdata dissemination challenge: Balancing data protection and data utility. Grant number: TF 015043, Project Number P094376. As well as from United Kingdom - DFID funding to the World Bank Multi-Donor Trust Fund - International Household Survey and Accelerated Data Program – TF071804/TF011722.





---

## Bibliography

---

- [DuBo10] Dupriez, O., & Boyko, E. (2010). **Dissemination of Microdata Files; Principles, Procedures and Practices**. International Household Survey Network (IHSN).
- [HDFG12] Hundepool, A., Domingo-Ferrer, J., Franconi, L., Giessing, S., Nordholt, E. S., Spicer, K., et al. (2012). **Statistical Disclosure Control**. Chichester, UK: John Wiley & Sons Ltd.
- [DuBo10] Dupriez, O., & Boyko, E. (2010). **Dissemination of Microdata Files; Principles, Procedures and Practices**. International Household Survey Network (IHSN).
- [DoTo03] Domingo-Ferrer, J., & Torra, V. (2003). **Disclosure Risk Assessment in Statistical Microdata Protection via Advanced Record Linkage**. *Statistics and Computing* 13 (4), 343-354
- [ElMa03] Elliot, M. J., & Manning, A. M. (2003). **Using DIS to Modify the Classification of Special Uniques**. Invited Paper. Joint ECE/Eurostat Work Session on Statistical Data Confidentiality. Luxembourg 2-9 April 2003.
- [ElMF02] Elliot, M. J., Manning, A. M., & Ford, R. W. (2002). **A Computational Algorithm for Handling the Special Uniques Problem**. *International Journal of Uncertainty, Fuzziness and Knowledge Based Systems*, 10 (5), 493-509.
- [ELMP10] Elliot, M. J., Lomax, S., Mackey, E. & Purdam, K. (2010) **Data Environment Analysis and the Key Variable Mapping System**. In *Privacy in Databases, PSD 2010* (ed. Domingo-Ferrer, J. & Magkos, E.), vol. 6344 of *Lecture Notes in Computer Science*, pp. 138-145. Berlin/Heidelberg: Springer
- [ElSD98] Elliot, M.J., Skinner, C.J. & Dale, A. (1998) **Special Uniques, Random Uniques, and Sticky Populations: Some Counterintuitive Effects of Geographical Detail on Disclosure Risk** *Research in Official Statistics* 1(2), pp. 53-67
- [EMMG05] Elliot, M. J., Manning, A., Mayes, K., Gurd, J., & Bane, M. (2005). **SUDA: A Program for Detecting Special Uniques**. Joint UNECE/Eurostat Work Session on Statistical Data Confidentiality. Geneva.
- [HDFG12] Hundepool, A., Domingo-Ferrer, J., Franconi, L., Giessing, S., Nordholt, E. S., Spicer, K., et al. (2012). **Statistical Disclosure Control**. Chichester, UK: John Wiley & Sons Ltd.
- [Lamb93] Lambert, D. (1993). **Measures of Disclosure Risk and Harm**. *Journal of Official Statistics*, 9 (2), 313-331.
- [MaHK08] Manning, A. M., Haglin, D. J., & Keane, J. A. (2008). **A Recursive Search Algorithm for Statistical Disclosure Assessment**. *Data Mining and Knowledge Discovery*, 16 (2), 165-196.

- [MKGv07] Machanavajjhala, A., Kifer, D., Gehrke, J., & Venkitasubramaniam, M. (2007). **L-diversity: Privacy Beyond K-anonymity**. *ACM Trans. Knowl. Discov. Data*, 1 (Article 3) (1556-4681).
- [TeMe08] Templ, M. & Meindl, B. (2008) **Robust Statistics Meets SDC: New Disclosure Risk Measures for Continuous Microdata Masking**. In *Privacy in Statistical Databases, PSD 2008* (eds. Domingo-Ferrer J. and Saygin Y.), vol. 5262 of *Lecture Notes in Computer Science*, pp. 177-189. Berlin/Heidelberg: Springer.
- [TeMK14] Templ, M., Meindl, B., & Kowarik, A. (2014, August). **Tutorial for SDCMicroGUI**. Retrieved from International Household Survey Network (IHSN): <http://www.ihsn.org/home/software/disclosure-control-toolbox>
- [TMKC14] Templ, M., Meindl, B., Kowarik, A., & Chen, S. (2014, August 1). **Introduction to Statistical Disclosure Control (SDC)**. Retrieved November 13, 2014, from <http://www.ihsn.org/home/software/disclosure-control-toolbox>.
- [BCRZ13] Burgert, C. R., Colston, J., Roy, T., & Zachary, B. (2013). **Geographic Displacement Procedure and Georeferenced Data Release Policy for the Demographic and Health Surveys**. DHS Spatial Analysis Report No. 7.
- [Bran02] Brand, R. (2002). **Microdata Protection through Noise Addition**. In J. Domingo-Ferrer (Ed.), *Inference Control in Statistical Databases - From Theory to Practice* (Vol. *Lecture Notes in Computer Science Series Volume 2316*, pp. 97-116). Berlin Heidelberg, Germany: Springer.
- [DMOT02] Domingo-Ferrer, J., Mateo-Sanz, J.M., Oganian, A. & Torres, A. **On the Security of Microaggregation with Individual Ranking: Analytics Attacks**. *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems* 10(5), pp. 477-492.
- [DoTo01a] Domingo-Ferrer, J., & Torra, V. (2001). **A Quantitative Comparison of Disclosure Control Methods for Microdata**. In P. Doyle, J. Lane, J. Theeuwes, & L. Zayatz (Eds.), *Confidentiality, Disclosure and Data Access: Theory and Practical Applications for Statistical Agencies* (pp. 111-133). Amsterdam, North-Holland: Elsevier Science.
- [DoTo05] Domingo-Ferrer, J., & Torra, V. (2005). **Ordinal, Continuous and Heterogeneous :math:'k'-anonymity through Microaggregation** *Data Mining and Knowledge Discovery* 11(2), pp. 195-212.
- [Drec11] Drechsler, J. (2011). **Synthetic Datasets for Statistical Disclosure Control**. Heidelberg/Berlin: Springer.
- [HDFG12] Hundepool, A., Domingo-Ferrer, J., Franconi, L., Giessing, S., Nordholt, E. S., Spicer, K., et al. (2012). **Statistical Disclosure Control**. Chichester, UK: John Wiley & Sons Ltd.
- [HuDr15] Hu, J., & Drechsler, J. (2015). **Generating synthetic geocoding information for public release**. NTTS - Conferences on New Techniques and Technologies for Statistics. Brussels.
- [KiWi03] Kim, J. J., & Winkler, W. W. (2003, April 17). **Multiplicative Noise for Masking Continuous Data**. Research Report Series.
- [MuSa06] Muralidhar, K., & Sarathy, R. (2006). **Data Shuffling- A New Masking Approach for Numerical Data**. *Management Science*, 658-670.
- [TeMK14] Templ, M., Meindl, B., & Kowarik, A. (2014, August). **Tutorial for SDCMicroGUI**. Retrieved from International Household Survey Network (IHSN): <http://www.ihsn.org/home/software/disclosure-control-toolbox>
- [TMKC14] Templ, M., Meindl, B., Kowarik, A., & Chen, S. (2014, August 1). **Introduction to Statistical Disclosure Control (SDC)**. Retrieved July 9, 2018, from <http://www.ihsn.org/home/software/disclosure-control-toolbox>.
- [Wolf15] de Wolf, P.-P. (2015). **Public Use Files of EU-SILC and EU-LFS data**.

- [DoTo01b] Domingo-Ferrer, J., & Torra, V. (2001). **Disclosure Protection Methods and Information Loss for Microdata**. In P. Doyle, J. Lane, J. Theeuwes, & Z. L., Theory and Practical Applications for Statistical Agencies (pp. 91-110). Amsterdam.
- [HDFG12] Hundepool, A., Domingo-Ferrer, J., Franconi, L., Giessing, S., Nordholt, E. S., Spicer, K., et al. (2012). **Statistical Disclosure Control**. Chichester, UK: John Wiley & Sons Ltd.
- [TMKC14] Templ, M., Meindl, B., Kowarik, A., & Chen, S. (2014, August 1). **Introduction to Statistical Disclosure Control (SDC)**. Retrieved July 9, 2018, from <http://www.ihsn.org/home/software/disclosure-control-toolbox>.
- [YaWC02] Yancey, W. W., Winkler, W. E., & Creecy, R. H. (2002). **Disclosure Risk Assessment in Perturbative Microdata Protection**. Research Report Series , Statistics 2002-01.
- [HDFG12] Hundepool, A., Domingo-Ferrer, J., Franconi, L., Giessing, S., Nordholt, E. S., Spicer, K., et al. (2012). **Statistical Disclosure Control**. Chichester, UK: John Wiley & Sons Ltd.
- [DuBo10] Dupriez, O., & Boyko, E. (2010). **Dissemination of Microdata Files; Principles, Procedures and Practices**. International Household Survey Network (IHSN).