

# Data Anonymization for Open Science

useR! 2024

Jiří Novák<sup>1,2</sup> Marko Miletic<sup>3</sup> Oscar Thees<sup>2</sup> Alžběta Beranová<sup>4</sup>

<sup>1</sup>University of Zurich <sup>2</sup>University of Applied Sciences Northwestern Switzerland

<sup>3</sup>Bern University of Applied Sciences <sup>4</sup>Czech Statistical Office

July 8, 2024



Jiří Novák CC BY-NC-ND (2024)



This license enables reusers to copy and distribute the material in any medium or format in unadapted form only, for noncommercial purposes only, and only so long as attribution is given to the creator.

# About Speakers

## **Jiří Novák**

- ▶ Ph.D. in Statistics with topic Simulation of Synthetic Microdata
- ▶ Statistical confidentiality of Czech Census 2021

# About Speakers

## **Jiří Novák**

- ▶ Ph.D. in Statistics with topic Simulation of Synthetic Microdata
- ▶ Statistical confidentiality of Czech Census 2021

## **Marko Miletić**

- ▶ XXX
- ▶ XXX

# About Speakers

## Jiří Novák

- ▶ Ph.D. in Statistics with topic Simulation of Synthetic Microdata
- ▶ Statistical confidentiality of Czech Census 2021

## Marko Miletić

- ▶ XXX
- ▶ XXX

## Oscar Thees

- ▶ XXX
- ▶ XXX

# About Speakers

## Jiří Novák

- ▶ Ph.D. in Statistics with topic Simulation of Synthetic Microdata
- ▶ Statistical confidentiality of Czech Census 2021

## Marko Miletić

- ▶ XXX
- ▶ XXX

## Oscar Thees

- ▶ XXX
- ▶ XXX

## Alžběta Beranová

- ▶ XXX
- ▶ XXX

Statistical Disclosure Control (SDC) or Statistical Disclosure Limitation (SDL) seeks to protect statistical data in such a way that they can be released without giving away confidential information that can be linked to specific individuals or entities.

# Why is confidentiality protection important?

There are three main reasons:

1. **Principle** It is a fundamental principle for Official Statistics that the statistical records of individual persons, businesses or events used to produce Official Statistics are strictly confidential, and are to be used only for statistical purposes.
2. **Legal** Legislation places a legal obligation to protect individual business and personal data. Legal frameworks regulate what is allowed and what is not allowed with regard to publication of private information.
3. **Quality** Respondents need confidence in the preservation of the confidentiality of individual information.
4. **Ethical:** It is unethical to disclose information that can be linked to specific individuals or entities.



- ▶ Legal frameworks regulate what is allowed and what is not allowed with regard to publication of private information.
- ▶ Before sensitive statistical databases can be made available to universities for research, confidentiality must be guaranteed.
- ▶ Users of statistical outputs should be aware of the reasoning and methodology behind statistical disclosure control.

Different outputs require different approaches to SDC and different mixtures of tools.

- ▶ Macrodata (Tabular data)
- ▶ Microdata
- ▶ Dynamic databases
- ▶ Statistical analyses

A disclosure occurs when a person or an organisation recognises or learns something that they did not know already about another person or organisation, via released data.

Types of disclosure risk:

- (1) Identity disclosure and
- (2) Attribute disclosure.
- (3) Inferential disclosure

SDC seeks to optimise the trade-off between the disclosure risk and the utility of the protected released data.

- ▶ Risk: the probability of a disclosure event occurring.
- ▶ Utility: the usefulness of the data for the intended purpose.
- ▶ The goal is to find a balance between risk and utility.

# Risk-utility trade-off

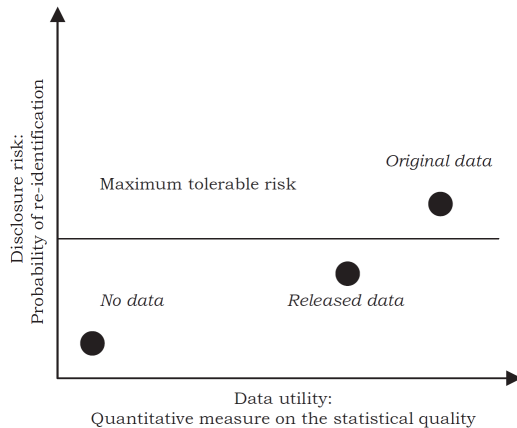


Figure 1: R-U confidentiality map (Duncan et al., 2001)

A unit is at risk of disclosure when it cannot be confused with several other units

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

# Variables

1. Identifiers
2. Quasi-identifiers or key variables
3. Confidential outcome variables
4. Non-confidential outcome variables



# Disclosure control methods

1. Masking original data
  - i. Non-perturbative masking.
  - ii. Perturbative masking.
2. Generating synthetic data

# Non-perturbation methods

Non-perturbative masking does not rely on distortion of the original data but on partial suppressions or reductions of detail.

Table 1: Non-perturbative methods vs. data types.

Method	Continuous data	Categorical data
Sampling		X
Global recoding	X	X
Top and bottom coding	X	X
Local suppression		X

- ▶ sdcMicro is an R package for statistical disclosure control.
- ▶ <https://cran.r-project.org/web/packages/sdcMicro/index.html>
- ▶ <https://github.com/sdcTools/sdcMicro>

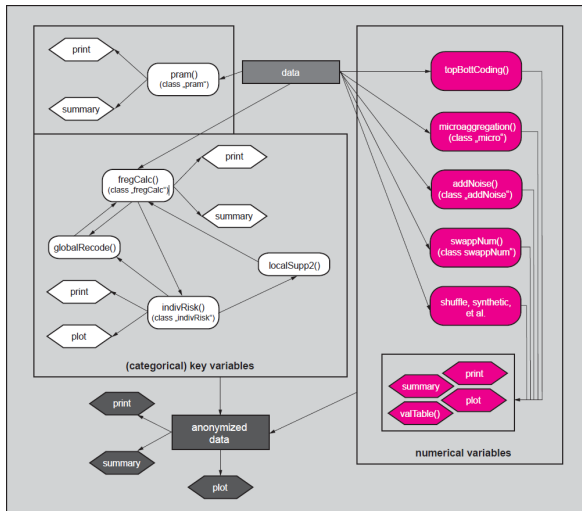


Figure 2: Certain procedures in package *sdcMicro* and their relationship

▶ ▶ sdcMicro

# Synthetic methods

- ▶ ▶ synthpop
- ▶ ▶ simPop
- ▶ ▶ GANs