



DDI Data Description Statistics Protection Software

Sebastian Kočar
Social Science Data Archives
University of Ljubljana

Privacy in Statistical Databases 2014,
Eivissa, September 18th 2014

The content of the presentation

- What DDI Data Description Statistics Protection Software **is**?
- What DDI Data Description Statistics Protection Software **is not**?
- The **background** of developing the software
- The **purpose** of developing the software
- DDI Data Description Statistics Protection Software **features**
- Quick **demonstration** of using the tool
- Discussion about the official statistics univariate statistics protection

What is DDI Data Description Statistics Protection Software?

- It is a relatively simple tool for protecting **aggregated data**
- It is a tool which works with **XML files**
- It protects aggregated data by automatically changing XML code
- It protects **DDI** Data Description displayed **univariate statistics**
- It enables various kinds of data protection – using techniques similar to microdata and tabular data protection ones




What DDI Data Description Statistics Protection Software is not?

- It is **not a microdata protection** tool (such as sdcMicro, μ -ARGUS)
- Datasets could be imported, but only to calculate certain univariate statistics
- It is **not a tabular data protection** tool (such as τ -ARGUS)
- It does not work with bivariate categorical data
- It is not meant to be used for protection of regularly published official statistics research results

The background / purpose of developing the software

- **Data without Boundaries** project's initiatives
- **Collaboration** between Social Science Data Archives and National Statistics Institutes
- Promotion of using official statistics data for **scientific research** purposes
- **Distribution** of aggregated official statistics data (with no restrictions)
- **Protection** of univariate statistics to minimize the disclosure risk

XML DDI Data Description Display



Social Science Data Archives

ads

Data
Content fields
Series
Catalogues

User support
About data
Training
Depositing
About arch

Labour force survey, 2011: Annual data - no confidentiality restrictions

Study description Data description Materials and publications Download data

Data file description

ADS11 - Labour Force Survey, 2011 [data file]

File ID: ADS11_P1_SL_V1_R1

*.txt - TEXT

- Number of variables: 214
- Number of units: 61888

Version: August 2012

Variable list

val Round

Value	Frequency
1	17797
2	13491
3	11580
4	9892
5	9128
-2 Missing value	0

Valid cases	Invalid cases	Minimum	Maximum	Arithmetic mean	Standard deviation
61888	0				

Valid range from 1 to 5

starost Age

Value	Frequency
-2 Missing value	

Valid cases	Invalid cases	Minimum	Maximum	Arithmetic mean	Standard deviation
61888	0	0		42.034	21.877

Valid range from 0 to Protected value

```
1 <dataDescr>
2 <var ID="V5" name="val" files="F1" dcml="0" intrvl="discrete">
3   <location width="1"/><labl>Round</labl>
4   <valrng><range min="1" max="5"/></valrng>
5   <invalrng><item VALUE="-2"/></invalrng>
6   <sumStat type="vald">61888</sumStat>
7   <sumStat type="invd">0</sumStat>
8   <catgry><catValu>1</catValu>
9     <labl>1.</labl>
10    <catStat type="freq">17797</catStat>
11  </catgry>
12 <catgry><catValu>2</catValu>
13   <labl>2.</labl>
14   <catStat type="freq">13491</catStat>
15  </catgry>
16 <catgry><catValu>3</catValu>
17   <labl>3.</labl>
18   <catStat type="freq">11580</catStat>
19  </catgry>
20 <catgry><catValu>4</catValu>
21   <labl>4.</labl>
22   <catStat type="freq">9892</catStat>
23  </catgry>
24 <catgry><catValu>5</catValu>
25   <labl>5.</labl>
26   <catStat type="freq">9128</catStat>
27  </catgry><catgry missing="Y">
28   <catValu>-2</catValu>
29   <labl>Missing value</labl>
30   <catStat type="freq">0</catStat>
31  </catgry>
32   <varFormat type="numeric" schema="other"/>
33 </var>
34 <var ID="V8" name="starost" files="F1" dcml="0" intrvl="contin">
35   <location width="2"/>
36   <labl>Age</labl>
37   <valrng>
38     <range min="0" max="Protected value"/>
39   </valrng>
40   <invalrng><item VALUE="-2"/></invalrng>
41   <sumStat type="vald">61888</sumStat>
42   <sumStat type="invd">0</sumStat>
43   <sumStat type="min">0</sumStat>
44   <sumStat type="max">Protected value</sumStat>
45   <sumStat type="mean">42.034</sumStat>
46   <sumStat type="stdev">21.877</sumStat>
47   <catgry missing="Y"><catValu>-2</catValu>
48   <labl>Missing value</labl>
49  </catgry>
50   <varFormat type="numeric" schema="other"/>
51 </var>
52 </dataDescr>
```

Software features and quick demonstration of using them

- Deleting variables/variable information
- Protecting descriptive statistics (min, max, mean, SD)
- Calculating replacement univariate statistics
- Top- and bottom- coding
- Bracketing/recoding
- Minimum frequency protection
- Displaying DDI Data Description statistics (importing XML)
- Guidelines and data protection report

Official statistics univariate statistics protection

- How much should univariate statistics (aggregated data) be protected to be publically distributed?
- Using the minimum frequency rule, what should be the threshold for one dimension table data - frequencies for values of one variable only?
- Which data protection functions should be added to the presented tool and which should be considered unnecessary?

Thank you for your attention!

Sebastian Kočar

sebastian.kocar@fdv.uni-lj.si

http://www.adp.fdv.uni-lj.si/

Privacy in Statistical Databases 2014,
Eivissa, September 18th 2014