# Statistical disclosure control and output checking

Remote access and on-site



#### Contents

- Settings
- Problems
- Guidelines
- Practical issues

### **Settings**

Remote acces and on-site

- Data remain at NSI
- Researcher analyses data
- Results
  - Intermediate
  - Final





#### **Settings**

Check results that leave the NSI

- On-site: intermediate and final
- Remote access: (intermediate and) final

#### **Settings**

Output of OS/RA can be everything...

Transforming data, combining data, ...

Set of rules that cover every possible output?

mission impossible?



#### **Problems**

What to check?

- Traditional output (tables)
- Descriptive statistics
- Output of many analyses
  - Regression
- Graphics

**Problems (traditional)** 

Tabular output

- Frequency tables
- Magnitude tables



You now know all about that

#### **Problems (descriptives)**

Extreme values (max, min)

- Often extremes belong to identifiable respondents
- Extremes are very visible

#### **Problems (descriptives)**

Mean

- Single respondent?
- Two respondents?
- Dominating respondent?
- Group disclosure?

Observations  $(x_{1i}, x_{2i}, x_{pi}, y_i)$ , i = 1, ..., nLinear regression:  $y = X\beta + u$ 

with 
$$\boldsymbol{X} = \begin{pmatrix} 1 & x_{11} & \cdots & x_{p1} \\ \vdots & \vdots & & \vdots \\ 1 & x_{1n} & \cdots & x_{pn} \end{pmatrix}$$
,  $\boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ \vdots \\ \beta_p \end{pmatrix}$ 

and  $\boldsymbol{u}$  is *n*-vector with error-terms.

Estimator

$$\widehat{y} = X\widehat{\beta} = X(X^tX)^{-1}X^ty$$

Simple case: p = 1, dummy variable  $x_{1i} = \begin{cases} 1 & \text{if } i = 1 \\ 0 & \text{else} \end{cases}$  hence  $X = \begin{pmatrix} 1 & 1 \\ 1 & 0 \\ \vdots & \vdots \\ 1 & 0 \end{pmatrix}$ 

#### Then, calculation of estimator yields

$$\hat{y}_1 = \cdots = y_1$$

## Differencing two regressions, with deletion/addition of single observation

Subtract two regression coefficients

$$\widehat{\boldsymbol{\beta}} - \widehat{\boldsymbol{\beta}}_0$$

May lead to information about individual values

- Known means
- Binary explanatory variables
- Binary dependent variable

#### **Problems (graphics)**

- Possible to identify individual respondent?
- Outliers present?
- Is underlying data safe?
- Does picture have "data embedded"?

#### **Problems (graphics)**

#### Operating returns (10<sup>6</sup> Euros) in manufacturing, 2010



#### SDC Output checking

16

Started as deliverable of ESSnet on common tools and harmonised methodology for SDC in the ESS

included in Wiley book Statistical Disclosure Control

currently improved upon and extended in FP7 Data without Boundaries

Classify output as either

*generally safe* or *generally unsafe* 

Generally safe:

researcher can expect to have output cleared with no or minimal changes

Generally unsafe:

not to be cleared unless researcher can demonstrate non-disclosiveness

Types of error with output checking:

- Releasing disclosive output
- Not releasing safe output

Two model approach:



(simple rules that make a first distinction between

safe output and potentially unsafe output)



(further check potentially unsafe output)

Overall rules-of-thumb



- 10 units
- 10 degrees of freedom
- No group disclosure
- No dominance

### **Descriptive statistics**

<u>Max, min</u>

- : Not released (refer to single unit)
- : Rules for magnitude tables apply (can be released if not associated with single unit)

#### **Descriptive statistics**

Means, indices, ratios, indicators

should come from  $\geq$  10 units (unweighted) largest contributor not more than 50%

indices considered more closely

Complex indices like Fisher Price Index:

$$\sqrt{\frac{\sum_{j=1}^{m} p_{1,j} q_{0,j}}{\sum_{j=1}^{m} p_{0,j} q_{0,j}}} \frac{\sum_{j=1}^{m} p_{1,j} q_{1,j}}{\sum_{j=1}^{m} p_{0,j} q_{1,j}}$$

Simple indices like:

$$\frac{\sum_{i=1}^{n} X_i}{\sum_{i=1}^{N} Y_i}$$

Complexity yields some protection

Simplicity: magnitude tables ( $X_i$  quantitative) frequency tables ( $X_i$  dichotomous)

#### Linear regression coefficients

Solution of the second state of the second

: check on degrees of freedom, not only categorical variables, not on a single unit

#### Graphs

I graphs are not allowed, unless based on undisclosive datapoints

: data point not identified with units, no significant outliers, no "embedded" data (e.g. .jpg, .bmp)

#### **Guidelines (overview)**

- Frequency tables
- Magnitude tables
- Max, min, percentiles (incl. median)
- Mode
- Means, indices, ratios, indicators
- Concentration ratios
- Higher moments like (co)variance, kurtosis, skewness, ...
- Graphics

#### **Guidelines (overview)**

- Linear regression coefficients
- Non-linear regression coefficients
- Estimation residuals
- Summary and test-statistics ( $R^2$ ,  $\chi^2$ , ...)
- Correlation coefficients
- Factor analysis
- Correspondence analysis

Missing something? Contact us or Expert Group SDC.

#### **Practical issues**

- SDC awareness of researcher
- Joint responsibility of researcher and NSI
- Output checking is context specific
  - Check by SDC-expert and contextexpert
- Consistency between checkers
- Costs and time-limits

#### **Practical issues**

Output checking is labour intensive!

Possibilities:

- Researcher facilitates NSI for checking
- Check sample of outputs
- Put price on checking of output
- Facilitate writing of reports at safe setting