

National guidelines on data editing; the foundation for building a solution for the future

ASLAUG HURLEN FOSS OG ANE SEIERSTAD



Statistisk sentralbyrå
Statistics Norway

Technological modernisation – cloud solution

- Storage in Google cloud buckets
- Access management
- Transfer services
- Automatic processing of data from source data to input data.
- Programming languages: Python and R.



Status of data editing in Statistics Norway

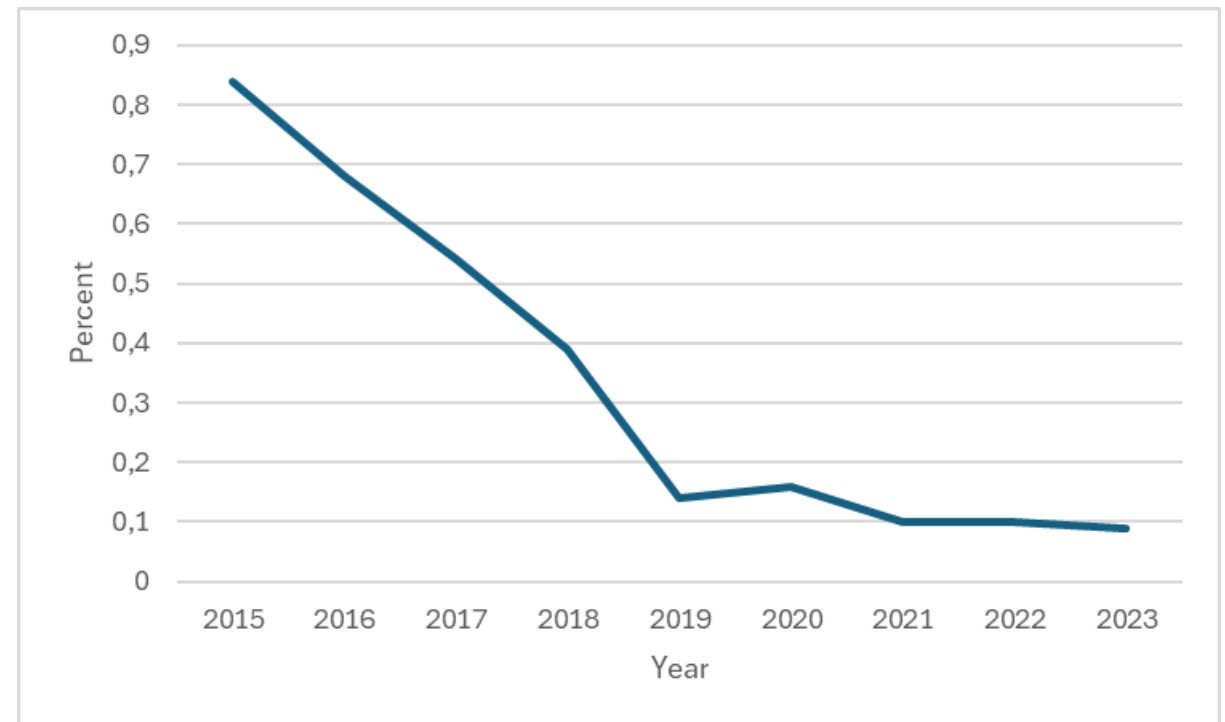
System of data editing:

- Dynarev - built 2005 – 150 surveys
- Driller – built 2011 – 25 surveys

Main functionality:

- Management of controls (edit rules)
- Execution of controls
- Rule-based and manual corrections
- Reports – analyzing

Proportion of manually edited values in percent, Dynarev



Lack of common understanding

- Data editing process
- Terminology
- Human interaction



Striving for a unified understanding

- We made a draft of the principles and guidelines
- Subject experts: feedback
- Formal consultation to all statistics departments
- Approved by the director's meeting and by the director general



Geir Axelsen Director general of Statistics Norway



Basis for principles and guidelines

- Apply international methodology
- Adapt to Norwegian condition
- Specification of lower level

modernstats

Generic Statistical Data Editing Model **GSDEM**

(Version 2.0, June 2019)

Methodology for data validation 1.0

Revised edition June 2016

Essnet Validat Foundation

Marco Di Zio, Nadežda Fursova, Tjalling Gelsema, Sarah Gießing, Ugo Guarnera, Jūratė Petrauskienė, Lucas Quensel-von Kalben, Mauro Scanu, K.O. ten Bosch, Mark van der Loo, Katrin Walsdorfer



Statistisk sentralbyrå
Statistics Norway

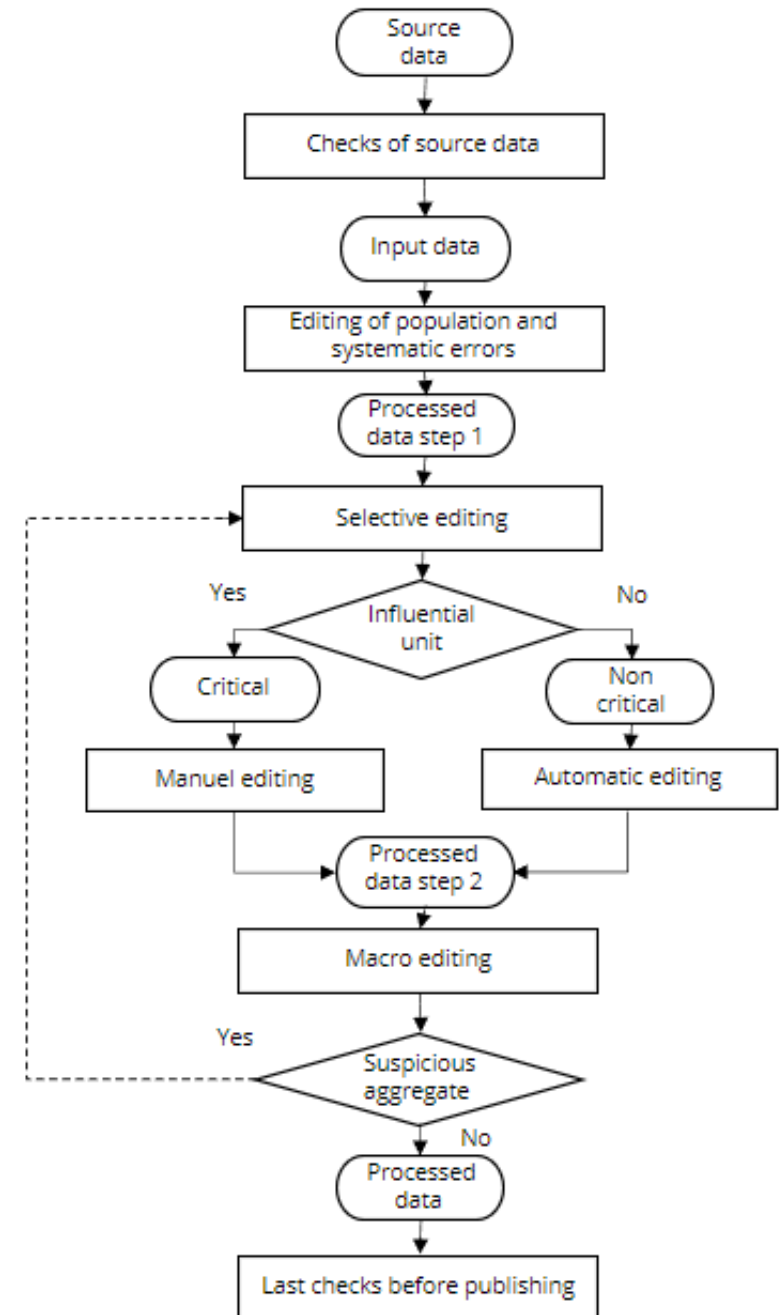
Process model the source for writing guidelines

Added two processes :

- Checks of source data
- Last checks before publishing

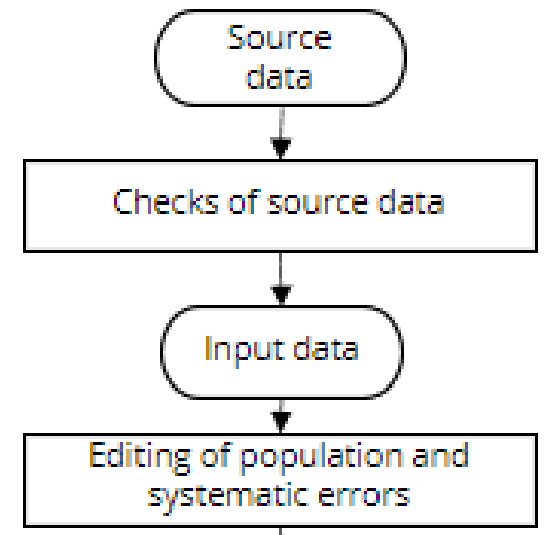
Sometimes separated:

- Social statistics
- Business statistics



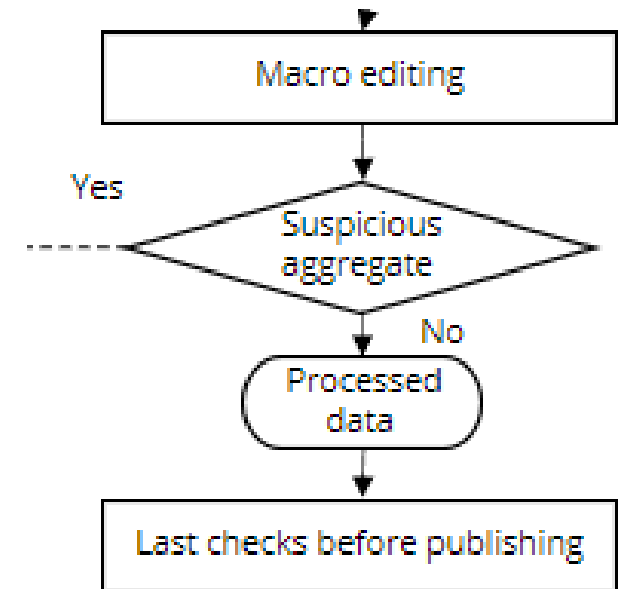
Example: Editing of Personal Identification number

- Control of personal identity numbers is done by cross-referencing against the latest SNR catalogue, containing all personal identity numbers that have ever been issued.
- Expired identity numbers are replaced with the latest identity number to maintain a stable identification.
- Document invalid identity numbers using standardized error codes.



Example: Macro editing

- Assess whether the aggregates are plausible given **historical trends**
- Evaluate the plausibility of the aggregate concerning **derived quantities**, such as ratios.
- Investigate the dataset for **potential errors or explanations** when changes occur beyond what is considered normal.
- Compare aggregates with those of **comparable countries** for assessment.
- Aggregates may undergo **editing** when they are part of a system where certain relationships must be maintained.



Example: Control of Statbank tables

- Statistical producers should perform controls on the tables before submission to Statbank
- Those receiving the tables should control compliance with Statbank rules.

StatBank Norway

StatBank contains detailed tables with time series. You can create your own selections and save these in different file formats. We also offer an API for StatBank.

[How to use StatBank Norway](#)

Our release time for all statistics is 08:00 CET.

[Changes to tables in Statbank](#)

At 05:00 and 11:30 StatBank's metadata are updated, and the tables can be temporarily unavailable for up to five minutes. Published figures which are being revised are shown as '0' or '.' between 05:00 and 08:00.

[Public APIs](#)

Choose topic and table

Agriculture, forestry, hunting and fishing



Banking and financial markets



Construction, housing and property



Culture and recreation



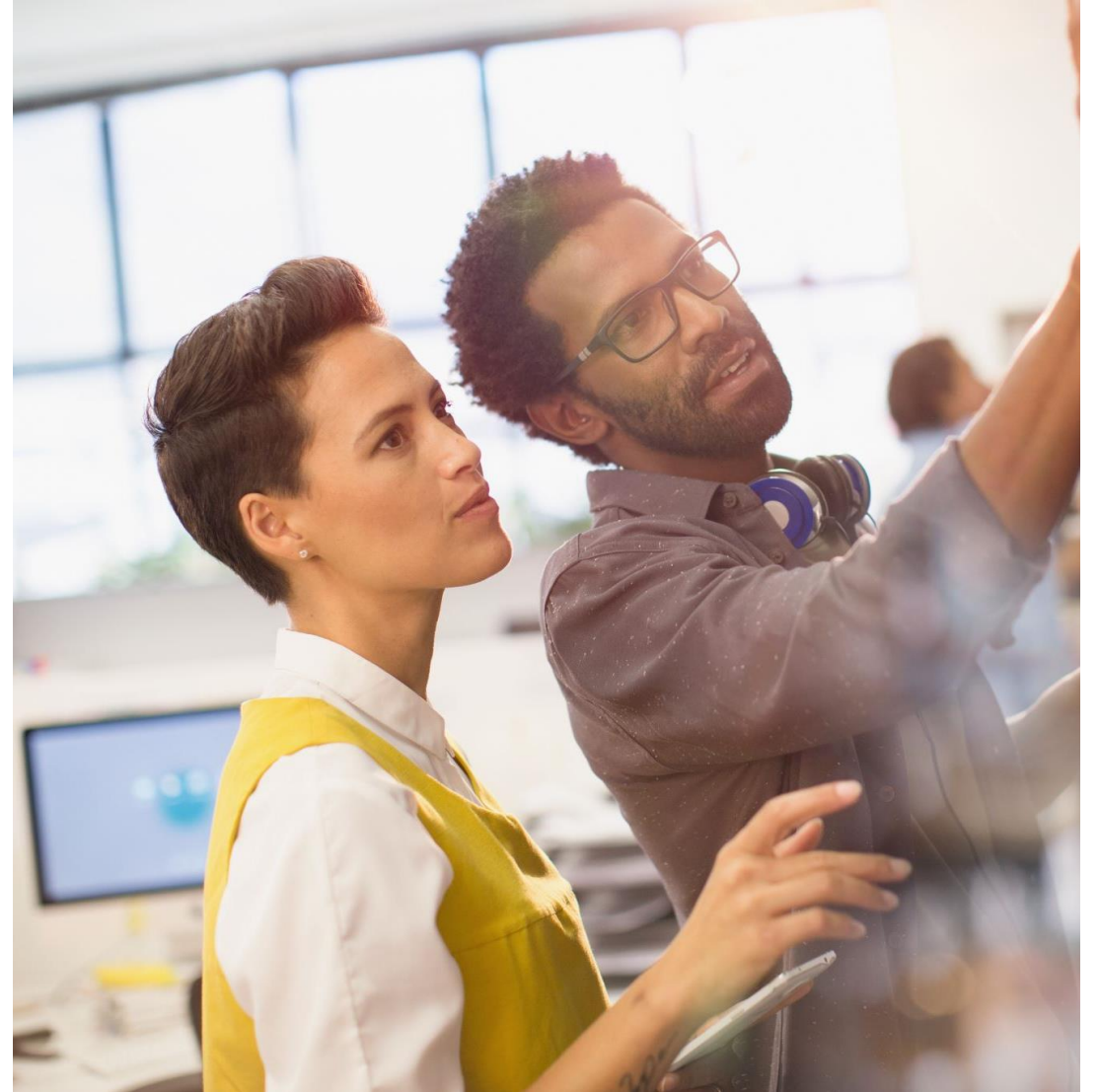
Principles of data editing

- Principle 1: Good knowledge about the subject area of statistics and the background of source data is the basis for a good editing process.
- Principle 2: The objective of data editing should be clearly defined.
- Principle 3: High-quality data input is the best
- Principle 4: Always control the data.
- Principle 5: The earlier, the better.
- Principle 6: The controls, control effects and changes made must be well-documented.
- Principle 7: Automate the editing process as much as possible.
- Principle 8: Streamline the editing work.
- Principle 9: Data editing should be evaluated.



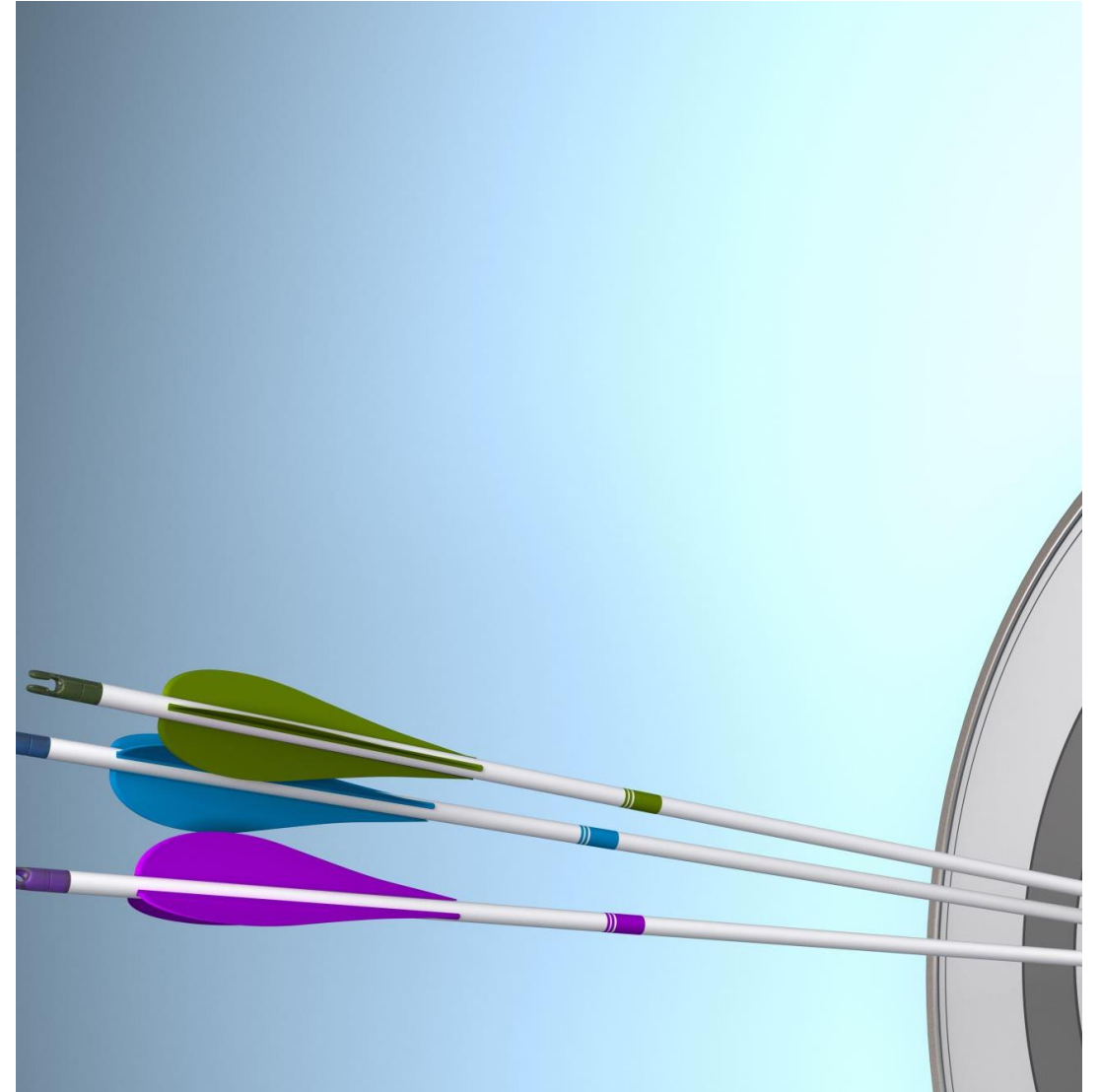
Knowledge

Principle 1: Good knowledge about the subject area of statistics and the background of source data is the basis for a good editing process.



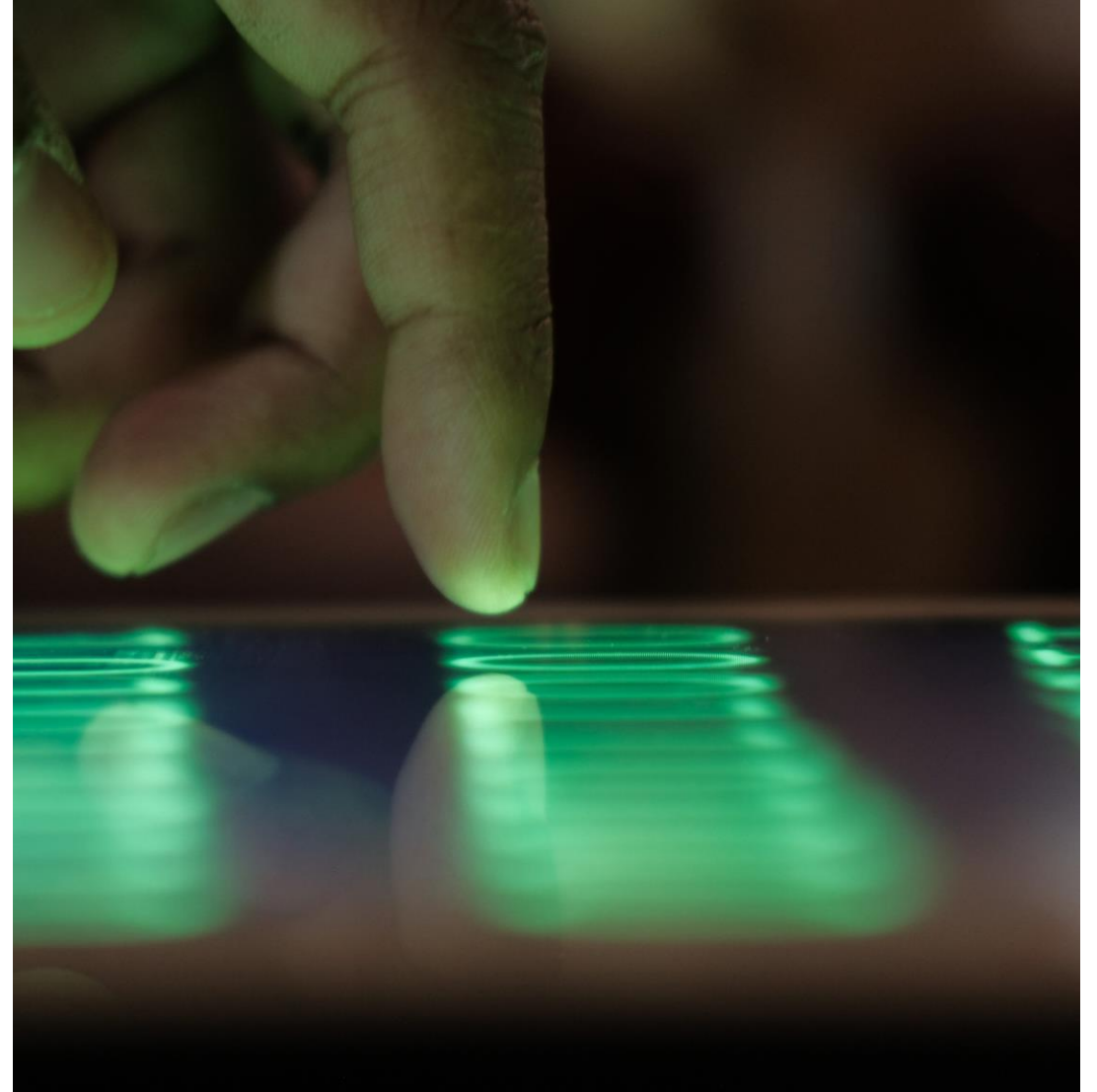
Objective

Principle 2: The objective of data editing should be clearly defined.



High-quality data

- Principle 3: High-quality data input is the best



Control

Principle 4: Always control the data.



Place

- Principle 5: The earlier, the better.



Well-documented

- Principle 6: The controls, control effects and changes made must be well-documented.



Automate

Principle 7: Automate the editing process as much as possible.

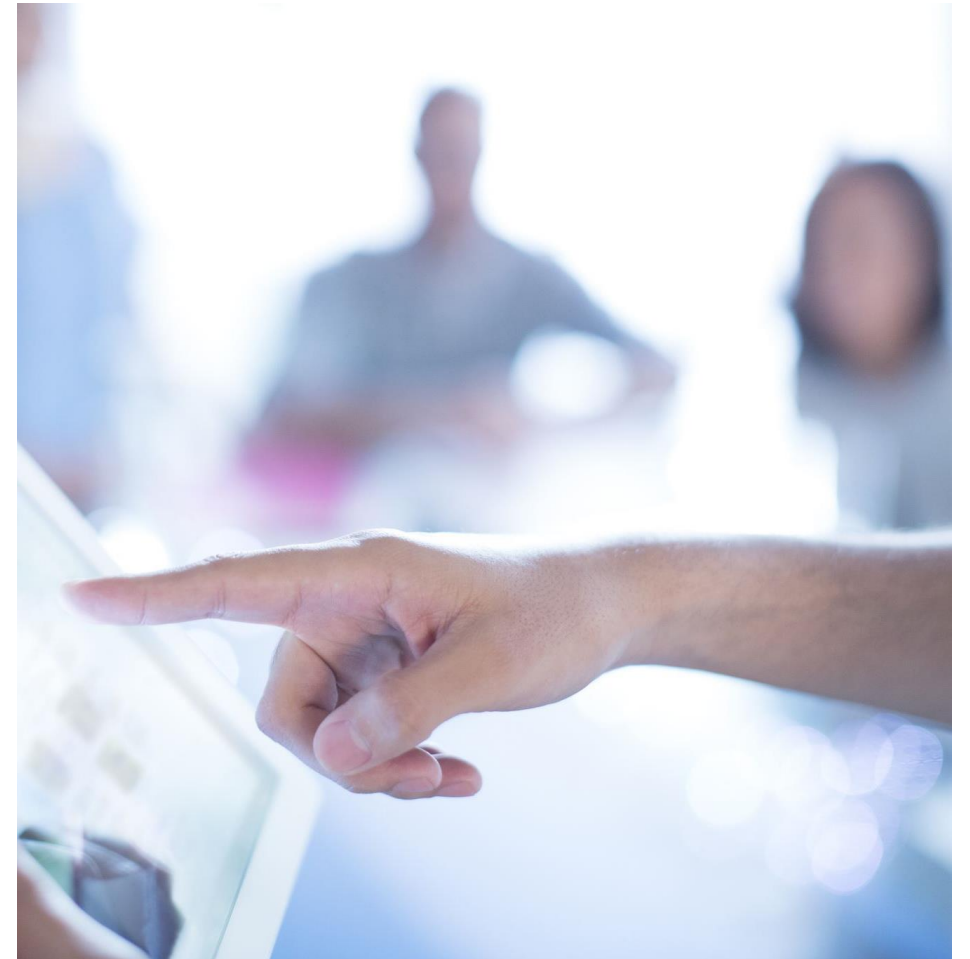


Streamline

Principle 8: Streamline the editing work.

Improve human interaction:

- Good interfaces
- Visualization
- Selective editing
- Drilling in data between macro and micro



Evaluate

Principle 9: Data editing should be evaluated.

- Create quality indicators
- Analyze the indicators
- Take action



Application under construction

- Build in python with some function in R
- Independent of the underlying storage technology
- Possible to manually edit values
- Module based
 - Aggregation and drilling in data
 - Controls – edit rules
 - Influential observation and Hidiroglou-Berthelot function from R
- Process data and quality indicators



Thanks to international community for:

Generic Statistical Data Editing Model: UNECE

Metodology for data validation: EUROSTAT

