

# Integration of simple data validation tools into the UNECE production tables

Michael Nagy

13<sup>th</sup> Session of the Joint Task Force on environmental  
Statistics and Indicators, 29-30 June 2017, Geneva



**UNECE**

# Content

1. Why this presentation
2. What is data validation?
3. Are there internationally agreed processes for data validation?
4. Data validation examples from other International Organisations
5. Recommendations for the UNECE production templates

# 1. Why this presentation

- Decision taken at the 11<sup>th</sup> Session of the JTF
- See following paragraph 31 of the report:

31. The need for specific guidelines for quality assurance (QA) of environment statistics and metadata as part of the SEIS self-assessment was highlighted. The Joint Task Force called for further action towards producing such guidelines, which should complement existing QA frameworks used in countries and by international bodies. Countries asked for case studies and examples of good practices in quality assurance of environment statistics.

The integration of simple validation tools into the UNECE production table (similar to validation tools already used by Eurostat and UNSD) could be a first practical step to support quality assurance.

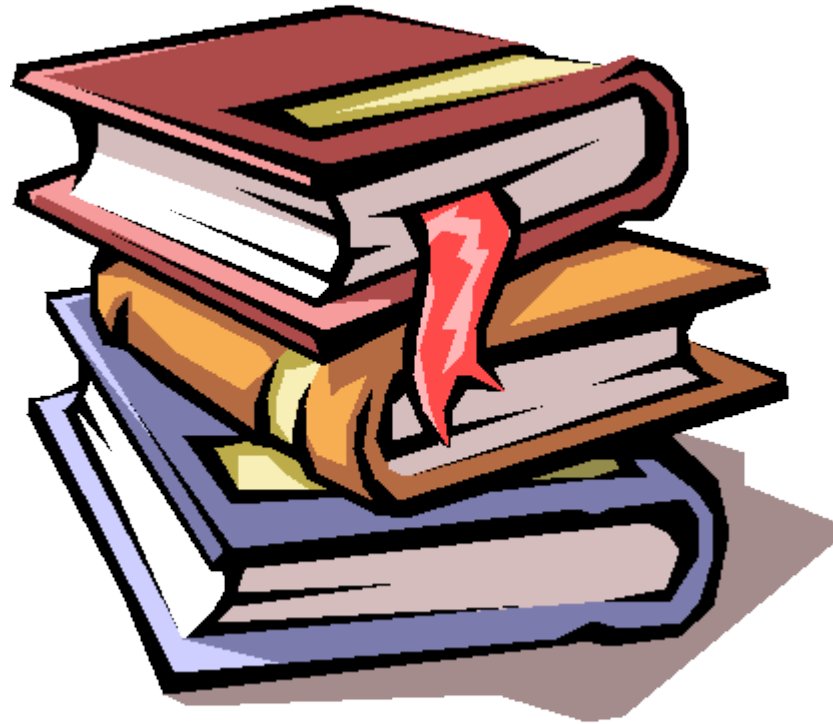
## 2. What is data validation?



# Definition of data validation

***An activity aimed at verifying whether the value of a data item comes from the given set of acceptable values.*** (UNECE glossary on statistical data editing, Eurostat and OECD use similar definitions)

### 3. Are there internationally agreed processes for data validation?



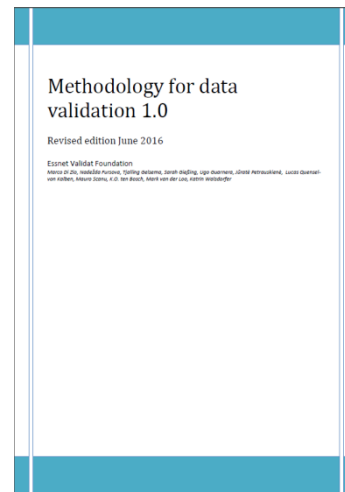
# Procedures and approaches for data validation have never been systematised

Data validation is carried out by all NSOs, but:

- procedures and approaches have never been systematised
- most of them are ad-hoc.

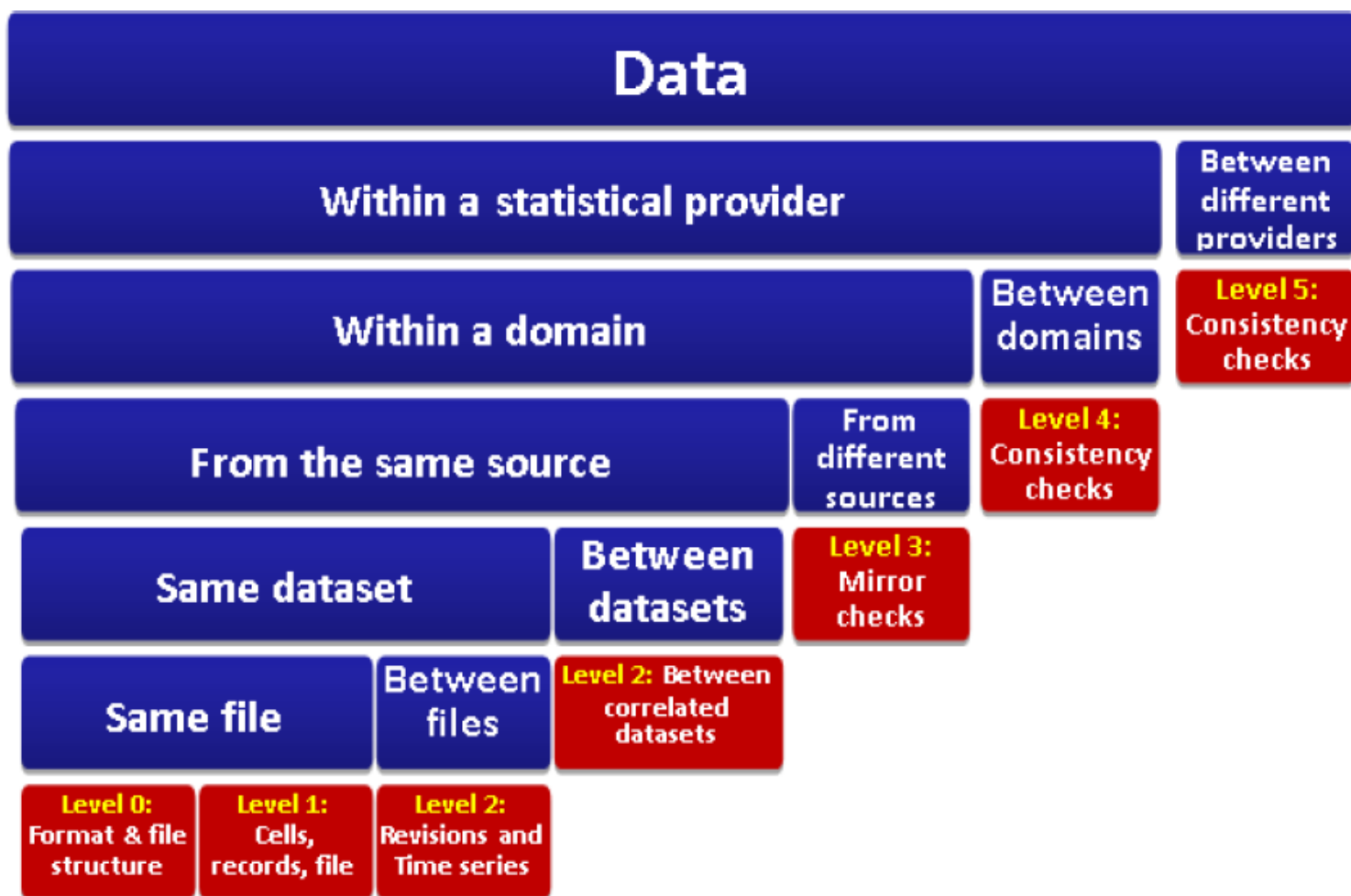
In June 2016 Eurostat has published a first document “Methodology for data validation 1.0” as an output of the Essnet project Validat foundation, primarily financed by Eurostat, involving CBS, Destatis, Istat and Statistics Lithuania.

[https://ec.europa.eu/eurostat/cros/system/files/methodology\\_for\\_data\\_validation\\_v1.0\\_rev-2016-06\\_final.pdf](https://ec.europa.eu/eurostat/cros/system/files/methodology_for_data_validation_v1.0_rev-2016-06_final.pdf)



# Eurostat methodology for data validation 1.0

## Data validation levels





# Eurostat methodology for data validation 1.0

Validation rules can be categorised as follows:

## A – integrity of the data file

- Formal validity of entries
- Presence of an entry
- No duplicate entries
- Code lists (if applicable) are used

## B - Logical validation and consistency

Table 1: Categories of a 2-way typology for validation rules for logical validation and consistency

Typology dimension	Types of checks	
1	Identity checks	Range checks <ul style="list-style-type: none"><li>• bounds fixed</li><li>• bounds depending on entries in other fields</li></ul>
2	Simple checks, based directly on the entry of a target field	More “complex” checks, combining more than one field by functions (like sums, differences, ratios)

+ conditional checks

# Eurostat methodology for data validation 1.0

## Examples for logical validation and consistency

- **Conditional check:** If a country is landlocked, then the discharge to the sea has to be 0
- **Identity check:** The total water use has to be equal to water available for use minus losses
- **Range check with fixed bounds:** Water use of households per capita should be between 100 l/day and 400 l/day
- **Range check with bounds depending on entries in other fields:** water losses must be smaller than water available for use

## 4. Data validation examples from other International Organisations





# UNSD/UNEP Environment Questionnaire

## Automatic validation by UNSD:

- Coherence:
  - E.g. total waste generation is not equal to the sum of waste generated by economic activities and households – value turns red (note: on purpose the questionnaires of UNSD do not calculate sums automatically)
  - Comparison with data from outside sources (e.g. World Bank, FAO Aquastat etc.)
- Time series
  - Significant changes will turn red in the grey data validation template (e.g. total amount of municipal waste collected increases by >25% within one year).

# International Energy Agency Monthly Oil Questionnaire

## MONTHLY OIL QUESTIONNAIRE

Country (select country)  
 Month of data January 2013  
 Status of data (final/provisional) Final  
 Prepared by \_\_\_\_\_  
 Date of transmission \_\_\_\_\_

ALL DATA SHOULD BE ENTERED  
 IN THOUSAND METRIC TONS  
 WITH NO DECIMAL PLACE

Table 1: SUPPLY OF CRUDE OIL, NGL, REFINERY FEEDSTOCKS, ADDITIVES AND OTHER HYDROCARBONS

Unit: Thousand metric tons

		Crude oil	Natural gas liquids	Refinery feedstocks	Additives / oxygenates	Of which Biofuels	Other hydrocarbons	Total (A to F, excl. E)	Probate
		A	B	C	D	E	F	G	H
+ Indigenous production	1	10						10	
+ Receipts from other sources	2							0	
+ Backflows <sup>1</sup>	3							0	
+ Products transferred <sup>2</sup>	4							0	
+ Imports (Balance) <sup>3</sup>	5	50						50	
- Exports (Balance) <sup>4</sup>	6	80						80	
- Direct use <sup>5</sup>	7	1						1	
- Stock changes <sup>6</sup>	8	0						0	
= Refinery intake (Calculated)	9	-21	0	0	0	0	0	-21	
- Statistical difference	10	-46	0	0	0	0	0	-46	
= Refinery intake (Observed)	11	25						25	
Memo Item: Refinery losses	12							0	

- Total of Backflows (reference G3) should correspond to Total products of Backflows to refineries in Table 2 (reference Z19).
- Total of Products transferred (reference G4) should correspond to Total products of Products transferred in Table 2 (reference Z9).
- Imports (Balance) should correspond to Total imports for these products in Table 3 (references A105 to F105).
- Exports (Balance) should correspond to Total exports for these products in Table 4 (references A96 to F96).
- Direct use should be carried over to Primary product receipts in Table 2 (line 1).
- Stock changes should correspond to Closing level minus Opening level of Stocks on national territory, Table 5 (line 2 minus line 1).

OK  
 OK  
 OK

Total Imports (Balance) do not equal Total Imports in Table 3

Total Exports do not equal Total Exports in Table 4

Total Direct use figure does not equal Total products of Primary product receipts in Table 2

### 3. Recommendations for the UNECE production templates



# General recommendations for automatised validation rules in the UNECE data production tables

## **Main goals of automatised validation:**

- Identify wrong use of units of measurement
- Identify typos (e.g. wrong decimals)
- Identify unexpected changes in time series
- Identify if data are balanced (when there should be a balance)

## **Validation should indicate unexpected values, but**

- Not block data entry
- Invite data providers to double-check, correct if necessary, and provide more information in footnotes

**Validation will be based on data within the same file**

**Logical validation and consistency checks only**

**Test with energy indicators and biodiversity indicators, then gradual implementation**



# Side condition

**No formulas for automatic calculation of sums (e.g. total water use is calculated automatically as the sum of water use by economic activities + households)**

- Sometimes only the total is known, but not the disaggregated values.
- A sum would be calculated automatically, even if not all disaggregated values are available.
- Automatic formulas for the data entry would outsmart parts of the data validation

# Recommended validation rules to be implemented into the UNECE data production templates

## Simple conditional checks:

- e.g. landlocked countries cannot have discharge to the sea or marine protected areas

## Identity checks:

- E.g. the total primary energy supply be equal to production + imports – exports – bunkers + stock changes
- E.g. the total waste generated has to be the sum of waste generated by economic activities + households

## Range check with fixed bounds

- E.g. Water use of households per capita should be between 100 l/day and 400 l/day

## Range check with bounds depending on entries in other fields:

- e.g. tolerance values for changes over time
- e.g. water losses must be smaller than water available for use

### 3. Questions, comments?

- Does the JTF agree with the recommended validation rules?
- Should the secretariat test it with the revised templates for energy and biodiversity indicators?
- Any other comments?

