

# Validation in R using metadata from SDMX registries

Olav ten Bosch, Mark van der Loo  
*Statistics Netherlands*

8<sup>th</sup> SDMX Global Conference, 27-30 September 2021

# Contents

- International data validation
- SDMX registries
- Data cleaning with R
- Connecting R-validate to SDMX
- Wrap up



# International data validation (1)

- ***Invalid*** data may lead to ***costly*** retransmissions or reprocessing (data ping pong)
- To guarantee overall data ***quality*** and ***efficiency***, the European Statistical System (ESS) is moving towards more harmonised validation activities
- International validation rules are agreed in domain specific ***statistical working groups***
- Data producer (NSIs) and data consumers (international organisations) ***validate*** data against the ***same rules***

## International data validation (2)

**ESSnet Validat Foundation 2015-2016**

**ESSnet Validat Integration, 2017 (DE,**

- Handbook on validation
- Validation principles
- Main types of rules

### ***Validation principles:***

1. *The sooner, the better*
2. *Trust but verify*
3. *Well-documented and appropriately communicated validation rules*
4. *Well-documented and appropriately communicated validation errors*
5. *Comply or explain*
6. *Good enough is the new perfect*

[https://ec.europa.eu/eurostat/cros/content/data-validation-overview\\_en](https://ec.europa.eu/eurostat/cros/content/data-validation-overview_en)

**ValidatFOSS & ValidatFOSS2, 2018-2021**

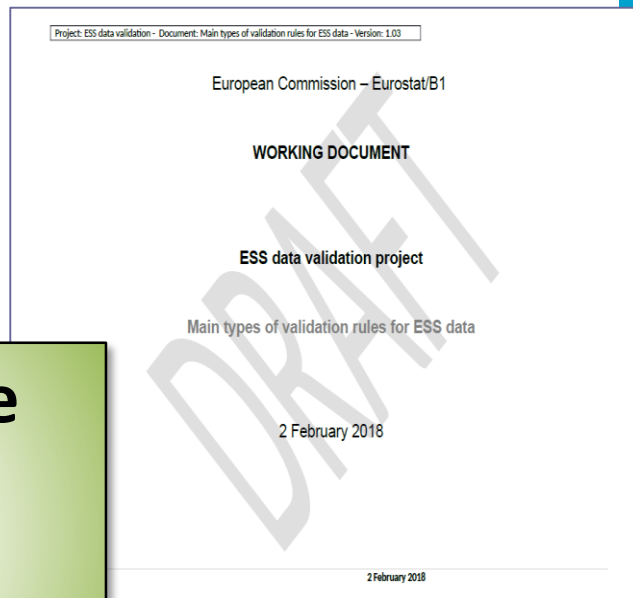
- Implementing validation in Free and open source software (R)

# International data validation (3)

## *'Main types of validation rules'*

- FDT: Field Type
- FDL: Field Length
- FDM: Field is Mandatory or empty
- COV: COdes are Valid
- RWD: Records are Without Duplicate id-keys
- REP: Records Expected are Provided
- RTS: Records are all present for Time Series
- RNR: Records' Number is in a Range
- COC: COdes are Consistent
- VIR: Values are In a Range
- VCO: Values are COnsistent
- VAD: Valueas for Aggregates are consistent with Details
- VSA: Values for Seasonally Adjusted data are plausible

**Can we easily use internationally agreed metadata for data validation?**



Eurostat, 2018

Implemented in R-package GenericValidationRules:

<https://github.com/SNStatComp/GenericValidationRules>

# SDMX registries (1)



[registry.sdmx.org](http://registry.sdmx.org)



[webgate.ec.europa.eu/sdmxregistry](http://webgate.ec.europa.eu/sdmxregistry)



[sdmxcentral.imf.org/overview.html](http://sdmxcentral.imf.org/overview.html)



[sdmx.data.unicef.org](http://sdmx.data.unicef.org)

Other & Internal  
SDMX registries



# SDMX registries (2)

SDMX Global Registry  
Version 9.8.11

Convert  
Validate and Convert Datasets

Data Registration  
Register the location of your datasets and re-register any dataset changes

Browse Data  
Browse all data that has been made available to SDMX Global Registry

## Codestats

| SDMX | Id              | Name                                 | State |
|------|-----------------|--------------------------------------|-------|
| SDMX | CL_AGE          | Age                                  | Final |
| SDMX | CL_AREA         | Reference area code list             | Final |
| SDMX | CL_CIVIL_STATUS | Civil (or Marital) Status            | Final |
| SDMX | CL_COFOG_1999   | Classification of the Functions o... | Final |
| SDMX | CL_COICOP_1999  | Classification of Individual Cons... | Final |
| SDMX | CL_CONF_STATUS  | Confidentiality Status               | Final |
| SDMX | CL_COPNI_1999   | Classification of the Purposes o...  | Final |
| SDMX | CL_COPP_1999    | Classification of the Outlays of ... | Final |
| SDMX | CL_DECIMALS     | Decimals                             | Final |

Viewing: Civil (or Marital) Status [1.0]

| Position | Id | Name  |
|----------|----|---|
| 1        | S  | Single person   |
| 2        | M  | Married person  |
| 3        | W  | Widowed person  |
| 4        | D  | Divorced person   |
| 5        | L  | Legally separated person  |
| 6        | P  | Person in registered partnership  |
| 7        | Q  | Person whose registered partnership ended with the death of the partner |
| 8        | E  | Person whose registered partnership was legally dissolved               |

- Access to many metadata resources
- Versioning, ownership

# SDMX registries (3)



Other & Internal  
SDMX registries

- Programmatic access via SDMX 2.1 REST API
- Excellent cheat sheet
- Note: changes in SDMX 3.0

<https://github.com/sdmx-twg/sdmx-rest>

## SDMX 2.1 RESTful web services cheat sheet, v1.5.0

| Structural metadata queries: <a href="https://ws-entry-point/resource/AgencyID/resourceID/version/itemID?queryStringParameters">https://ws-entry-point/resource/AgencyID/resourceID/version/itemID?queryStringParameters</a>   |  | Default     |
|--|--|-------------|
| <b>resource</b>  | The type of metadata to be returned. Values: <code>datastructure</code> , <code>metadatatstructure</code> , <code>categoryscheme</code> , <code>conceptscheme</code> , <code>odelist</code> , <code>hierarchicalodelist</code> , <code>organisationscheme</code> , <code>agencyscheme</code> , <code>dataconsumerscheme</code> , <code>organisationunitscheme</code> , <code>dataflow</code> , <code>metadatatflow</code> , <code>reportingtaxyonomy</code> , <code>provisionagreement</code> , <code>structureset</code> , <code>process</code> , <code>categorisation</code> , <code>contentconstraint</code> , <code>attachmentconstraint</code> , <code>actualconstraint</code> , <code>allowedconstraint</code> , <code>structure</code> , <code>transformationscheme</code> , <code>rulesetscheme</code> , <code>userdefinedoperatorscheme</code> , <code>customtypescheme</code> , <code>namepersonalisationscheme</code> , <code>vtlnappingscheme</code> |             |
| <b>agencyID</b>  | Agency maintaining the artefact (e.g.: SDMX)   | all         |
| <b>resourceID</b>  | Artefact ID (e.g.: CL_FREQ)  | all         |
| <b>version</b>   | Artefact version (e.g.: 1.0)   | latest      |
| <b>itemID</b>  | ID of the item (for item schemes) or hierarchy (for hierarchical codelists) to be returned.  | all         |
| <b>detail</b>  | Desired amount of information ( <code>allstubs</code> , <code>referencestubs</code> , <code>allcompletetubs</code> , <code>referencecompletetubs</code> , <code>referencepartial</code> , <code>full</code> )  | full        |
| <b>references</b>  | References to be returned with the artefact ( <code>none</code> , <code>parents</code> , <code>parentsandsiblings</code> , <code>children</code> , <code>descendants</code> , <code>all</code> , a resource type)  | none        |
| Data query: <a href="https://ws-entry-point/resource/FlowRef/key/providerRef?queryStringParameters">https://ws-entry-point/resource/FlowRef/key/providerRef?queryStringParameters</a>  |  | Default     |
| <b>resource</b>  | <code>data</code> or <code>metadada</code>   |             |
| <b>FlowRef</b>   | Dataflow ref (e.g. ECB_EXR_latest)   |             |
| <b>key</b>   | Key of the series to be returned (e.g. D_NOK_EUR.SP00.A). Wildcarding (e.g. D_ .EUR.SP00.A) and OR (e.g. D_NOK+RUB_ .EUR.SP00.A) supported.  | all         |
| <b>providerRef</b>   | Data provider (e.g.: ECB)  | all         |
| <b>startPeriod</b>   | Start period (inclusive). ISO8601 (e.g. 2014-01) or SDMX reporting period (e.g. 2014-Q3).  |             |
| <b>endPeriod</b>   | End period (inclusive). ISO8601 (e.g. 2014-01-01) or SDMX reporting period (e.g. 2014-W53).  |             |
| <b>updatedAfter</b>  | Last time the query was performed. Used to retrieve deltas. Must be percent-encoded (e.g. 2009-05-15T14%3A15%3A00%2B01%3A00)   |             |
| <b>firstNObservations</b>  | Maximum number of observations starting from the first observation   |             |
| <b>lastNObservations</b>   | Maximum number of observations counting back from the most recent observation  |             |
| <b>dimensionAtObservation</b>  | Id for the dimension attached at the observation level   | TIME_PERIOD |
| <b>detail</b>  | Desired amount of information to be returned. Values: <code>full</code> , <code>dataonly</code> , <code>serieskeyonly</code> , <code>nodata</code>   | full        |
| <b>includeHistory</b>  | Whether to return vintages   | false       |
| <b>Legend: Mandatory path parameter / Optional path parameter / Query string parameter (all optional)</b><br>For available data queries, cf. <a href="https://github.com/sdmx-twg/sdmx-rest/wiki/Data-Availability">https://github.com/sdmx-twg/sdmx-rest/wiki/Data-Availability</a> |  |             |



## SDMX registries (4)

In ValidatFOSS2 we retrieved metadata from SDMX registries using their SDMX API's. Some experiences:

- There is no central list of SDMX endpoints / registries
- Some REST behaviours differ, but we came to one script for all
- For R: rsdmx ([cran.r-project.org/package=rsdmx](https://cran.r-project.org/package=rsdmx)) makes life easy
- Some thoughts about registry content:
  - does a statistician know where to look for which content?
  - is there content overlap among registries?
  - are registry contents consistent?

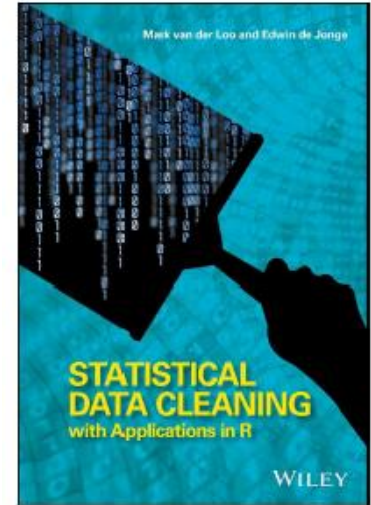
<https://github.com/SNStatComp/validatesdmx> (demo notebooks Python, R)



# Data cleaning with R (1)

MPJ van der Loo and E de Jonge (2018)  
*Statistical data cleaning with applications in R*  
John Wiley & Sons, NY.

- R data cleaning ecosystem
  - **validate**: check data based on validation rules
  - **dcmodify**: change data based on 'if-this-then-that' rules
  - **errorlocate**: locate errors based on validation rules and mark them for correction
  - **simputation**: many different imputation methods
  - **rspa**: adapt numerical records to fit (in)equality restrictions
  - **deductive**: solve errors based on control rules
  - **validatetools**: find inconsistencies and redundancies



# Data cleaning with R (2)

## Rules

```
# Range limits:  
Age >= 0  
Age <= 120  
Working_hours >= 0  
Working_hours <= 100  
  
# Some checks between variables:  
if (Married > 0) Age > 18  
if (Working_hours > 0) Employed > 0  
  
# Such a rule depends on country legislation:  
if (Age > 65) Working_hours = 0  
  
# ID must be unique  
any(duplicated(ID)) == FALSE
```

## Data

| ID | Age | Marital status | Status in employment | Working hours per week |
|----|-----|----------------|----------------------|------------------------|
| 1  | 36  | 0              | 1                    | 40                     |
| 2  | 40  | 1              | 1                    | 40                     |
| 3  | 25  | 0              | 0                    | 0                      |
| 4  | 31  | 0              | 1                    | 20                     |
| 5  | 62  | 1              | 1                    | 43                     |
| 6  | 55  | 1              | 1                    | 41                     |
| 7  | 34  | 1              | 1                    | 40                     |

## Summary

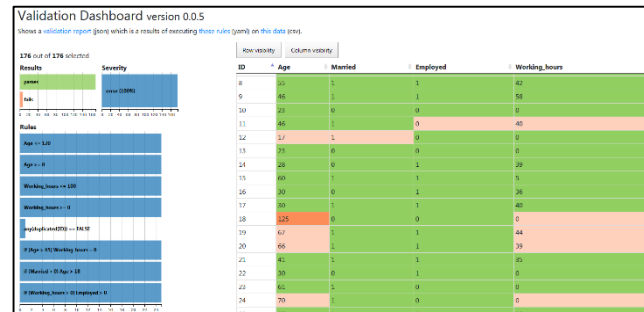
```
> summary(validation)  
name items passes fails nNA error warning expression  
1 V1 25 25 0 0 FALSE FALSE (Age - 0) >= -1e-08  
2 V2 25 24 1 0 FALSE FALSE (Age - 120) <= 1e-08  
3 V3 25 25 0 0 FALSE FALSE (working_hours - 0) >= -1e-08  
4 V4 25 25 0 0 FALSE FALSE (working_hours - 100) <= 1e-08  
5 V5 25 24 1 0 FALSE FALSE !(Married > 0) | (Age > 18)  
6 V6 25 24 1 0 FALSE FALSE !(working_hours > 0) | (Employed > 0)  
7 V7 25 21 4 0 FALSE FALSE !(Age > 65) | (working_hours = 0)  
8 V8 1 0 1 0 FALSE FALSE any(duplicated(ID)) == FALSE
```

## Per rule



confront

## Dashboard: data & results



## Data cleaning with R (3)

***R-validate*** supports the majority of checks needed in today's data validation, such as:

- Variable types, missingness, field length, format, range checks, codelists
- Availability, uniqueness, gaps in time series
- Multivariate checks, balance equalities
- Statistical & groupwise checks

Documented in online ***cookbook***:



[data-cleaning.github.io/validate](https://data-cleaning.github.io/validate)

# Connecting R-validate to SDMX (1)

***R-validate*** ( $\geq 1.1.0$ ) now supports:

- Rules based on ***any codelists***  
from ***any registry***
- ***Caching*** of registry results within a session
- Convenience functions for global and ESTAT registry
- Result can be used in a natural way in R expressions:

```
Activity %in% global_codelist(agency_id="ESTAT", resource_id="CL_ACTIVITY")
```

- New chapter in cookbook on SDMX:  
[data-cleaning.github.io/validate/sect-sdmxrules.html](https://data-cleaning.github.io/validate/sect-sdmxrules.html)

| function           | what it does                                      |
|--------------------|---|
| sdmx_endpoint      | retrieve URL for SDMX endpoint                    |
| sdmx_codelist      | retrieve sdmx codelist                            |
| estat_codelist     | retrieve codelist from Eurostat SDMX registry     |
| global_codelist    | retrieve codelist from Global SDMX registry       |
| validator_from_dsd | derive validation rules from DSD in SDMX registry |



# Connecting R-validate to SDMX (2)

- Deriving all rules from a DSD:

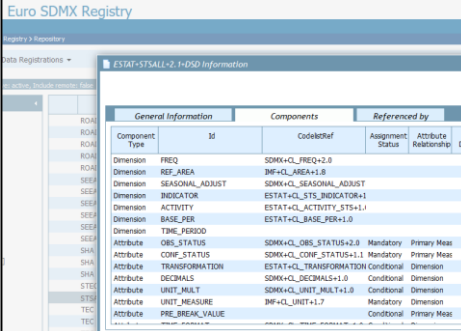
```
# import data
my_data <- read.csv("mydata.csv")

# derive all rules from a DSD
rules <- validator_from_dsd(endpoint = sdmx_endpoint("ESTAT")
  , agency_id = "ESTAT", resource_id = "STSALL", version="latest")

# confront data with rules
out <- confront(my_data, rules)

# plot results
plot(out)
```

- Generates and checks multiple codelist rules derived from the ESTAT registry
- Easy to integrate in statistical processes



| General Information |                 | Components                | Referenced by     |                        |
|---------------------|-----------------|---------------------------|-------------------|------------------------|
| Component Type      | Id              | Code/Ref                  | Assignment Status | Attribute Relationship |
| Dimension           | FREQ            | SDMX+CL_FREQ+2.0          |                   |                        |
| Dimension           | REF_AREA        | IMF+CL_AREA+1.0           |                   |                        |
| Dimension           | SEASONAL_ADJUST | SDMX+CL_SEASONAL_ADJUST   |                   |                        |
| Dimension           | INDICATOR       | ESTAT+CL_STS_INDICATOR+1  |                   |                        |
| Dimension           | ACTIVITY        | ESTAT+CL_ACTIVITY_STS+1.1 |                   |                        |
| Dimension           | BASE_PER        | ESTAT+CL_BASE_PER+1.0     |                   |                        |
| Dimension           | TIME_PERIOD     |                           |                   |                        |
| Attribute           | OBS_STATUS      | SDMX+CL_OBS_STATUS+2.0    | Mandatory         | Primary Measure        |
| Attribute           | CONF_STATUS     | SDMX+CL_CONF_STATUS+1.1   | Mandatory         | Primary Measure        |
| Attribute           | TRANSFORMATION  | ESTAT+CL_TRANSFORMATION   | Conditional       | Dimension              |
| Attribute           | DECIMALS        | SDMX+CL_DECIMALS+1.0      | Conditional       | Dimension              |
| Attribute           | UNIT_MULT       | SDMX+CL_UNIT_MULT+1.0     | Conditional       | Dimension              |
| Attribute           | UNIT_MEASURE    | IMF+CL_UNIT+1.7           | Mandatory         | Dimension              |
| Attribute           | PRE_BREAK_VALUE |                           | Conditional       | Primary Measure        |

# Wrap-up

- The ESS works on improving **international data validation** i.a. handbook, principles, main types of rules
- R has a rich **data-cleaning ecosystem**. **R-validate** covers many of today's validation needs in data validation
- SDMX metadata contained in **registries** is necessary input for international validation rules
- R-validate now supports **checks** against **any** SDMX registry in high level functions.
- Some **experiences** working with SDMX registries and their APIs have been shared



# Questions, ideas, suggestions



Olav ten Bosch

o.tenbosch@cbs.nl

@olavtenbosch

Mark van der Loo

mpj.vanderloo@cbs.nl

@markvdloo

and keep an eye on:

[awesomeofficialstatistics.org](https://awesomeofficialstatistics.org)



Star

183



Fork

47

