# SDMX in a FAIR official statistics landscape

Statistics Netherlands

Olav ten Bosch, Edwin de Jonge, Henk Laloli

12th SDMX Experts workshop, 7-11 Oct. Amsterdam

# Contents

- The official statistics landscape

- Open source software for access; awesome list

- Features commonly offered: towards FAIRness

- And what about linked data?

- Wrap-up

# Statistics Netherlands output



**Layered Open Data architecture**
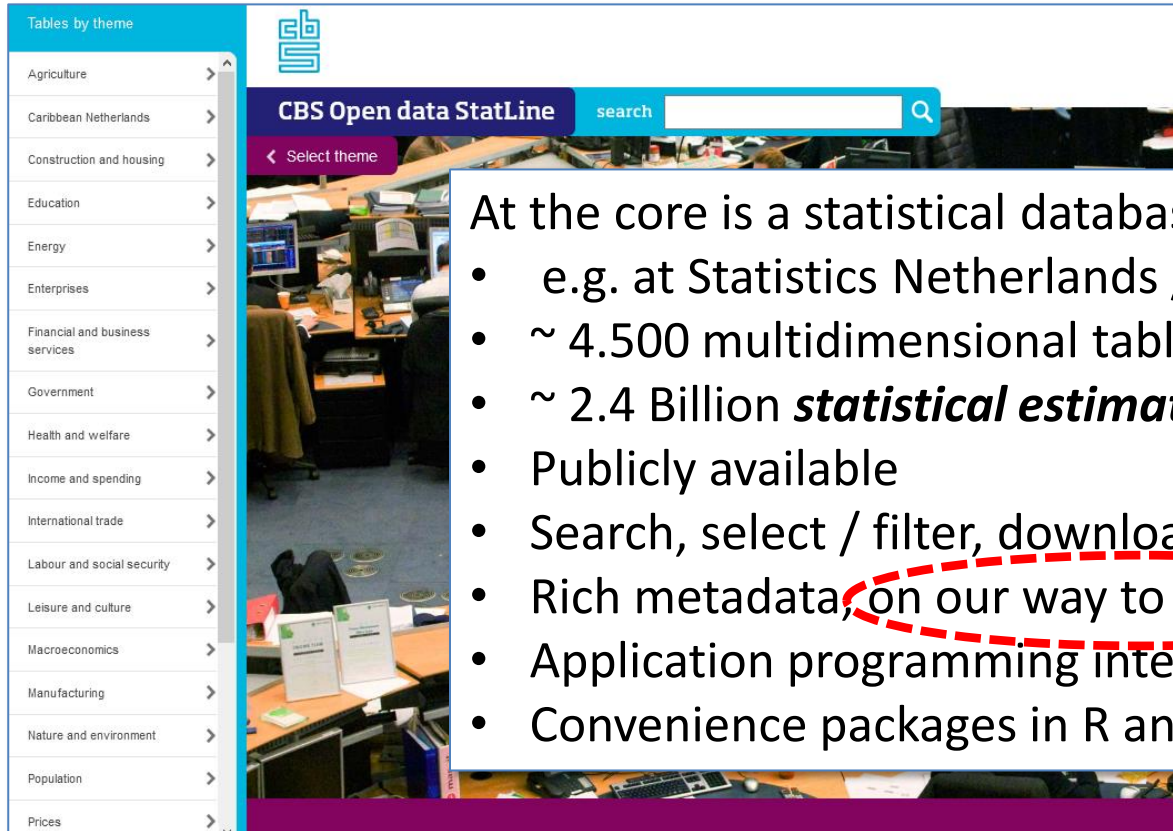
News articles

Thematic pages & visualisations

Selected data slices

Statistical Estimates

3

# Statistical dissemination database

Tables by theme

- Agriculture
- Caribbean Netherlands
- Construction and housing
- Education
- Energy
- Enterprises
- Financial and business services
- Government
- Health and welfare
- Income and spending
- International trade
- Labour and social security
- Leisure and culture
- Macroeconomics
- Manufacturing
- Nature and environment
- Population
- Prices

**CBS Open data StatLine**   search

‹ Select theme

At the core is a statistical database:
- e.g. at Statistics Netherlands / CBS:
- ~ 4.500 multidimensional tables
- ~ 2.4 Billion *statistical estimates* (active tables)
- Publicly available
- Search, select / filter, download
- Rich metadata, on our way to SDMX
- Application programming interface (API): Odata
- Convenience packages in R and Python
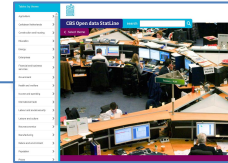
# Official Statistics landscape (OSL)



OSL = collective output of all official statistics organizations

DK

NL

DE

IT
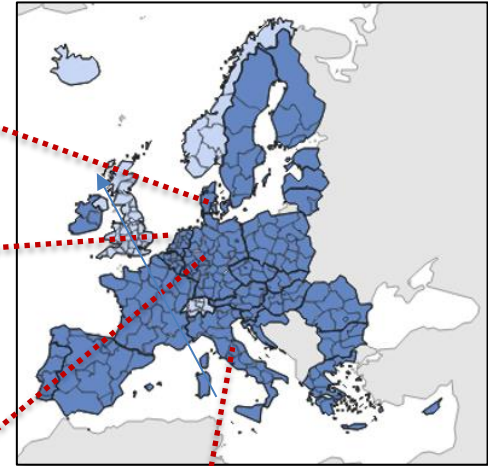
User

# Official Statistics landscape (OSL)

OSL = collective output of all official statistics organizations

Happy user?

DK

NL

DE

IT

# Awesome list of statistical software

- A **community approach** to knowledge management
- To **collectively remember useful software** in official statistics
- Started UNECE SDE 2017
- Maintained by **statistical community**
- How:
  - Using **awesome concept**
  - **awesomeofficialstatistics.org**
  - *Largest category:*
    *"Access to official statistics"*

uRos17-23

# awesome packages by GSBPM

Over 30 software packages, giving access to > 60 dataproviders, majority are R-packages

**Statistical disclosure control (GSBPM 6.4)**
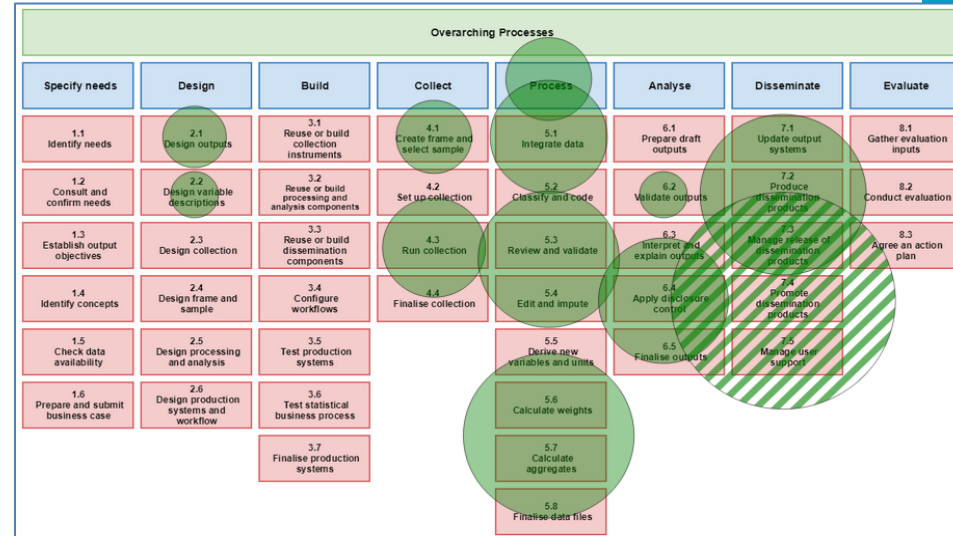
- `GitHub v5.1.7b3` `last commit march` `license EUPL-1.2`

  Java and C++ application Mu-ARGUS. Tool to create safe micro-data files. S

- `GitHub v4.2.4.2` `license EUPL-1.2`

  Java C++ Fortran and Delphi application T-ARGUS. Tool to protect statistica

- `CRAN 5.7.6 – 2 months ago` `license GPL-2`

  R package sdcMicro. Disclosure control for statistical microdata.

- `CRAN 0.32.6 – 4 months ago`

  R package sdcTable. Disclosure control for tabulated data.

**Sampling (GSBPM 4.1)**

- `CRAN 2.10 – a month ago` `license GPL (>= 2)`

  R package sampling. Several algorithms for drawing survey samples, including
  sampling designs (high entropy, systematic, Rao-Sampford, etc.), and calibratin

- `CRAN 4.0 – 4 years ago` `license GPL (>= 2)`

  R package surveyplanning. Tools for sample survey planning, including sample
  expected precision for the estimates of totals, and calculation of optimal samp

**Data integration and record linkage (GSBPM 5.1)**

- `CRAN 0.3.4 – 5 months ago` `license GPL-3`

  R package reclin2. Functions to assist in performing probabilistic record linkage a
  pairs, comparing records, em-algorithm for estimating m- and u-probabilities, fo
  also be used for pre- and post-processing for machine learning methods for reco

- `CRAN 0.4-12.4 – a year ago` `license GPL (>= 2)`

  R package RecordLinkage. Implementation of the Fellegi-Sunter method for reco

- `CRAN 1.4.1 – 2 years ago` `license GPL (>= 2)`

  R package StatMatch. Statistical Matching or Data Fusion

- `CRAN 0.6.1 – 24 days ago` `license GPL (>= 3)`

  R package fastLink. Implements a Fellegi-Sunter probabilistic record linkage mod

**Access to official statistics (GSBPM 7.4)**

- `CRAN 0.6-3 – 7 months ago` `license GPL (>= 2)`

  R package rsdmx. Access to data or metadata from statistical organisations that support SDM
  The package contains a list of SDMX access points of various national and international statisti

- `CRAN 0.3.1 – 7 months ago` `license GPL-3`

  R package readsdmx. Read SDMX into dataframes from local SDMX-ML file or web-service. Pa
  OECD.

- `GitHub v2.14.0` `last commit last wednesday` `license Apache-2.0`

  Python sdmx. Python library that implements SDMX 2.1 to explore data from SDMX data provi
  data and metadata and convert it into Pandas objects.

- `CRAN 0.4.3 – 7 months ago` `license MIT + file LICENSE`

  R package rjstat. Read and write data sets in the JSON-stat format.

- `PyPI v2.4.0` `license Apache License 2.0`

  Python pyjstat. Read and write JSON-stat.

- `GitHub v0.2.8` `last commit march 2023` `license MIT`

  Java application json-stat.java. Read and write JSON-stat. By Statistics Norway.

- `CRAN 0.2.5 – 2 years ago` `license CC0`

  R package oecd. Search and Extract Data from the OECD

- `CRAN 0.8.21 – 7 months ago` `license BSD_2_clause + file LICENSE`

  R package sorvi. Finnish Open Government Data Toolkit

- `CRAN 4.0.0 – 3 months ago` `license BSD_2_clause + file LICENSE`

  R package eurostat. Tools to download data from the Eurostat database together with search
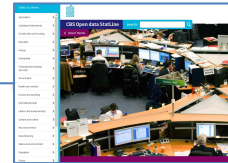
# Software layer surrounding OSL

> 30 software packages
> 60 dataproviders
> majority R-packages
> Some use SDMX



OSL
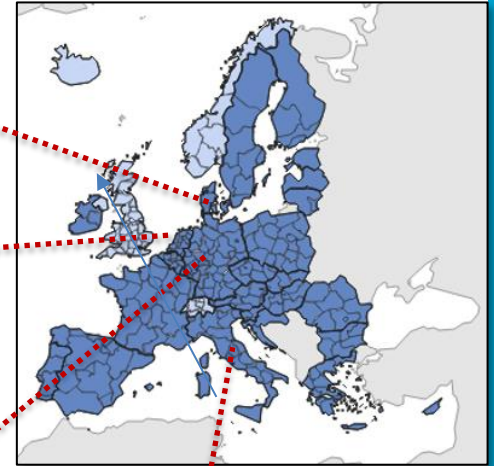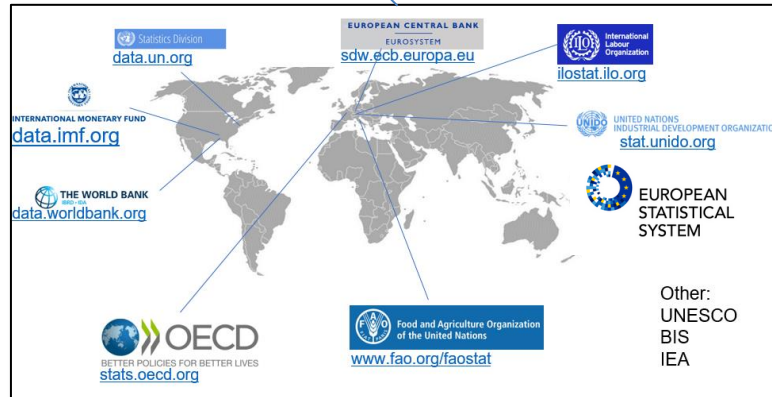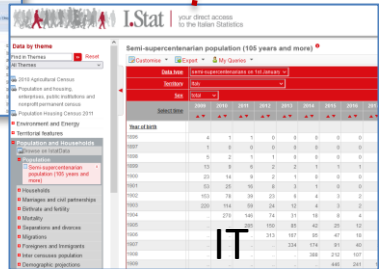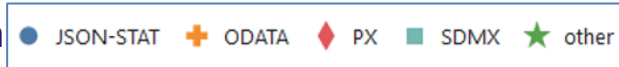
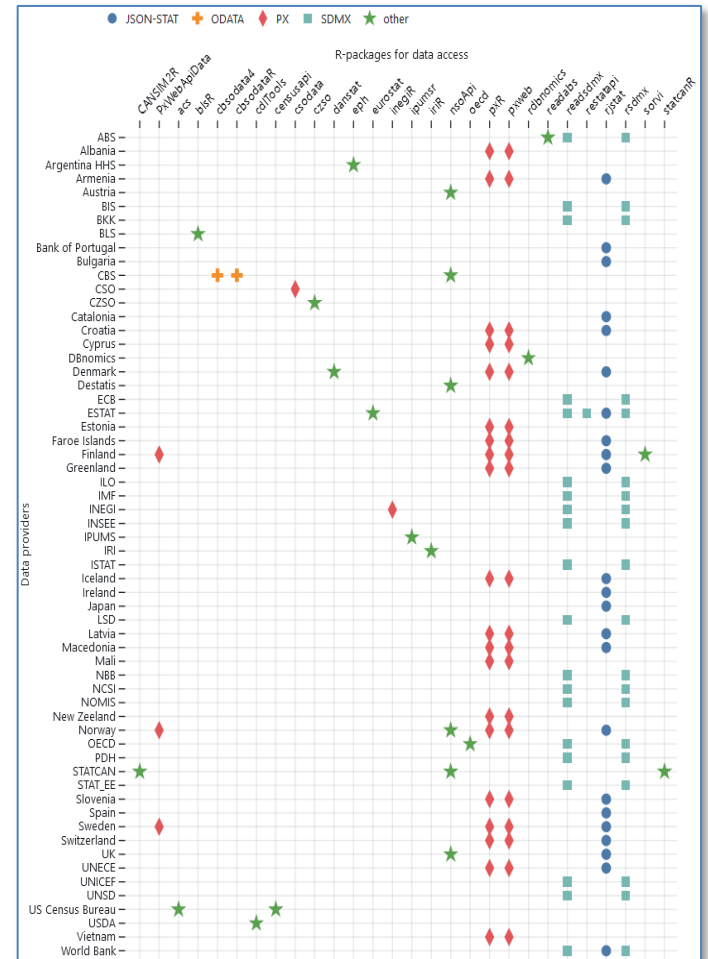DK

NL

DE

IT

Data user

9

# "access to official statistics" software landscape

- Matrix derived from documentation and links to web pages and packages execution:

packages (38) x data

- Standards:



https://observablehq.com/@olavtenbosch/access_to_official_statistics

# "access to official statistics" software landscape



- Matrix from docs, links to web pages and packages execution:

  packages (38) x dataproviders (60) x standards (5)

- Standards: 

- JSON-STAT/PX v. SDMX: almost disjunct worlds

- ODATA: CBS only

https://observablehq.com/@olavtenbosch/access_to_official_statistics

# "access to official statistics" software landscape



- Standards-oriented packages:

  *rsdmx, readsdmx, rjstat, px\**

- Data provider-centric packages:

  *inegiR, readabs, statcanR, eurostat*

- Official statistics aggregator sites:

  *rdbnomics*: economic data

  *ipumsr*: census & survey data
  
  time&space harmonised

https://observablehq.com/@olavtenbosch/access_to_official_statistics

# Features commonly offered

- **endpoint hiding**: wrapping the preconfigured endpoint(s) in a function
- **catalogue retrieval**: to list the availability datasets on the endpoint(s)
- **search**: to search for datasets or within datasets
- **endpoint queries**: query for subsets / slices on the endpoint(s) side
- **local queries**: the ability to easily slice or filter on the client
- **caching**: preventing unnecessary roundtrips
- **cartographic queries**: retrieve a geo data or a map with the data
- **registry access**: access to coordinated metadata in registries

Could we build a "one-for-all" solution with all features?

Software features supporting FAIR principles

| Software feature | Findability | Accessibility | Interoperability | Reusability |
|---|---|---|---|---|
| endpoint hiding | yes | yes | | |
| catalogue retrieval | yes | yes | | |
| search | yes | | | |
| endpoint queries | | yes | | |
| local queries | | yes | | |
| caching | | yes | | |
| cartographic queries | yes | | yes | yes |
| registry access | yes | yes | yes | |

# Features powerful in SDMX (1)



- **_endpoint hiding_**: wrapping the preconfigured endpoint(s) in a function
- *catalogue retrieval: to list the availability datasets on the endpoint(s)*
- *search: to search for datasets or within datasets*
- *endpoint queries: query for subsets / slices on the endpoint(s) side*
- *local queries: the ability to easily slice or filter on the client*
- *caching: preventing unnecessary roundtrips*
- *cartographic queries: retrieve a geo data or a map with the data*
- *registry access: access to coordinated metadata in registries*

khaeru.github.io/sdmx/

# Features powerful in SDMX (2)

- ***endpoint hiding***: wrapping the preconfigured endpoint(s) in a function
- ***catalogue retrieval***: to list the availability datasets on the endpoint(s)
- ***search***: to search for datasets or within datasets
- ***endpoint queries***: query for subsets / slices on the endpoint(s) side
- ***local queries***: the ability to easily slice or filter on the client
- ***caching***: preventing unnecessary roundtrips
- ***cartographic queries***: retrieve a geo data or a map with the data
- ***registry access***: access to coordinated metadata in registries



Rich API &
cheat sheets

# Features powerful in SDMX (3)

- ***endpoint hiding***: wrapping the preconfigured endpoint(s) in a function
- *catalogue retrieval*: to list the availability datasets on the endpoint(s)
- *search*: to search for datasets or within datasets
- ***endpoint queries***: query for subsets / slices on the endpoint(s) side
- *local queries*: the ability to easily slice or filter on the client
- *caching*: preventing unnecessary roundtrips
- *cartographic queries*: retrieve a geo data or a map with the data
- ***registry access***: access to coordinated metadata in registries

2021



SDMX registries

Validating data in R using SDMX registry

# And what about linked data? (1)

- Multiple organisations offer (meta)data as Linked Data (LD)

- URIs, SKOS, XKOS, schema.org, SPARQL



CBS Knowledge graph

op.europa.eu/en/web/eu-vocabularies/eurostat

vocabs.cbs.nl/taxonomie

statistics.gov.scot

# And what about linked data? (2)

- Statistical metadata as LD helps modeling *changes in metadata & linkability*



Linking within ESS

Linking outside ESS

vocabs.cbs.nl/cbs_geo/en

# Reflections on linked data (LD) in offstats

- LD is well suited to *link across ESS* and to *non-ESS communities*
- LD suited to model *complex statistical metadata dependencies*
  - Example: a complete historical graph of Dutch geo-changes on municipality level for over two centuries
- Challenges:
  - LD can be perceived *complex* to end-users => need for *data stories*
  - *Flexibility* of LD model gives room for different implementations among ESS => need for *harmonization*

- Food for thought: with the current *AI training data hunger*: could statistical LD content more easily picked up than ESS-specific standards like SDMX?

19

# Becoming (more) FAIR…



- *# standards*
- *Identifiability*
- *Note the software layer!*
- *Work with communities*

> A1: (Meta)data are retrievable by their identifier using a standardised communication protocol

20

# Wrap-up

- Official Statistics Landscape grows towards *standardisation: SDMX, JSON-stat/PX*
- Data user can choose from over *30 open source software packages*, each offering specific functionality to over *60 data sources*
- Many NSIs offer *targeted (R-)packages* for data and metadata access
- *Common features* in software layer identifie
- *SDMX* scores well on 3 of the features
- There is *no software package* that provides access to *all* official statistics
- Dream:
  - *one generic (R-)package* to access all official statistics content from all dataproviders supporting maximum FAIRness*
  - offstats community and the open source / data science community have to work together to reach that goal

* FAIR: Findable, Accesible, Interoperable, Reusable

Feedback: Olav ten Bosch o.tenbosch@cbs.nl  pubs

awesomeofficialstatistics.org Please ☆ Star 280