# Access to official statistics from R: part II

## Statistics Netherlands

Olav ten Bosch, Edwin de Jonge

**uRos2024,** Greece, 27-29 November 2024

# Contents

- What is the awesome list?
    - History, concept, working
    - What happened from uRos2023 -> uRos 2024?

- Category "access to official statistics"
    - Packages, data providers, standards, *features*
    - A "one for all" package?

- Wrap-up

# What is the awesome list?

# Awesome list of official statistics software

- Started at **UNECE SDE conference** 2017 (The Hague)
- a **community approach** to **remember useful software**
- A **public** list, clear and simple **criteria**
- Majority is **R software**
- **awesomeofficialstatistics.org**

**2018**

**2019**

**2023**

Data integration and record linkage (GSBPM 5.1)

- CRAN `0.5.0 – 9 months ago` license `GPL-3`

  R package reclin2. Functions to assist in performing
  pairs, comparing records, em-algorithm for estimati
  also be used for pre- and post-processing for mach

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

  R package RecordLinkage. Implementation of the Fe

- CRAN `1.4.2 – 6 months ago` license `GPL (>= 2)`

  R package StatMatch. Statistical Matching or Data F

- CRAN `0.6.1 – a year ago` license `GPL (>= 3)`

  R package fastLink. Implements a Fellegi-Sunter pro
  and the inclusion of auxiliary information. Documen

Statistical disclosure control (GSBPM 6.4)

- GitHub `v5.1.7b4` last commit `march` license `EUPL-1.2`

  Java and C++ application Mu-ARGUS. Tool to create safe micro-d

- GitHub `v4.2.5.2` last commit `august` license `EUPL-1.2`

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

- CRAN `5.7.8 – 8 months ago` license `GPL-2`

  R package sdcMicro. Disclosure control for statistical microdata.

- CRAN `0.32.6 – a year ago` license `GPL (>= 2)`

  R package sdcTable. Disclosure control for tabulated data.

- CRAN `1.0.7 – 2 years ago` license `Apache License 2.0 | file LICENSE`

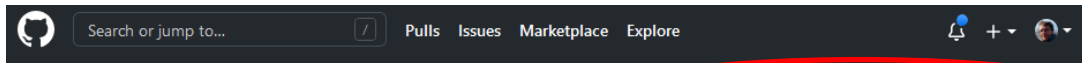  R package easySdcTable. Provides an interface to the package sdcTable.

4

# How does it work?
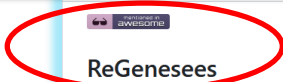


Curated list of software for official statistics

awesome

www.awesomeofficialstatistics.org

Social interactions

The right to wear the badge

Criteria

Working together

Open license

# uRos2023 -> uRos2024



Growth of the awesome list:

137
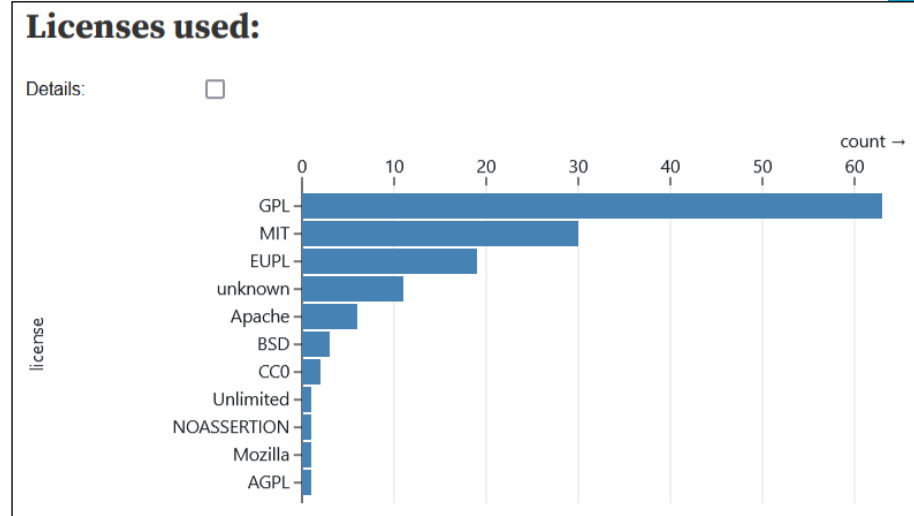
- Quite stable: growing to 137 entries (~+10)

- UNECE HLG-MOS open source project:

    - ESS principles on OSS (derived from awesome list) adopted

    - Charter, awesome list case study

    - License stats



Packages by GSBPM:



Licenses used:

https://observablehq.com/@olavtenbosch/visualizing-awesomeofficialstatistics-org

6

# Category "access to official statistics"

# Over 30 software packages, giving access to > 80 dataproviders, majority are R-packages

## 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 SDMX webservices. The package contains a list of SDMX access points of various national and international statistical institutes.

- `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. Parts in C++. By 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 providers, parse 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 and manipulation utilities.
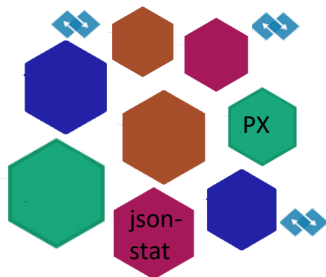
# Software layer surrounding Official Statistics

> 40 software packages
> 80 dataproviders
> majority R-packages
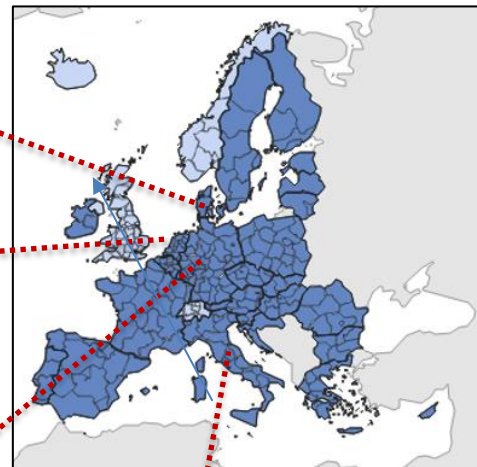> Some use standards (SDMX, PX, JSON-STAT)



**?**

User

PX

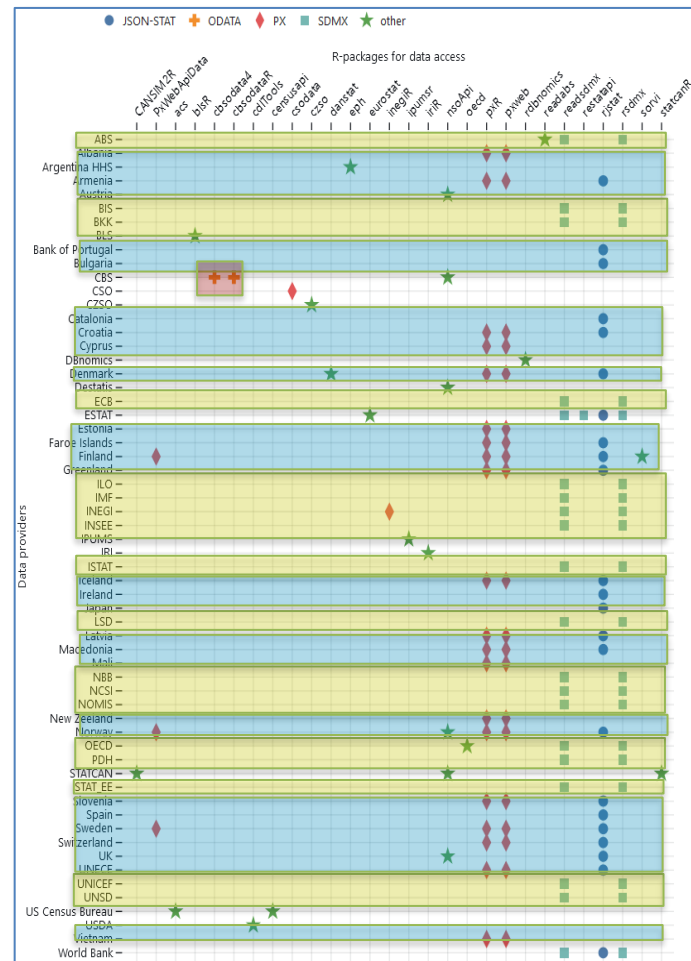json-stat

DK

NL

DE

IT

9

# "access to official statistics" software landscape



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

packages (40) x dataproviders (80) x standards (5)

- Standards:



- JSON-STAT/PX v. SDMX: almost disjunct worlds
- ODATA: CBS only, moving to SDMX

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

# uRos2023 -> uRos2024

- New packages:
  - TEMPO (Romania), cancensus (STATCAN), insee (INSEE), nomisr (UK)
- New dataproviders per standard:
  - JSON-STAT: Kenya;
  - PX: Kosovo, Liechtenstein, North Macedonia, Jordan, Philippines, Ghana;
  - SDMX: Lithuania, Estonia, UNESCO, Chile, Cambodia, El Salvador, FAO, UAE, Luxembourg, Maldiven, Malta, Thailand, UNESCAP, Uruguay, Greece.

# Features commonly offered

- **to_df:** getting (selections of) the data in a dataframe
- **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

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 | |

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

# One FAIR R-package for all data providers?



13

# Features x R-packages

# Features x R-packages

R software for data access

| Feature | CANSIM2R | OECD | PxWebApiData | TEMPO | acs | blsR | cancensus | cbsodata4 | cbsodataR | cdlTools | censusapi | csodata | czso | danstat | eph | eurostat | inegiR | insee | ipumsr | nomisr | pxR | pxweb | rdbnomics | readabs | readsdmx | restatapi | restatis | rjstat | rsdmx | sorvi | statcanR |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| to_df | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ |
| endpoint_hiding | ○ | ○ | | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | | | ○ | | ○ | ○ | | ○ | | ○ | ○ |
| catalogue_retrieval | | ○ | | | ○ | ○ | ○ | ○ | | ○ | | ○ | ○ | ○ | ○ | | ○ | | ○ | ○ | ○ | | | ○ | ○ | | ○ | ○ | | ○ | |
| search | ○ | | | | | ○ | | ○ | | | ○ | | | ○ | | ○ | | ○ | | ○ | | | | ○ | ○ | ○ | | | ○ | | |
| endpoint_queries | | ○ | | | | | | ○ | | | | | | | | | | | | | ○ | | | | | | ○ | | | | |
| local_queries | | | | | | | | | | | | | | | | | | | | | | | | | ○ | | | | | | |
| caching | | | | | | ○ | | ○ | | | ○ | | | ○ | | | | | | | | | | | ○ | | | | | | |
| cartographic queries | | | | | ○ | | ○ | | ○ | ○ | ○ | ○ | | | | ○ | ○ | | | | | | | | | | | | ○ | | |
| registry_access | | | | | | | | | | | | | | | | | | | | | | | | | ○ | | | | | | |

NL: will be re-implemented for .STAT/SDMX

Caching surprisingly often supported

Registry support only in SDMX

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

15

# Remarkable: scraping functions

- If not supported by API, then ***scraped*** and offered as a function in an R-package

- API developers: watch the R community for what's ***really needed*** ☺

Restatis

| gen_update_evas | *gen_update_evas* |
|---|---|

**Description**

Function to web scrape the EVAS numbers from the EVAS website and save them as a .rda file. Takes no parameters.

readabs

| scrape_abs_catalogues | *Helper function for* download_abs_data_cube *to s...* *able catalogues from the ABS website.* |
|---|---|

**Description**

This function downloads a new version of the lookup table used by s...

**Usage**

scrape_abs_catalogues()

censusapi

| listCensusApis | *Get useful dataset metadata on all available APIs as a data frame* |
|---|---|

**Description**

Scrapes https://api.census.gov/data.json and returns a dataframe that includes: title, description, name, vintage, url, dataset type, and other useful fields.
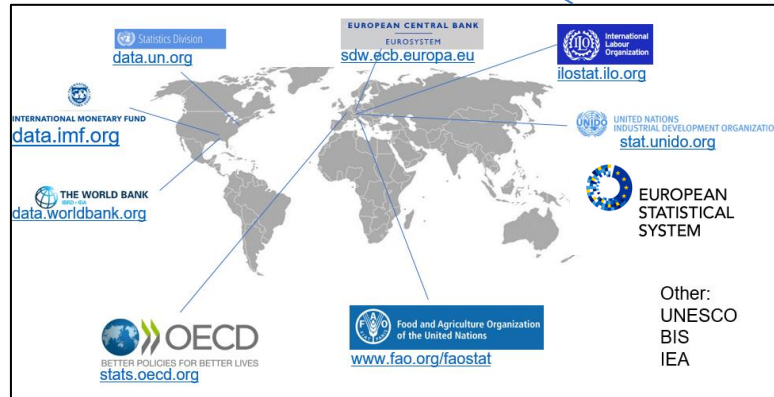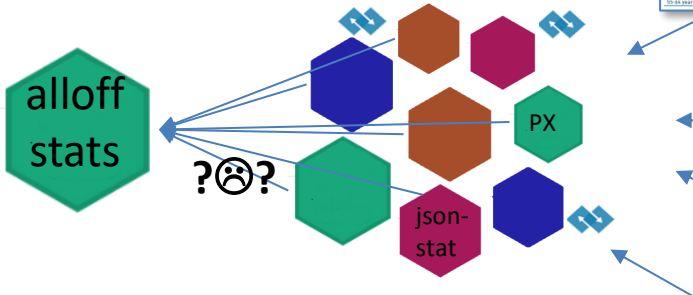
# One R-package for all data providers?

- Only for basic features:
  to_df, endpoint_hiding, catalogue_retrieval
- does it make sense?
- Would the 'user' use it?



DK

NL

DE

IT

alloff stats

???

json-stat

PX

User

data.un.org
sdw.ecb.europa.eu
ilostat.ilo.org
data.imf.org
stat.unido.org
data.worldbank.org

EUROPEAN STATISTICAL SYSTEM

Other:
UNESCO
BIS
IEA

stats.oecd.org
www.fao.org/faostat

# Observations

- Currently **> 30 R-packages** for access to **> 80 dataproviders**
- Further standardization on**: SDMX, JSON-stat, PX**
- Large **variety** but **common** features **widely** supported
- **GEO features** often supported, **caching** useful
- Some R-packages **scrape** metadata not covered in API -> added value
- Statistics Netherlands will redesign **cbsOdataR** for .Stat (SDMX)

- **one** package for all official statistics with **all** features -> too hard
- **'alloffstats'** package for all official statistics with **common** features -> ☺?

Common features: to_df, endpoint_hiding, catalogue_retrieval

# Wrap-up

# Wrap-up

- [www.awesomeofficialstatistics.org](http://www.awesomeofficialstatistics.org) Please ⭐ Starred 290
  - Spread the word and help maintain!

- Study on category "access to offstats":
  - Features R-packages reflect data science needs *in practice*
  - *Common features* broadly supported, but different patterns
  - '*alloffstats*' package with common features new goal???

Paper [Cosmos 2024](https://olavtenbosch.github.io/pdf/2024_COSMOS2024_ten_Bosch_PuhlishedRef.pdf) conference :
[https://olavtenbosch.github.io/pdf/2024_COSMOS2024_ten_Bosch_PuhlishedRef.pdf](https://olavtenbosch.github.io/pdf/2024_COSMOS2024_ten_Bosch_PuhlishedRef.pdf)

Olav ten Bosch          o.tenbosch@cbs.nl
Edwin de Jonge          e.dejonge@cbs.nl