# Access to official statistics from R: an overview

Statistics Netherlands

Olav ten Bosch, Edwin de Jonge
**uRos2023** 12-14 December 2023

# Contents



2018



2019

- ## What is the awesome list?
  - History, concept, status

- ## Category "access to official statistics"
  - Packages, data providers, standards, features
  - What are potential improvements in this software landscape?
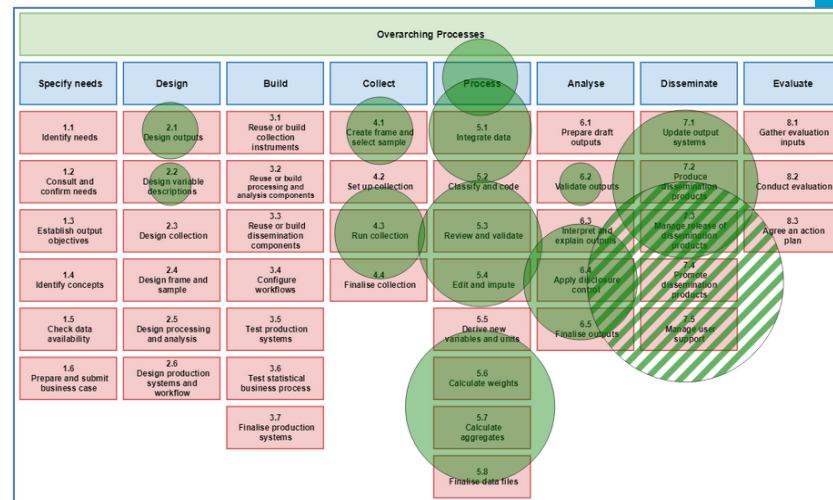
- ## Wrap-up

# Awesome list of official statistics software

- When: started at **UNECE SDE conference** april 2017 (The Hague)
- Why: to **collectively remember useful software** in official statistics
- Who: maintained by **statistical community**
- What: a **community approach** to knowledge management
- How:
  - Using **awesome concept**
  - A **public** list of awesome things
  - Clear and simple **criteria**
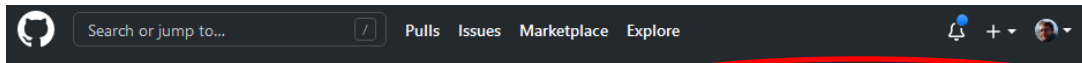  - **awesomeofficialstatistics.org**

# awesome packages by GSBPM

# What is the awesome list?

Social interactions

The right to wear the badge

Criteria

1. it is free, open source, and available for download and
2. it is confirmed to be used in the production of official statistics by at least one institute or it provides access to official statistics publications.

Working together

Open license

## Design frame and sample (GSBPM 2.1)

- `CRAN` `1.5-4 – a year ago` `license` `GPL (>= 2)`

  R package SamplingStrata. Optimal Stratification of Sampling Frames for Multipurpose Sampling Surveys.

- `CRAN` `1.0.5 – 7 months ago` `license` `EUPL`

  R package R2BEAT. Multistage Sampling Allocation and PS...

## Design variable descriptions (GSBPM 2.2)

- `GitLab` `no releases found` `last commit` `november` `license` `MIT Licens...`

  Excel SDMX_Matrix_Generator. Excel-based visual SDMX a...

## 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 cr...

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

  Java C++ Fortran and Delphi application T-ARG...

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

  R package sdcMicro. Disclosure control for stati...

- `CRAN` `0.32.6 – 4 months ago`

  R package sdcTable. Disclosure control for tabul...

## Sampling (GSBPM 4.1)

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

  R package sampling. Several algorithms for drawing survey samples, including a variety of unequal probabiltiy sampling designs (high entropy, systematic, Rao-Sampford, etc.), and calibrating design weights.

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

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

- `CRAN` `1.4.2 – 6 days ago` `license` `GPL-3`

  R package PracTools. Functions and datasets related to Valliant, Dever, and Kreuter (2018 2nd ed), *Practical Tools for Designing and Weighting Survey Samples*.

- `CRAN` `0.3.0 – 9 months ago` `license` `MIT + file LICENSE`

  R package prnsamplr. Coordinated stratified sampling using permanent random numbers (PRN's). Supports simple random sampling and probability-proportional-to-size sampling and includes a function for transforming...

## 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 l... pairs, comparing records, em-algorithm for estimating m- and u-probab... also be used for pre- and post-processing for machine learning methods

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

  R package RecordLinkage. Implementation of the Fellegi-Sunter method

- `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 linka... and the inclusion of auxiliary information. Documentation

- `CRAN` `0.9.12 – 13 days ago` `license` `GPL-3`

  R packages stringdist. Approximate string matching. Supports various st... Hamming, Levenshtein, optimal sting alignment), qgrams (q- gram, cosin... (Jaro, Jaro-Winkler). An implementation of soundex is provided as well.

- `CRAN` `0.1.6 – 4 years ago` `license` `MIT + file LICENSE`

  ... or similar matches. Allo...

... represented in data sou...

...Information from XBRL.

## Access to official statistics (GSBPM 7.4)

- R package rsdmx. Easy access to data from statistical organisations that support SDMX webservices. The package contains a list of SDMX access points of various national and international statistical institutes.
- R package and C++ readsdmx. Read SDMX into dataframes from local SDMX-ML file or web-service. By OECD.
- Python pandaSDMX. Python interface to SDMX that facilitates the acquisition and analysis of SDMX-2.1 compliant data and metadata.
- R package rjstat. Read and write data sets in the JSON-stat format.
- Python package pyjstat. Read and write JSON-stat.
- Java module json-stat.java Read and write JSON-stat. By Statistics Norway.
- R package oecd Search and Extract Data from the OECD
- R package sorvi Finnish Open Government Data Toolkit
- R package eurostat Tools to download data from the Eurostat database together with search and manipulation utilities.
- R package acs Download, Manipulate, and Present American Community Survey and Decennial Data from the US Census.
- R package inegiR Access to data published by INEGI, Mexico's official statistics agency.
- R package cbsodataR. Access to Statistics Netherlands' (CBS) open data API from R.
- Node.js package cbsodata.js. Access to Statistics Netherlands' (CBS) open data API from js.
- Python package cbsodata.py. Access to Statistics Netherlands' (CBS) open data API from Python.
- R package censusapi A wrapper for the U.S. Census Bureau APIs that returns data frames of Census data and metadata.
- R package nsoApi builds on other packages to access data from official statistics and tries to harmonize the API.
- R package CANSIM2R. Extract CANSIM (Statistics Canada) tables and transform them into readily usable data.
- Python package pyscbwrapper. Access to the open data API of the Swedish Statistical Institute
- R package pxweb. Generic interface for the PX-Web/PC-Axis API used by many National Statistical Agencies.
- R package PxWebApiData. Easy API access to e.g. Statistics Norway, Statistics Sweden and Statistics Finland.
- R package rdbnomics. Access to the DB.nomics database which provide macroeconomic data from 38 official providers such as INSEE, Eurostat, Wolrd bank, etc.
- R package readabs Download data from the Australian Bureau of Statistics.
- R package destatiscleanr. Clean csv files from Genesis, the database of the Federal Statistical Office of Germany (Destatis) and its regional outlets.
- R package statcanR. An R connection to Statistics Canada's Web Data Service. Open economic data (formerly CANSIM tables) are accessible as a data frame in the R environment.
- R package cdlTools. Downloads USDA National Agricultural Statistics Service (NASS) cropscape data for a specified state.
- Java package SDMX Connectors. Browse SDMX data providers, build your queries and get data directly in your favourite tool (R, SAS, Matlab, Stata and Excel). By Banca d'Italia.
- Node.js package sdmx-rest. This library allows to easily create and execute SDMX REST queries from a JavaScript client application.
- R package csodata Download data from Central Statistics Office (CSO) of Ireland.
- R package iriR. Client for the EU Industrial Research and Innovation Scoreboard.

5

# Category "access to official statistics"

- Over 30 software packages, 28 R-packages
- Cran standard documentation used as starting point

rsdmx: Tools for Reading SDMX Data and Metadata

OECD: Search and Extract Data from the OECD

czso: Use Open Data from the Czech Statistical Office in R

readsdmx: Read SDMX-XML Data

rjstat: Handle 'JSON-stat' Format in R

readabs: Download and Tidy Time Series Data from the Australian Bureau of Statistics

sorvi: Functions for Finnish Open Data
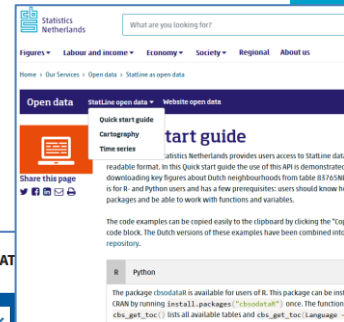
eurostat: Tools for Eurostat Open Data

restatapi: Search and Retrieve Data from Eurostat Database

inegiR: Integrate INEGI's (Mexican Stats Office) API with R

pxweb: R Interface to PXWEB APIs

csodata: Download Data from the CSO 'PxStat' API

danstat: R Client for the Statistics Denmark Databank API

# "access to official statistics" R-software landscape (1)

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

R-packages (28) x dataproviders (60) x standards (5)

- Standards:



https://observablehq.com/@olavtenbosch/access-to-official-statistics-from-r

# "access to official statistics" R-software landscape (2)

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

R-packages (28) 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-from-r

# "access to official statistics" R-software landscape (3)

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

R-packages (28) x dataproviders (60) x standards (5)

- Standards:

JSON-STAT   ODATA   PX   SDMX   other

- Few data providers support multiple standards: **ESTAT, WB**

https://observablehq.com/@olavtenbosch/access-to-official-statistics-from-r

# "access to official statistics" R-software landscape (4)

- 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-from-r

# Features commonly offered

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

Future research: Packages x endpoints x standards x features

# Considerations

- Offstats landscape grows towards *standardisation: SDMX, JSON-stat, PX*
- R user can choose from (at least) *28 packages*, each offering specific functionality to *60 data sources* in total
- Many NSIs offer *targeted R-packages* for data and metadata access
- There is *no 'one-for-all' R-package* that provides access to all official statistics data providers
- *Common features* identified

- Dream: *one generic R-packages* to all official statistics from all dataproviders supporting maximum FAIRness

FAIR: Findable, Accesible, Interoperable, Reusable

# International

## 71st CES



EU OSS strategy

- IE, NL, NO, PL, UK, CA



ESS principles on OSS (OS4OS group)

1. OSS by default

2. Work in the open

3. Improve and give back

4. Think general statistical building blocks

5. Test, package and document

6. Choose permissive

7. Promote

# Wrap-up

- [www.awesomeofficialstatistics.org](www.awesomeofficialstatistics.org) 👓 Please ☆ Star 239 !
  - Spread the word and help maintain!
- Meta analysis on awesome list "Access to offstats":
  - *standardisation*: SDMX, JSON-stat, PX
  - *Common features* identified; often reimplemented
  - *No* 'one-for-all' R-package for access to *all* offstats data
  - Lets work *together* to improve access to offstats from R

Extended abstract: [https://olavtenbosch.github.io/pdf/2023_Uros_tenBosch_Access.pdf](https://olavtenbosch.github.io/pdf/2023_Uros_tenBosch_Access.pdf)

Follow-up paper to be published at: [http://cosmos-conference.org/2024/](http://cosmos-conference.org/2024/)

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