# Open source statistical software at the statistical office

Mark van der Loo
Statistics Netherlands
m.vanderloo@cbs.nl

CREATED WITH FREE SOFTWARE

ISI 2017
MARRAKECH
61° WORLD STATISTICS CONGRESS

# What is open source software?

## Free and Open Source Software

- ▶ Use
- ▶ Study
- ▶ Change
- ▶ Redistribute

# FOSS is driving modern stats and data science



And much, much more. . .
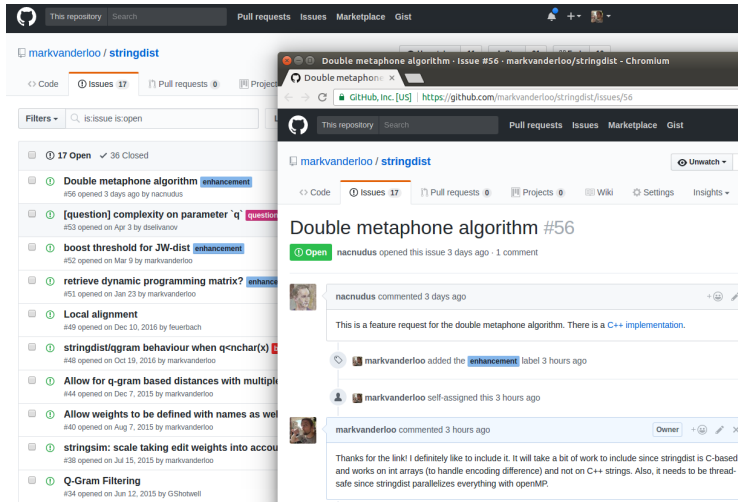
# Communities (1) Social Coding with github

# Communities (2) Q&A with stackoverflow

# Communities (3) news & discussions on Twitter

# Role of commercial parties, foundations



And many,many more, . . .
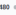
# Motivations

## Wikipedians

*Wikipedians enjoy a sense of accomplishment, collectivism, and benevolence while working with exceptional freedom and ease. The values of reputation, community, reciprocity, altruism and autonomy are fostered by both the people and the technology[...].*
(Kutzetsnov, 2006)

## Commecial parties

*[Companies] expect to benefit from their expertise in some segment whose demand is boosted by the success of a complementary open source program.* (Learner, 2000)

# Motivations for Official Statistics

## Use

- ▶ Economic (its free!)
- ▶ New hires
- ▶ Supporting community

## Contribute

- ▶ Solving shared problems
- ▶ Many eyes make all bugs shallow
- ▶ Influence, reputation
- ▶ Built with tax payer's money

# FOSS for official statistics

Awesome list 😎

▶ Community effort

▶ Curated

▶ 50+ Software packages

▶ Covering 14 GSBPM areas

▶ Growing
  – 75 commits
  – 5 PR's
  – 6 contributors



www.awesomeofficialstatistics.org

# When is it awesome?

You may be awesome when. . .

▶ Free, open source, available for download

▶ Used in at least one statistical institute for production or, offers access to official statistics

▶ Relatively easy to install and use (for non-dev's)

▶ Actively maintained

▶ At least one stable release

# What's on the Awesome List?

**Statistical data editing and imputation (GSBPM 5.3 | 5.4)**

- R package validate. Rule ma
- R package errorlocate. Error
  - Uses validate rule definiti
  - supports categorical and
  - supports linear equalities
  - Configurable backend fo
- R package VIM. Visualisation
  - Advanced visualisation o
  - Imputation using (robust
  - Imputation using severa
- R package VIMGUI. Graphic
- R package simputation. Simp
  - Allows to easily combin
  - Supports regression (sta
    randomForest, EM-base
    user-defined methods a
- R package SeleMix. Detection of outliers and influential err
- R pac
- R pac

**Estimation and weighting (GSBPM 5.6 | 5.7)**

- R package survey. Weighting and estimation for complex
  estimator variance. See also R package srvyr for integra
- R package hbsae. Small area estimation based on hiera
- R package rsae. Small area estimation based on (robust
- R package calibrateSSB. Calculate weighs and estimate

**Time series and seasonal adjustment (GSBPM 5.6 | 5.7)**

- X-13ARIMA-SEATS Seasonal adjustment software prod
- R package seasonal. Interface to the `x13-ARIMA-SEATS`
- R package x12. Alternative interface to the `x13-ARIMA-`
  series.
- JDemetra+ The seasonal adjustment software officially r

**Process (GSBPM 5)**

- Java application Java-VTL. A partial implementation
  draft specification. By Statistics Norway.

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

- R package RecordLinkage. Implementation of the F
- R packages stringdist and fuzzyjoin allow for matchi
- R package XBRL. Extraction of Business Financial

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

- R package rsdmx. Easy access to data fr
  contains a list of SDMX access points of
- R package oecd Search and Extract Data
- R package sorvi Finnish Open Governme
- R package eurostat Tools to download da
- R package acs Download, Manipulate, an
  Census.
- R package inegiR Access to data publish
- R package cbsodataR. Access to Statisti
- npm package cbsodata.js. Access to Sta

**Sampling (GSBPM 4.1)**

- R package sampling. Several algorithm

**Scraping for Statistics (GSBPM 4.3)**

- Java application URLSearcher. An appl
- Java application URLScorer. Gives a ru
- node.js tool RobotTool A tool for check
- node.js package S4Si
  functionalities for stati

**Output validation (GSBPM 6.2)**

- R package validate. Rule management an

**Statistical disclosure control (GSBPM 6.4)**

- Argus and SDC Tools. Tools like Tau-Argu
  and the Statistical disclosure control netw
- R package sdcMicro. Disclosure control fo
- R package sdcTable. Disclosure control fo
- R package simPop. Simulation of syntheti

**Statistical Dissemination (GSBPM 7.2)**

- SDMX Converter. Converter between diffe
- SDMX-RI. Framework for disseminating d
- R package rsdmx. Writing SDMX from R.
- StatMiner, Experimental visualization fran
- SDMX-JSON. JSON variant of SDMX. Th
- JSON-Stat. Lightweight JSON based mes

# In the works...

# FOSS policy at Statistics Netherlands (in short)

## Usage

Selection and introduction follows the same procedure as for COTS
(commercial off-the-shelve).

- ▶ R, Python, `node.js`

## Contributing

When relevant to Statistics Netherlands, with positive business case.

- ▶ R, `node.js`

# Deployment of R and Python at Statistics Netherlands

- ▶ Central read-only folder for executables.
- ▶ All users have access to the same version with curated list of libraries installed.
- ▶ Scripts can be prepared and integrated for non-developers with ease.
- ▶ Repositories (CRAN, Anaconda) on internal website, updated frequently.
- ▶ Old versions stay available for some time so existing applications stay working.

# Example using R in data editing: tools and roles



Users

IT dpt — Microsoft .net — Graphical user interface

Rules, metadata, control parameters...

TXT

Methodology — R → R → R → R → R → R

IT dpt — Microsoft SQL Server — Database

# Example contributions

▶ R packages
  — data editing: `validate`, `simputation`, `errorlocate`, `deductive`, `dcmodify`
  — data logging: `lumberjack`
  — small area estimation: `hbsae`
  — datavis: `tabplot`, `tabplotd3`, `tmap`, . . .

▶ node.js packages
  — scraping: `RobotTool`, `S4Robo`
  — dashboard: `StatMine`

▶ . . .

# So you want to contribute?

## Here are some options

1. Use it (& send a thumbs-up!)
2. Advocate
   - Tell your friends and colleagues
   - Write blog posts / articles / presentations
   - Social media (twitter...)
3. File bug reports, suggestions
4. Add code to an existing project
5. Start your own project

# FREE TIP!



What a brilliant idea! Here, eat my shoe.

your e cards
someecards.com

9GAG.COM/GAG/2927719

## Don't work alone

▶ Join your local community
  − meetups, news letters, hackatons
▶ Set up a community in your institute
  − Local wiki, user meetings, hackatons, ask-the-expert