



Why and for what purpose R in Official Statistics?

- uRos 2024 -

Sandra Barragán



- Statistics Spain (INE) -

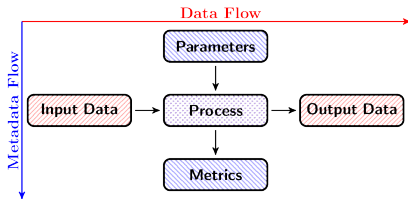


INē


Instituto Nacional de Estadística


Conclusions: Take home messages

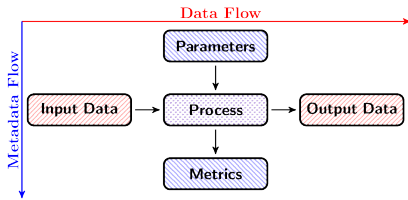
- ▶  is open source, user friendly (easy to learn), constantly evolving. R has a bunch of packages implementing statistical methods and a wide community of users.
- ▶  facilitates a proper **implementation** and **automation** of **modular processes** in the production of Official Statistics.



Conclusions: Take home messages

-  is open source, user friendly (easy to learn), constantly evolving. R has a bunch of packages implementing statistical methods and a wide community of users.

-  facilitates **HOW?** in and **automation** of **modular processes** in the production of Official Statistics.



- ▶ Introduction
- ▶ Industrialization of Official Statistical production
- ▶ Implementation
- ▶ Use Cases implemented in R
- ▶ Conclusions

Introduction



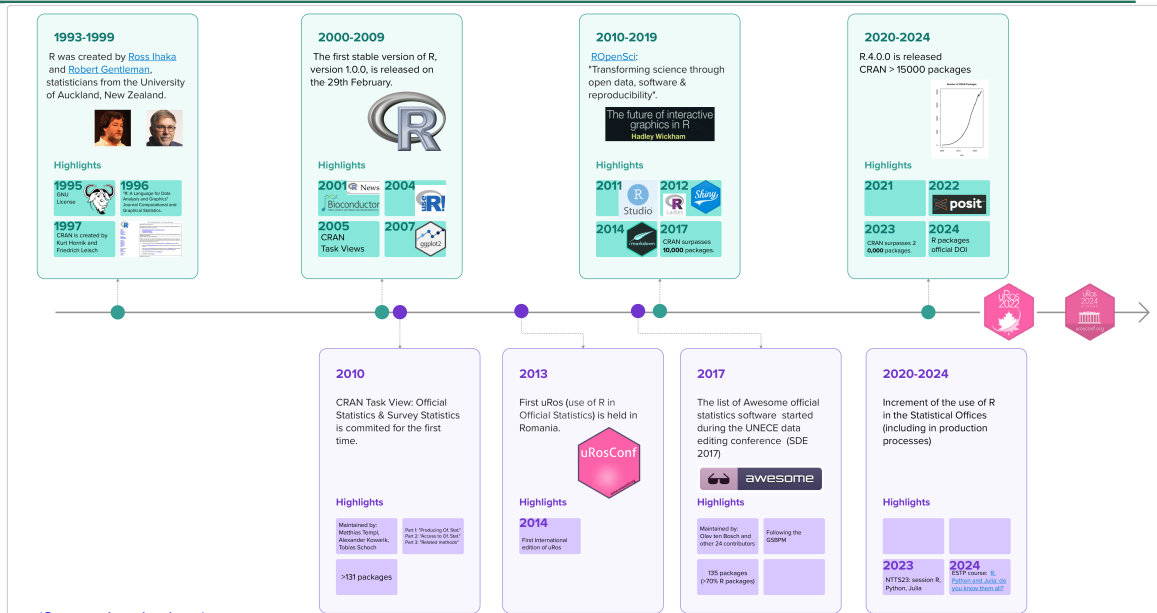
Run, or he's going to tell us about
again!

R

Some people are afraid of R language because it is something new (for them).

NEW???

... not so new...



► CRAN Task View: Official Statistics & Survey Statistics

<https://cran.r-project.org/web/views/OfficialStatistics.html>

CRAN Task View: Official Statistics & Survey Statistics

Maintainer: Matthias Templ, Alexander Kowarik, Tobias Schoch

Contact: matthias.templ at gmail.com

Version: 2023-02-19

URL: <https://CRAN.R-project.org/view=OfficialStatistics>

Source: <https://github.com/cran-task-views/OfficialStatistics/>

Contributions: Suggestions and improvements for this task view are very welcome and can be made through issues or pull requests on GitHub or via e-mail to the [guide](#).

Citation: Matthias Templ, Alexander Kowarik, Tobias Schoch (2023). CRAN Task View: Official Statistics & Survey Statistics. Version 2023-02-19. URL

Installation: The packages from this task view can be installed automatically using the [ctv](#) package. For example, `ctv::install.views("OfficialStatistics")`, `ctv::update.views("OfficialStatistics")` installs all packages that are not yet installed and up-to-date. See the [CRAN Task View Initiative](#) for

This CRAN Task View contains a list of packages with methods typically used in official statistics and survey statistics. Many packages provide functions for more strict categorization and packages may be listed more than once.

The task view is split into several parts

- First part: [“Producing Official Statistics”](#). This first part is targeted at people working at national statistical institutes, national banks, international organizations and using methods from survey statistics. It is loosely aligned to the [“Generic Statistical Business Process Model”](#).
- Second part: [“Access to Official Statistics”](#). This second part’s target audience is everyone interested to use official statistics results directly from within R.
- Third part: [“Related Methods”](#) shows packages that are important in official and survey statistics, but do not directly fit into the production of official statistics. It also includes a collection of packages that are loosely linked to official statistics or that provide limited complements to official statistics and survey methods.

First Part: Production of Official Statistics

1 Preparations/ Management/ Planning (questionnaire design, etc.)

- [questionr](#) package contains a set of functions to make the processing and analysis of surveys easier. It provides interactive shiny apps and addins for data retrieval and analysis.
- [surveydata](#) makes it easy to keep track of metadata from surveys, and to easily extract columns with specific questions.

► Awesome official statistics software

<https://github.com/SNStatComp/awesome-official-statistics-software>

Awesome official statistics software 


An awesome list of open source software for official statistics.

An item on this list is awesome because it is

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

We prefer software that is easy to install and use, has at least one stable version, and is actively maintained.

Design frame and sample (GSBPM 2.1)

- CRAN 1.5-4 – 5 months ago 

R package [SamplingStrata](#). Optimal Stratification of Sampling Frames for Multipurpose Sampling Surveys.

- CRAN 1.0.4 – a year ago 

R package [R2BEAT](#). Multistage Sampling Allocation and PSU Selection.

Design variable descriptions (GSBPM 2.2)

- ▶ **Eurostat: Working group on Open Source Software for the use in Statistics**
- ▶ **European Comission: Open source software strategy**
https://commission.europa.eu/about-european-commission/departments-and-executive-agencies/informatics/open-source-software-strategy_en
- ▶ **European Comission: Open Source Observatory (OSOR)**
<https://joinup.ec.europa.eu/collection/open-source-observatory-osor>



The 12th International Conference Use of R in Official Statistics

uRos2024

27-29 November 2024, Piraeus, Greece

November 29, Friday

8³⁰ – 9⁰⁰ Walk In

9⁰⁰ – 10⁰⁰



Keynote Speaker 2: Romain Lesur, INSEE France

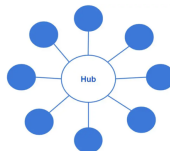
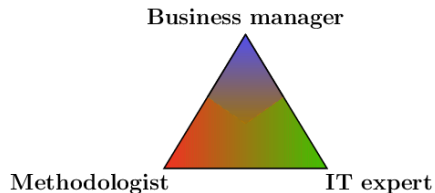
Presentation - Moving away from SAS: an opportunity to modernise the practices of statisticians



Industrialization of Official Statistical production

Industrialization of official statistical production

- Standardization
- Automation
- Modularity



- ▶ Standardization
- ▶ Automation
- ▶ Modularity

Advantages:

- + Reuse of processes.
- + Easier to detect bugs and errors.
- + Maintenance.
- + Resource optimisation.
- + Evolution.

Focused on the process instead of the product, but flexible enough to be adapted to the particularities of each product.

- ▶ **Parameterisation:** Same process can be executed for different products without changes in the code, just changing the parameters.
- ▶ **Reproducibility:** consistent execution when it is repeated.
- ▶ **Documentation:** every part of the process is documented.
- ▶ **Quality** control.
- ▶ Continuous integration and deployment: **CI/CD**.

See paper in JOS: *Data organisation and process design based on functional modularity for a standard production process*, Salgado et al. [2018].

Implementing manual processes to be executed by a machine.
Scheduling task to be executed based on a calendar.

Ex1: Generate files from collection to a repository at a time everyday.

Ex2: Trigger the execution of selective editing once all the files are available and a date is reached.

Scheduling R processes in the system from R directly:

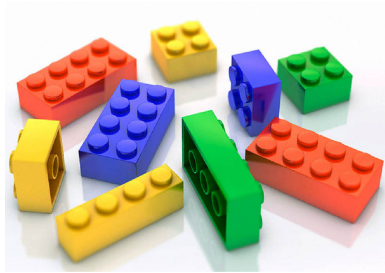
- ▶ `cronR` in Linux:
 - ▶ Create a job: `cron_add()`.
 - ▶ Remove a job: `cron_rm()`.
 - ▶ List all jobs scheduled: `cron_ls()`.
- ▶ `taskscheduleR` in Windows:
 - ▶ Create a job: `taskscheduler_create()`.
 - ▶ Remove a job: `taskscheduler_delete()`.
 - ▶ List all jobs scheduled: `taskscheduler_ls()`.



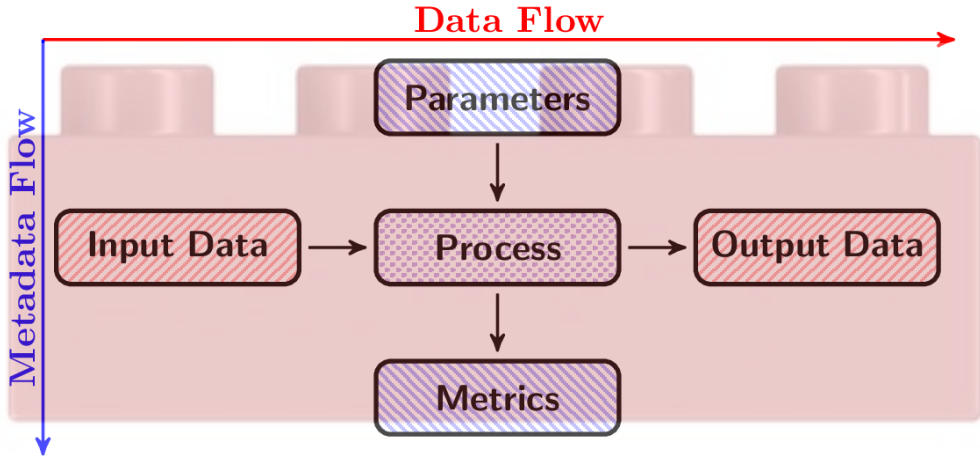
Different from the Background jobs of RStudio:

<https://docs.posit.co/ide/user/ide/guide/tools/jobs.html>.

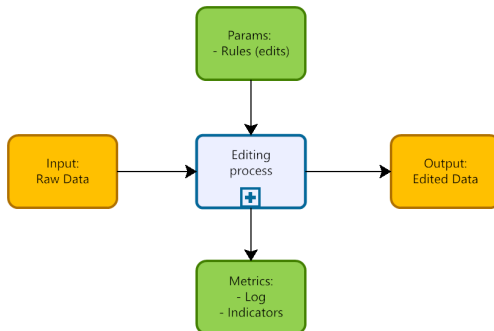
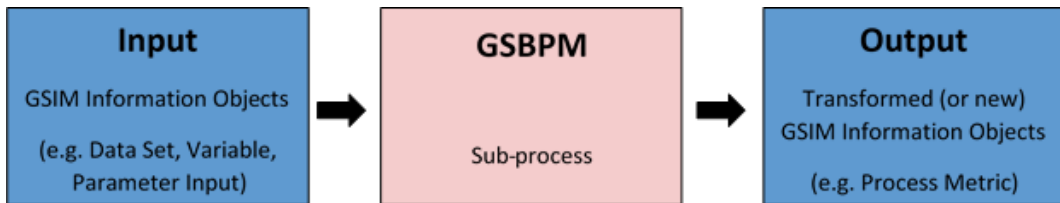
Industrialization of official statistical production: Modularity



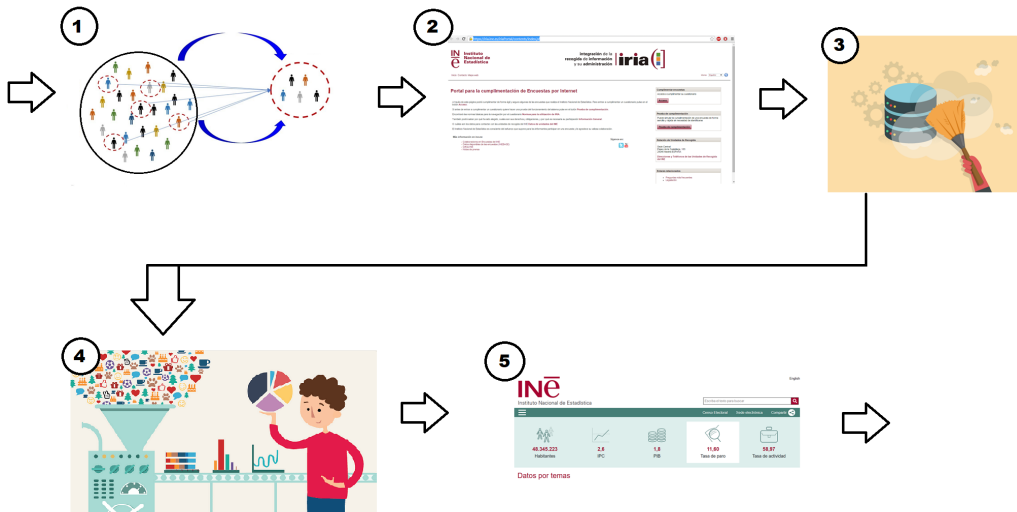
Industrialization of official statistical production: Modularity



Industrialization of official statistics: Modularity

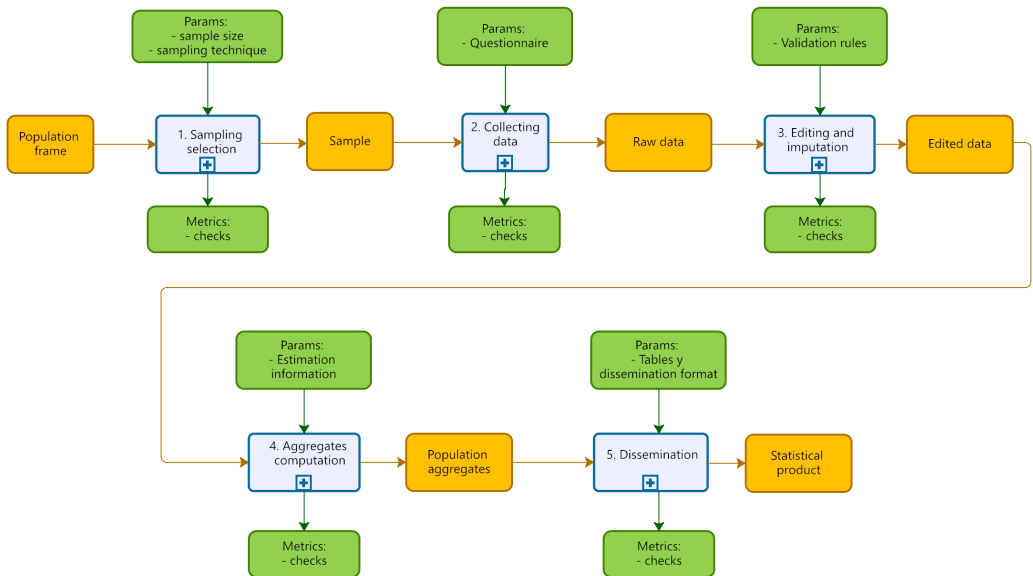


Industrialization of official statistics: Modularity



Industrialization of official statistical production: Modularity

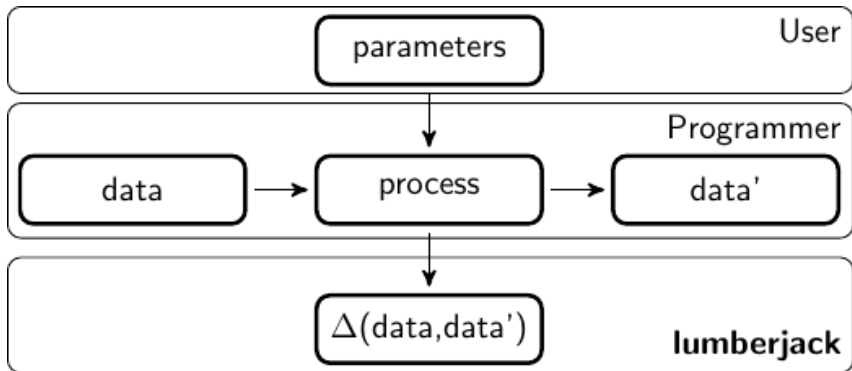
Modular statistical production process example



Modularity in R

Tip!

R package `lumberjack`: <https://github.com/markvanderloo/lumberjack>



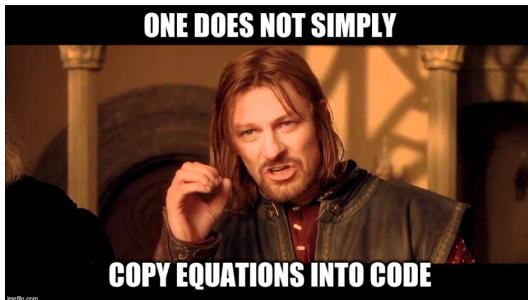
Implementation

*If you can't describe what you are doing as a process,
you don't know what you're doing.*

W. Edwards Deming

Implementation

Implementing methods is not trivial. See van der Loo [2020].



Principles¹:

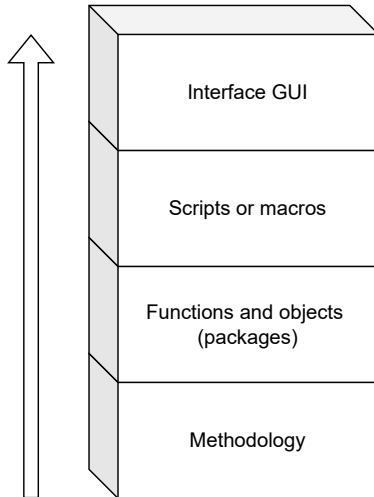
- ▶ Funcionality
- ▶ Object-oriented



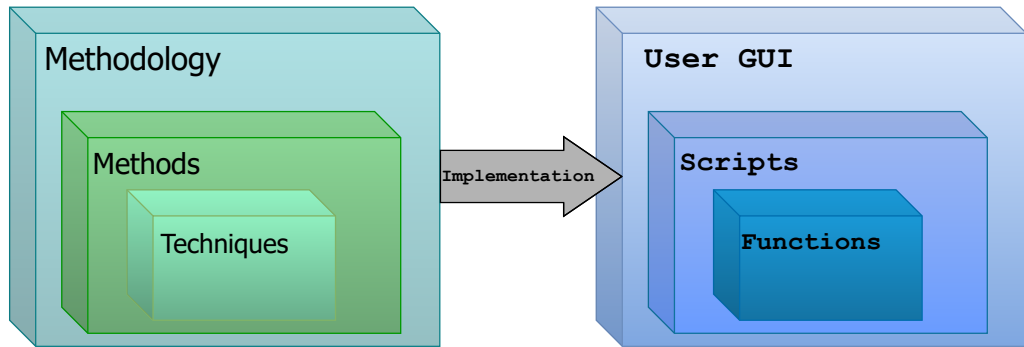
¹ See Tucker and Noonan [2007]

² See Matloff [2011]

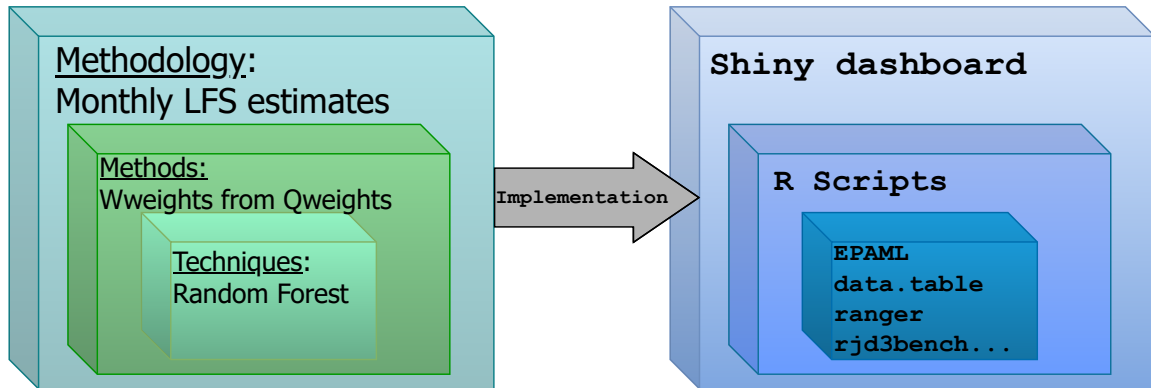
Traditional levels of implementation



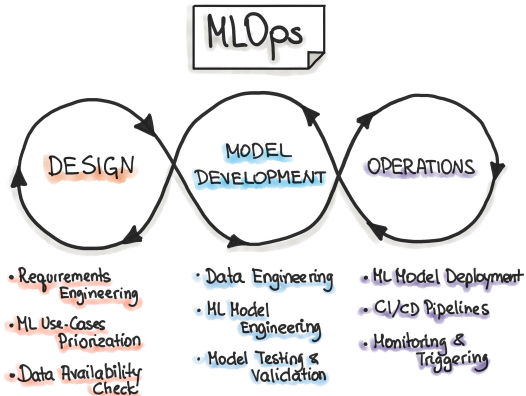
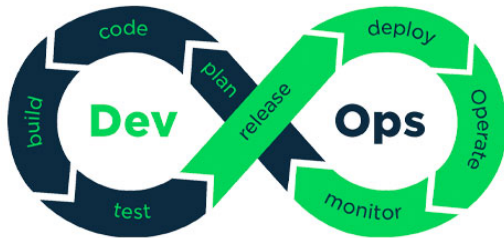
Implementation



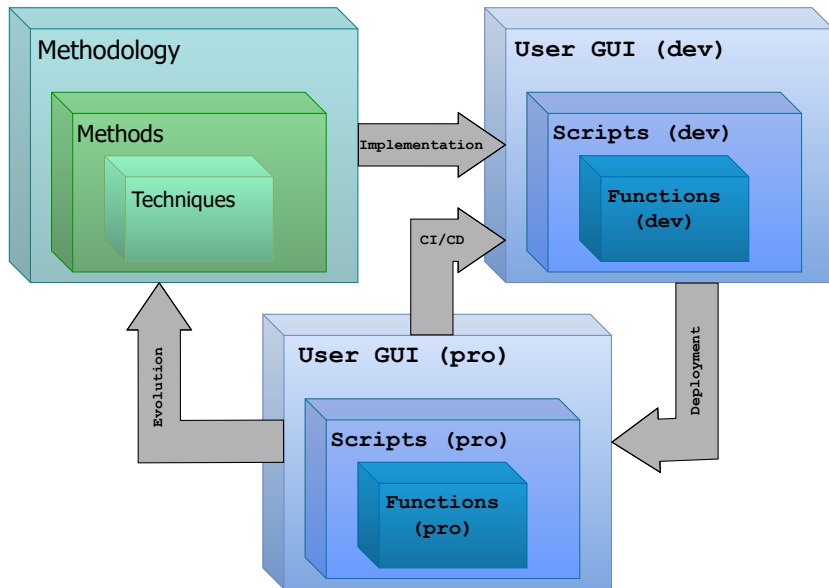
Implementation: example



Implementation for production: DevOps and MLOps



Implementation for production: StatOps



- ▶ **P**roduction vs. **p**roduction
- ▶ Project oriented (usethis, here)
- ▶ Environments: production, test, development.

More details in: [R in production - Hadley Wickham](#)

Relevant considerations:

1. Running code on another machine
 - ▶ Someone else. R package [pointblank](#).
 - ▶ Debugging, Logging
 - ▶ Configuration (Authentication, package installation)
2. Running code multiple times (Changes over time)
 - ▶ Data schema
 - ▶ Package versions: R package [pak](#).
 - ▶ System libraries
 - ▶ Configuration (OS, Requirements...)
3. Shared responsibility
 - ▶ Parquet format: R packages [arrow](#) and [nanoparquet](#).
 - ▶ Code review: R packages [cyclocomp](#) and [covr](#).
 - ▶ Git

Implementation: Some good practices

- ▶ Version control: Git (github, gitlab)
- ▶ Code organization:
 - ▶ Naming: variables names.
 - ▶ Indentation: margins, spaces.
 - ▶ Well structured documentation (comments).
 - ▶ Functions and modularity.
 - ▶ Check arguments validity.
- ▶ Files and folders:
 - ▶ Naming: files names.
 - ▶ Well structured organization.
- ▶ Good habits:
 - ▶ Release memory.
 - ▶ Executing from the beginning.

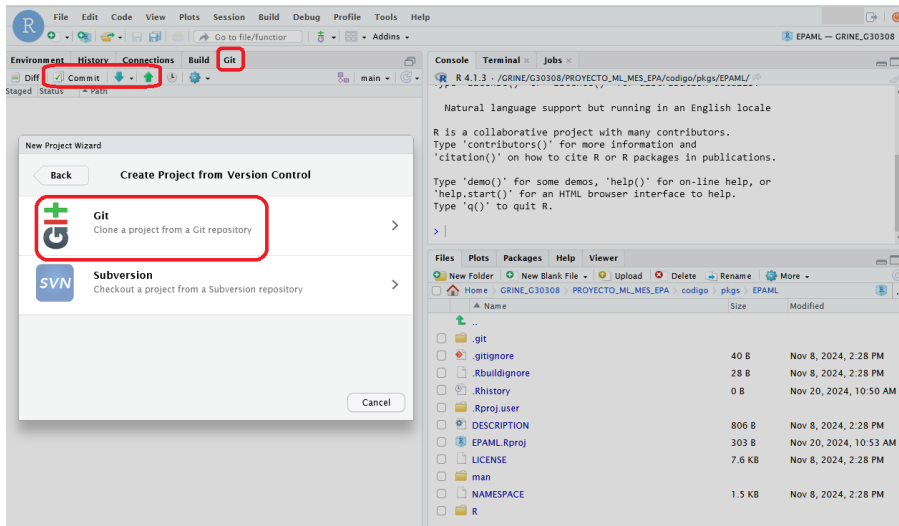


RStudio eases your work

```
##### .....  
##### PROCESO  
## * Leer microdatos  
schema_raw.stSchema <- StxlSxToSchema(  
  xlsName = file.path(path_param, schema_raw_filename),  
  sheetToRead = schema_raw_sheetname  
)  
  
microdata_raw.dt <- fread_fwf(  
  filename = file.path(path_data, data_raw_filename),  
  StfwfSchema = schema_raw.stSchema,  
  outFormat = 'data.table', perl = TRUE,  
  validate = validate_data, convert = convert_types  
)  
  
microdata_recontacto.dt <- fread_fwf(  
  filename = file.path(path_data, data_recontacto_filename),  
  StfwfSchema = schema_raw.stSchema,  
  outFormat = 'data.table', perl = TRUE,  
  validate = validate_data, convert = convert_types  
)  
  
## * Leer estructura de la cnae09  
cnae09.dt <- as.data.table(read.xlsx(  
  xlsFile = file.path(path_param, cnae09_estructura_filename),  
  sheet = cnae09_estructura_sheetname))  
  
cnae09_clases.dt <- cnae09.dt[nchar(COD_CNAE2009) == 4]  
setnames(cnae09_clases.dt, 'COD_CNAE2009', 'cnae09_clase')  
cnae09_grupos.dt <- cnae09.dt[nchar(COD_CNAE2009) == 3]  
setnames(cnae09_grupos.dt, 'COD_CNAE2009', 'cnae09_grupo')  
cnae09_divisiones.dt <- cnae09.dt[nchar(COD_CNAE2009) == 2]  
setnames(cnae09_divisiones.dt, 'COD_CNAE2009', 'cnae09_division')
```

```
#####  
#####  
#####  
PAQUETES  
.....  
RUTAS  
.....  
FUNCIONES  
.....  
PARAMETROS  
* Parámetros para lectura  
* Parámetros generales  
* Parámetros para depuración DOS  
* Parámetros para depuración selectiva  
* Parámetros para depuración automática  
.....  
PROCESO  
* Leer microdatos  
* Leer estructura de la cnae09  
* Generar variables derivadas  
:.....;  
DOS  
* Leer edits  
* Aplicar edits y tratar errores  
** sample-wise  
** record-wise  
*** norden_unico  
*** clase_longitud  
*** (variable)_formato  
*** (variable)_recorrido  
*** cifraNeg(, Nac, Ext)_NA -- balance  
*** ocupados(, Fijos, Event)_NA -- balance  
** imputación NA  
:.....;  
Depuración selectiva  
:.....;  
Depuración interactiva  
:.....;  
Depuración automática  
** Localización de errores  
** Reemplazamiento de errores por NA  
** Imputar NA en edits de balance activos  
:.....;
```

RStudio eases your work



Implementation: Some good practices

More details in:

- ▶ [R Style Guide de Google](#).
- ▶ [Advanced R by Hadley Wickham](#) (Wickham [2019]).
- ▶ [Writing Better R Code workshop](#).
- ▶ Boswell and Foucher [2011].
- ▶ R packages [goodpractice](#).





This is not just code...

program
reuse version
test document
metadata share
model process
automate
?

How are resources (mis)used:
Duplication? Redundancy?
Cost?

No thanks!

We are
too busy



➤ but also **culture & management**



"I want you to find a bold and innovative way
to do everything exactly the same way
it's been done for 25 years!"



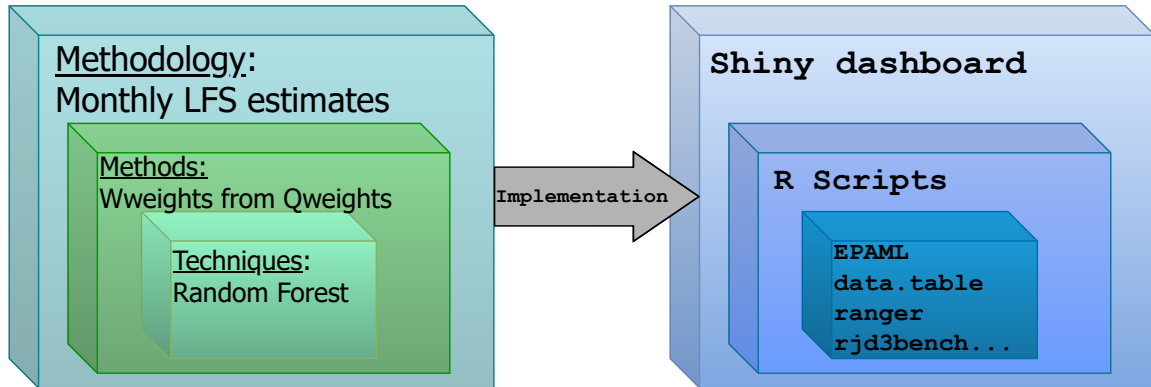
"What if we don't change at all ...
and something magical just happens?"

Use Cases implemented in R

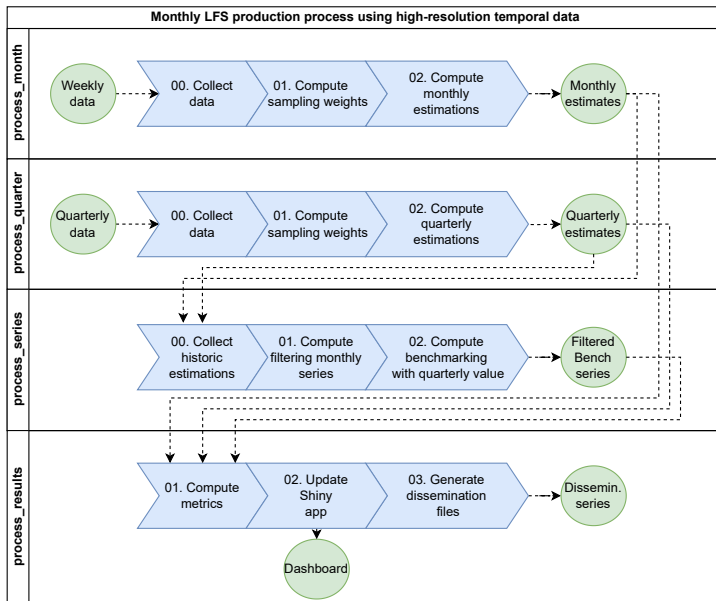
Use Case 1: Monthly LFS production process using high-resolution temporal data

- ▶ Novel end-to-end statistical production process that combines machine learning techniques, time series filtering, and benchmarking.
- ▶ Monthly (un)employment statistics:
 - ▶ Gender: male, female.
 - ▶ Age: 15-24, 25-74.
- ▶ LFS in Spain. Quarterly.
 - ▶ Weekly interviews.
 - ▶ Survey data: Quarterly sample.
- ▶ Significant improvement in time granularity.

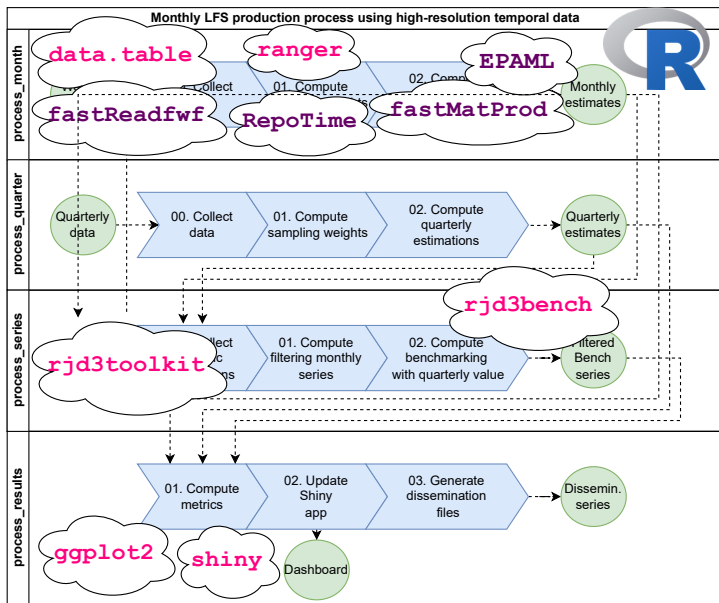
R as facilitator to implement an end-to-end statistical production process



R as facilitator to implement an end-to-end statistical production process



R as facilitator to implement an end-to-end statistical production process



High-resolution temporal data. Weekly Sampling Weights.

From quarterly to weekly.

1. Intermediate weights:

$$\omega_k^{o[W]} = \frac{N_{U_h}^{[M]}}{\sum_{k \in r_h^{[W]}} \frac{d_k^{[Q]}}{\pi_k^{[W \rightarrow Q]}} n_k} \frac{d_k^{[Q]}}{\pi_k^{[W \rightarrow Q]}}.$$

01. Compute sampling weights

where $\pi_k^{[W \rightarrow Q]}$ is the conditional inclusion probability for unit k to be interviewed in week W conditioned on its quarterly sample member.

$$\hat{\pi}_k^{[W \rightarrow Q]} = \sum_{r=1}^6 \mathbb{P}_k [(W, r) | S, P, R],$$

where r are the rotation group, S is the *stratum*, R is the *region* (NUTS2) and P is the *province* (NUTS3). This conditional probability is estimated with a Random Forest model (see [Murphy, 2012]) using R package [ranger](#).

2. Calibrated weights: $\omega_k^{[W]}$ computed using R package [fastMatProd](#) (in-house dev based on [calib](#) R package).

High-resolution temporal data. Aggregates and variances.

- ▶ Weekly aggregate for week w : $\hat{Y}_w^{[W]} = \sum_{k \in s_w} \omega_k^{[W]} y_k$.
- ▶ Monthly aggregate for month m based on weekly aggregates:

$$\hat{Y}_m^{[M]} = \frac{1}{n_m^W} \sum_{w \in W_m} \hat{Y}_w^{[W]}.$$

where W_m is the subset of weeks in month m , with cardinal n_m^W .

- ▶ Filtering time series with ARIMA-model-based decomposition: $Y = T + S_{13} + I$. ARIMA models for T , S_{13} and I are given by the canonical decomposition for the Y model. Filtering is done with the R package [rjd3toolkit](#).
- ▶ Benchmarking is used once the quarterly estimates are computed to correct the provisional monthly estimates to be coherent with the corresponding quarterly estimate. Benchmarking is done with the R package [rjd3bench](#).

High-resolution temporal data. Aggregates and variances.

- ▶ Weekly aggregate for week w : $\hat{Y}_w^{[W]} = \sum_{k \in s_w} \omega_k^{[W]} y_k$.
- ▶ Monthly aggregate for month m based on weekly aggregates:

$$\hat{Y}_m^{[M]} = \frac{1}{n_m^W} \sum_{w \in W_m} \hat{Y}_w^{[W]}.$$

02. Compute
monthly
estimations

where W_m is the subset of weeks in month m , with cardinal n_m^W .

- ▶ Filtering time series with ARIMA models for T , S_{13} and the Y model. Filtering is done by the R package [rjd3toolkit](#).
01. Compute filtering monthly series
- ▶ Benchmarking is used once the quarterly estimates are computed to correct the provisional monthly estimates with the corresponding quarterly estimate. Benchmarking is done by the R package [rjd3bench](#).
02. Compute benchmarking with quarterly value

R as facilitator to implement an end-to-end statistical production process

LFS-ML Dashboard Intro Totals CV Volatility Turning points Sampling weights

Select time limits
2016-01-07 (01) 2022-06-12 (19) 2023-12-28 (52)
2016-01-07 (01) 2017-10-26 (43) 2019-08-15 (32) 2021-06-03 (22) 2023-03-23 (12) 2024-11-01 (44)

Select Time Frequency
☐ Quarter
☐ Month
☒ Week

Select AOI
☒ Employed
☒ Unemployed
☐ Inactive

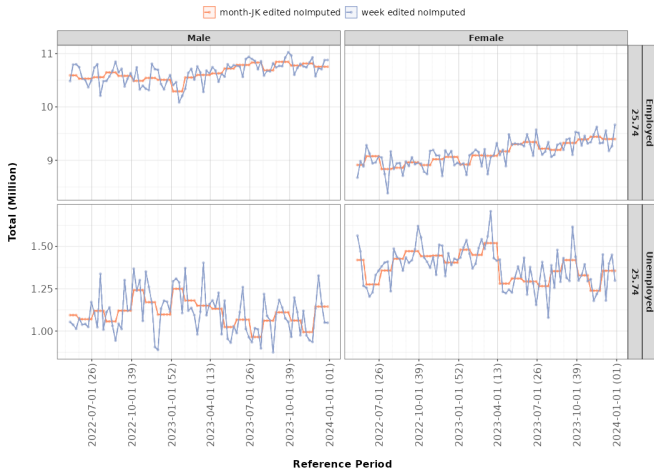
Select Age Group
☐ 16-24
☒ 25-74
☐ 75-

Select Sex Group
☒ Male
☒ Female
☐ y-scale begins at 0
☒ Plot Conf Interval

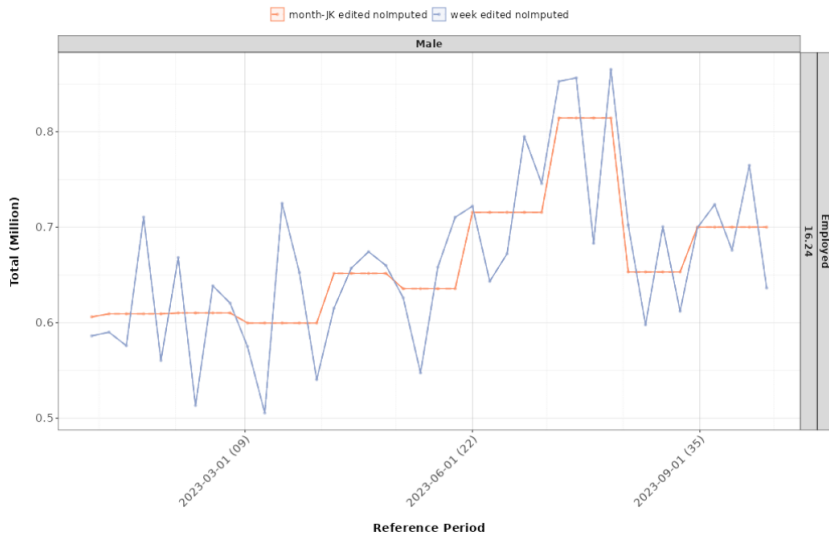
Select Conf Level
☐ 99%
☒ 95%
☐ 90%

Scenario: edited, not imputed
☐ Quarter-ORIG
☐ Quarter
☐ Month
☒ Month-JK
☐ Month-JKfiltered-prov
☐ Month-JKfiltered-def
☐ Month-JKbench
☒ Week
☐ Eurostat

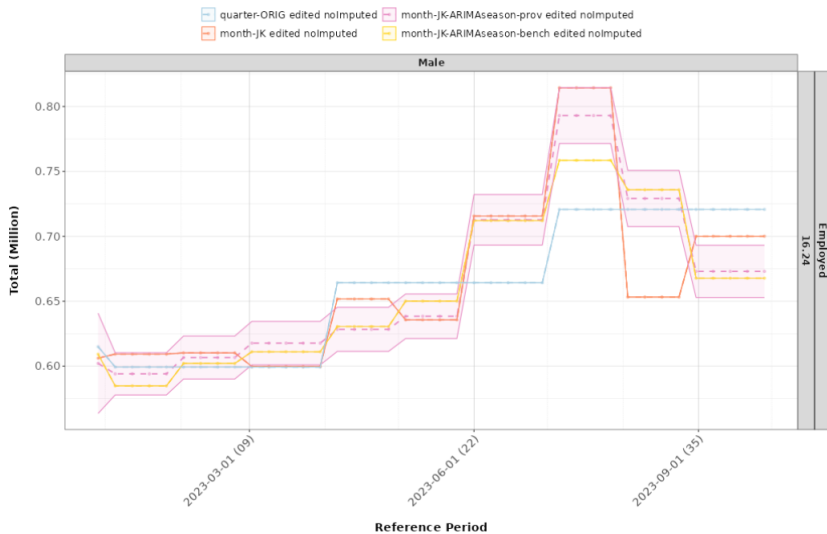
02. Update Shiny app



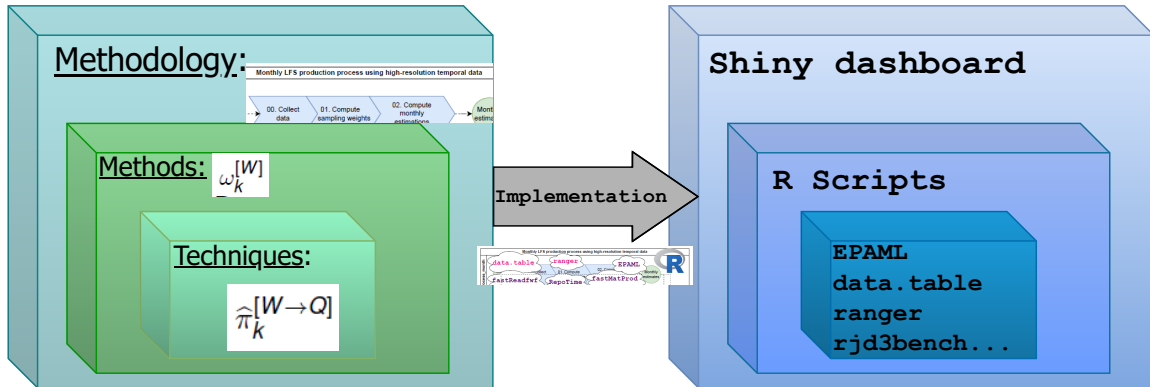
R as facilitator to implement an end-to-end statistical production process



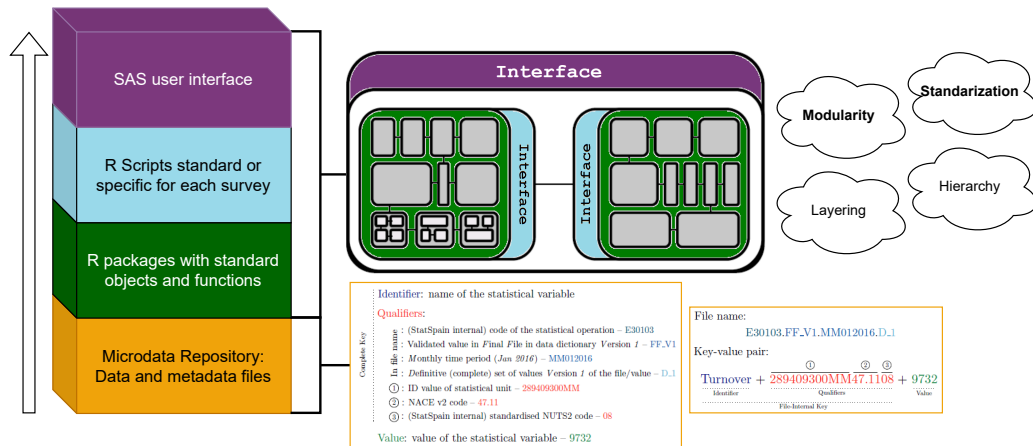
R as facilitator to implement an end-to-end statistical production process



R as facilitator to implement an end-to-end statistical production process



Use Case 2: Ecosystem of R packages to interact with the Microdata Repository



R as facilitator to connect with updated published data

Use Case 3: ineapir

R package `ineapir`: <https://github.com/es-ine/ineapir>

INE data extraction with `ineapir`: : CHEAT SHEET

The **ineapir** package allows to extract open data and metadata published by the **INE** (Spain). The data is obtained using calls to the INE API JSON service which access via URL requests to the data required by introducing the ID of the serie/tabla desired.



1 How to obtain ID's Go to **INE** website and find a table/ series with the desired data (no need to be filtered). Depending the type of data hosted on the INE database we can distinguish several types of URL's:

FROM A TABLE

CASE 1: tempus type
The ID is the t parameter, 50902

`ine.es/jaxi/T3/Tabla.htm?t=50902`

CASE 2: pc-axis
The ID is the concatenation of path and file, `t20/e245/p08/i0/i01001.px`

`ine.es/jaxi/Tabla.htm?path=t20/e245/p08/i0&file=i01001.px`

CASE 3: tpx
The ID is the tpx parameter, 33387

`https://www.ine.es/jaxi/Tabla.htm?tpx=33387&L=0`

FROM A SERIES

CASE 4: series

1. Browse a table of interest
2. Filter the selected values
3. Click on the corresponding value cell
4. Take the **code** that appears on the plot

2 Main functions

OBTAINING DATA

- `get_data_table(idTable, filter, nlast, det, tip, lang, validate, verbose, unnest, metanames, metacodes)`
It returns the data of the idTable specified according to the filter and the other arguments.
- `get_data_series(codSeries, nlast, dateStart, dateEnd, det, tip, lang, validate, verbose, unnest)`
It returns the data of the codSeries specified according to the parameters.
- `get_data_series_filter(operation, filter, periodicity, nlast, det,...)`
It returns the data of the operation specified according to the filter and the other parameters.

Auxiliar functions for Tables

- `get_metadata_table_groups/values(idTable, (idGroup),...)`
It returns all available groups and values for the specified table.
- `get_metadata_table_varval(idTable, ...)`
Get metadata information about the variables and values for a given table
- `get_metadata_series_table(idTable, filter,...)`
Get all the series for a given table
- `get_metadata_operation_table(idTable,...)`
It returns the operation for the specified table.

3 Arguments

OPTION	TYPE	DEFAULT	EFFECTS
<code>idTable</code>	<code>int</code>		Id of the table

3a filter argument

Data from **tables** and **operations** can be filtered with a list according to the variables/values they contain. Let's see how to construct the filter:

FOR TABLES

See `get_metadata_table_varval()` to get all the values at once. There are different approaches to build the filter depending on the table type:

1. **tempus**: The filter is based on ids, with format: `list(id_variable1 = id_value1, id_variable2 = id_value2)`

Use Case 4: Common interface for Machine Learning algorithms

13⁰⁰ – 14⁰⁰

Scientific Session *Machine Learning and AI*

Chairman: Bogdan Oancea, University of Bucharest



1. *Supervised statistical (machine) learning for domain estimation with business survey data*

Vasilis Chasiotis, Department of Statistics, Athens University of Economics and Business, Athens, Greece



2. *The Synergy of R and Generative AI in Statistics*

Vytas Vaiciulis, Central Statistics Office, Ireland



3. *Automating Data Validation on SQL Server Using R and the Machine Learning package*

Paula Hartung, Statistics Iceland

4. Lightning talk: *RMLUtils: an Official Statistics oriented common interface for Machine Learning*

Luis Sanguiao Sande, Statistics Spain (INE)

Carlos Sáez Calvo, Statistics Spain (INE)

Sandra Barragán Andrés, Statistics Spain (INE)

Ester Puerto Sanz, Statistics Spain (INE)



María Novás Filgueira, Statistics Spain (INE)

Javier Villaescusa Almagro, Statistics Spain (INE)

Sergio Pardina Quirós, Statistics Spain (INE)

Álvaro García Tenorio, Statistics Spain (INE)

Use Case 5: an R package for time series model identification

15³⁰ - 16³⁰ Scientific Session *NLP in Official Statistics and Time series*
Chairman: Mark Van Der Loo, Statistics Netherlands (CBS)

1. *A Neural Network Approach to Text Classification for International Standardized Codes*



Nina Niederhametner, Statistics Austria
Alexander Kowarik, Statistics Austria
Johannes Gussenbauer, Statistics Austria

2. *Semantic address matching using Keras for R*



Paula Cruz, Statistics Portugal; NOVA IMS
Leonardo Vanneschi, NOVA IMS
Marco Painho, NOVA IMS
Filipa Ribeiro, Statistics Portugal



3. *TEAM: an R package for time series model identification*

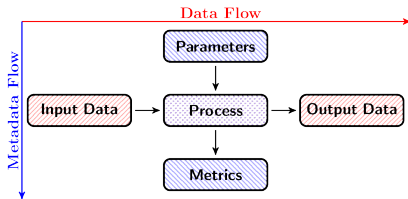


C. Sáez Calvo, S.G. for Methodology and Sampling Design, Statistics Spain
L. Sanguiao Sande, S.G. for Methodology and Sampling Design, Statistics Spain
Félix Aparicio Pérez, S.G. for Methodology and Sampling Design, Statistics Spain
María Teresa Vázquez Gutiérrez, S.G. for Information Technologies and Communications, Statistics Spain
José Fernando Arranz Arauzo, S.G. for Information Technologies and Communications, Statistics Spain

Conclusions

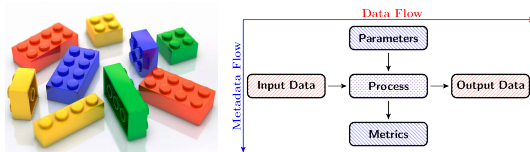
Conclusions: Take home messages

- ▶  is open source, user friendly (easy to learn), constantly evolving. R has a bunch of packages implementing statistical methods and a wide community of users.
- ▶  facilitates a proper **implementation** and **automation** of **modular processes** in the production of Official Statistics.

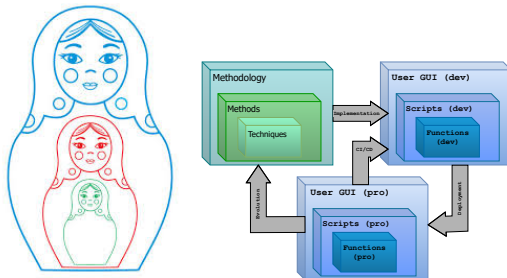


Conclusions: Take home messages

- Design the main process and subprocesses: **Modularity**



- Implement each subprocess: **Levels of implementation**





Run, or he's going to tell us about
again!

R



Thank you for your inspiration...



... and thank all of you here for your attention.



Questions?



sandra.barragan.andres@ine.es

References

- Dustin Boswell and Trevor Foucher. *The Art of Readable Code: Simple and Practical Techniques for Writing Better Code*. " O'Reilly Media, Inc.", 2011.
- Norman Matloff. *The art of R programming: A tour of statistical software design*. No Starch Press, 2011.
- K.P. Murphy. *Machine Learning: A Probabilistic Perspective*. MIT Press, New York, 2012.
- D. Salgado, M E. Esteban, M. Novás, S. Saldaña, and L. Sanguiao. Data organisation and process design based on functional modularity for a standard production process. *Journal of official statistics*, 34(4):811–833, 2018.
- Allen B Tucker and Robert Noonan. *Programming languages: principles and paradigms*. McGraw-Hill, 2007.
- Mark PJ van der Loo. A method for deriving information from running r code. *arXiv preprint arXiv:2002.07472*, 2020.