# Between R and Excel - The ETER Data Validation and Quality Process

*Daniel Wagner-Schuster*

# *European Tertiary Education Register*

- Collects and provides data on higher education institutions in Europe

  - at the institutional level

  - 2013-2019
    - Tender for another 3 years

  - 37 countries
    - ~ 3,000 institutions per year

  - currently data for 6 academic years
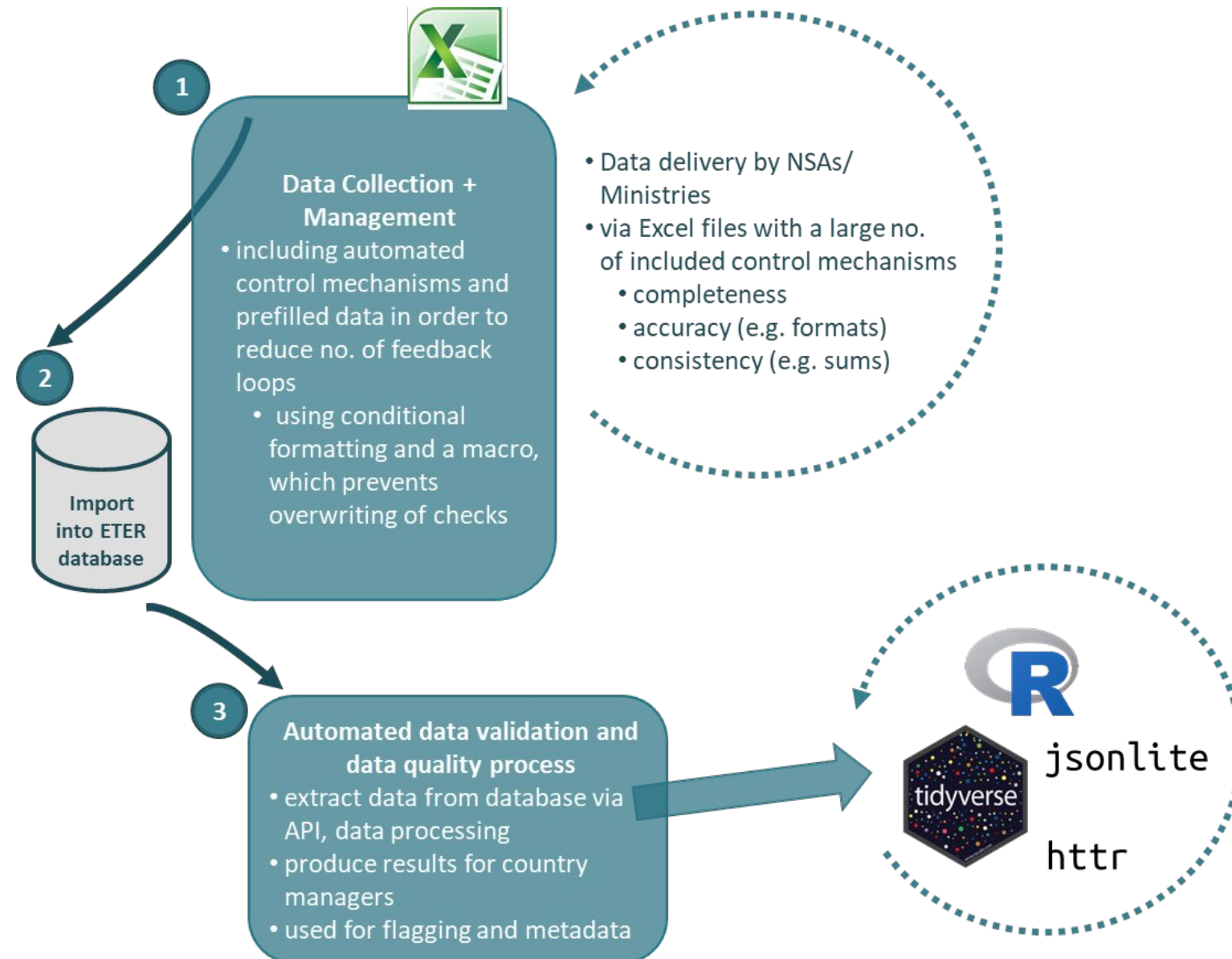    - 17,285 observations for 618 variables

  - eter-project.com

# *Some challenges*

- Collect data from 37+ different sources (countries)

- Ensure tidy data (data validation)

  - Find missing variables, inconsistencies, wrong formats

- Ensure high data quality

  - Multiannual checks

  - Outlier detection

  - Consistency with external data sources (e.g. Eurostat national aggregates)

- Constant interaction between data deliverer, country managers and data collection managers

# Challenge #1: Data Collection



**1** 

**Data Collection + Management**
- including automated control mechanisms and prefilled data in order to reduce no. of feedback loops
  - using conditional formatting and a macro, which prevents overwriting of checks

**2** **Import into ETER database**

- Data delivery by NSAs/ Ministries
- via Excel files with a large no. of included control mechanisms
  - completeness
  - accuracy (e.g. formats)
  - consistency (e.g. sums)

**3** **Automated data validation and data quality process**
- extract data from database via API, data processing
- produce results for country managers
- used for flagging and metadata

R
jsonlite
tidyverse
httr

# *Challenge #2: Tidy data*

**3a**

**Data validation**
- check for completeness, accuracy and consistency
- produce one report for each country and year

R > PDF

kableExtra  rmarkdown

**Snippet of the Austria data validation report for 2016**

Country: AT

Year: 2016

Test date: February 08, 2019

| ETER-ID | Variable | Type of error | Flag | Remarks |
|---------|----------|---------------|------|---------|
| AT0035 | STA.TOTALFTE | accuracy problem - breakdowns numeric but total is not | i | only academic staff available |
| AT0036 | STA.TOTALFTE | accuracy problem - breakdowns numeric but total is not | i | only academic staff available |
| AT0037 | STA.TOTALFTE | accuracy problem - breakdowns numeric but total is not | i | only academic staff available |
| AT0038 | STA.TOTALFTE | accuracy problem - breakdowns numeric but total is not | i | only academic staff available |
| AT0039 | STA.TOTALFTE | accuracy problem - breakdowns numeric but total is not | i | only academic staff available |
| AT0040 | STA.TOTALFTE | accuracy problem - breakdowns numeric but total is not | i | only academic staff available |
| AT0041 | STA.TOTALFTE | accuracy problem - breakdowns numeric but total is not | i | only academic staff available |
| AT0042 | STA.TOTALFTE | accuracy problem - breakdowns numeric but total is not | i | only academic staff available |
| AT0043 | STA.TOTALFTE | accuracy problem - breakdowns numeric but total is not | i | only academic staff available |

# *Challenge #3: Data quality*



**3b**

**Data quality**
- Internal data quality
  - Multiannual checks
  - Outlier detection
    - determine thresholds
    - calculate test values
    - generate MS-Excel file for country managers
      - use function **conditionalFormatting** of package openxlsx for highlighting already flagged cases

- External data quality
  - e.g. comparison with Eurostat national aggregates

openxlsx

tidyverse

# Challenge #3: Data quality (output)

**Snippet of the sheet containing positive multiannual test cases**
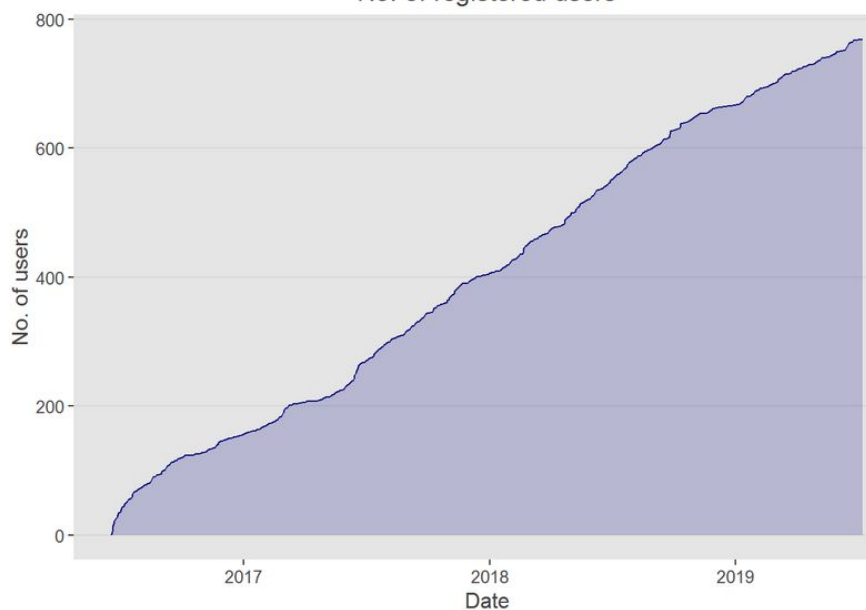
# *Paper*

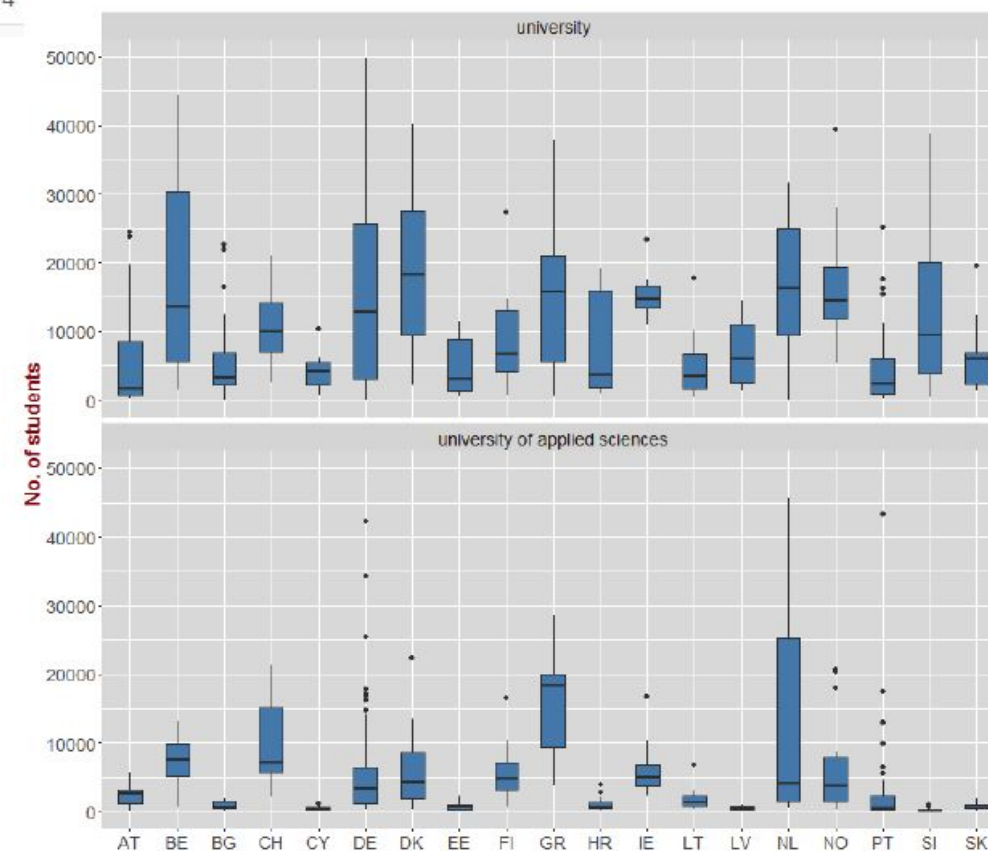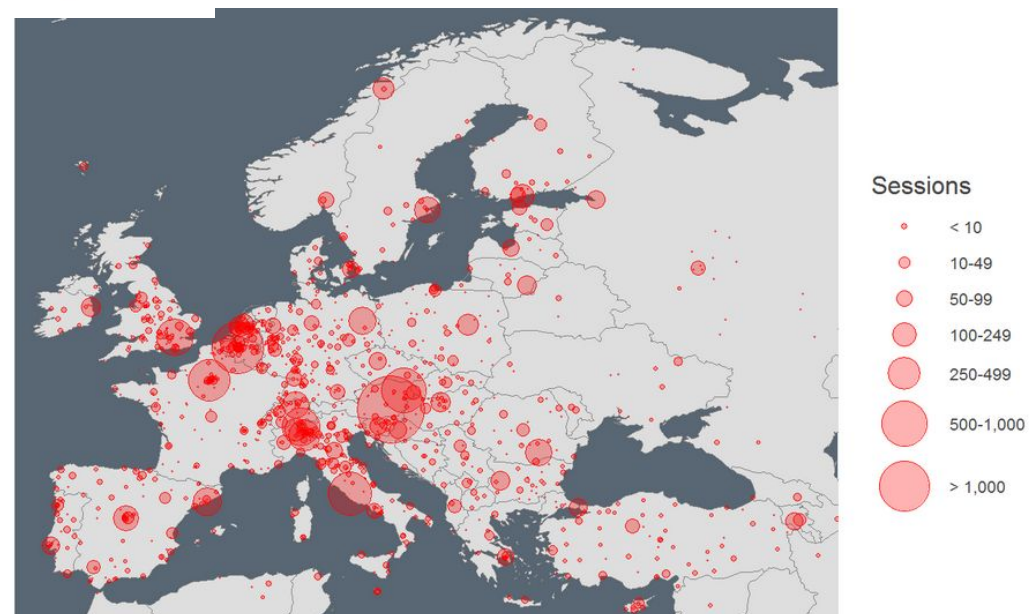A Tailor-made Data Quality Approach for Higher Educational Data

# *Other R usage in ETER*

No. of registered users

| | No. of Sessions | No. of new Users | Share of new users (%) | Share of returning users (%) |
|---|---|---|---|---|
| **2016** | | | | |
| August | 371 | 278 | 74.93 | 25.07 |
| September | 683 | 497 | 72.77 | 27.23 |
| October | 738 | 568 | 76.96 | 23.04 |
| November | 1058 | 669 | 63.23 | 36.77 |
| December | 737 | 537 | 72.86 | 27.14 |

THE INNOVATION COMPANY

JOANNEUM
RESEARCH
POLICIES

www.joanneum.at/policies