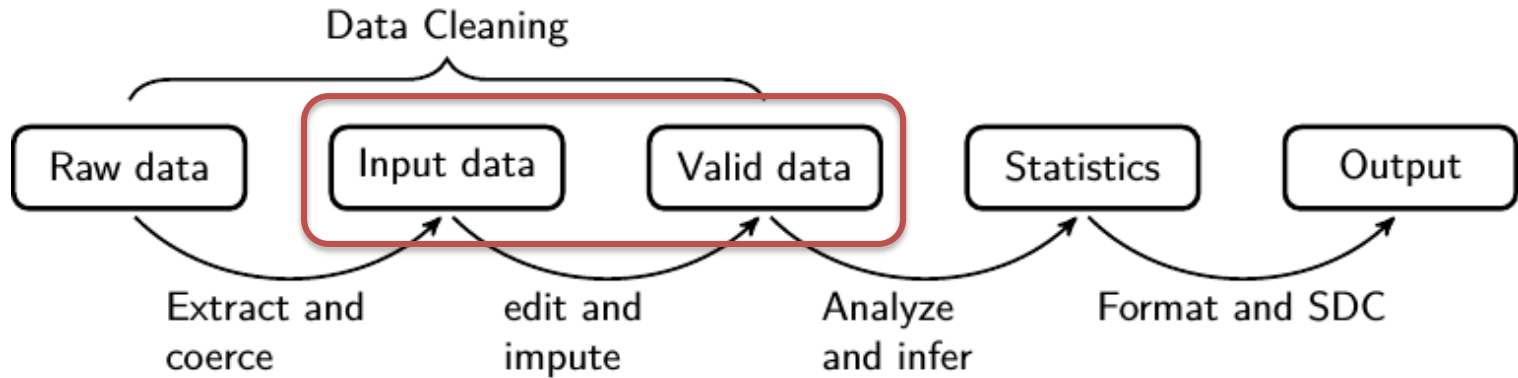# Systematic Data Cleaning using R

Mark van der Loo
Dpt of Methodology
NTTS 13-03-2019
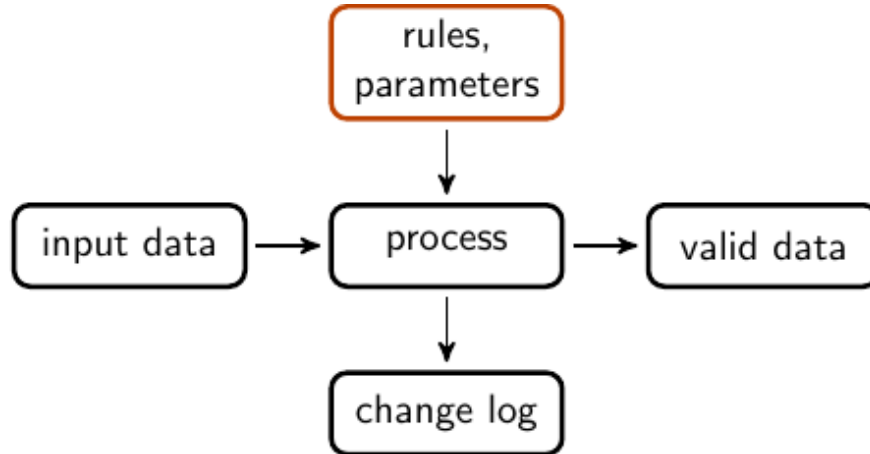
# Data Cleaning in the Statistical Value Chain

# Ultimate goal

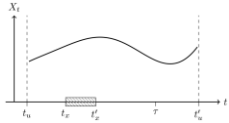- User specifies what is 'valid data'.
- Everything else is automated.

# Approach

- Think deeply about the nature of 'data validation'
- Design tools
  - In **open source**
  - Following the **Unix philosopy** (small, powerful, combinable)
  - With **humans** and modern **standards** in mind
  - While **re-using** existing tools where possible

# Results, ESS collaboration

# R-based tools, available on CRAN

**R> dcmodify**
*User-defined data cleaning*

**R> rspa**
*Alter data to pass validation rules*

**R> simputation**
*Missing data imputation*

**R> deductive**
*Use validation rules to repair data*

**R> errorlocate**
*Localize erroneous fields*

**R> lumberjack**
*Track changes in data*

**R> validate**
*Define rules, measure data quality*
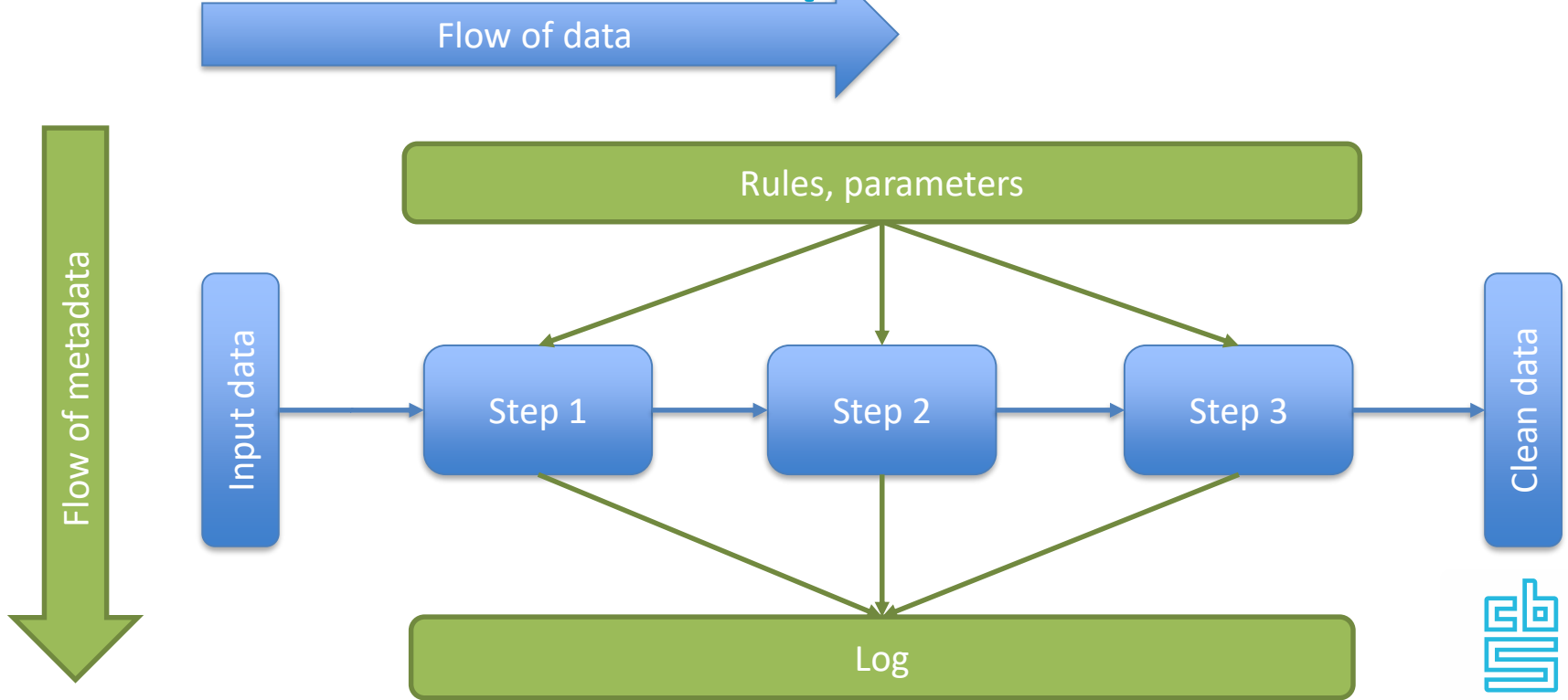
**R> validatetools**
*Maintain and investigate rules*

**R> validatereport**
*Validation reports in ESS standard*

# Flow of data, metadata, paradata

# Implementation (1): preparation

```
dat    <- read.csv( "SBS2000.csv" )
rules  <- validate::validator( .file = "rules.R" )
logger <- validate::lbj_rules( rules )
```

| | id | staff | turnover | other.rev | total.rev | total.costs | profit |
|---|---|---|---|---|---|---|---|
| 1 | RET01 | 75 | NA | NA | 1130 | 18915 | 20045 |
| 2 | RET02 | 9 | 1607 | NA | 1607 | 1544 | 63 |
| 3 | RET03 | NA | 6886 | -33 | 6919 | 6493 | 426 |
| 4 | RET04 | NA | 3861 | 13 | 3874 | 3600 | 274 |
| 5 | RET05 | NA | NA | 37 | 5602 | 5530 | 72 |
| 6 | RET06 | 1 | 25 | NA | 25 | 22 | 3 |
| 7 | RET07 | 5 | NA | NA | 1335 | 136 | 1 |
| 8 | RET08 | | 404 | 13 | 417 | 342 | 75 |
| | | | | NA | 2596 | 2486 | 110 |
| | | | | NA | NA | NA | NA |
| | | | | NA | 645 | 636 | 9 |

**SBS2000.csv**

```
1
2  # range resitrictions
3  staff >= 0
4  turnover >= 0
5  other.rev >= 0
6
7
8  # balance restrictions
9  turnover + other.rev == total.rev
10 total.rev - total.costs == profit
11 profit <= 0.6 * total.rev
12
```

**rules.R**

# Implementation (2): execution

```
dat %L>%
    lumberjack::start_log(logger) %L>%
    errorlocate::replace_errors(rules) %L>%
    tag_missing() %>%
    simputation::impute_mf(. - id ~ . - id) %L>%
    rspa::match_restrictions(rules, eps=1E-8) %L>%
    dump_log() ->
    clean_data
```

# Implementation (2): execution

```
dat %L>%
    lumberjack::start_log(logger) %L>%
    errorlocate::replace_errors(rules) %L>%
    tag_missing() %>%
    simputation::impute_mf(       - id) %L>%
    rspa::match_restrictions(       eps=1E-8) %L>%
    dump_log() ->
    clean_data
```
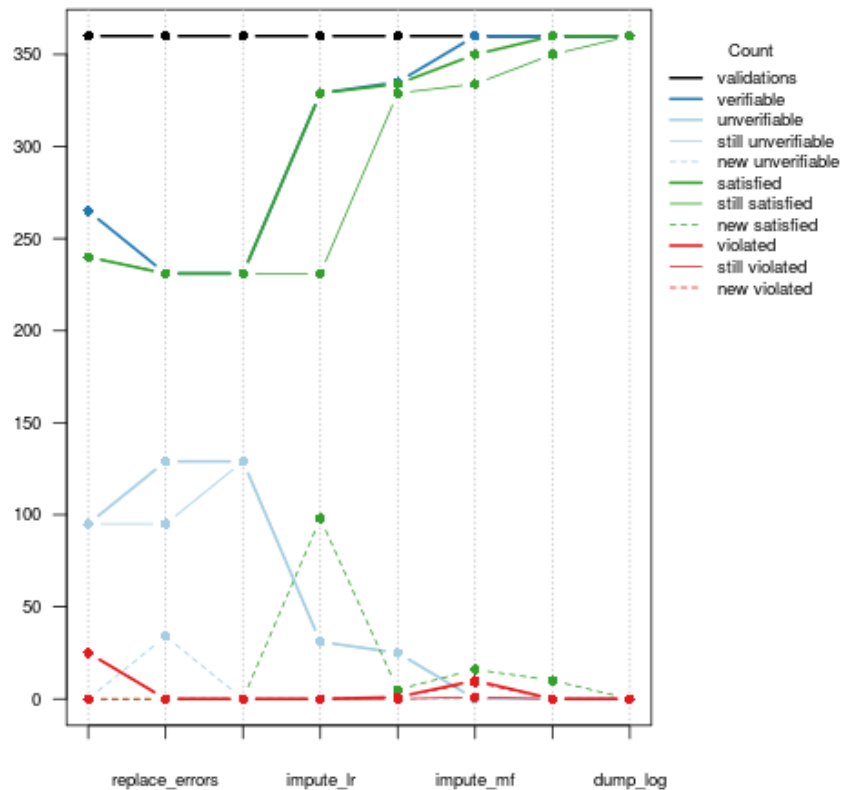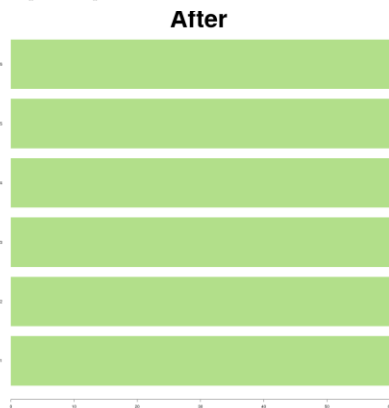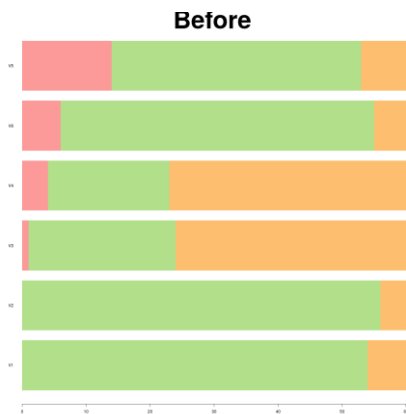
Process block

Pipe operator, controlling flow of data and meta-data

User Specification

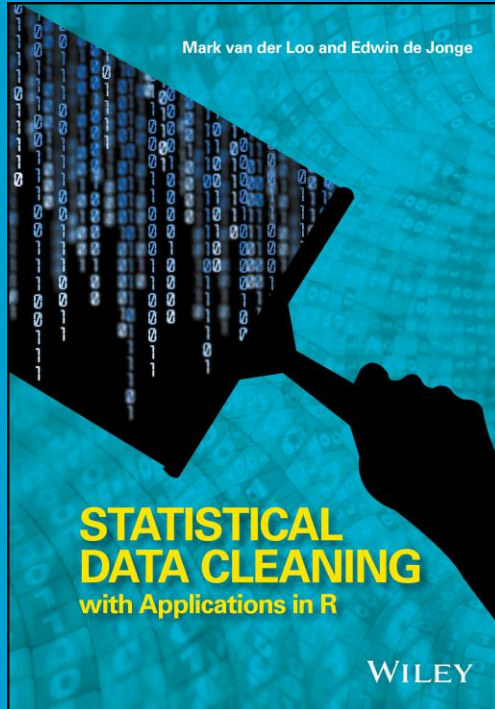# Implementation (3) process monitoring

# Implementation in CBS

– Used in production, e.g.

- Health care institutes

- Energy (currently under redesign)

- Imputation of population registers

- …

– How?

- Internal courses (CBS academy)

- Redesign of production systems

# More information



*MPJ van der Loo and E de Jonge (2018)*
*John Wiley & Sons.*

Upcoming tutorials:
- uRos2019 conference (Bucharest)
- useR2019 conference (Toulouse)
- ENBES2019 conference (BilBao)

More FOSS for official statistics:
www.awesomeofficialstatistics.org