David Klein

# How to save a scientist's career with data classes?

deRSE    Potsdam    5.6.2019

# "I can't replicate your results."

# "I can't replicate your results."

# "You made an Excel error."

P I K

https://science.sciencemag.org/content/314/5807/1856
https://www.bloomberg.com/news/articles/2013-04-18/faq-reinhart-rogoff-and-the-excel-error-that-changed-history

# Increase robustness of research software



Performance    Robustness

# Increase robustness of research software

Performance    Robustness

# Increase robustness of research software



Performance     Robustness

# 1. Standardized, generic structure

```
#>      AFR.IndexA CPA.IndexA EUR.IndexA AFR.IndexB CPA.IndexB EUR.IndexB
#> y2000          1          4          7         10         13         16
#> y2001          2          5          8         11         14         17
#> y2002          3          6          9         12         15         18
```

```
#>    Cell Region Year  Data1 Value
#> 1    NA    AFR 2000 IndexA     1
#> 2    NA    CPA 2000 IndexA     4
#> 3    NA    EUR 2000 IndexA     7
#> 4    NA    AFR 2001 IndexA     2
#> 5    NA    CPA 2001 IndexA     5
#> 6    NA    EUR 2001 IndexA     8
#> 7    NA    AFR 2002 IndexA     3
#> 8    NA    CPA 2002 IndexA     6
#> 9    NA    EUR 2002 IndexA     9
#> 10   NA    AFR 2000 IndexB    10
#> 11   NA    CPA 2000 IndexB    13
#> 12   NA    EUR 2000 IndexB    16
#> 13   NA    AFR 2001 IndexB    11
#> 14   NA    CPA 2001 IndexB    14
#> 15   NA    EUR 2001 IndexB    17
#> 16   NA    AFR 2002 IndexB    12
#> 17   NA    CPA 2002 IndexB    15
#> 18   NA    EUR 2002 IndexB    18
```

- Avoid mistakes by using a standardized data structure
- Various input formats are transferred into the same structure.
- Data class is flexible with regard to numbers of dimensions (columns)
- It detects dimensions in the raw data and adds them to the object (as columns)

# 2. Name matching

- Automatically matches entries when performing operations
- Columns can not be mixed up

| | GDP | | | / | | Population | | | = | | GDP per capita | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|

|  GDP | | | | Population | | | | GDP per capita | | |
|---|---|---|---|---|---|---|---|---|---|---|
| region | period | value | | region | period | value | | region | period | value |
| EUR | 2005 | 13386 | | EUR | 2005 | 496 | | EUR | 2005 | 26.99 |
| LAM | 2005 | 4984 | | EUR | 2010 | 505 | | LAM | 2005 | 8.88 |
| USA | 2005 | 12584 | | LAM | 2005 | 561 | | USA | 2005 | 42.51 |
| EUR | 2010 | 14024 | | LAM | 2010 | 597 | | EUR | 2010 | 27.77 |
| LAM | 2010 | 5908 | | USA | 2005 | 296 | | LAM | 2010 | 9.90 |
| USA | 2010 | 13291 | | USA | 2010 | 309 | | USA | 2010 | 43.01 |

Ordered by year                    Ordered by region

# 3. Units

```
> a <- 2
> a <- set_units(a,"m")
> b <- 4
> b <- set_units(b,"km")

> a+b
4002
> units(a+b)
[m]

> time <- 5
> time <- set_units(time,"s")
> velo <- a/time

> velo
0.4
> units(velo)
[m*s^-1]

> units(velo) <- "km*h^-1"    ⟵    Manually changing the unit
> velo
1.44                          ⟵    Auto conversion of values
```

- Variables have units attached
- Applies conversion factors when combining variables
- Reports an error if conversion is unknown

# 4. History

- Automatically **logs operations** performed on the data
- Allows following the operations and **detecting mistakes ex post**
- even if routine that performed the operations is not available

```
c = a + b
d = b + a
e = c + d
```

History of „e"

```
$calcHistory

1 c + d
2  ¦--a + b
3  °--b + a
```

# 4. History

```
$calcHistory

1   calcOutput("TauTotal")
2    °--toolAggregate(x = x$x, weight = x$weight, rel = reg_rel)
3       |--toolAggregate(x * weight, rel, from = from, to = to, dim = dim, partrel = partrel) * weight2
4       |    |--x * weight
5       |    |    °--readSource("Tau", "paper")
6       |    |         |--toolAggregate(tau, rel = iso_cell, weight = collapseNames(xref))
7       |    |         |    °--toolAggregate(x * weight, rel, from = from, to = to, dim = dim, partrel = partrel) * weight2
8       |    |         |         |--x * weight
9       |    |         |         °--1/(toolAggregate(weight, rel, from = from, to = to, dim = dim, partrel = partrel, verbosity = 10) + 10^-100)
10      |    |         °--toolAggregate(xref, rel = iso_cell)
11      |    °--1/(toolAggregate(weight, rel, from = from, to = to, dim = dim, partrel = partrel, verbosity = 10) + 10^-100)
12      |         °--readSource("Tau", "paper")
13      |              |--toolAggregate(tau, rel = iso_cell, weight = collapseNames(xref))
14      |              |    °--toolAggregate(x * weight, rel, from = from, to = to, dim = dim, partrel = partrel) * weight2
15      |              |         |--x * weight
16      |              |         °--1/(toolAggregate(weight, rel, from = from, to = to, dim = dim, partrel = partrel, verbosity = 10) + 10^-100)
17      |              °--toolAggregate(xref, rel = iso_cell)
18   °--readSource("Tau", "paper")
19        |--toolAggregate(tau, rel = iso_cell, weight = collapseNames(xref))
20        |    °--toolAggregate(x * weight, rel, from = from, to = to, dim = dim, partrel = partrel) * weight2
21        |         |--x * weight
22        |         °--1/(toolAggregate(weight, rel, from = from, to = to, dim = dim, partrel = partrel, verbosity = 10) + 10^-100)
23        °--toolAggregate(xref, rel = iso_cell)
```

# 5. Metadata: Documentation

```
> getMetadata(population)
$unit
1e+06 [people]

$user
[1] "dklein"

$date
[1] "2019-06-04 18:07:20"

$calcHistory
1 c + d
2  ¦--a + b
3  °--b + a

$source
@TechReport{,
  title = {SRES Population scenarios},
  author = {{IPCC}},
  institution = {IPCC},
  year = {2000},
  url = {https://www.ipcc.ch/report/emissions-
scenarios/},
}

$description
[1] "Regional popuation data for SRES scenarios"

$note
[1] "Ahhh, too many people"
```

- Documentation is attached to the data
- Helps user assessing whether it's the data he/she expects

https://github.com/pik-piam/magclass

rse@pik-potsdam.de