



Make your paper figures professionally: Scientific data analysis and visualization in R

Julien Gagneur

Gagneur lab - Computational biology

gagneurlab.in.tum.de

To understand the genetic basis of gene regulation and its implication in diseases



Uhrenturm der TUM

You've been building algos and software processing data...



... now sail!



Caution: different skills required!



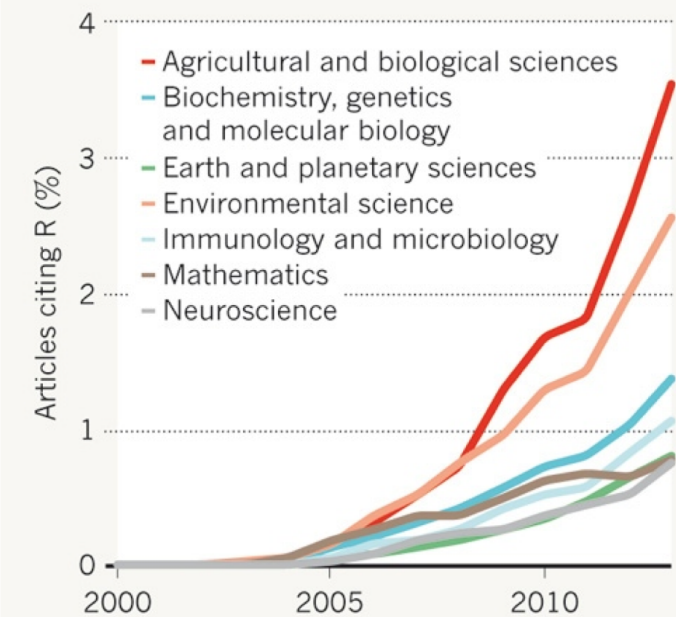
R: A data analyst-oriented language

- Written by statisticians for statisticians
- GNU, Open source
- Widely used
- Thousands of packages available for specific statistical analysis and application domains.

R usage in Science

A RISING TIDE OF R

An increasing proportion of research articles explicitly reference R or an R package.



https://www.youtube.com/watch?v=TR2bHSJ_eck

R: A data analyst-oriented language

- R is used **interactively** to explore data. The outcome of one analysis session are turned into **scripts**.
- Most valuable time is the human thinking time, not machine computing time. Data analysts spend most of your time in exploring the data and thinking about it, not in coding.
 - ➡ High-level language, loosely typed.
e.g. variables are not declared
 - ➡ Running time can be slow because the language takes care of what you don't need to care of.

Two sides of R programming

1. As **data analyst**
scripting. Few functions and data structures
2. As **package developer**
robust code that copes with flexibility of the language.

Our teaching team has experience with both.

This lecture focuses on 1.

Reproducibility: Reports with R Markdown

Analyze. Share. Reproduce.

Your data tells a story. Tell it with R Markdown.

Turn your analyses into high quality documents, reports, presentations and dashboards.



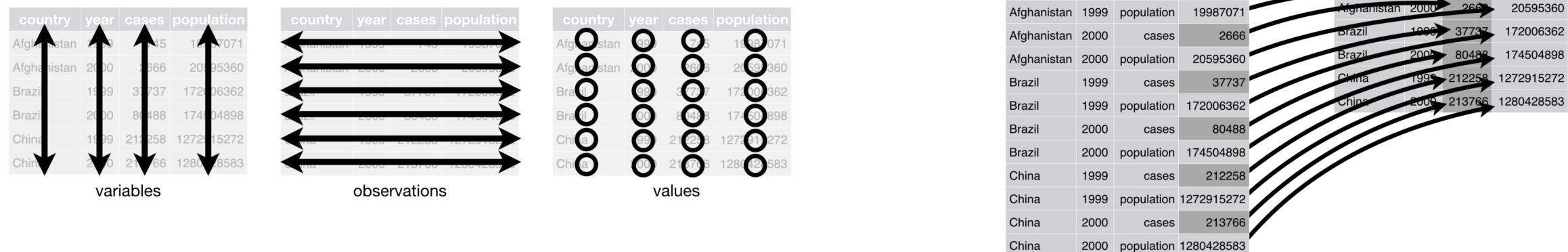
Get: reading and manipulating data

Data import

flat files, excel, XML, JSON, relational database

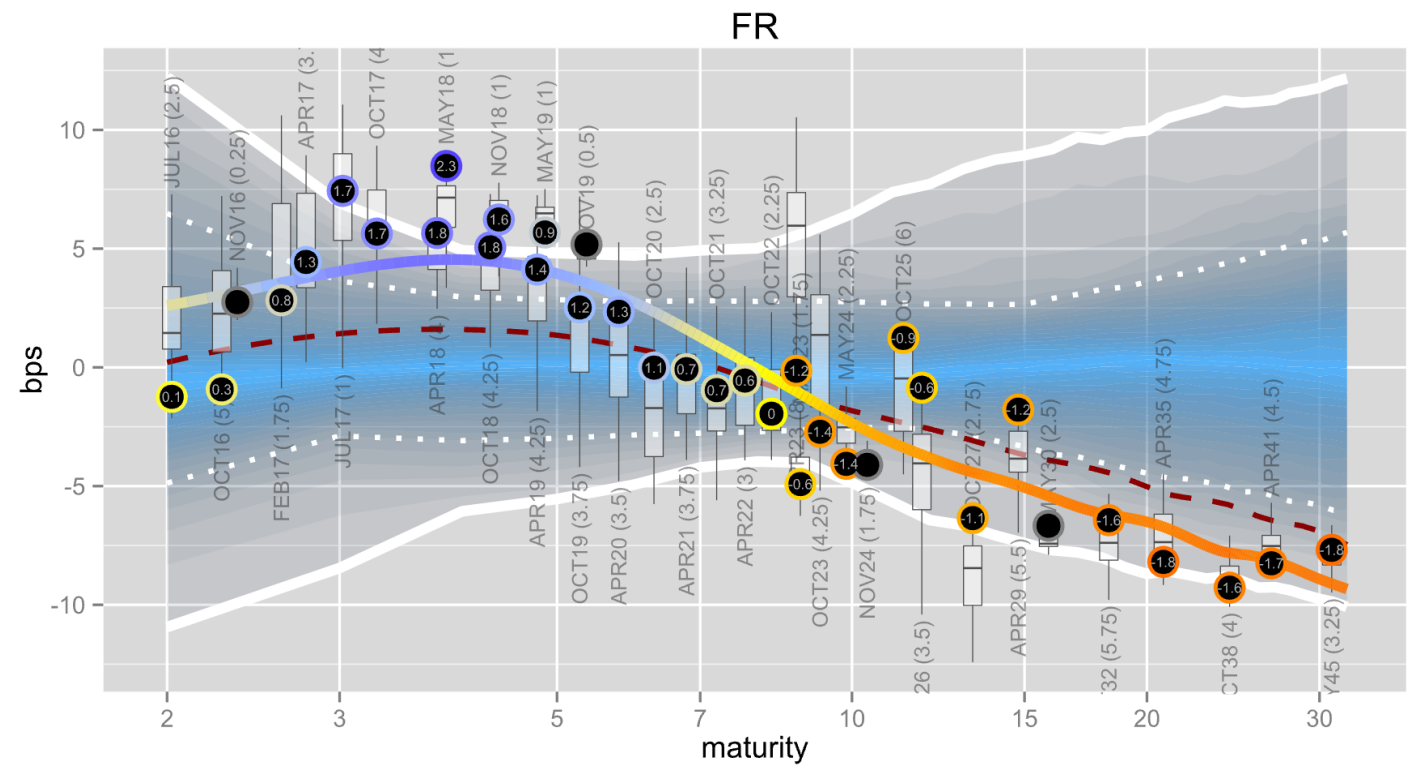
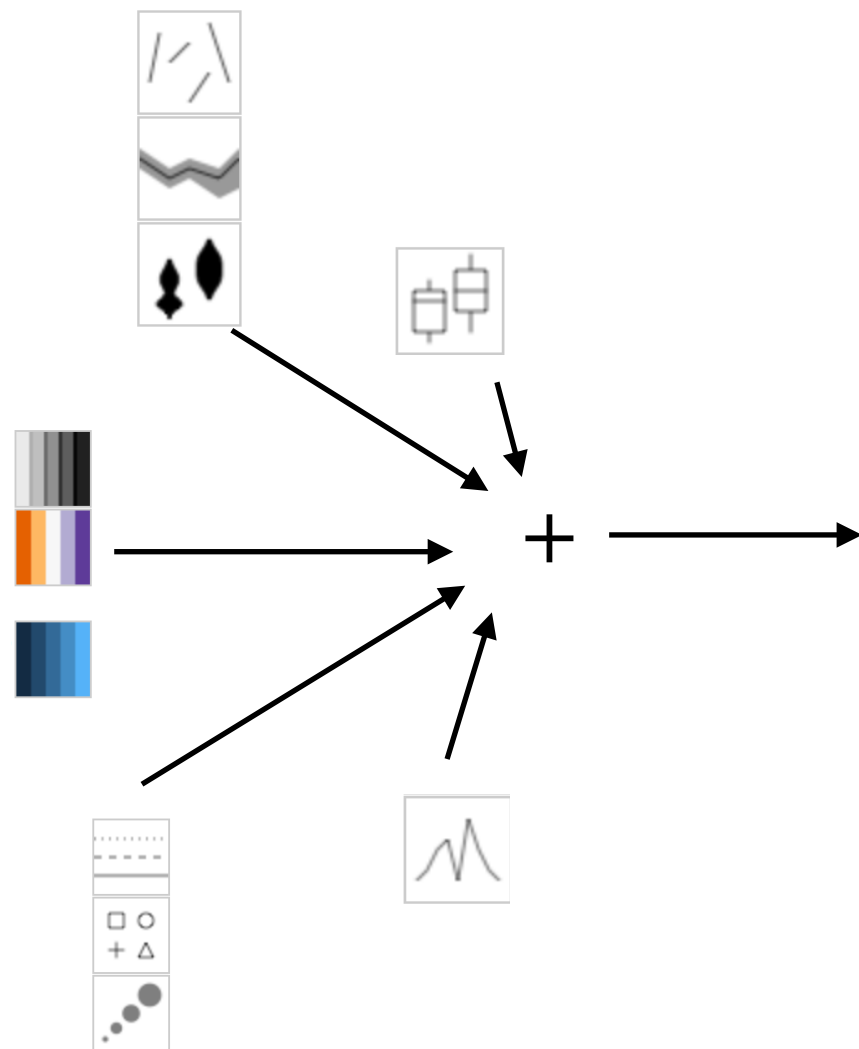
Tidy data

A standard way of structuring dataset for statistical analysis
data.table: Efficient in memory storage and operations



Look: Grammar of graphics with ggplot2

Flexible plotting using a set of independent complements that can be composed in many different ways



<http://stackoverflow.com/questions/24828341>

Conclude: Drawing robust conclusions

Not covered in the CEDOSIA module. See our master module “Data analysis and visualization in R”

Data exploration leads to **hypotheses** made on the data:

“Does smoking significantly associates with increased lung cancer?”

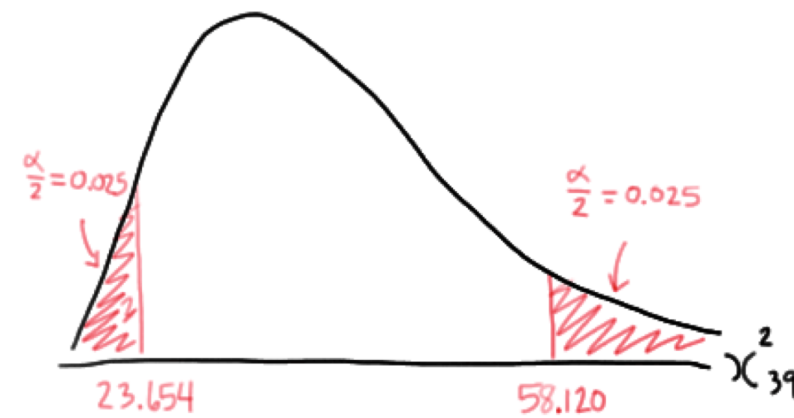
“Does smoking significantly associates with higher academic grades...

.... when I control for student age?

➔ **Statistical testing** (empirical and theoretical)

Classical tests: T-test, Wilcoxon, Fisher test

Generic resampling approaches



Syllabus

13.05.19 R basics (optional)

20.05.19 Introduction

2' flash presentations

Grammar of graphics

Data table

Tidy data I

27.05.19 Tidy data II

Plot types

Advanced plots

Improve one of your own plots (in groups)

03.06.19 10' presentation group I

10' presentation group II



@home: prepare a
dataset and code
producing one plot



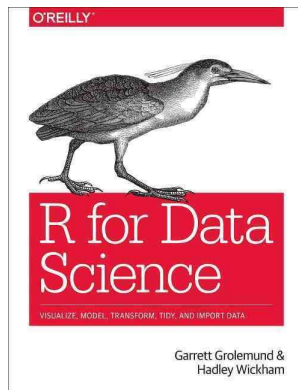
@home: prepare a 10'
presentation & code
producing plots

Evaluation

Presence (signatures)

Honest attempt at last day presentation.

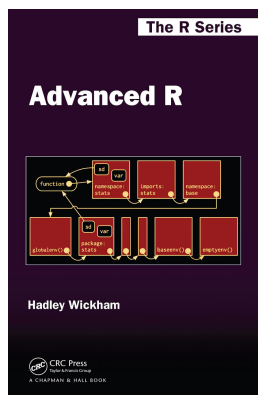
Recommended reading



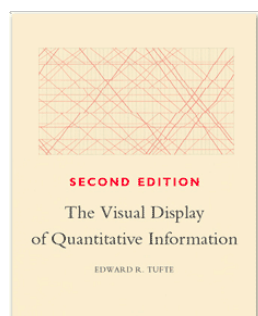
R for Data Science *, by Garrett Grolemund and Hadley Wickham



Modern Dive *, by Chester Ismay and Albert Y. Kim
(under construction)



Advanced R *, by Hadley Wickham



The Visual Display of Quantitative Information, Edward Tufte

* have a free online web version

Conclusion

- Data analysis is complementary to statistical methods and software development
- Get, look, (and conclude). Statistics and visualisation are both necessary.
- Reproducible analyses (Scripted reports).