# Data analysis and machine learning with R: an overview

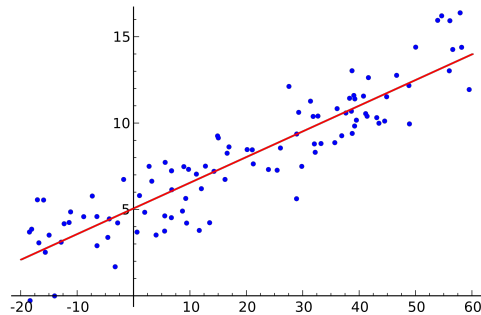Enrique Garcia Ceja
May 18 2018

R is a programming language for statistical computing and graphics.

Is an implementation of the S programming language and was created by Ross Ihaka and Robert Gentleman.

Stable beta version released in 2000.

- Linear and nonlinear modelling
- Classical statistical tests
- Time-series analysis
- Classification
- Clustering
- Deep learning
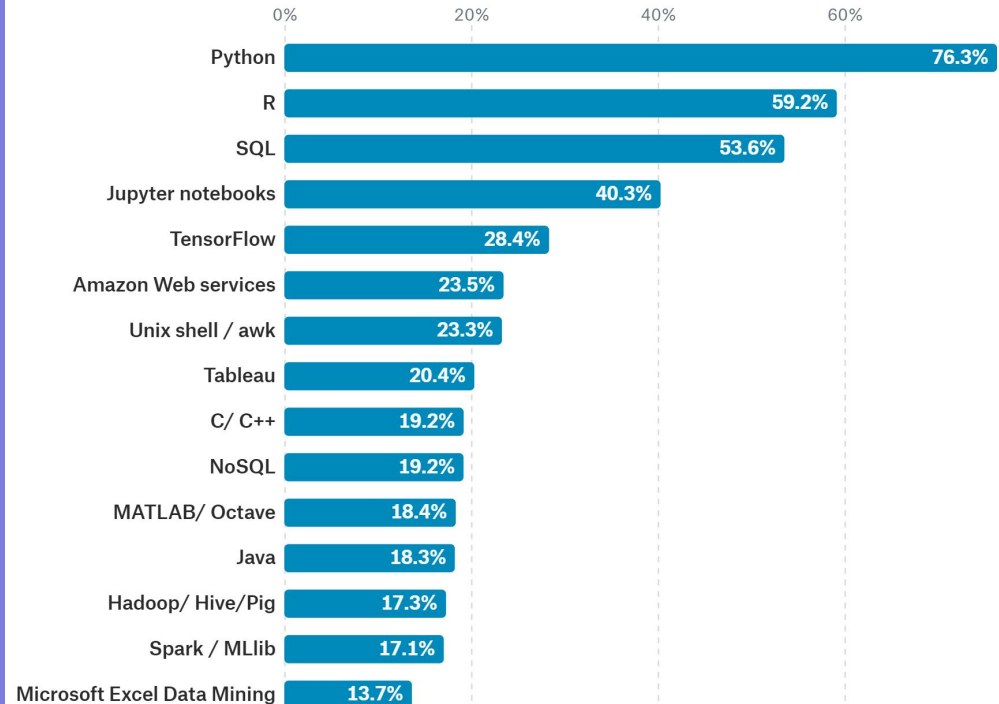- Visualization
- Spatial data

# Some stats



What tools are used at work?

Python was the most commonly used data analysis tool across employed data scientists overall, but more Statisticians are still loyal to R.
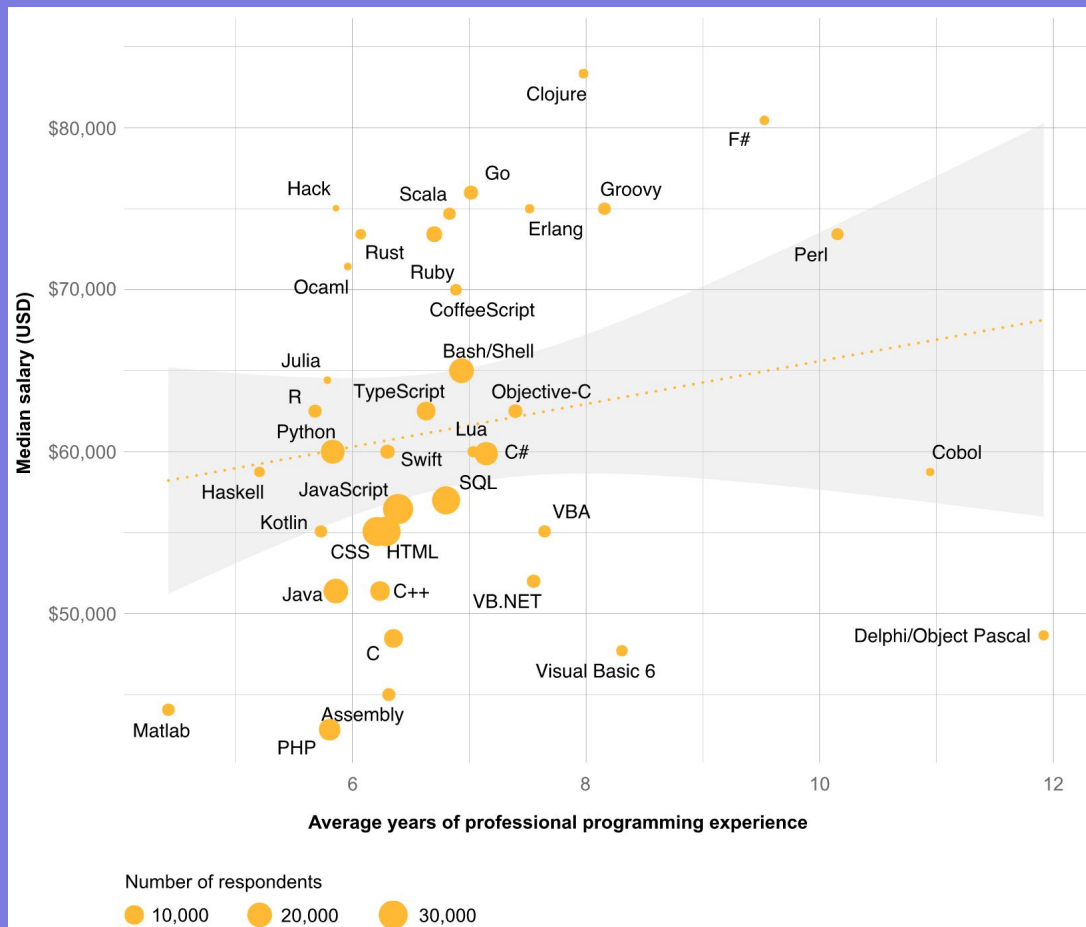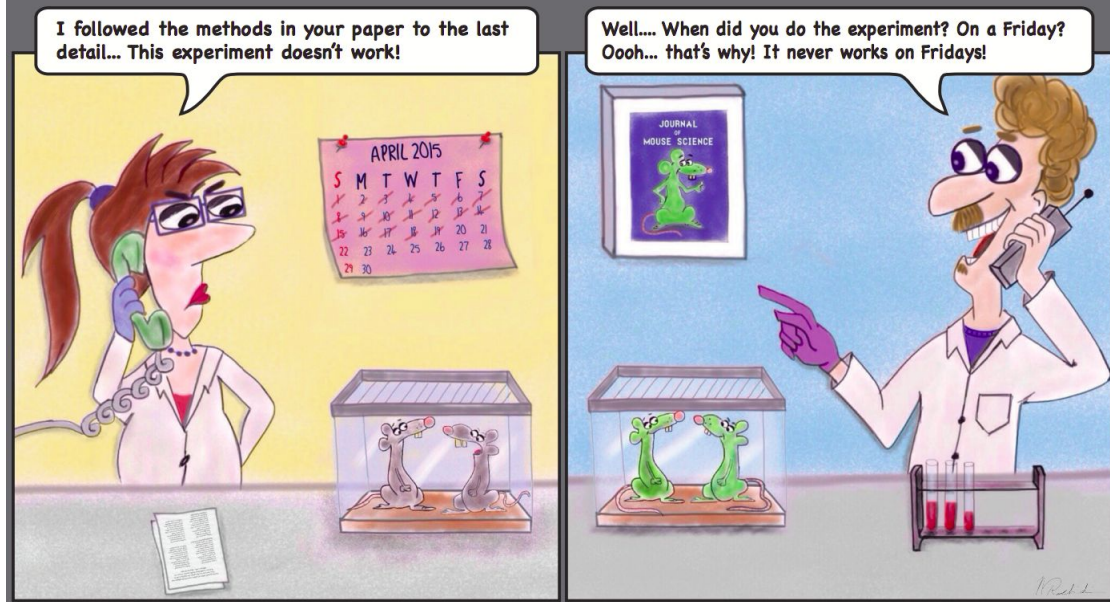
Company Size ▼    Industry ▼    Job Title ▼

| | 0% | 20% | 40% | 60% |
|---|---|---|---|---|
| Python | | | | 76.3% |
| R | | | 59.2% | |
| SQL | | | 53.6% | |
| Jupyter notebooks | | 40.3% | | |
| TensorFlow | 28.4% | | | |
| Amazon Web services | 23.5% | | | |
| Unix shell / awk | 23.3% | | | |
| Tableau | 20.4% | | | |
| C/ C++ | 19.2% | | | |
| NoSQL | 19.2% | | | |
| MATLAB/ Octave | 18.4% | | | |
| Java | 18.3% | | | |
| Hadoop/ Hive/Pig | 17.3% | | | |
| Spark / MLlib | 17.1% | | | |
| Microsoft Excel Data Mining | 13.7% | | | |

7,955 responses

Kaggle survey 2017
https://www.kaggle.com/surveys/2017

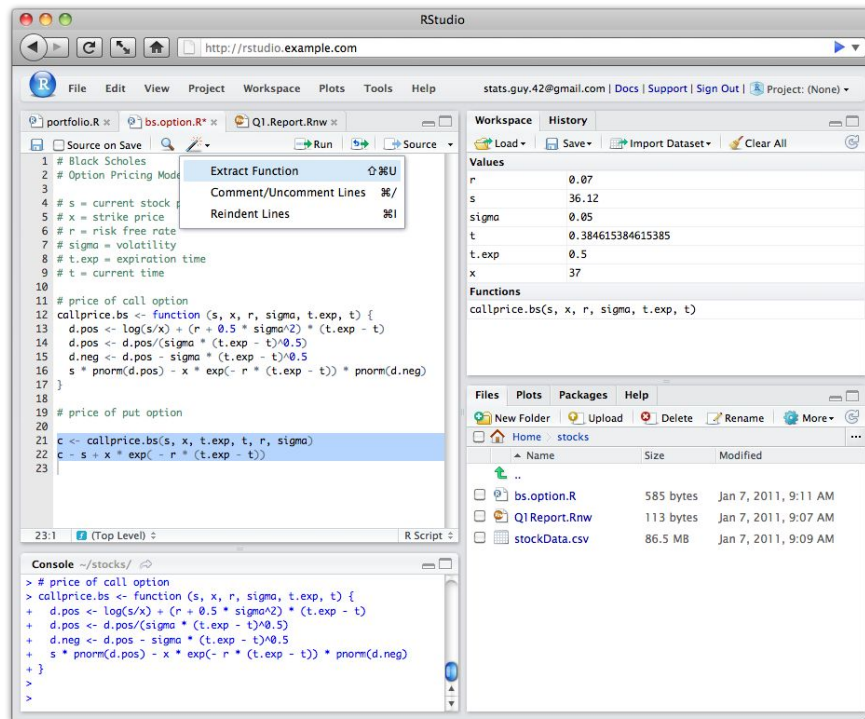# Some stats

# REPRODUCIBILITY



More than 70% of researchers have tried and failed to reproduce another scientist's experiments, and more than half have failed to reproduce their own experiments (Monya Baker, 2016).

The R ecosystem provides tools for reproducible research such as Markdown documents, interactive Notebooks, etc.

# RStudio

One of the most popular IDE.

# OUTLINE

- Basics
- R Notebooks
- Time series visualization
- Data frames
- Interactive presentations
- Machine learning
- Deep Learning with Keras

# Data frames

A **data frame** is the most common data structure in R. It can be thought of as a table.

Columns can be of different types (numeric, string, date, boolean, etc).

| | len | supp | dose |
|---|---|---|---|
| 1 | 4.2 | VC | 0.5 |
| 2 | 11.5 | VC | 0.5 |
| 3 | 7.3 | VC | 0.5 |
| 4 | 5.8 | VC | 0.5 |
| 5 | 6.4 | VC | 0.5 |
| 6 | 10.0 | VC | 0.5 |
| 7 | 11.2 | VC | 0.5 |
| 8 | 11.2 | VC | 0.5 |
| 9 | 5.2 | VC | 0.5 |
| 10 | 7.0 | VC | 0.5 |

# CODE

Code used in this presentation

https://github.com/enriquegit/data-analysis-r

# References

- Monya Baker, "1,500 scientists lift the lid on reproducibility", Nature 533, 452–454 (26 May 2016) doi:10.1038/533452a