

Computational Thinking for Social Scientists

Jae Yeon Kim

2021-12-13

Contents

Chapter 1

Hello World

```
print("Hello, World!")
```

```
## [1] "Hello, World!"
```

Make simple things simple, and complex things possible. - [Alan Kay](#)

This is the website for *Computational Thinking for Social Scientists*. This open-access book intends to help social scientists think computationally and develop proficiency with computational tools and techniques to research computational social science. Mastering these tools and techniques not only enables social scientists to collect, wrangle, analyze, and interpret data with less pain and more fun, but it also let them work on research projects that would previously seem impossible.

Horace Mann, the first great American advocate of public education, claimed that “Education, then, beyond all other divides of human origin, is a great equalizer of conditions of men—the balance wheel of the social machinery.’ I believe in this potential of education; however, I also fully acknowledge that quality education is not accessible equally. Often, the gap between education and technology is greater among disadvantaged groups. As an educator, this book is my small contribution to making this democratic vision possible, at least in the emerging field of computational social science.

That said, this book is not intended to be a comprehensive guide for computational social science or any particular programming language, computational tool, or technique. If you are interested in a general introduction to computational social science, I highly recommend [Matthew Salganik’s Bit By Bit \(2017\)](#). Salganik’s book is comprehensive, accessible, and pedagogically friendly.

The book comprises two main subjects (fundamentals and applications) and eight main sessions.

1.1 Part I Fundamentals

1. Why computational thinking
2. Best practices in data and code management using Git and Bash
3. How to wrangle, model, and visualize data easier and faster
4. How to use functional programming to automate repeated things
5. How to develop data products (e.g., packages and shiny apps)

1.2 Part II Applications

6. How to collect and parse semi-structured data at scale (e.g., using APIs and web scraping)
7. How to analyze high-dimensional data (e.g., text) using machine learning
8. How to access, query, and manage big data using SQL

The book teaches how to do all of these, mostly in **R**, and sometimes in **bash** and **Python**.

- Why R? R is free, easy to learn (thanks to **tidyverse** and **RStudio**), fast (thanks to **Rcpp**), runs everywhere (Mac/Windows/Linux), open (16,000+ packages; counting only ones **available at CRAN**), and has a growing, large, and inclusive community (**#rstats**).
- Why R + Python + bash?

“For R and Python, Python is first and foremost a programming language. And that has a lot of good features, but it tends to mean, that if you are going to do data science in Python, you have to first learn how to program in Python. Whereas I think you are going to get up and running faster with R, than with Python because there’s just a bunch more stuff built in and you don’t have to learn as many programming concepts. You can focus on being a great political scientist or whatever you do and learning enough R that you don’t have to become an expert programmer as well to get stuff done.” - Hadley Wickham

- However, this feature of the R community also raises a challenge.

Compared to other programming languages, the R community tends to be more focused on results instead of processes. Knowledge of software engineering best practices is patchy: for instance, not enough R programmers