# Data Science for Social Scientists

Daniel E. Ehrlich     Anna Tremblay

2022-08-19

# Contents

# Preface

This is my very first bookdown book. I hope you like it.

# Syllabus Details

## Course Aims

Provide students with relevant, hands on, methodological training in data literacy and visualization.

## Course Description

This course is broken down into 3, 5-week units. Unit 1 focuses on familiarizing yourself with R and the IPUMS dataset. In Unit 2, each week will showcase a method/analysis using preselected variables. In class, students will walk through a given problem set and produce a lab report by the end of class. In Unit 3, students will work towards answering a research question that they pose, creating a research paper with literature review, data analysis, conclusion, and data outputs.

## LEARNING OUTCOMES

LEARNING OUTCOMES: * Understand the depth of the IPUMS database and the variables it has to
offer * Compose R code to analyze the IPUMS data * Produce visually pleasing data outputs in R * Synthesize the information in a written report * Present the analysis in a poster format for other students

## Guiding Principles

- phenomenon-based learning
    - try to start the class with a **question** or **problem**
    - *why* does the data look the way it does

  – structure class so students work towards solving the problem

- RELEVANT examples

  – try to touch on 2 or more disciplines (eg, economics, demography)

## 0.1  Weekly Schedule

This syllabus is initially envisioned as 3 5-week sections. However, compilation and content are intended to be modular with templates for instructors to include their own specialties.

The basic structure of this course is:

### 0.1.1  Unit 1 Understanding and Testing Data

Intro to data/ simple analysis

Students will be able to:

Technical:

- Download R and RStudio
- Read data into R and
- Write (save) data out of R.
- Summarize data visually

  – Using base R
  – Using ggplot (tidyverse)

- Summarize data tabularly

  – Using base R
  – Using gttable / tidyverse

- Formally state and test assumptions of data

  – *EG:* t-test, anova, (maybe) correlations

Conceptual:

- Understand main types of data

  – *EG:* logical, numeric, character, etc
  – R specic vs general terms

- Recgonize various data distributions

  – *EG:* normal, poisson, etc

- Know which types statistical tests are appropriate for a given set of data.

## 0.1.2 Unit 2 Finding Data and Asking Questions

Here we demonstrate two **different** approaches to conducting research. Students become familar writing up short lab reports detailing their findings. For unit 0.1.2.1, we/instructor provides students with simple datasets from IPUMS (or other real-world data). Students will learn exploratory data analysis techniques and how to create lab reports to summarize key findings.

For unit 0.1.2.2, students will learn to develop their own simple research questions or social-science hypotheses. They will seek out data to answer these questions, learning to navigate ipums.org, and create **data extracts**, as well as hypothesis-testing statistical methods. Again, lab reports to summarize findings.

### 0.1.2.1 Exploratory Analysis

If you've just collected a survey, or other raw data, you may not know what you're looking for. This is perfectly ok but goes against *the scientific method* most people learned in grade school (More on that to follow(***include_link***)).

This unit begins by presenting data/distributions and asking students to begin interpreting the data . visual exploration is encouraged and basic of data manipulation are taught * *EG:* how to subset data, how to reshape data, how to recode data, how to convert from one `data type` to another.

Example lab exercise:

Students given a data set (xls, csv, etc) * load data, perform manipulations, basic summaries + cross tabs + group means by a covariate * inspect data visually + *DESCRIBE* the distribution - is it normal? significant? * *FIND* aquestion in the spread of the data + how can you test this (maybe small group work) * write up/ present results + think on confounding factors / biases

### 0.1.2.2 Hypotethsis Driven

If, on the other hand you have an a pre-exisiting idea you want to test. We can follow the traditional *scientific method*. With a question in mind, the first question is: where to look. What better place than IPUMS!

Begin introducing navigation of web resources - mainly IPUMS international

Students should become comfortable working through lab exercises: * Define a question (or be presented with one) * Download variables from IPUMS (course downloads possible) * Perform a basic analysis (discussed in Unit 1) * Generate a **visual argument** for your analysis + Include explanation/interpretation/reflection on the question at hand, and the data used + Any obvious biases + Any obvious confounding factors

### 0.1.3   Unit 3 Discussing Data and Student Research

Students will select their own research question that can be answered with the IPUMS data set and will spend five weeks producing a research paper complete with data analysis, visualization, and interpretation.

In this section we encourage the instructor to provide ample time for independent student/small-group research. Some class time should be devoted to modelling healthy discussion and critique of methods.

We provide some examples here but encourage instructors (or students) to bring in recent journal/popular articles that do (or do not) apply data science methods well.

# DEV NOTES

## TO DO

- discuss style
  - key terms section for each chapter?
  - key terms in **bold**
  - italics for *emphasis*
  - are we pro-hyphens, or are they pedantic?

## MISC IDEAS

- Application forward
- Present research/ analysis/results FIRST, then explain the mathematical principals behind it
- daily/weekly "i'm stuck on…"
  - Students send in questions (night before class) and instructor spends 10-15 mins talking through (or collaboratively working through with class) solutions
  - Alternatively, once a month maybe a longer class covering "common problems asked this month" daily/weekly "recent research"
- pick out a recent article with good visualization (or bad) and spend 5-10 mins discussing what makes it good (or bad)
  - Encourage students to find articles for extra credit

## Documentation

This function grabs any packages in your project and adds them to a local list that can be referenced using `R-pacakgename` * **NOTE** in practice, that needs to be wrapped in markdown syntax, eg: `[@R-bookdown]` * See help files for more info - might be able to create/add a `citation` file

# Chapter 1

# Introduction

Welcome to our open-source and cross-listed course on Data Science for kids who can't read data but want to learn to read data and other things too. in R.

We hope your knowledge of R increases, just like the graph below! * (This is just making up reasons for the sample plot)

```
par(mar = c(4, 4, .1, .1))
plot(pressure, type = 'b', pch = 19)
```
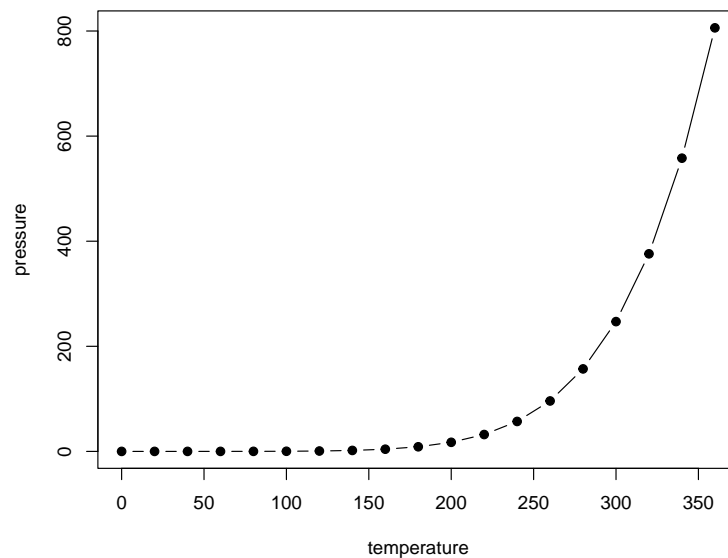


Figure 1.1: Here is a nice figure!

Table 1.1: Here is a nice table!

| Sepal.Length | Sepal.Width | Petal.Length | Petal.Width | Species |
|---:|---:|---:|---:|:---|
| 5.1 | 3.5 | 1.4 | 0.2 | setosa |
| 4.9 | 3.0 | 1.4 | 0.2 | setosa |
| 4.7 | 3.2 | 1.3 | 0.2 | setosa |
| 4.6 | 3.1 | 1.5 | 0.2 | setosa |
| 5.0 | 3.6 | 1.4 | 0.2 | setosa |
| 5.4 | 3.9 | 1.7 | 0.4 | setosa |
| 4.6 | 3.4 | 1.4 | 0.3 | setosa |
| 5.0 | 3.4 | 1.5 | 0.2 | setosa |
| 4.4 | 2.9 | 1.4 | 0.2 | setosa |
| 4.9 | 3.1 | 1.5 | 0.1 | setosa |
| 5.4 | 3.7 | 1.5 | 0.2 | setosa |
| 4.8 | 3.4 | 1.6 | 0.2 | setosa |
| 4.8 | 3.0 | 1.4 | 0.1 | setosa |
| 4.3 | 3.0 | 1.1 | 0.1 | setosa |
| 5.8 | 4.0 | 1.2 | 0.2 | setosa |
| 5.7 | 4.4 | 1.5 | 0.4 | setosa |
| 5.4 | 3.9 | 1.3 | 0.4 | setosa |
| 5.1 | 3.5 | 1.4 | 0.3 | setosa |
| 5.7 | 3.8 | 1.7 | 0.3 | setosa |
| 5.1 | 3.8 | 1.5 | 0.3 | setosa |

Reference a figure by its code chunk label with the `fig:` prefix, e.g., see Figure 1.1. Similarly, you can reference tables generated from `knitr::kable()`, e.g., see Table 1.1.

```
knitr::kable(
  head(iris, 20), caption = 'Here is a nice table!',
  booktabs = TRUE
)
```

You can write citations, too. For example, we are using the **bookdown** package (Xie, 2022) in this sample book, which was built on top of R Markdown and **knitr** (Xie, 2015).

Each `.Rmd` creates a unique **CHAPTER**. * Organizing by content probably makes the most sense * Gives more flexibility to adapt 2/3 class days. * *other thoughts??*

# Chapter 2

# The Basics

For now, I have 3 main chapters for each of the main sections: * Basics of data science / R 2 * Applications/critiques using IPUMS data 3 * Student-driven projects 4

Each of these **Chapters** contains multiple sections. We'll likely want to break these sections out into their own `.Rmd` files as they get fleshed out. For now, I'll try to keep the abundance of files limited.

**NOTE:** As these actually get filled out, we will probably want to insert different `part`s to the book (EG, the content of Unit 1 is covered in `Part I`). * Declare parts with `# (PART) Part I {-}` immediately before the first chapter `#` it contains.

**Topics to include:** * What is data? * Everything can be data * How do we interpret data * Tables * Plots * Univariate distributions * What can they tell us * Multi-modality in distributions * Categorical vs continuous data * Don't need to get ahead of this yet * Add in a grouping category - multi state/multi-national dataset * Ttest / anova

**Type of Data:** Age distributions Specifically generate a dataset with old/young folks over-represented to highlight a bimodal distribution

Start with single state/country Add a second state/country to demo ttest Add more to demo anova

Alternatively, income by education level - may be more interesting/relevant to college students (or depressing)

## 2.1   Intro to R/RStudio

## 2.2   Reading Data / Distributions

### 2.2.1   What *is* a normal distribution

#### 2.2.1.1   How normal is it?

show increasingly unclear examples of normal vs not

introduce tests of normality

#### 2.2.1.2   Measuring normality - single sample

reinforce [concept of statistical] **normality**

is a value from a sample? - one way ttest something about tails

#### 2.2.1.3   comparing normality - two saples

standard / two-way t test

#### 2.2.1.4   comparing more than two - ANOVA

# Chapter 3

# IPUMS

Some text to break up the sub-section headers

## 3.1 Intro to IPUMS website

### 3.1.1 background on ipums

### 3.1.2 navigating website

Find certain (very common) variables to answer (common) social science questions.

We describe our methods in this chapter.

Math can be added in body using usual syntax as follows. This may be useful, particularly for explaining the math side of things.

## 3.2 math example

$p$ is unknown but expected to be around 1/3. Standard error will be approximated

$$SE = \sqrt{(\frac{p(1-p)}{n})} \approx \sqrt{\frac{1/3(1-1/3)}{300}} = 0.027$$

You can also use math in footnotes like this[1]. Footnotes are helpful because they re-link to where you left off.

---

[1] where we mention $p = \frac{a}{b}$

We will approximate standard error to $0.027^2$

The `longnote` footnote seems particularly useful.

---

[2] $p$ is unknown but expected to be around 1/3. Standard error will be approximated

$$SE = \sqrt{(\frac{p(1-p)}{n})} \approx \sqrt{\frac{1/3(1-1/3)}{300}} = 0.027$$

# Chapter 4

# Independent Research

By this point, students should be familiar with basic concepts from Chapter 2. These include:

- Basic Coding
  - read/write data in/out of R
  - basic manipulations
- Theoretical Basis
  - looking at data distributions
  - formal assessment of distributions

Students will also be familiar with how these concepts are applied from Chapter 3. Hopefully students will be able to:

- Come up with a social science question they are interested in
  - Critically think about target variable(s) of interest. Any *a priori* covariates? confounders?
  - Acquire relevant data from IPUMS
  - Analyze, Summarize, Visualize Data
    * scope and complexity at student/teach discretion
  - Present research to class
    * **potentially** critically discuss/evaluate each others work.
    * **science is collaborative** everyone should be out to do their best work and represent the data as best we can. We all have conscious and unconscious biases, and the best way to confront them is share and receive (respectful) feedback.

During this Unit, we suggest giving ample class time for independent student research, peer-to-peer collaboration, and basic R/stats troubleshooting. This would also be a great time to model how to give respectful criticism by discussing recent research papers. * We could maybe come up with 1-2 seed examples, with a few talking points

## 4.1   Example one

## 4.2   Example two

# Chapter 5

# Example RMD code

For now, this chapter is a bit of a placeholder. I'm not sure what/how the `references.Rmd` file actually fits in to the code/construction (it looks automatic) so I want to keep that in place and need a section to note that.

I also want a more centralized reference point to put any example code I find helpful while working in R/bookdown. This section could get really unrully really fast, but oh well.

## 5.1 Core

`index.Rmd` is required and treated as file `00`. Chapters *should* be numbered for ease of sorting but custom orders are possible by specifying filenames somewhere **in this file**

Remember each Rmd file contains one and only one chapter, and a chapter is defined by the first-level heading `#`. + **IE** beyond the YAML header this file functions as a normal chapter since it starts with a top level header. + Note that `index.Rmd` has its own YMAL in addition to the various .yml files...not sure exactly how these relate.

Reference a figure by its code chunk label with the `fig:` prefix, e.g., see Figure 1.1. Similarly, you can reference tables generated from `knitr::kable()`, e.g., see Table 1.1. * Again, this prints an auto-generated numeral * also leaving this in the context of the plots in Chapter 1

You can write citations, too. See `knitr::write_bib()` for more on this. Quick example from demo/index (may not work without write_bib() though): we are using the **bookdown** package (Xie, 2022) in this sample book, which was built on top of R Markdown and **knitr** (Xie, 2015). * If included, "Refernces" section gets added to each chapter. * Not exactly sure where

Embed html renders (EG, fancy tables (IPUMS_var_desc), or any shiny app) with `webshot` R package and `phantomJS`.

```
install.packages("webshot")
webshot::install_phantomjs()
```

## 5.2   Tips

\***Autonumber sections** Note the `{-}` used to indicate "do not number this section" eg: preface.

**LABEL EVERYTHING** you'll likely want to reference it later \* code chunks that produce figures can be referenced via `@\ref(fig:[LABEL])`

You can label chapter and section titles using `{#label}` after them, e.g., we can reference Chapter 1. If you do not manually label them, there will be automatic labels anyway, e.g., Chapter 2. \* No idea how the automatic references work, so always be sure to declare them. \* **NOTE** these display as the relevant Chapter `numeral`.

## 5.3   Syntax

*italics* or *italics* (can handle spaces) **bold** `code` *equations*

### 5.3.1   Math

Randal Pruim features an extensive list of common math expression on their github page. Here are some quick notes:

In-line equations can be written within `$` and will be displayed right there: $a^2 + b^2 = c^2$. In contrast, you can also add equation chunks by using `$$`

This can be coded in-line,

$$\sum_{n=1}^{10} n^2$$

, but will result in a page break.

Alternatively, a more "classic" equation chunk:

$$ Plain text doesnt get spaces

how

very

odd

$$

To compile this example to PDF, you need XeLaTeX. You are recommended to install TinyTeX (which includes XeLaTeX): https://yihui.name/tinytex/.

# Bibliography

Xie, Y. (2015). *Dynamic Documents with R and knitr*. Chapman and Hall/CRC, Boca Raton, Florida, 2nd edition. ISBN 978-1498716963.

Xie, Y. (2022). *bookdown: Authoring Books and Technical Documents with R Markdown*. R package version 0.27.