

Introduction to Data Science Using R

Version: April 4, 2018

Prof. Darin Christensen¹

¹ 6341 Public Policy
darinc@luskin.ucla.edu
(310) 825-7196

1 Course Information

- **Class Location & Schedule:** Public Affairs 2250; TR, 5-6:15p
- **Section Location & Schedule:** Public Affairs 4357; W, 3-3:50p
- **Office Hours:** W, 8:00-10:30a; F, by appointment.
- **Course Website:** <https://piazza.com/ucla/spring2018/pp290dis3/home>
- **Prerequisites:** None.

2 Course Overview

There is a growing demand that public policy be “data-driven.” The White House, to take a prominent example, appointed its first-ever Chief Data Scientist in 2015. While ten weeks isn’t enough time to prepare you for that job, this class introduces you to a set of common steps that data scientists use to transform a spreadsheet into a useful description or illustration. We cover acquiring, cleaning, merging, managing, summarizing, and visualizing quantitative data.² While these tasks can be accomplished using different software, this course demonstrates how to use R, a free and **popular** statistical program.

² If you’re interested in effective data visualization consider taking Prof. Reber’s *Arguing with Data* course.

3 Learning Goals

There are several goals for this course. By the end of the quarter, you should be able to:

1. Write down a recipe (i.e., a set of ordered steps) for transforming a dataset into a table or figure that describes a pattern;
2. Identify and apply functions/packages in R for accomplishing the steps of this recipe;³
3. Combine data from multiple (unrelated) sources; and
4. Apply best practices in writing code, including using clear naming conventions, commenting steps, and avoiding repetition. These conventions help ensure that code is error-free and legible to our colleagues and future selves.

³ You won’t be expected to memorize every R function, but rather identify sources for useful code snippets and understand how to adapt those snippets.

An important disclaimer: like any new language, R is challenging to learn, and you'll spend time fighting with the program. In those darker moments, try to remember two things. First, you are not alone; your peers and I can both relate and also offer help when you're hitting a wall.⁴ Second, remember that the first goal of this course has nothing to do with R. Writing down the recipe that transforms data into insights requires careful, ordered reasoning (i.e., "algorithmic thinking"), not a specific piece of software.

⁴ Our Piazza message board is a great place to post coding questions and get quick responses.

4 Assignments

- **Participation** (10%): There are several ways to participate in and outside of class. I encourage you to raise questions and make comments during our class meetings. In addition, I will reward students for assisting peers during in-class group work, as well as posting questions or responses to our Piazza discussion board.
- **Weekly Problem Sets** (30%): These assignments serve two purposes. First, they allow you to apply the reasoning and coding skills we cover in class. Second, some problems will help you get started on your final project. While there will not be much reading, you will need to set aside time for problem sets.

All problem sets must be completed in R Markdown (discussed below). These assignments will be evaluated according to the following scale:⁵

- ✓ : complete
- ✓ – : submitted, but incomplete
- 0 : not submitted or less than half complete

- **Midterm Exam** (35%): The midterm exam offers an opportunity to assess your progress on the first two learning goals. As such, part of the exam will test your problem solving skills; another component will ask you to write R code.

The mid-term will be conducted during class time. However, you can take the exam from the location of your choice. You will need R to complete the exam.

- **Final Project / Hackathon** (25%): As a group of up to four students, you will write a quantitative report based on one of the datasets listed below. As with other "hackathons," this is a friendly competition between groups to squeeze insights out of the same data.

This project will allow you to practice the real-world tasks of manipulating, and analyzing quantitative data and clearly conveying

⁵ I will post detailed solutions online but cannot grade every problem. Please come to office hours or schedule an appointment if you have any questions.

your approach and results. In addition to the written report, you will (time permitting) briefly present your findings during the last week of class. More detail about the format and rubric for the final project will be provided.

5 Course Policies

5.1 Late Submissions

I don't want you to fall behind, so I will not accept late homework assignments. However, I understand that there are forces beyond your control, so I will drop your lowest homework grade.

Final projects will be docked a letter grade for every day that they are late. For example, an A project submitted in the first 24 hours after the deadline will receive an B. If that same project is submitted 24 to 48 hours after the deadline, the grade drops to a C.

5.2 Academic Accommodations

Academic Accommodations. Students needing academic accommodations based on a disability should contact the [Center for Accessible Education](#) (CAE). When possible, students should contact the CAE within the first two weeks of the term as reasonable notice is needed to coordinate accommodations. For more information visit.

5.3 Academic Integrity

Please review [UCLA's rules related to academic integrity](#). If you're feeling overwhelmed or are unsure about the collaboration policy, please speak with me; don't risk violating the honor code. Unfortunately, past MPP students have been investigated and punished by the Office of the Dean of Students.

Collaboration Policy:

- Problem Sets: You can work with peers (in groups of four or less). However, you must write every line of code and text that appears in any assignment that you submit. In short, never copy code or use other students' words to describe your approach or results.
- Midterm Exam: This is an individual assessment that helps us identify any knowledge gaps. I want to know what you know, so there's no collaboration.
- Final Project: You will work with a partner on your final project. You are responsible for ensuring that no written work is copied or summarized from other sources without appropriate attribution.

6 Resources

6.1 Mental Health

Let Public Policy's Student Affairs Officer (SAO) [Annie Kim](#) know if you are having academic or personal problems that affect your ability to participate fully in your education. If you feel comfortable, let me know as well.

[UCLA Counseling and Psychological Services](#) offers services and programs in a confidential environment to promote mental health and wellness.

6.2 Immigration

The [Bruin Resource Center's \(BRC\) Undocumented Student Program](#) offers caring and personalized support to undergraduate and graduate undocumented students.

Even if you are not undocumented, you may be able to get legal help for a family member. The USP office provides immigration legal services to students and their family members through a partnership with the [UC Undocumented Legal Services Center](#).

6.3 Software

You will need to acquire some software *before* the first day of class:

- You will need to download and install R from this website: <https://cran.r-project.org/>.
- Use RStudio, which can be downloaded here: <https://www.rstudio.com/products/RStudio/>. RStudio provides a nice interface for writing and executing R code and composing R Markdown documents.
- You will use R Markdown (<http://rmarkdown.rstudio.com/>) to write your homework assignments and final project. R Markdown allows you to create well-formatted documents that embed R code and outputs.⁶
- To compile R Markdown to PDF, you'll need to install [MacTeX](#) (Mac) or [MiKTeX](#) (Windows).

⁶ I won't devote class time to teaching Markdown, but I can answer questions and point you to useful online resources.

6.4 Using R

The resource we will use most is *R for Data Science* by Wickham and Golemund. An online version of the text is available here (<http://r4ds.had.co.nz/>) or you can purchase the book once it's published (<http://goo.gl/EUUNeh>).

There are many online and print resources for learning R. I've included a few below, but a Google search will turn up many others.

Online:

- Interactive Tutorials: [Try R](#)
- Examples: [UCLA's Institute for Digital Research and Education](#)
- Examples: [Cookbook for R](#)
- Crowd-sourced Solutions: [Stack Overflow](#)

Print:

- *R for Dummies* (<http://goo.gl/mel0fo>)
- *R Cookbook* (<http://goo.gl/Hh1REr>)

We are going to rely heavily on two R packages: `dplyr` for data manipulation (<https://goo.gl/Nc5gIV>), and `ggplot2` for visualization (<http://docs.ggplot2.org/current/>).

6.5 *Data Sources for Final Project / Hackathon*

I recognize that there are diverse policy interests in the course. In an effort to accommodate this, I'll permit groups to choose one of two datasets for their final project:

1. **World Bank Geocoded Research Release:** "This geocoded dataset release represents all World Bank projects in the IBRD and IDA lending lines approved from 1995-2014. This dataset tracks 5684 geocoded projects across 61243 locations."
2. **Database on Ideology, Money in Politics, and Elections:** "The resulting database contains over 130 million political contributions made by individuals and organizations to local, state, and federal elections spanning a period from 1979 to 2014."

7 Course Schedule

Legend: ★: Required; ● Suggested.

Caveat: This schedule is approximate, and we may drop or reorder sections as the quarter proceeds.

0 Before the First Class

- ★ Complete the [Try R](#) interactive tutorials; and
 - ★ Download, install, and open R, RStudio, and MiKTeX (Windows) or MacTeX (Mac).
 - ★ Open and compile this [R Markdown file](#). If this compiles to a PDF, then you're ready for class.
 - ★ Sign up for the course website at piazza.com/ucla/spring2018/pp290dis3.
-

1 Introductions (4/3)

- ★ Chapters 1 and 2: *Introduction*. Garrett Golemund and Hadley Wickham. *R for Data Science*. O'Reilly
- Mike Loukides. What is data science?, June 2010. URL <https://www.oreilly.com/ideas/what-is-data-science>
- Sharp Sight Labs. Why you should learn R first for data science, January 2015. URL <http://www.r-bloggers.com/why-you-should-learn-r-first-for-data-science/>
- Paul Curzon. A recipe for programming, 2014. URL https://teachinglondoncomputing.files.wordpress.com/2014/07/cs4fnissue16_pr_6_1.pdf

2 Description

2.1 Visualizing Data (4/5, 4/10)

- ★ Chapter 3: *Data Visualisation*. Golemund and Wickham.
- ★ Jonathan A Schwabish. An Economist's Guide to Visualizing Data. *The Journal of Economic Perspectives*, 28(1):209–34, February 2014
- H Wickham. A layered grammar of graphics. *Journal of Computational and Graphical Statistics*, 19(1), 2010

2.2 Manipulating Data Frames (4/12, 4/17)

- ★ Chapter 5: *Data Transformation*. Golemund and Wickham.

3 Wrangling

3.1 Importing Data (4/19)

- ★ Chapter 11: *Data Import*. Golemund and Wickham.

3.2 Cleaning (or “Tidying”) Data (4/24)

- ★ Chapter 12: *Tidy Data*. Golemund and Wickham.

3.3 *Working with Dates and Strings (4/26, 5/8)*

- ★ Chapter 16: *Dates and Times*. Grolemund and Wickham.
- ★ Chapter 14: *Strings*. Grolemund and Wickham.

[Midterm Exam (5/3)]

3.4 *Merging, Joining, and Appending (5/10)*

- ★ Chapter 13: *Relational Data*. Grolemund and Wickham.

4 *Coding*

4.1 *Functions (5/15)*

- ★ Chapter 19: *Functions*. Grolemund and Wickham.

4.2 *Iteration (5/17, 5/22)*

- ★ Chapter 21: *Iteration*. Grolemund and Wickham.
- ★ Chapter 20: *Vectors*. Grolemund and Wickham.

5 *Special Topics (Time Permitting)*

5.1 *GIS in R (5/24)*

- Using the `sf` Package. Here's a [vignette](#) that demonstrates how to use the package.

5.2 *Basic Models (5/29)*

- ★ Chapter 22-23: *Introduction & Model Basics*. Grolemund and Wickham.

6 *Skill Consolidation (5/31, 6/5)*

7 *Final Presentations*

[Final Project Due (11 June 2018 by 11:59 p.m.)]