

Williams College ECON 370:
Data Science for Economic Analysis

Topic 0: Getting Started

Professor: Pamela Jakiela

Outline

- Introductions
- Syllabus
- Next steps

ECON 370: Data Science for Economic Analysis

- Class meetings: Wednesdays and Fridays from 8:30 to 9:45 AM in 129 Schapiro
- Instructor: Professor Pamela Jakiela (pj5@williams.edu)
- Office hours: Mondays from 1:30 to 3:30 in 339 Schapiro or by appointment
- Prerequisites: ECON 255 or STAT 346, ideally one other economics elective using data
- Course websites:
 - ▶ <https://pjakiela.github.io/ECON370/> (course materials)
 - ▶ <https://www.gradescope.com/courses/854937> (submitting assignments)

Learning Objectives

- Developing advanced data skills, with a particular focus on identifying data new sources, data cleaning and wrangling, exploratory data analysis, and data visualization
- Learning to implement new statistical methods (e.g. machine learning, text analysis) which are increasingly being incorporated into the empirical economist's toolbox, and learning how these tools are used by economists for prediction and causal inference
- Honing the ability to master new empirical tools such as programming languages (e.g. R, Python) and statistical methods (e.g. machine learning, text analysis) as they emerge

TLDR: how to be a super RA and be ready for grad school in quantitative social science

What This Course Is About

- Working with data: finding data, cleaning/wrangling data, visualization/EDA
- Machine/statistical learning: shrinkage methods, random/causal forests, clustering
- Text analysis: working with text, web scraping, word frequencies, topic models

What This Course Is NOT About

- Big data
- Computation
- Database management
- Data science methods that are focused exclusively on prediction
- Causal inference, statistical techniques that are not used by economists

Stata vs. R vs. Python

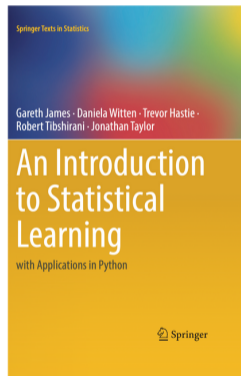
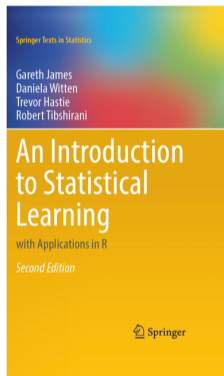
ECON 255 is taught in Stata, and Stata is still by far the most widely used in by economists

- Many of the techniques we'll cover in this class can't be done in Stata (at present)
- Economics is shifting away from Stata, albeit slowly; not yet clear where we are heading
- Being able to get yourself up-to-speed in a new language is a (data science) skill in itself
- Familiarity with R or Python is not a prerequisite (as long as you know Stata)

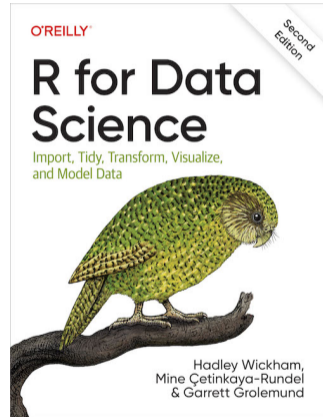
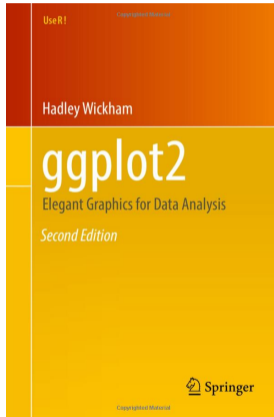
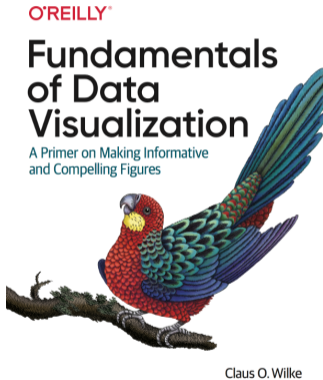
I plan to teach the course primarily in R, to the extent that I focus on a particular language

- Lab 1 is in both R and Python (you will submit two versions)
- Labs 2 through 15 can (almost all?) be done in R or Python
- I would like you to end the course (somewhat) familiar with both languages
- Absolutely fine to focus your energy on your preferred language

Readings: Download This

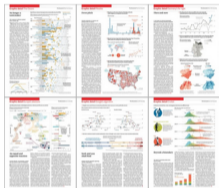


Readings: Bookmark These



Topic-Specific Readings Will Be Posted Online

ECON 370



source: The Economist

Instructor:
Pamela Jakiela

[home](#)
[syllabus](#)
[schedule](#)

1 Data

Readings

These are useful references as opposed to required readings.

R for Data Science: 7, 5, and 19

Intro to Data Science: 2, 4, 6

An Intro to Stats with Python: 2.1, 2.2, and 3.1

Lecture

Slides from Lecture 1

Lab

Assignments: Labs

```
## ECON 370 LAB 1: DATA CLEANING AND PREP
## NAME:
## DATE:

# preliminaries -----

## install packages / load libraries as needed and load libraries
## you'll need at least tidyverse for R, numpy and pandas for Python

## specify your file path by defining mypath as the path to your working directory for this lab

## R suggestion: mypath <- "[YOUR FILE PATH IN QUOTES]"
## Python suggestion: mypath = "[YOUR FILE PATH IN QUOTES]"

# get the barro-lee data on educational attainment -----

## load barro-lee educational attainment data set from a csv file
## the barro lee website is here: http://barrolee.com/
## download the csv for the data set on Education Attainment for Population Aged 15 and Over (Total Population)
## the file should be called: BL2013_MF1599_v2.2.csv
## load the data as a tibble or pandas dataframe called bldata

## R suggestion: read_csv()
## Python suggestion: mypath = pd.read_csv

## how many observations are in the data set? how many variables?
## which variables are string (or character) variables?

## R suggestions: dim(), spec(), head(), summary()
## Python suggestions: type(), df.info, df.shape, df.columns, df.dtypes

## define a new tibble called bl2010 which only contains BL data from 2010
## keep only the columns: BLcode, country, yr_sch, MBcode, region_code
## rename yr_sch as sean_edu, MBcode as isocode, and region_code as sb_region

## R suggestion: use the pipe and the filter, select, and rename functions
## Python suggestion: df[df['x1'] == VALUE] selects a subset of the rows of dataframe df
## Python suggestion: df[['x1', 'x2', 'x3']] selects a subset of the columns of dataframe df
## Python suggestion: rename columns with
## df = df.rename(columns={
##     'old_var_name': 'new_var_name',
## })

## make a table/tibble summarizing the average rate of educational attainment by region

## R suggestion: use group_by() and pipe to summarize()
## Python suggestion: use df.groupby().agg()
```

Assignments: Projects/Presentations

- Data visualization project using provisions data
- Using machine learning to predict infant mortality in the Demographic and Health Surveys
- Final project on a topic of your choice (text? unsupervised learning? new data?)

Presentations and Other Important Dates

Date	Description
9/20	Guest Speaker: Bilal Zia (coffee at 9:00, brownbag at 12:00)
9/25	Data Visualization Group Meetings
9/27	Data Visualization Project Presentations (possibly at Provisions)
10/16 or 10/18	Predicting Infant Mortality Group Meetings
10/23, 10/25	Predicting Infant Mortality Presentations
11/22	Final Project Group Meetings
12/4, 12/6	Final Project Presentations

Assignments and Grading

Lab Assignments	45 points
Provisions Data Visualization Project	12 points
DHS Predicting Infant Mortality Project	12 points
Final Empirical Project	21 points
Class Participation in Lecture and Lab	9 points
Getting-to-Know-You Survey	1 point

ChatGPT

- The words, analysis tables, and figures/images you produce should be your own
 - ▶ If the English words/sentences, tables, or figures that you submit can be found in published articles, on the internet, or anywhere else, you do not receive credit for that lab/project (as if you submitted an article to a journal but that article was already published elsewhere)
 - ▶ You must always submit replication code that generates your tables and figures
- Generative AI is a powerful coding tool, and you should learn to use it
 - ▶ I recommend handing it small coding tasks (e.g. turning a few lines of R code into Python)
 - ▶ Make sure you understand each step in your final script/program (I may ask)
 - ▶ It is **always** your responsibility to make sure that your code is correct, and that it runs

Next Steps

1. Install R and RStudio
2. Install Python
3. Create a gradescope account
4. Complete the getting-to-know-you survey
5. Create a DHS account and request data for Kenya + one other African country