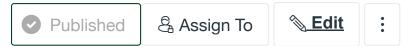
# 1.1 Prepare Datasets



Download this <u>step-by-step PDF</u> ⇒

(http://static.trogu.com/documents/classes/523/2025/fall2025/1.1\_prepare\_datasets\_step-by-step\_F2025.pdf) to complete this assignment. The PDF also includes links to videos. Some links in the PDF have changed but we will go over the changes in class as we do the exercise together.

Below is some important software information and the list of files to upload at the end of the assignment.

## Project 1 is a 4-week project composed of three parts:

This page is Part 1: 1.1 Prepare datasets

2 clean comma-separated value files (CSV)

3 plots (PDF)

1 R file (R)

Part 2: 1.2 Generate graphs (line, bar, scatterplot)

(<u>https://sfsu.instructure.com/courses/61861/assignments/697651</u>)

Part 3: 1.3 Clean up and format final graphs in Illustrator

(<u>https://sfsu.instructure.com/courses/61861/assignments/697652</u>)

**Note:** for spreadsheets, you must use Excel, not Google Sheets, because by editing in a web browser, HTML can introduce "dirty" characters like "/" that can cause errors later in the visualizations. The clean dataset file should always be checked with a text-only editor. Download all the necessary programs, including **Excel**  $\Longrightarrow$ 

(https://its.sfsu.edu/service/office365students), BBEdit ⇒

(<u>https://www.barebones.com/products/bbedit/</u>) (Mac) or <u>Notepad++</u> ⇒ (<u>https://notepad-plus-plus.org/downloads/</u>) (PC), and <u>R</u> ⇒ (<u>https://www.r-project.org/</u>).

Some notes on the plots and info for dowloading and installing R and RStudio.

The goal of this project is to visualize the disparity in breast cancer mortality between

White females and Black females, and show that the rate is almost 50% higher for Black, despite them having a lower incidence rate (fewer cases per 100K). We will not look at causes in this assignment, but it turns out that while lack of health care plays a big role, genetics also does.

After creating the two dataset files, plot a matrix for each using R, showing every possible combination of pairs of variables (columns). There is a bit of coding, just a few short lines. Plot also a scatterplot of just the white and black death rates.

Note: although R (or any other program) won't be required as the only program to use in the class, this little exercise will show that it's very reliable to just get a base graph plotted, and therefore I highly recommend it. Illustrator can then be used to clean up the base graph.

Use the links below to download and install all the necessary software for this assignment:

## BBEdit (Mac)

https://www.barebones.com/products/bbedit/download.html (https://www.barebones.com/products/bbedit/download.html)

Notepad++ (PC)

<u>https://notepad-plus-plus.org/downloads/</u> <u>□→ (https://notepad-plus-plus.org/downloads/</u>)

# **Excel (MS Office)**

• Install **Excel** (https://its.sfsu.edu/service/office365students), it's free from the university and part of the full Microsoft Office 365 suite:

https://its.sfsu.edu/service/office365students ⊟

## (https://its.sfsu.edu/service/office365students)

Note; work in Google Sheets will not be accepted, you must install Microsoft Excel!

### R

• Download and install R (R Project: <a href="https://www.r-project.org">https://www.r-project.org</a> (<a href="https://www.r-project.org">https://www.r-project.org</a>

From one of the servers: <a href="https://cran.r-project.org/mirrors.html">https://cran.r-project.org/mirrors.html</a> (<a href="https://cran.r-project.org/mirrors.html">https://cran.r-project.org/mirrors.html</a>)

For example Iowa State University: <a href="https://mirror.las.iastate.edu/CRAN/">https://mirror.las.iastate.edu/CRAN/</a>)

The R website is very old school looking — just make sure you download the correct version for your system, Mac or PC.

Once on the mirror site, download the latest version for Mac. Check you OS for compatibility but the two versions below (Mac or PC) should work.

<u>https://cran.r-project.org/bin/macosx/</u> <u>⇒ (https://cran.r-project.org/bin/macosx/)</u>

Latest release: R-4.5.1-arm64.pkg ⇒ (https://cran.r-project.org/bin/macosx/big-sur-arm64/base/R-4.5.1-arm64.pkg) (Please note that this version only works on silicon Macs: M1,2,...)

For older Intel-based mac: R-4.5.1-x86\_64.pkg (https://cran.r-project.org/bin/macosx/big-sur-x86\_64/base/R-4.5.1-x86\_64.pkg)

Or the latest version for Windows: <a href="https://cran.r-project.org/bin/windows/">https://cran.r-project.org/bin/windows/</a>)

Click on <u>base</u> <u>⇒ (https://cran.r-project.org/bin/windows/base/)</u>, then click on:

**Download R-4.3.1 for Windows** ⇒

(https://mirror.las.iastate.edu/CRAN/bin/windows/base/R-4.5.1-win.exe)

#### **RStudio**

After installing R (RStudio will remind you) download and install also **RStudio** (https://posit.co/download/rstudio-desktop/), which is a graphic interface that runs on

top of R. Download the free RStudio Desktop version:

<u>https://posit.co/download/rstudio-desktop/</u> <u>□→ (https://posit.co/download/rstudio-desktop/</u>)

Go to the <u>step-by-step PDF</u> ⇒

(<a href="http://static.trogu.com/documents/classes/523/2025/fall2025/1.1\_prepare\_datasets\_step-by-step\_F2025.pdf">http://static.trogu.com/documents/classes/523/2025/fall2025/1.1\_prepare\_datasets\_step-by-step\_F2025.pdf</a>) for file and plot instructions. (See also note at top of this page). Upload 6 files by the deadline:

- 1. lastName brecan 75 17.csv
- 2. lastName\_brecan\_wb\_2019.csv
- 3. lastName\_7517\_plot1.pdf
- 4. lastName\_wb\_19\_plot2.pdf
- 5. lastName\_wb\_19\_scatterplot3.pdf
- 6. lastName\_brecan.R

Points 10

Submitting a file upload

Due	For	Available from	Until
Aug 27 at 8:30am	Everyone	-	-

+ Rubric