

Homework 2

Due September 16, 2020

2020-08-25

For each assignment, turn in by the due date/time. Late assignments must be arranged prior to submission. In every case, assignments are to be typed neatly using proper English in Markdown.

In Lectures 2 and 3, we spoke about Reproducible Research, R and version control, getting, cleaning, munging and ‘tidying’ data and finally, summarizing data. Again, we are focusing on Reproducible Analysis which, for us, is accomplished by mixing code, figures and text into a cohesive document that fully describes both the process we took to go from data to results and the rationale behind our data driven conclusions. In this exercise we begin creating tidy data sets. While others have proposed standards for sharing data with statisticians, as practicing data scientists, we realize the often onerous task of getting, cleaning and formatting data is usually in our hands. From here on out, we will use GitHub to retrieve and turn in the homework assignments.

Problem 1

Work through Primers titled “Work with Data” and “Tidy Your Data”.

For a different take on how to use R, try swirl:

```
install.packages("swirl")
install_course("Getting_and_Cleaning_Data")
install_course("Exploratory_Data_Analysis")
swirl()
```

These can be a little cheezy, but are overall a pretty decent intro. Each one takes 5-10 min. Pick and choose as you like. Don’t worry about submitting anything at the end, these were used for credits in a Coursera Data Science Course.

Problem 2

Create a new R Markdown file within your local GitHub repo folder (file->new->R Markdown->save as).

The filename should be: HW2_pid, i.e. for me it would be HW2_rsettag

You will use this new R Markdown file to solve problems 3-5.

Problem 3

In the lecture, there were two links to StackOverflow questions on why one should use version control. In your own words, summarize your thoughts (2-3 sentences) on version control in your future work. No penalties here if you say, useless!

Problem 4

In this exercise, you will import, munge, clean and summarize datasets from Wu and Hamada’s *Experiments: Planning, Design and Analysis* book you will use in the Spring. For each dataset, you should perform the cleaning 2x: first with base R functions (ie no dplyr, piping, etc), second using tidyverse function. Make sure you weave your code and text into a complete description of the process and end by creating a tidy dataset

describing the variables, create a summary table of the data (summary, NOT full listing), note issues with the data, and include an informative plot.

- a. Sensory data from five operators. – see video, I am doing this one
<http://www2.isye.gatech.edu/~jeffwu/wuhamadabook/data/Sensory.dat>
- b. Gold Medal performance for Olympic Men's Long Jump, year is coded as 1900=0.
<http://www2.isye.gatech.edu/~jeffwu/wuhamadabook/data/LongJumpData.dat>
- c. Brain weight (g) and body weight (kg) for 62 species.
<http://www2.isye.gatech.edu/~jeffwu/wuhamadabook/data/BrainandBodyWeight.dat>
- d. Triplicate measurements of tomato yield for two varieties of tomatoes at three planting densities.
<http://www2.isye.gatech.edu/~jeffwu/wuhamadabook/data/tomato.dat>

Problem 5

Finish this homework by pushing your changes to your repo. In general, your workflow for this should be:

1. git pull – to make sure you have the most recent repo
2. In R: do some work
3. git add – this tells git to track new files
4. git commit – make message INFORMATIVE and USEFUL
5. git push – this pushes your local changes to the repo

If you have difficulty with steps 1-5, git is not correctly or completely setup. See me for help.

Only submit the .Rmd and .pdf solution files. Names should be formatted HW2_lastname.Rmd and HW2_lastname.pdf

Optional preperation for next class:

TBD