

# Lab Session 1: Read and manipulate data

Ari Anisfeld

8/26/2021

We expect you to watch the `class 1` material, [here](#) prior to lab.

In this lab we'll use two data sets. You can download them in advance [here](#) and [here](#).

## Working with data and scripts (To be covered during Q&A session)

We recommend a file structure for coding lab. If you have your own preferred way of organizing code feel free to follow it.

### Setting up working directory and coding environment

1. Do you have a folder on your computer for coding lab material?  
If not, create one and make sure you know the path to the folder.
2. We recommend creating a `problem_set` folder inside your coding lab folder.
3. Make folder called `data` inside the `problem_set` folder.

### Putting your files in place

4. Create a new R script. Save your script in the `problem_set` folder. From now on, when you start a script or Rmd save it there.
5. Download the first data set [here](#) and put the data in your `data` folder.<sup>1</sup>

### Tell R where to find files

- Local paths are like addresses on your computer. Use `getwd()` to see how your computer
- In lab0 we downloaded data form a URL which is an address on the internet.<sup>2</sup>

```
# compare url
covid_data <- read_csv("https://data.cdc.gov/api/views/qfhf-uhaa/rows.csv?")
# compare path
wid_data <-
  read_xlsx(
    "~/coding-lab/harris-coding-lab.github.io/data/world_wealth_inequality.xlsx"
  )
```

<sup>1</sup>The data source is the World Inequality Database where you can find data about the distribution of income and wealth in several contries over time. Outside of lab time, check out [wid.world](#) for more information.

<sup>2</sup>Fun fact: `read_xlsx()` cannot download and open data from the internet. I googled “Error: `path` does not exist: downloading xlsx from github” and found this on [github](#)

6. Add a line to your script where you `setwd()` to the `data` folder.

## Working with the files

7. Finally, we are using data in an excel format. We need the package `readxl` to process data of this type. In the console, run `install.packages("readxl")`.
8. Add code to load the `tidyverse`.
9. If you followed the set-up from above, you should be able to run the following code with no error.

```
wid_data <- read_xlsx("world_wealth_inequality.xlsx")
```

## Warm-up (Lab 1)

1. Your stats 1 partner comes to you and says they can't get data to load after restarting R. You see the code:

```
install.packages("haven")
awesome_data <- read_dta("awesome_data.dta")

Error in read_dta("awesome_data.dta") :
could not find function "read_dta"
```

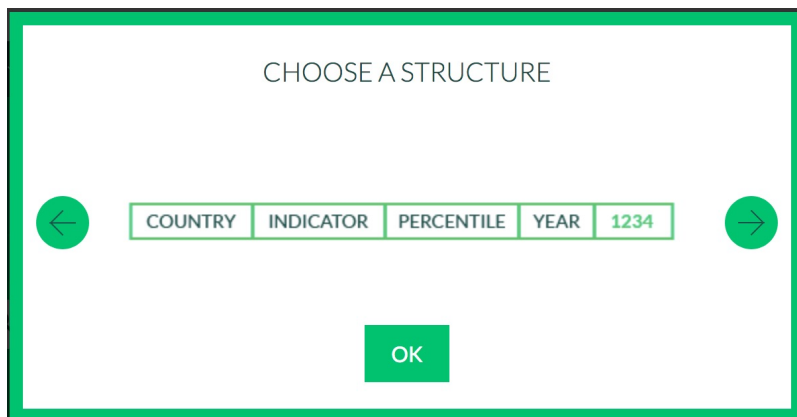
- a. Diagnose the problem.

Note: If they say the code worked before, it's likely they had loaded `haven` in the console or perhaps in an earlier script. R packages will stay attached as long as the R session is live.

2. In general, once you have successfully used `install.packages(pkg)` for a "pkg", you won't need to do it again. Install `haven` using the console.
3. In your script, load `haven`. Notice that if you had to restart R right now. You could reproduce the entire warm-up by running the script. We strive for reproducibility by keeping the code we want organized in scripts or `Rmds`.
4. It's good practice when starting a new project to clear your R environment. This helps you make sure you are not relying on data or functions you wrote in another project. After you `library()` statements add the following code `rm(list = ls())`.
5. `rm()` is short for remove. Find the examples in `?rm` and run them in the console.

## Examining 'wid\_data

1. Look at the data. What is the main problem here?
2. We don't have columns headers. The World Inequality Database says the "structure" of the download is as shown in the image below.



So we can create our own header in `read_xlsx`.

```
wid_data_raw <-
  read_xlsx("world_wealth_inequality.xlsx",
            col_names = c("country", "indicator", "percentile", "year", "value"))
```

Now when we look at the second column. It's a mess. We can separate it based on where the `\n` (or on some Windows computers where `\r\n`) are and then deal with the data later.

```
wid_data_raw <-
  read_xlsx("world_wealth_inequality.xlsx",
            col_names = c("country", "indicator", "percentile", "year", "value")) %>%
  separate(indicator, sep = "[\\r]?\\n", into = c("row_tag", "type", "notes"))
```

NOTE: We want a clean reproducible script so you should just have one block of code reading the data: that last one. The other code were building blocks. If you want to keep “extra” code temporarily in your script you can use `#` to comment out the code.

## manipulating world inequality data with `dplyr` (20 - 25 minutes)

Now we have some data and are ready to use `select()`, `filter()`, `mutate()`, `summarize()` and `arrange()` to explore it.

1. The data comes with some redundant columns that add clutter when we examine the data. What `dplyr` verb let's you choose what columns to see? Remove the unwanted column `row_tag` and move `notes` to the last column position and assign the output to the name `wid_data`.<sup>3</sup>
2. Let's start to dig into the data. We have two types of data: “Net personal wealth” and “National income”. Start by `filter()`ing the data so we only have “Net personal wealth” for France, name the resulting data `french_data` and then run the code below to visualize the data.

```
# replace each ... with relevant code
french_data <- wid_data %>%
  filter( ... , ...)
```

Note: When referring to words in the data, make sure they are in quotes “France”, “Net personal wealth”. When referring to columns, do not use quotes. We'll talk about data types in the next lecture.

```
french_data %>%
  ggplot(aes(y = value, x = year, color = percentile)) +
  geom_line()
```

<sup>3</sup>Hint: You can type all the column names or use the slicker `select(-notes, everything())`

Now we're getting somewhere! The plot shows the proportion of national wealth owned by different segments of French society overtime. For example in 2000, the top 1 percent owned roughly 28 percent of the wealth, while the bottom 50 percent owned about 7 percent.

1. Explain the gaps in the plot. Using `filter()`, look at `french_data` in the years between 1960 and 1970. Does what you see line up with what you guessed by looking at the graph?
2. Using `mutate()`, create a new column called `perc_national_wealth` that equals `value` multiplied by 100. Adjust the graph code so that the y axis shows `perc_national_wealth` instead of `value`.
3. Now following the same steps, explore data from the "Russian Federation".
4. The data for "Russian Federation" does not start in 1900, but our y-axis does. That's because we have a bunch of NAs. Let's filter out the NAs and remake the plot. You cannot test for NA using `==`. Run: `NA == NA`, we see the result is `NA`, but we need `TRUE` or `FALSE` to know what to filter! Instead we have a function called `is.na()`. (Try: `is.na(NA)` and `!is.na(NA)`).
5. Use two `dplyr` verbs to figure out what year the bottom 50 percent held the least wealth. First, choose the rows that cover the bottom 50 percent and then sort the data in descending order using `arrange()`.<sup>4</sup>.

```
# replace ... with relevant code
russian_data %>%
  filter(...) %>%
  arrange(...)
```

1. For both the Russian Federation and French data, calculate the average proportion of wealth owned by the top 10 percent over the period from 1995 to 2010. You'll have to filter and then summarize with `summarize()`.

```
# replace ... with relevant code
russian_data %>%
  filter(...) %>%
  summarize(top10 = mean(...))
```

## manipulating midwest demographic data with dplyr

1. Now we'll use midwestern demographic data which is [at this link](#). The dataset includes county level data for a single year. We call data this type of data "cross-sectional" since it gives a point-in-time cross-section of the counties of the midwest. (The world inequality data is "timeseries" data).
2. Save `midwest.dta` in your data folder and load it into R with the name `midwest`.<sup>5</sup>
3. Run the following code to get a sense of what the data looks like:

```
glimpse(midwest)
```

4. I wanted a tibble called `midwest_pop` that only had county identifiers and the 9 columns from `midwest` concerned with population counts. Replicate my work to create `midwest_pop` on your own.

```
names(midwest_pop)
```

```
## [1] "county"      "state"      "poptotal"   "popdensity"
## [5] "popwhite"    "popblack"   "popamerindian" "popasian"
## [9] "popother"    "popadults" "poppovertyknown"
```

<sup>4</sup>Hint: Look at the examples in `?arrange` or at your notes

<sup>5</sup>Hint: `read_dta()`

Hint 1: I went to `?select` and found a *selection helper* that allowed me to select those 9 columns without typing all their names.<sup>6</sup>

```
# replace ... with relevant code
midwest_pop <- midwest %>% select(county, state, ...)
```

5. From `midwest_pop` calculate the area of each county. Notice that  $\text{popdensity} = \frac{\text{poptotal}}{\text{area}}$ .<sup>7</sup> What's the largest county in the midwest? How about in Illinois?
6. From `midwest_pop` calculate percentage adults for each county. What county in the midwest has the highest proportion of adults? What's county in the midwest has the lowest proportion of adults?
7. How many people live in Michigan?

Want to improve this tutorial? Report any suggestions/bugs/improvements on [here!](#) We're interested in learning from you how we can make this tutorial better.

---

<sup>6</sup>Hint 2: notice that all the columns start with the same few letters.

<sup>7</sup>in `midwest` there's an `area` variable, but we want you to create your own using the information in `midwest_pop`.