

Data Visualization

First Steps: Getting Started

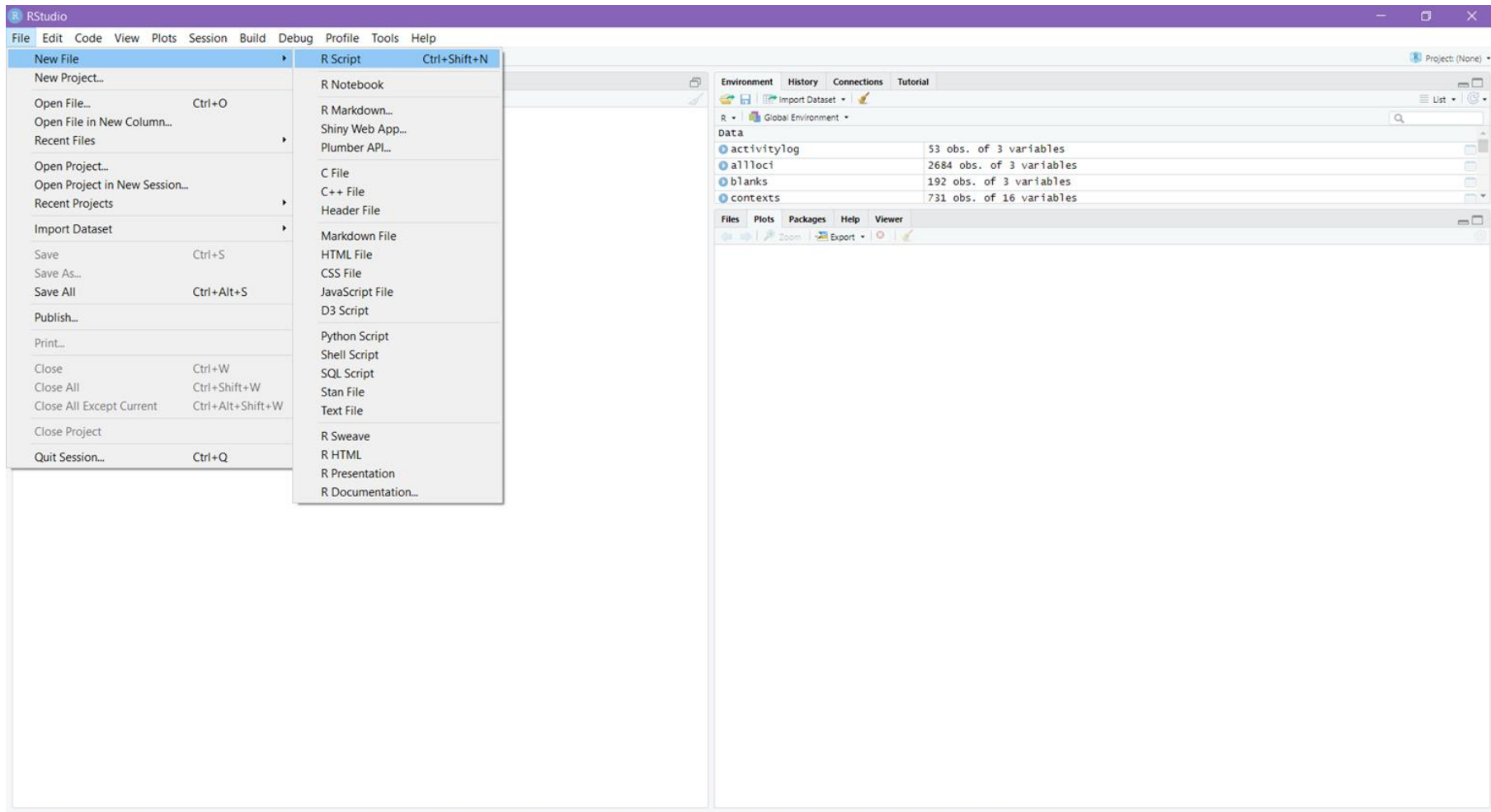
Ciara Zogheib

Data Sciences Institute, University of Toronto

In this lesson, we will...

- Work through practice code from Chapter 2 (*Getting Started*) of Healy, K. (2018). *Data Visualization: A Practical Introduction*. Princeton University Press. This code will let us:
 - Set up our workspace in RStudio
 - Get the data we want to visualize into RStudio
 - Make our first data viz with RStudio!

Setting up R for data visualization



The image shows the RStudio desktop environment. A large blue arrow points from the text 'This is our R Script where we type our code' to the 'Untitled1' script editor. Another blue arrow points from the text 'Plots and figures will appear in the 'Plots' tab' to the 'Plots' tab in the 'Environment' pane.

This is our R Script where we type our code

Plots and figures will appear in the 'Plots' tab

Environment

| Object | Class | Attributes |
|-------------|------------|--------------------------|
| activitylog | data.frame | 53 obs. of 3 variables |
| allicci | data.frame | 2684 obs. of 3 variables |
| blanks | data.frame | 192 obs. of 3 variables |
| contexts | data.frame | 333 obs. of 16 variables |

Files

Plots

Console

```
R version 4.0.3 (2021-03-31) -- "Snake and Throw"
Copyright (C) 2021 The R Foundation for Statistical Computing
Platform: x86_64-w64-mingw32/x64 (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

Natural language support but running in an English locale


R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

warning: namespace 'pool' is not available and has been replaced
by .GlobalEnv when processing object 'pool'
[workspace loaded from ~/.RData]

> |
```

We click 'Install' in the 'Packages' tab (next to 'Plots') to install the packages we need



The image shows the RStudio interface with the 'Packages' tab selected. The 'Global Environment' pane on the right lists installed packages: activitylog (53 obs. of 3 variables), allloci (2684 obs. of 3 variables), blanks (192 obs. of 3 variables), and contexts (731 obs. of 16 variables). The 'User Library' pane lists various installed packages, including abind, animation, and many others. The 'Install Packages' dialog box is open, showing the 'Repository (CRAN)' selected and the packages 'tidyverse' and 'tidyverse-brary' entered. The 'Install dependencies' checkbox is checked. The 'Install' button is highlighted.

```
by .GlobalEnv when processing object "pool"
[Workspace loaded from ~/.RData]

> library(tidyverse)
-- Attaching packages ----- tidyverse 1.3.1 --
v ggplot2 3.3.3    v purrr  0.3.4
v tibble  3.1.1    v dplyr  1.0.5
v tidyr   1.1.3    v stringr 1.4.0
v readr   1.4.0    v forcats 0.5.1
-- Conflicts ----- tidyverse_conflicts() --
x dplyr::filter() masks stats::filter()
x dplyr::lag()    masks stats::lag()
warning messages:
1: In function(kind = NULL, normal.kind = NULL, sample.kind = NULL) :
  non-uniform 'Rounding' sampler used
2: In function(kind = NULL, normal.kind = NULL, sample.kind = NULL) :
  non-uniform 'Rounding' sampler used
> library(socviz)
+
+
> library(socviz)
Error in library(socviz) : there is no package called 'socviz'
>
```

Load required packages

- Once our packages have been installed, type the following in your script:

```
library(tidyverse)  
library(socviz)  
library(ggplot2)
```

- ctrl+enter (Windows) or cmd+enter (Mac) over highlighted code to run it and load the libraries

Load required packages

- We will (later!) be using the ggplot library to make our visualizations

```
library(ggplot2)
```


Getting our data into R

For our purposes...

- In this class, we will be using already-prepared datasets from the socviz library rather than separate data files
- BUT when visualizing data for your own purposes, you will most likely need to read in data from flat files
- We do this using the `read_csv()` function from the readr package to assign our dataset to an object from our computer or from a URL:

```
ourdata <- read_csv(/location/of/file.csv)
ourdata <- read_csv(file="fakewebsite.com/file.csv")
```

Preparing our data for visualizing

- For graphing with ggplot (AKA what we want to do), our data should be in **long** format as opposed to **wide** format
- That is: **in long formatted data, every observation should be a row, and every variable should be a column**

Make a basic figure

Load our sample dataset

- We will use a preloaded sample dataset called 'gapminder', which can be loaded as follows:

```
install.packages("gapminder")  
  
library(gapminder)  
  
gapminder
```

- We can see that our dataset contains data about countries over several years

Make a scatterplot

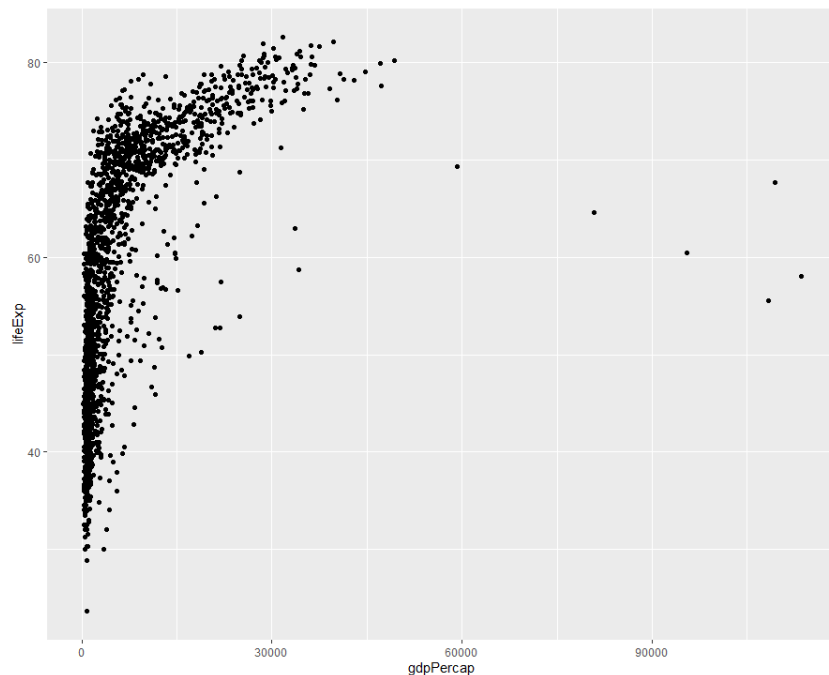
- To make some sense of what is happening in our dataset, we can make a simple scatterplot as follows:

```
p <- ggplot(data = gapminder,  
            mapping = aes(x = gdpPercap, y = lifeExp))  
p + geom_point()
```

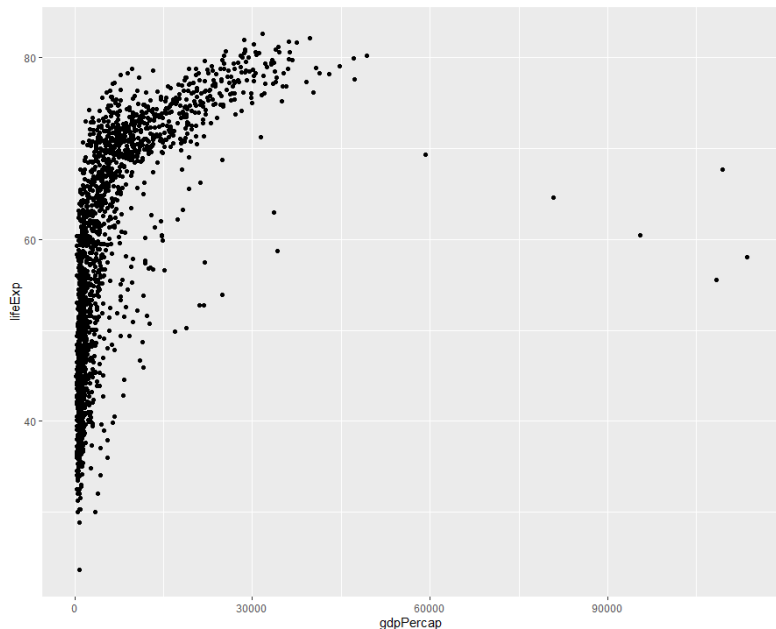
- (Don't worry about the code yet, we'll break it down later)

Make a scatterplot

- If we run the code, we should see an output in our 'plots' pane:



Assess our first plot



The Good

- Mostly legible
- Has axis labels
- Shows some kind of relationship between variables

The Not-So-Good

- Font size is tiny
- Not visually interesting or pleasing
- What are we trying to show?

Next...

- Making a plot
- Understanding the different components of the ggplot we made in this class