# 1.3.2: Data Wrangling

## (Asynchronous-Online)

## Session Objectives

1. To read in data.
2. To subset data.
3. To select, create, and modify variables.
4. To filter the observations according to a certain condition.
5. To explore and summarize data.
6. To run descriptive statistics.

Key points to cover:

1. Import data.
2. Introduce to tidyverse.
3. Summarize data and run descriptive statistics.
4. Introduce other resources (e.g., books, blogs, or websites) trainees can refer to.

---

## 0. Prework - Before You Begin

### Install Packages

Before you begin, please go ahead and install the following packages - these are all on CRAN, so you can install them using the RStudio Menu Tools/Install Packages interface:

- haven
- readr
- readxl
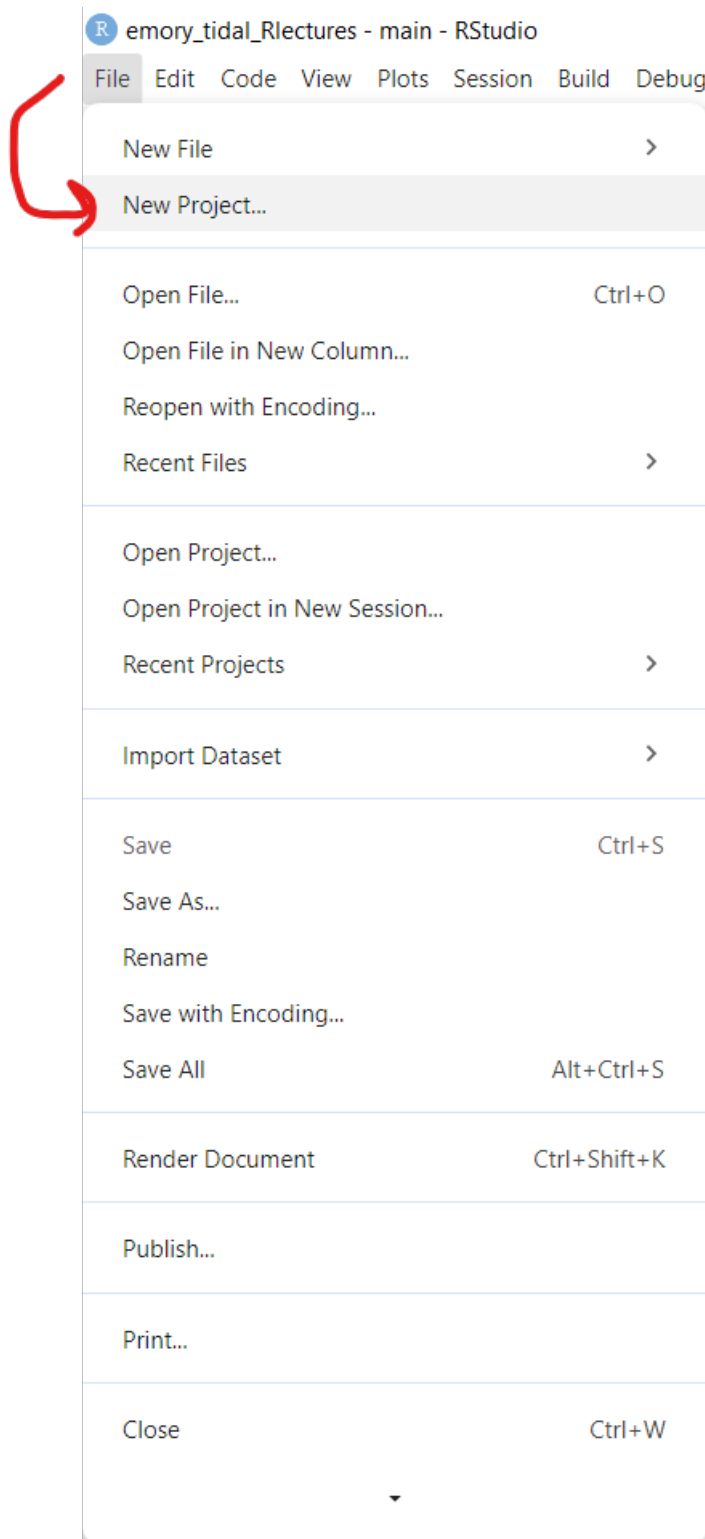- dplyr

See Module 1.3.1 on Installing Packages

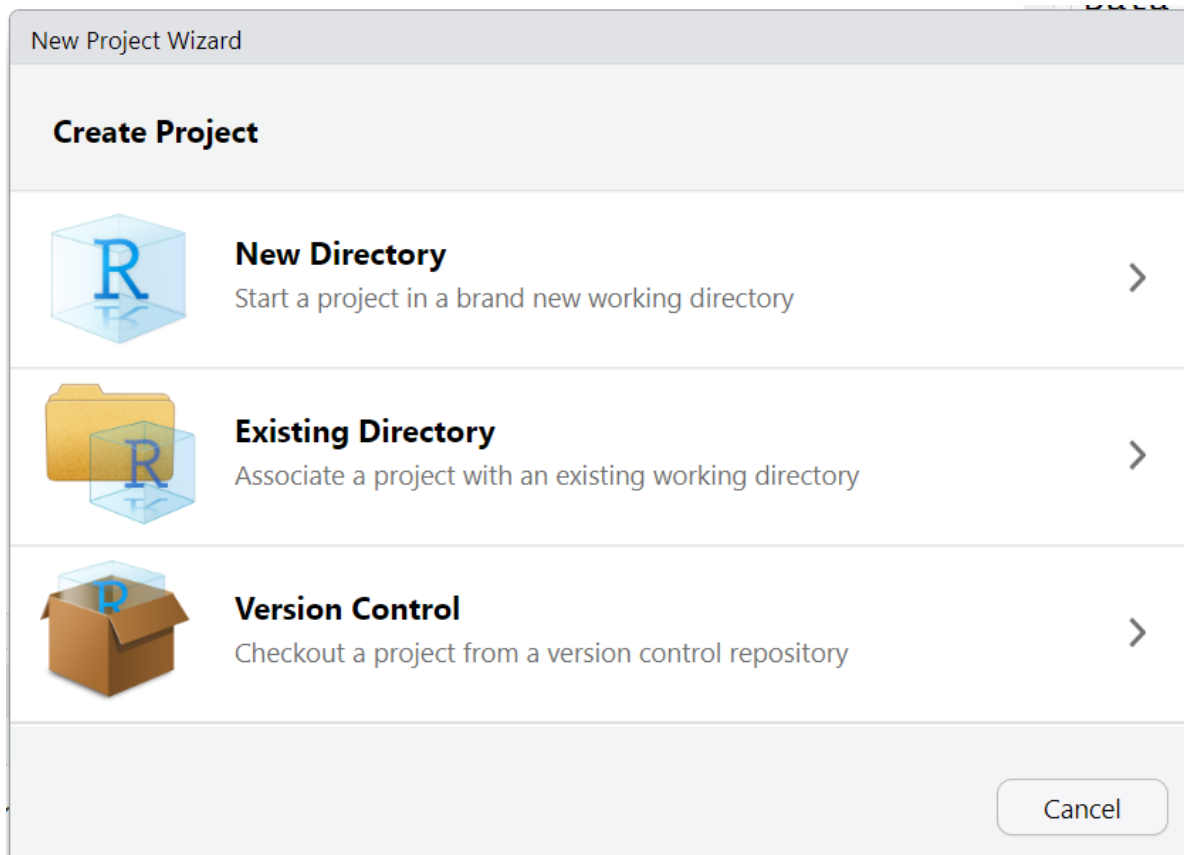## 1. To read in data.

**Begin with a NEW RStudio Project**

Let's begin with a new RStudio Project.

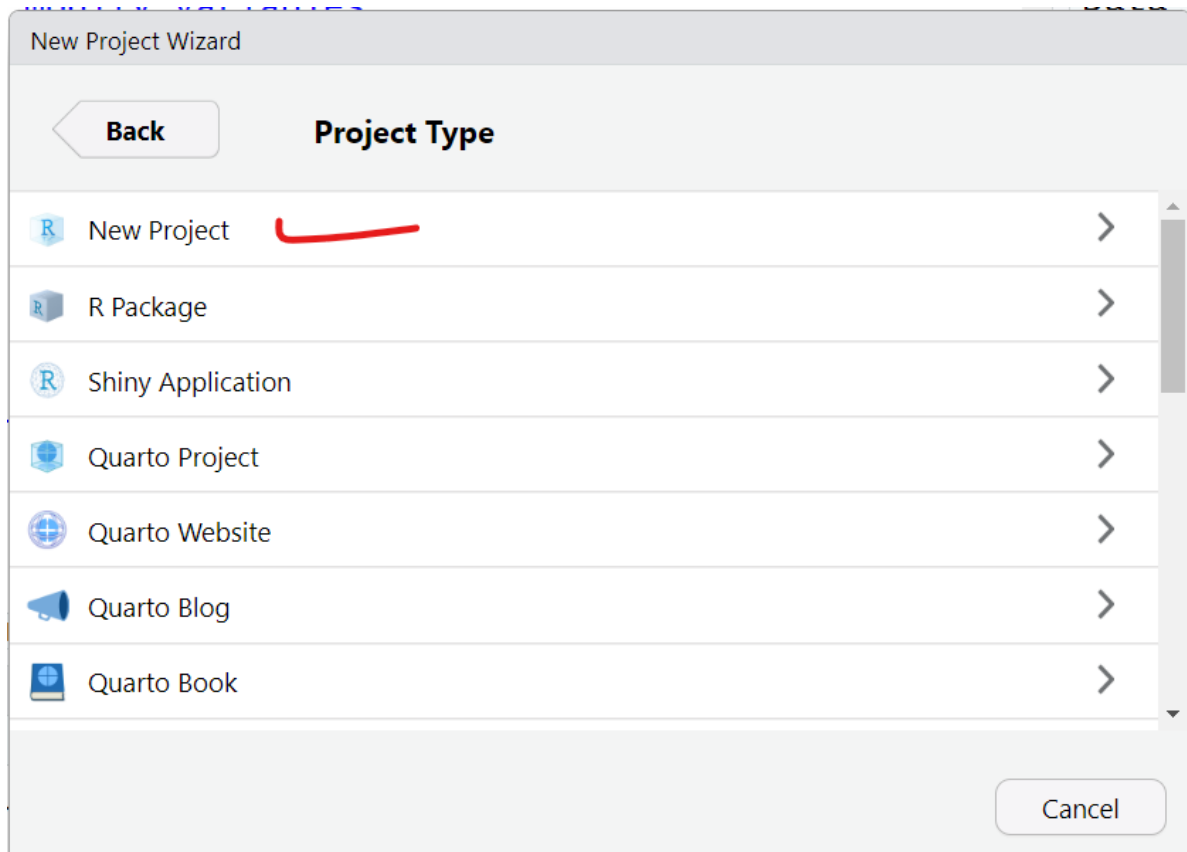1. First click on the menu at top for File/New Project:

2. Next choose either an "Existing Directory" or "New Directory" depending on whether you want to use a folder that already exists on your computer or you want to create a new folder.

**New Project Wizard**

**Create Project**

R | **New Directory** | >
Start a project in a brand new working directory

R | **Existing Directory** | >
Associate a project with an existing working directory

R | **Version Control** | >
Checkout a project from a version control repository

Cancel

3. For now, let's choose a "New Directory" and then select "New Project"

New Project Wizard

Back      **Project Type**

R  New Project  ⌄  >

R  R Package  >

R  Shiny Application  >

Quarto Project  >

Quarto Website  >

Quarto Blog  >

Quarto Book  >

Cancel

4. When the next window opens, as an example, I'm creating a new folder called `myfirstRproject` for my RStudio project under the parent directory, `C:\MyGithub`.

5. So, if I look back on my computer in my file manager (I'm on a computer with Windows 11 operating system) - I can now see this new folder on my computer for `C:\MyGithub\myfirstRproject`.
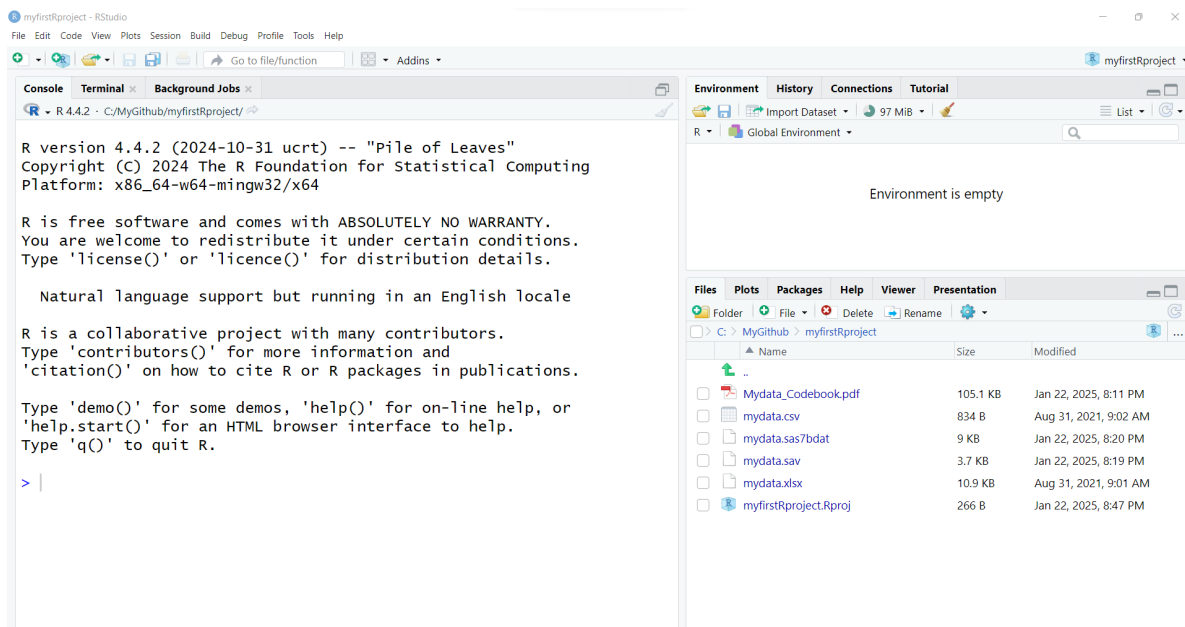


6. Now let's put some data into this folder. Feel free to move datasets of your own into this new RStudio project directory. But here are some test datasets you can download and place into this new directory on your computer - choose at least one to try out - right click on each link and use "Save As" to save the file on your computer.

- `mydata.csv` - CSV (comma separated value) formatted data
- `mydata.xlsx` - EXCEL file
- `mydata.sav` - SPSS Dataset
- `mydata.sas7bdat` - SAS Dataset
- `Mydata_Codebook.pdf` - Codebook Details on "mydata" dataset

7. After putting these files into your new RStudio project folder, you should see something like this now in your RStudio Files Listing (bottom right window pane):
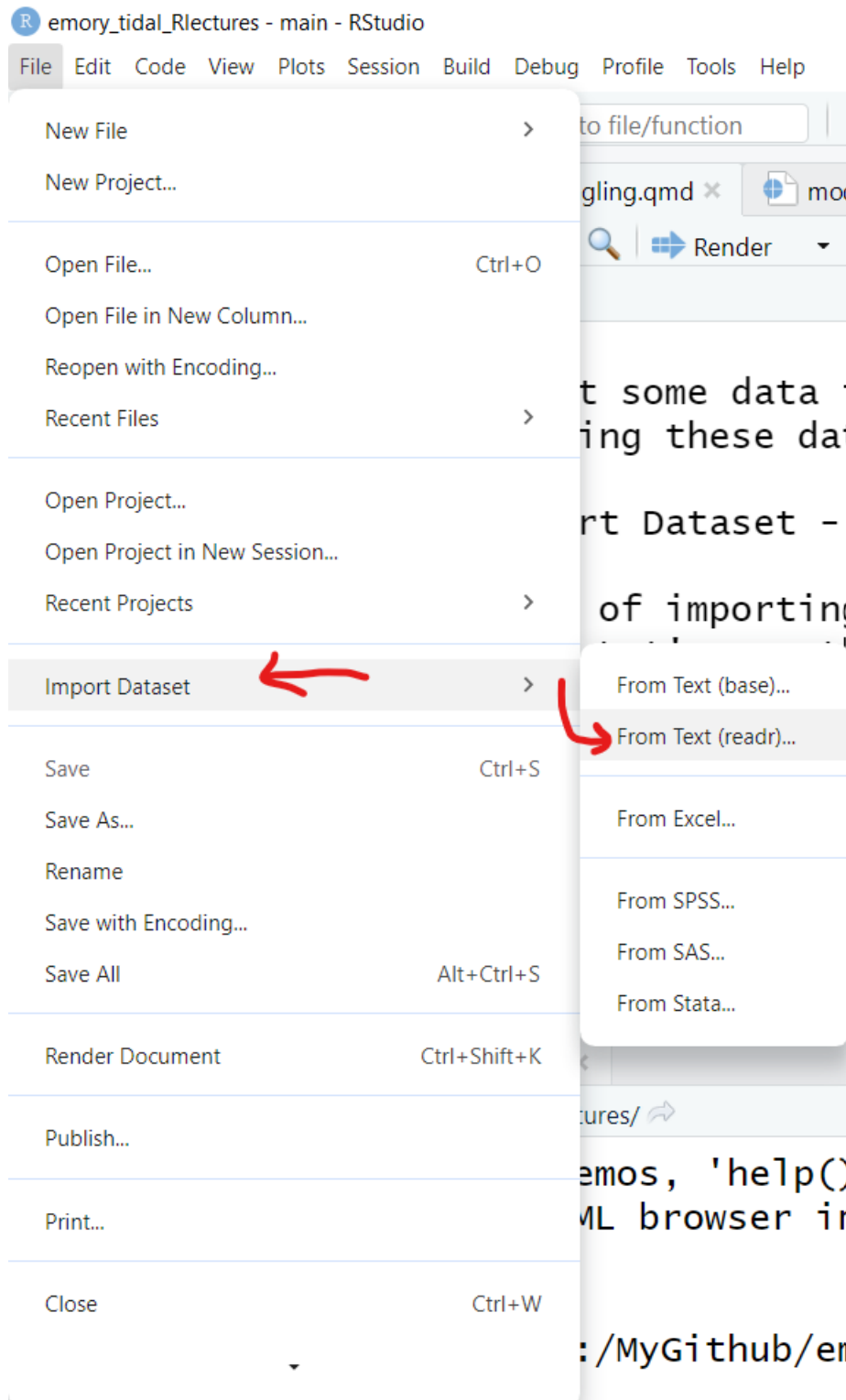


## Importing Data

Now that you've got some data in your RStudio project folder, let's look at options for importing these datasets into your RStudio computing session.

Click on File/Import Dataset - and then choose the file format you want.

## Import a CSV file

Here is an example of importing the `mydata.csv` - CSV formatted data. Let's use the `From Text (readr)` option.

Once this opens, click on "Browse" and choose the `mydata.csv` file. Assuming all goes well, this window will read the top of the datafile and show you a quick preview to check that the import will work. And on the bottom right, the R code commands needed to import this dataset will be shown to you. You can then click on the little "clipboard" on the bottom right to copy this R code to your "clipboard", *(the R code option will be explained below)*.

OR You can also just click "Import" and the R code will be executed for you and the dataset brought into your R computing session.



But the better way is to save the R code commands to import the data so you will be able to reproduce all steps in your data analysis workflow using code as opposed to non-reproducible point-and-click steps.

Once you copied the R code above to your clipboard, go to "File/New File/R Script" to open a script programming window:
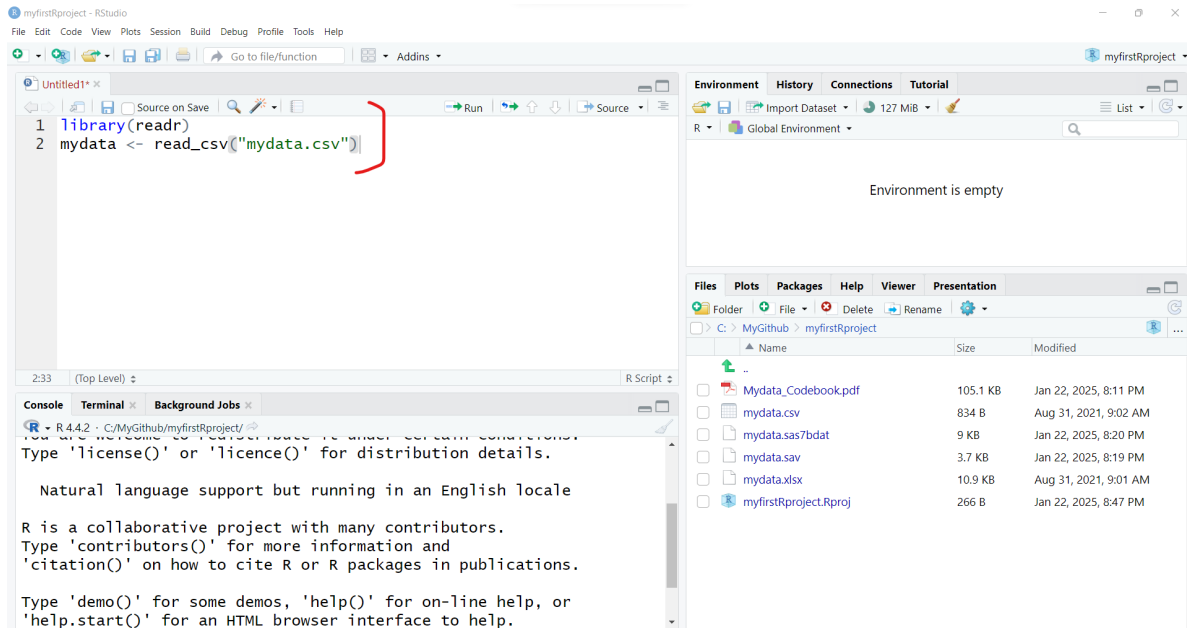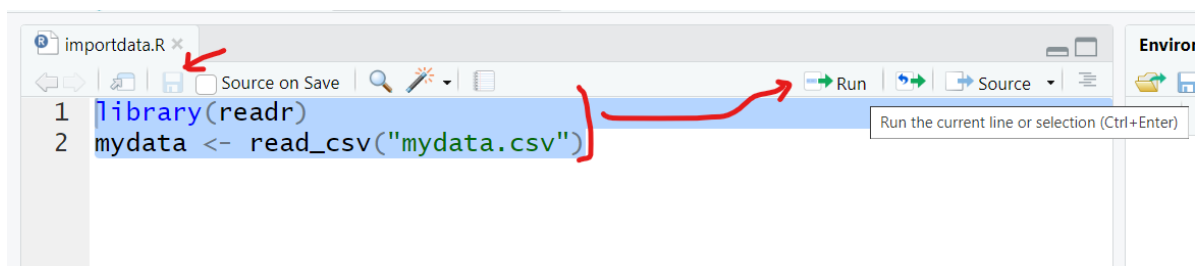
And then "paste" your R code into this window.

So as you can see importing the `mydata.csv` dataset, involves 2 steps:

1. Loading the `readr` package into your RStudio computing session, by running `library(readr)`
2. Running the `read_csv()` function from the `readr` package and then assigning `<-` this output into a new R data object called `mydata`.



To import the dataset, select these 2 lines of code and then click "Run" to run the R code. And be sure to click "Save" to save your first R program - for example "importdata.R".



After running these 2 lines of code, you should see something like this - the code messages in the bottom left "Console" window pane and a new R data object "mydata" in the top right "Global Environment" window pane.

## Import an EXCEL file

aaaaaaaaaaa

## exploring builtin datasets - see environment also - including data with packages

**2. To subset data.**

---

**3. To select, create, and modify variables.**

**4. To filter the observations according to a certain condition.**

---

## 5. To explore and summarize data.

**base r - summary stats**

**making tables**

**gtsummary package**

**other packages - arsenal, tableone and Rmarkdown formatting of tables... get inspired**

table competitions

---

**6. To run descriptive statistics.**

Module 1.3: Data Analytics Using R

## References

R Core Team. 2024. *R: A Language and Environment for Statistical Computing.* Vienna, Austria: R Foundation for Statistical Computing. https://www.R-project.org/.

## Other Helpful Resources

**Other Helpful Resources**