

# Day 1 Worksheet: Installation, R Basics, and R Markdown

Reference for Stat 415/615, Regression

Originally Prepared for Stat 412/612, Statistical Programming in R, 2022

## Contents

<b>1</b>	<b>Learning Objectives</b>	<b>1</b>
<b>2</b>	<b>References for RStudio and R Markdown</b>	<b>1</b>
<b>3</b>	<b>Installations</b>	<b>1</b>
<b>4</b>	<b>Getting Started</b>	<b>2</b>
<b>5</b>	<b>RStudio Options</b>	<b>3</b>
<b>6</b>	<b>R Basics</b>	<b>3</b>
<b>7</b>	<b>R Markdown</b>	<b>7</b>

## 1 Learning Objectives

- Install R and RStudio
- Understand the basics of R.
- R Markdown.
- Chapters 1, 2, 4, and 27 of [RDS](#).

## 2 References for RStudio and R Markdown

- [RStudio IDE Cheat Sheet](#).
- [R Markdown Cheat Sheet](#).
- [Introduction to R Markdown](#).
- [Markdown Basics](#).
- [RMarkdown for Scientists](#).

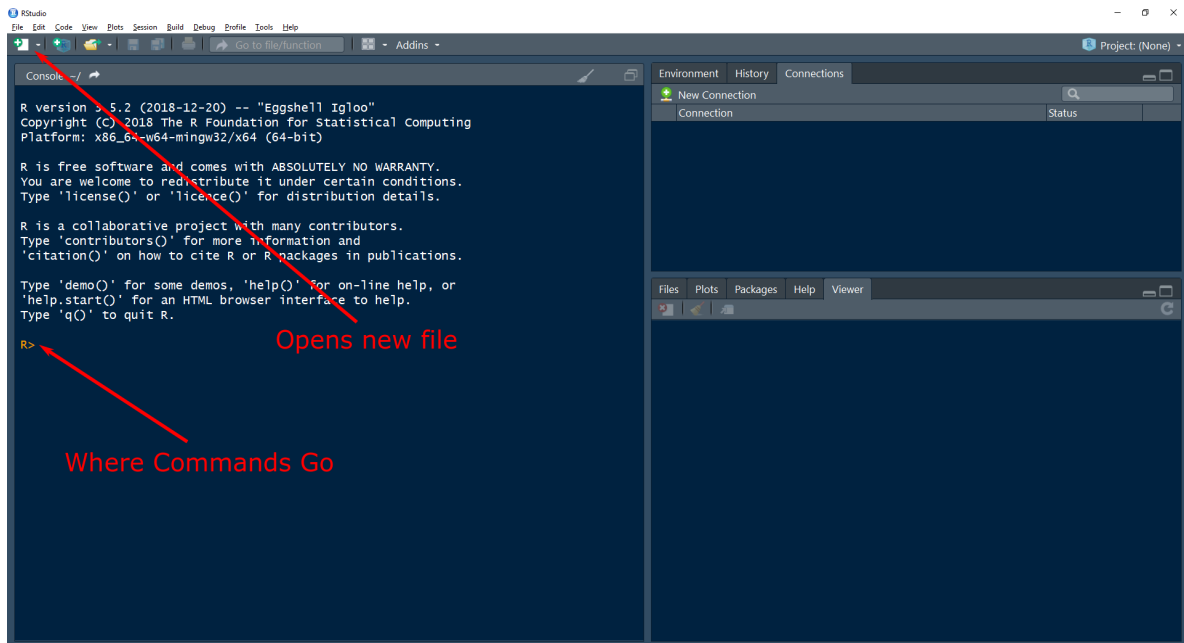
## 3 Installations

- Install R from the Comprehensive R Network (CRAN): <https://cran.r-project.org/>
  - You will have to choose the appropriate operating system. If you are using Ubuntu, there are special instructions if you want the newest version of R (v3.5.2 as of writing this document) as the default is to install an older version. This is not too important for this class, but some packages (mainly the ones in [Bioconductor](#)) require the newest version of R. If you want to use these packages, you should follow these instructions: <https://cran.r-project.org/bin/linux/ubuntu/README.html>
- Install RStudio: <https://www.rstudio.com/>

- RStudio is **different** from R. R is a statistical programming language. RStudio is an Integrated Development Environment (IDE), a software application used to write R code. It contains tools (such as shortcut buttons and syntax autocmpletion and highlighting) that make coding in R much easier. Most data scientists either use [RStudio](#), [Jupyter Notebook](#), or [Emacs Speaks Statistics](#) to interact with R.

## 4 Getting Started

- When you first open up RStudio, it should look something like this:



- The area right next to the carrot “>” is called a **prompt**. You insert commands into the prompt.
- You can use R as a powerful calculator. Try typing some of the following commands into the command prompt:

```
3 * 7
9 / 3
4 + 6
3 - 9
(3 + 5) * 6
3 ^ 2
4 ^ 2
```

- What do you think %% does? Try it out on a few examples:

```
1 %% 3
2 %% 3
3 %% 3
4 %% 3
9 %% 3
11 %% 3
```

- In general if you don’t know what something does, there are three things you can do:
  1. Google it. Always add “in R” at the end of your search term. E.g. “%% in R”.

2. Use the `help()` function to see if there is a help page on it. E.g., type `help("%")` in the command prompt. This is the most important function in R.
3. Play around with it, changing the inputs until you figure out what it does.

## 5 RStudio Options

- RStudio will automatically save a bunch of leftover stuff when you exit it. This is bad for reproducibility because you might be depending on code that you ran in a previous session (which you no longer can view). We'll get rid of these silly options.
- Open up “Tools > Global Options... > General” and
  - **Uncheck** “Restore .RData into workspace at startup”.
  - Change “Save workspace to .RData on exit” to “Never”.
- Other options in “Tools > Global Options... >”
  - “Code” has several options related to coding convenience and diagnostics.
  - “Pane layout” lets you change the layout.

## 6 R Basics

### 6.1 Variables

- R consists of two things: **variables** and **functions** (computer scientists would probably disagree with this categorization).
- A **variable** stores a value. This value can be:

1. A **numeric** (a number): 1, 2, 3, 4
2. A **string** (a word/character). These are always surrounded by quotation marks: "Hello", "world", "32".
3. A **logical** (either TRUE or FALSE).
4. A **factor** (categorical variable), more on this later.
5. A **vector** of numerics/strings/logicals/factors. A vector is a sequence of something of the same type. E.g.

```
## [1] 1 2 3 4
```

```
## [1] "Hello" "world" "32"
```

6. A **matrix** of numerics/strings/logicals/factors. A matrix is a rectangle of elements of the same type. E.g.

```
##      [,1] [,2] [,3]
```

```
## [1,]    1    3    5
```

```
## [2,]    2    4    6
```

```
##      [,1] [,2] [,3]
```

```
## [1,] "this" "a"   "of"
```

```
## [2,] "is"   "matrix" "strings"
```

7. A **list** of any sequence of objects.

```
## [[1]]
```

```
## [1] "hello"
```

```
##
```

```
## [[2]]
```

```
## [1] 3
##
## [[3]]
## [1] TRUE
```

8. A **data frame** where each column is a vector of a certain type (a vector of numerics or strings or logicals or factors).

- You **assign** variables values with `<-`. You can also use `=`, but most experienced users `<-` (or `->` in piping). For example, assign `x` the value of 7.

```
x <- 7
```

- If you type `x` now into the command prompt, it will return what is stored in `x`:

```
x
```

- Creating variables is useful because it allows you to reuse those values over and over again.
- You can apply arithmetic operations to variables as you would to numbers.

```
y <- 9
x + y
x * y
x / y
x - y
```

- You assign strings and logicals via:

```
z <- "I AM A STRING"
x <- TRUE
```

- Notice that we overwrote `x` to now be `TRUE` instead of 7. This is possible in R (in other languages the type is pre-defined and not changed easily).
- I've chosen horrible variable names so far. Generally variable names should be:
  - Nouns
  - Informative about what it is storing, yet concise.
  - Use either the ``camelCaseSystem`` or the ``snake_case_system``.
  - You are also able to use periods in variable names, but only noobs do that (it makes things more confusing in more advanced R).

- Variables are **case sensitive**, so `Variable` is different from `variable`. Try this out:

```
variable <- 1
Variable <- 2
variable
Variable
```

- When I write a variable inline, it will look like `this`.
- Exercise: spot the difference?

```
my_variable <- 10
my_variable
```

```
Error: object 'my__variable' not found
```

```
data0 <- 1
data0
```

```
Error: object 'dataO' not found
```

## 6.2 Functions

- Functions take objects (such as numbers or variables) as input and output new objects. Let's look at a simple function that takes the log of a number:

```
log(x = 4, base = 2)
```

- The inputs are called “arguments”. Generally, every function will be for the form:

```
function_name(arg1 = val1, arg2 = val2, ...)
```

- If you do not specify the name of the argument, R will assume you are assigning in their order.

```
log(4, 2)
```

- You can change the order of the arguments if you specify them.

```
log(base = 2, x = 4)
```

- To see the list of all possible arguments of a function, use the `help()` function:

```
help(log)
```

- In the help file, there are often **default** values for an argument. For example, the following indicates the the default value of `base` is `exp(1)`.

```
log(x, base = exp(1))
```

- This indicates that you can omit the `base` argument and R will assume that it should be `exp(1)`.

```
log(x = 4, base = exp(1))
```

```
## [1] 1.386294
```

```
log(x = 4)
```

```
## [1] 1.386294
```

- If an argument does not have a default, then it must be specified when calling a function.

- Type this:

```
log(x = 4,
```

- The “+” indicates that R is expecting more input (you forgot either a parentheses or a quotation mark). You can get back to the prompt by hitting the ESCAPE key.

- There are a *lot* of functions that come pre-installed with R. Don't try to memorize the list of all of them functions. The best course of action is to Google the functionality you need and you'll automatically pick up the names as you go.

- But here is a reference card with some basic functions: <https://cran.r-project.org/doc/contrib/Short-refcard.pdf>

- When I want to write the name of a function, I will write it like `this()`.

- Generally function names should be:

- Verbs
- Informative about what it is does, yet concise.
- ``useTheCamelCaseSystem()`` or ``use_the_snake_case_system()``.
- You are also able to use periods in function names, but only noobs do that (it makes things more confusing in more advanced R).
- But lots of base R functions have periods in them: ``is.na()``,

```
`which.max()`, `data.frame()`, etc.
```

## 6.3 Package

- A **package** is a collection of functions that don't come with R by default.
- There are **many many** packages available. If you need to do any data analysis, there is probably an R package for it.
- Using `install.packages()`, you can install packages that contain functions and datasets that are not available by default. Do this now with the tidyverse package:

```
install.packages("tidyverse")
```

- You will only need to install a package once per computer. Once it is installed you can gain access to all of the functions and datasets in a package by using the `library()` function.

```
library(tidyverse)
```

- You will need to run `library()` at the start of every R session if you want to use the functions in a package.

## 6.4 Vectors and Indexing

- In statistics we deal with data. To manage and analyze data in R, we need a way to store those data.
- The most basic and useful way to store data is in a vector.
- A vector is a sequence of data elements.
- It's tempting to use the word "list" to describe it, however a List is actually a very specific thing in R.
- Vectors are great because we can perform operations on every single value in a vector without iterating through it.
- For those of you who have used Python before, you might be used to writing for-loops to work on values in a list. That's not necessary here.
- Vectors are created in R using the `c()` function.
- "c" stands for concatenate, which means stick together.
- Create a of numbers now:

```
my_first_vector <- c(3, 6, 2, 5)
```

- **Exercise:** What does `my_first_vector + 4` do?
- **Exercise:** What does `my_first_vector * 4` do?
- You can refer to individual elements of a vector using brackets `[]`

```
my_first_vector[1]  
my_first_vector[2]  
my_first_vector[3]  
my_first_vector[4]
```

- You can refer to multiple elements in a vector by creating an indexing vector first.

```
indexing_vector <- c(2, 4)  
my_first_vector[indexing_vector]
```

- You can skip the intermediate step by placing `c(2, 4)` directly in the brackets.

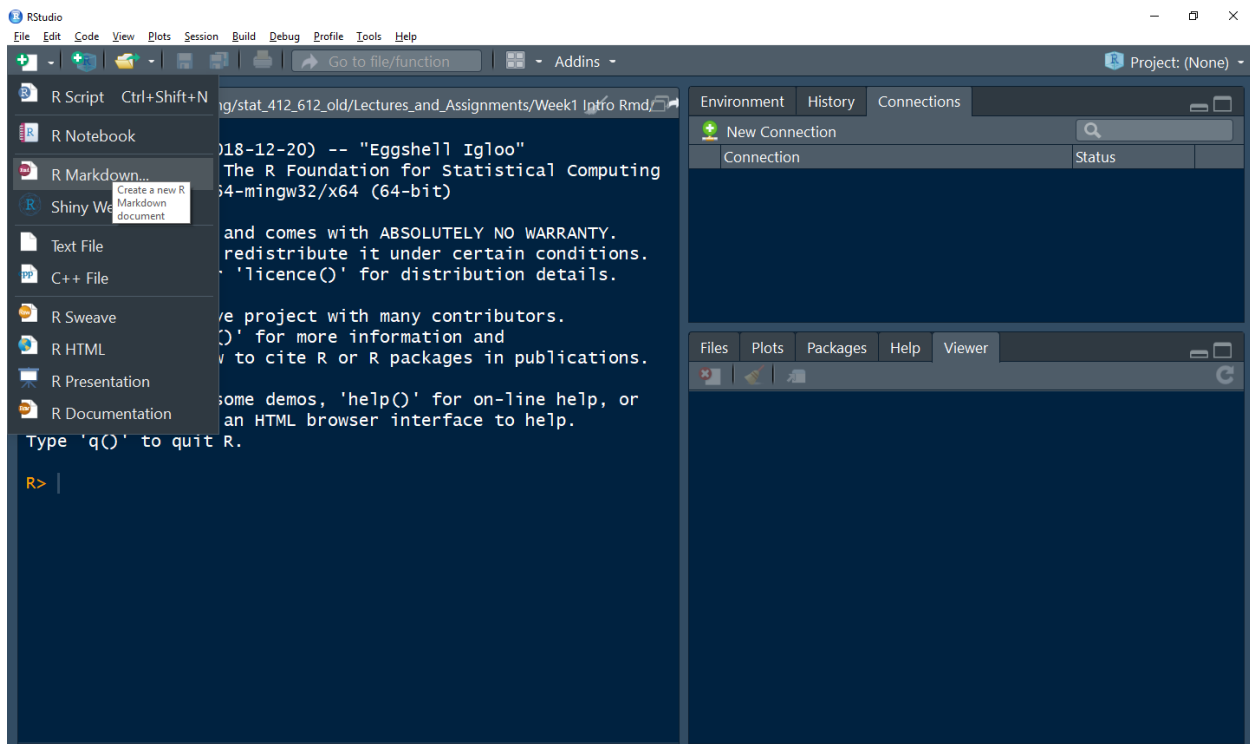
```
my_first_vector[c(2, 4)]
```

- You can create a vector in a sequence using the `seq()` function.
- **Exercise:** Read up on the `seq()` function to create a vector of even numbers from 52 to 1000. Then extract the 3rd and 29th elements from this vector.
- You can extract elements of a vector by using logicals. What does the following code do?

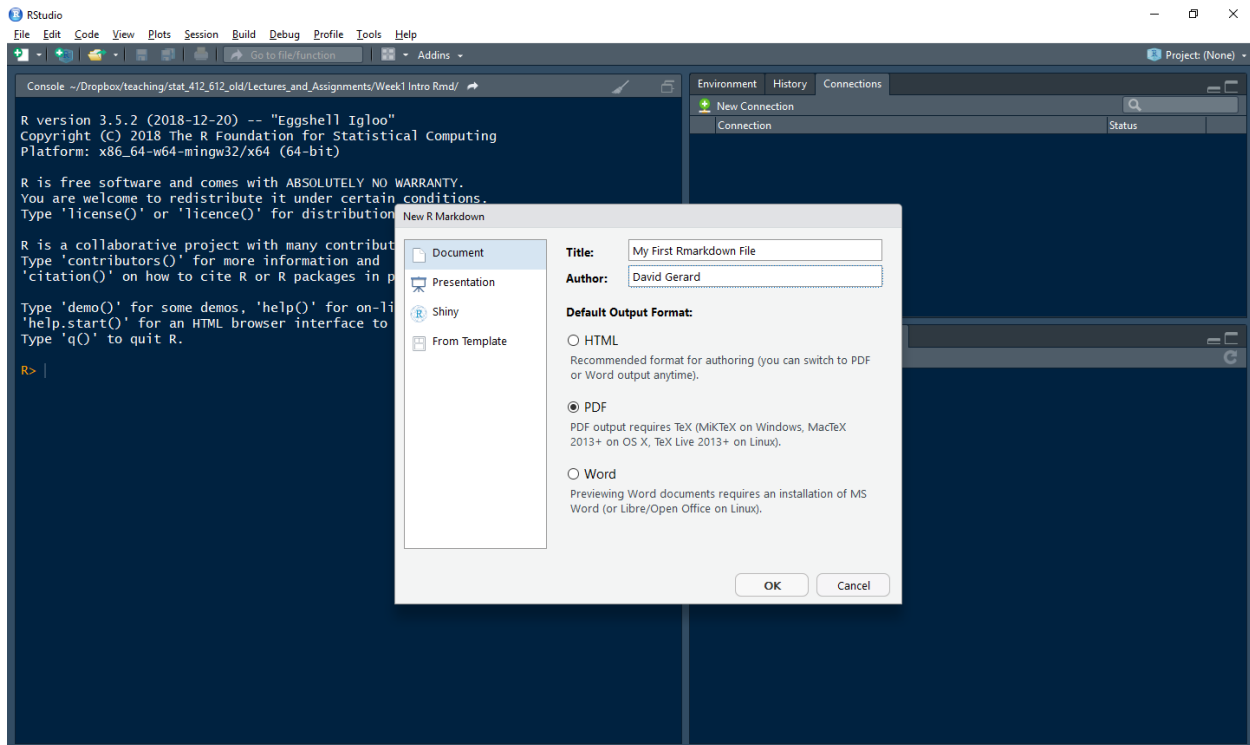
```
logvec <- c(TRUE, FALSE, FALSE, TRUE)
my_first_vector[logvec]
```

## 7 R Markdown

- R Markdown is a file format that is a combination of plain text and R code.
- A free book on R Markdown is available here: <https://bookdown.org/yihui/rmarkdown/>
- You write code and commentary of code in one file. You may then compile (RStudio calls this “Knitting”) the R Markdown file to many different kinds of output: pdf (including beamer presentations), html (including various presentation formats), Word, PowerPoint, etc.
- R Markdown is useful for:
  1. Communication of statistical results.
  2. Collaborating with other data scientists.
  3. Using it as a modern lab notebook to *do* data science.
- *You will use R Markdown for all assignments in this class.*
- Open up a new R Markdown file:



- Choose the options for the type of output you want



- You should now have a rudimentary R Markdown file.
- Save a copy of this file in your “analysis” folder in the “week1” project.
- R Markdown contains three things
  1. A YAML (Yet Another Markup Language) header that controls options for the R Markdown document. These are surrounded by `---`.
  2. Code **chunks** — bits of R code that that are surrounded by ````${r}````. Only valid R code should go in here.
  3. Plain text that contains simple formatting options.
- All of these are displayed in the default R Markdown file. You can compile this file by clicking the “Knit” button at the top of the screen or by typing `CONTROL + SHIFT + K`. Do this now.

## 7.1 YAML Header

- When open a new R Markdown file, you can choose the output file format. The YAML header will be generated automatically.
- If you are new to LaTeX (a document preparation system), set the output as Word document. After getting the Word document knitted by R Markdown, You can convert the Word file into PDF using Microsoft Office or similar software.

```
---
title: "Document Title"
author: "Author 1, Author 2"
date: "`r Sys.Date()` or other dates"
output: word_document
---
```

- If you have are familiar with Latex, you can set the output file as pdf document. You may need to install additional style files for your existing LaTeX (MiKTeX, etc.), or try TinyTex (see <https://yihui.org/tinytex/>).



```

---
title: "Document Title"
author: "Author 1, Author 2"
date: "`r Sys.Date()`"
output: pdf_document
urlcolor: "blue"
---

```

- All of those settings are fairly self-explanatory.
- The full list of outputs supported is here: <https://rmarkdown.rstudio.com/formats.html>
- The ``r Sys.Date()`` will insert the current date.

## 7.2 Code Chunks

- You can insert new code-chunks using CONTROL + ALT + I (or using the “Insert” button at the top of RStudio).
- You write all R code in chunks. You can send the current line of R code (the line where the cursor is) using CONTROL + ENTER (or the “Run” button at the top of RStudio).
- You can run all of the code in a chunk using CONTROL + ALT + C (or using the “Run” button at the top of RStudio).
- You can run all of the code in the next chunk using CONTROL + ALT + N (or using the “Run” button at the top of RStudio).

## 7.3 Inline R Code

- Sometimes, you want to write the output of some R code inline (rather than as the output of some chunk). You can do this by placing code within ``r R_code``. The YAML header uses such notation (``r Sys.Date()``) for automatically writing the date.
- When preparing a report, you might calculate the  $p$ -value from a linear regression, then write this  $p$ -value in the paragraph of a report. The code “The resulting  $p$ -value is ``r 1-pt(2.4, df=35)`` will result in”The resulting  $p$ -value is 0.0109249)

## 7.4 Formatting markdown

- Here is Hadley’s brief intro to formatting text in R Markdown:

Text formatting

-----

```

*italic*   or _italic_
**bold**   __bold__
`code`
superscript^2^ and subscript~2~

```

Headings

-----

```
# 1st Level Header
```

```
## 2nd Level Header
```

```
### 3rd Level Header
```

## Lists

---

- \* Bulleted list item 1
  - \* Item 2
    - \* Item 2a
    - \* Item 2b
1. Numbered list item 1
  1. Item 2. The numbers are incremented automatically in the output.

## Links and images

---

<<http://example.com>>

[linked phrase](<http://example.com>)

![optional caption text]([path/to/img.png](#))

## Tables

---

First Header	Second Header
Content Cell	Content Cell
Content Cell	Content Cell