

# My first R Markdown

Rudy Rupp

21/08/2020

## Data Scientist's Toolbox Week 1

**What is Data Science?** 0:43 - 0:56 They state that a data scientist is broadly defined as someone who combines the skills of software programmer, statistician, and storyteller/artists to extract the nuggets of gold hidden under mountains of data.

## Data Scientist's Toolbox Week 2

### R Studio tour

```
matrix(c(1, 2, 3, 4, 5, 6, 7, 8), nrow = 4, ncol = 2)
```

```
##      [,1] [,2]
## [1,]    1    5
## [2,]    2    6
## [3,]    3    7
## [4,]    4    8
```

### R Packages

#### Installing from CRAN

If you are installing from the CRAN repository, use the `install.packages()` function, with the name of the package you want to install in quotes between the parentheses (note: you can use either single or double quotes). For example:

```
install.packages("ggplot2")
```

To install multiple packages at once, you can do so by using a character vector:

```
install.packages(c("ggplot2", "devtools", "lme4"))
```

#### Installing from Bioconductor

The BioConductor repository uses their own method to install packages. First, to get the basic functions required to install through BioConductor, use:

```
source("https://bioconductor.org/biocLite.R")
```

This makes the main install function of BioConductor, `biocLite()`, available to you. Following this, you call the package you want to install in quotes, between the parentheses of the `biocLite` command, like so:

```
biocLite("GenomicFeatures")
```

## Installing from GitHub

You first must find the package you want on GitHub and take note of both the package name AND the author of the package. Check out this guide for installing from GitHub:

- `install.packages("devtools")` - only run this if you don't already have devtools installed.
- `library(devtools)` - loads the devtools package
- `install_github("author/package")` replacing "author" and "package" with their GitHub username and the name of the package.

## Loading packages

Installing a package does not make its functions immediately available to you. First you must load the package into R; to do so, use the `library()` function. For example:

```
library(ggplot2)
```

NOTE: Do not put the package name in quotes! Unlike when you are installing the packages, the `library()` command does not accept package names in quotes!

## Checking what packages you have installed

`installed.packages()` or `library()` with nothing between the parentheses to check.

## Updating packages

`old.packages()` identify packages that have been updated since last updated.

To update all packages, use `update.packages()`. If you only want to update a specific package, use `install.packages("packagename")`

## Unloading packages

To unload a given package you can use the `detach()` function. For example, `detach("package:ggplot2", unload=TRUE)`

## Uninstalling packages

If you no longer want to have a package installed, you can simply uninstall it using the function `remove.packages()`. For example, `remove.packages("ggplot2")`

## How do you know what version of R you have?

Type `version` into the console and it will output information on the R version you are running.

`sessionInfo()` - will tell you what version of R you are running along with a listing of all of the packages you have loaded. The output of this command is a great detail to include when posting a question to forums - it tells potential helpers a lot of information about your OS, R, and the packages that you are using.

## What commands are included in a package?

You can use the `help()` function to access a package's help files.

```
help(package = "ggplot2")
```

If you still have questions about what functions within a package are right for you or how to use them, many packages include "vignettes." These are extended help files, that include an overview of the package and its functions, but often they go the extra mile and include detailed examples of how to use the functions in plain words that you can follow along with to see how to use the package. To see the vignettes included in a package, you can use the `browseVignettes()` function.

Eg: `browseVignettes("ggplot2")`

You should see that there are two included vignettes: "Extending ggplot2" and "Aesthetic specifications." Exploring the Aesthetic specifications vignette is a great example of how vignettes can be helpful, clear instructions on how to use the included functions.

# Data Scientist's Toolbox Week 3

## Version control (Git & GitHub) Vocabulary

- **Repository** online folder that contains the whole, shared project Commit to save your changes locally, and add explanations of what has been changed.
- **Push** upload changes made locally to the repository, and merge them.
- **Pull** download the most recent changes from the repository and merge with your local folder.
- **Staging** to commit files separately and add separate explanations to each section edited, in order to simplify looking at the change log.
- **Branch** your local files and edits, before you push them.
- **Merge** independent branches are joined together into a single file. May need manual control if the same section of the same file was edited differently in both branches.
- **Conflict** when branches have been edited in the same section and file, making automatic merging hard.
- **Clone** make a local copy of a whole Git repository's files and all the change history.
- **Fork** make a personal copy of a repository that doesn't affect the original.

## Best Practices:

- Purposeful, single issue commits
- Informative commit messages
- Pull and Push often

Git account:

- \* Git username Rodrupp
- \* Git email rodrupp@gmail.com
- \* Git password: automatically generated on firefox.

**Linking an Existing Project with GitHub** Open Git Bash or Terminal and navigate to the directory containing your project files. Move around directories by typing `cd ~/dir/name/of/path/to/file` Once here, type `git init` followed by `git add .` -> this initializes (init) this directory as a git repository and adds all of the files in the directory (.) to your local repository. Commit these changes to the git repository using `git commit -m "Initial commit"`

Go to GitHub.com, and again, create a new repository:

1. Make sure the name is the exact same as your R project;
2. Do NOT initialize a README file, .gitignore, or license.

Find “...push an existing repository from the command line” and copy to clipboard. Paste these into the Terminal, then reload the repository page.

# Data Scientist's Toolbox Week 4

## R Markdown Example

This is an R Markdown document. Markdown is a simple formatting syntax for authoring HTML, PDF, and MS Word documents. For more details on using R Markdown see <http://rmarkdown.rstudio.com>.

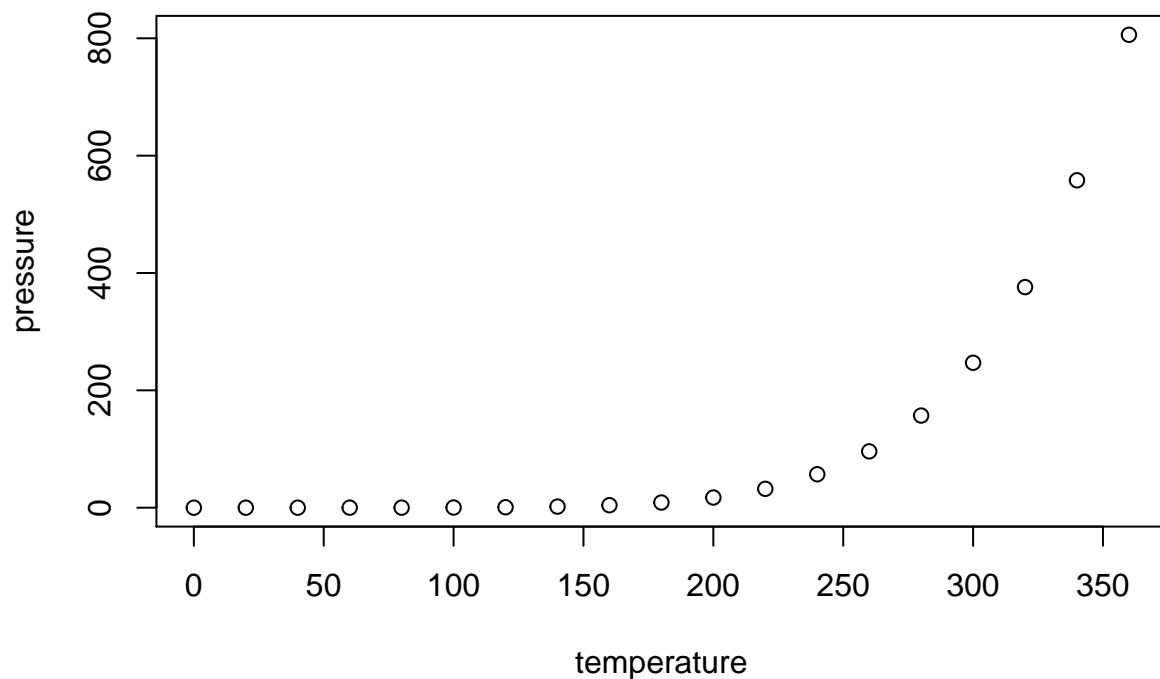
When you click the **Knit** button a document will be generated that includes both content as well as the output of any embedded R code chunks within the document. You can embed an R code chunk like this:

```
summary(cars)
```

```
##      speed      dist
## Min.   : 4.0    Min.   :  2.00
## 1st Qu.:12.0    1st Qu.: 26.00
## Median :15.0    Median : 36.00
## Mean   :15.4    Mean    : 42.98
## 3rd Qu.:19.0    3rd Qu.: 56.00
## Max.   :25.0    Max.    :120.00
```

## Including Plots

You can also embed plots, for example:



Note that the `echo = FALSE` parameter was added to the code chunk to prevent printing of the R code that generated the plot.

## R Markdown formatting

To bold text use double asterisks **bold**.

to make it italics use single asterisks *italics*.

to strikethrough text use double tildes ~~strikethrough~~

Use # for headings: - # Biggest Header

- ## Second biggest

- ### Third heading

Code boxes created with triple backdash “`

```
print("Testing some code here.")
```

```
## [1] "Testing some code here."
```

```
print("testing the keyboard shortcut ctrl+alt+i")
```

```
## [1] "testing the keyboard shortcut ctrl+alt+i"
```

```
print("testing the insert button on RStudio")
```

```
## [1] "testing the insert button on RStudio"
```

```
print("testing multiple lines of code.")
```

```
## [1] "testing multiple lines of code."
```

```
print("and here is the third line")
```

```
## [1] "and here is the third line"
```

Use Ctrl+enter to run a line of code in RStudio without having to knit the document.

Use Ctrl+shift+enter to run all lines in a chunk.

Making bulleted lists:

- bullets must finish with two spaces
- this is an example
- i wonder what happens if i don't include two spaces?
- ok, here goes:
- did it work?
- indent bullets with a + sign

This is a link to google

## Types of Data Analysis

- 1) Descriptive
- 2) Exploratory
- 3) Inferential
- 4) Predictive
- 5) Causal
- 6) Mechanistic

**Descriptive** Describes or summarises a set of data. Often the first analysis to perform. eg mean, standard deviation, etc.

**Exploratory Analysis** Examine data and find new relationships. Help formulate hypotheses and drive future studies.

**Inferential Analysis** Use a small sample of data do infer something about a larger population provides estimate and uncertainty. Depends on your sampling.

**Predictive Analysis** Use current data and historical data to make predictions. Dependent on measuring the right variables. More data and simple model is usually better.

**Causal Analysis** See what happens to one variable when we manipulate another. Gold standard in data analysis. Often applied to randomized studies designed to identify causation. Often challenging to obtain data.

**Mechanistic Analysis** Understand the exact changes in variables that lead to exact changes in other variables. Less common and more complex. Applied to simple situation or those that are nicely modeled by deterministic equations (eg physics or engineering).

## Experimental Design

Experimental design is organizing an experiment, so that you have the correct data and enough of it to clearly and effectively answer your data science question.

In this order:

- 1) Formulate your question before collecting any data
- 2) Design the experiment.
- 3) Identify problems.
- 4) Collect data.

And then a bunch of stuff I already know...

## Big Data

Volume, Velocity, Variety, Unstructured

Examples: Text files and documents, Websites and applications, Sensor data, Image Files, Audio files, Video files, Email Data, Social Media data.

**Before you can start looking for answers, you need to turn your unstructured data into a format that you can analyze.**