# R and RStudio

UNIVERSITY OF
SAN FRANCISCO

Abbie M Popa

BSDS 100 - Intro to Data Science with R

# Class Organization

- In Class: Mix of Lecture and "Coding Challenges"
  - Not "graded" (though attendance is worth 20% of your grade), but you will do better on the homework if you attempt the activities!
- 8 Homework assignments (40% of your grade)
  - Based on in-class activities, completed **individually** at home
  - 30 points each, 10 points deducted automatically if turned in late
  - Attempt all questions! Even wrong answer may receive partial credit

# Class Organization

- 2 Case Studies (20% of your grade)
  - Like homework assignments, but longer
  - Completed in pairs
  - You will receive **some** in-class time to work on these to help with scheduling issues
- Final Project (20% of your grade)
  - Similar to case studies, but more open-ended (you will pick the data and goals)
  - Completed in groups of 3 or 4

- The R Programming Language and RStudio

  - Comparison with other programming languages

  - Installation

  - Handy Shortcuts

# What is R ?

[From the R website]: "R is a language and environment for statistical computing and graphics. It is a GNU project which is similar to the S language and environment which was developed at Bell Laboratories by John Chambers and colleagues. R can be considered as a different implementation of S. There are some important differences, but much code written for S runs unaltered under R.

R provides a wide variety of statistical (linear and nonlinear modelling, classical statistical tests, time-series analysis, classification, clustering, etc.) and graphical techniques, and is highly extensible. The S language is often the vehicle of choice for research in statistical methodology, and R provides an Open Source"

- R is a user-friendly integrated development environment (IDE) that is open source.

- Has many packages and functions for statistical analyses

- You can even run R from a terminal window if you wish (using BASH scripting)

# Why Use R?

- Open source (free)

- Runs on just about any platform

- Great visualization capabilities (`ggplot2`)

- Read/write from/to various data sources

- Scripting language (interpreted)

- Massive library of data manipulation and statistical packages

# Excel is Great for Certain Things...

# ...but Not Everything

## Sample Data

- Six columns of data with ~ 1.05 million rows

- Column 5: `startDate`

- Column 6: `endDate`

- **Objective**: test to see if `endDate` < `startDate`

## RESULTS

- **Excel**: good luck...

- R: 33 min (poor coding technique)

- R: 58.5 sec (improved coding technique)

# ...but Not Everything

## Sample Data

- Six columns of data with ~ 1.05 million rows

- Column 5: `startDate`

- Column 6: `endDate`

- **Objective**: test to see if `endDate` < `startDate`

## RESULTS

- **Excel**: good luck...

- R: 33 min (poor coding technique)

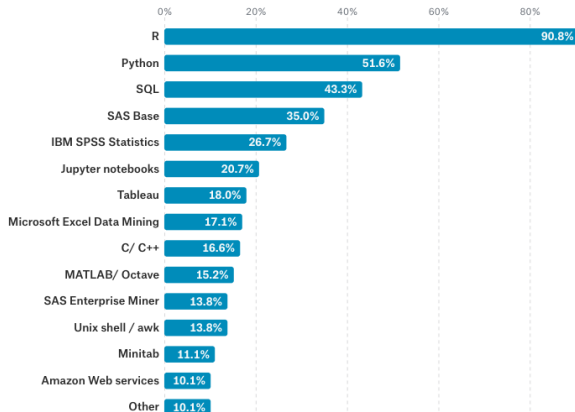- R: 58.5 sec (improved coding technique)

# Vectorization in R

- Vectorization: a style of coding where an operation is applied to all elements in an array, rather than looping

- Vectorized code saves time asking `type` questions

- There is an optimized engine—a basic linear algebra system (BLAS)—that is highly efficient at solving linear algebra problems

- A lot or R functions are written in C (or variants)

- MATLAB, Mathematica and the NumPy package for Python are also vectorized

    http://www.noamross.net/blog/2014/4/16/
        vectorization-in-r--why.html

R wins at statistics

- Download and install at this website:
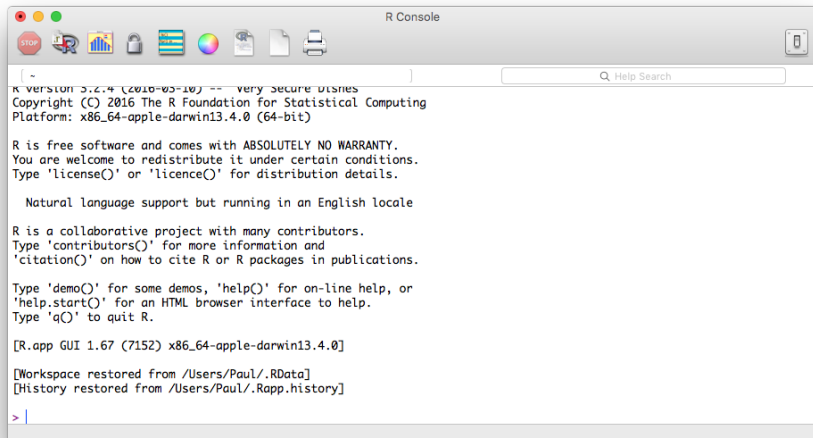
  https://www.r-project.org

- **Important**: You will have to re-install R from time-to-time to maintain the newest version so that code remains compatible! New versions generally come out every 4 - 6 months.

# Installing RStudio

- RStudio has a very nice graphical user interface (GUI) that is easier to use than base R

- We will be using this throughout the course

- Make sure that you have R installed first. Then, download and install at this website:
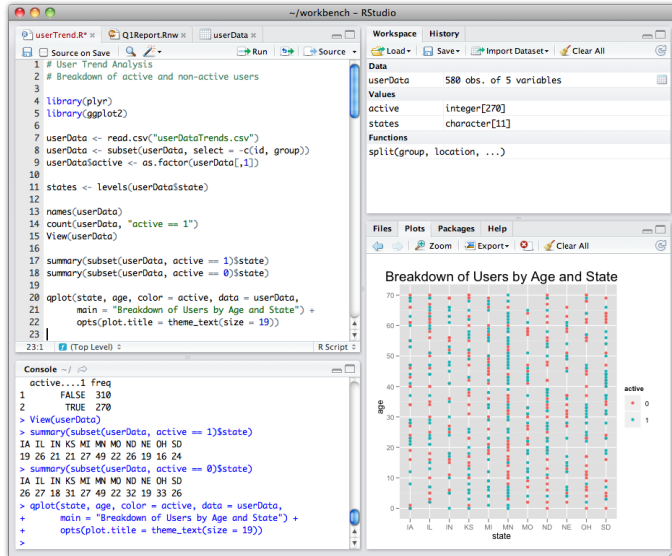
  https://www.rstudio.com/products/RStudio/

# The RStudio GUI

# RStudio

## RStudio has Four Panels

- Console (bottom left) - where all calculations are performed

- Scripting/Viewing (top left) - where writing of new functions / scripts should be done

- Files/Packages/Help/Plots (bottom right) - for easy analysis

- Variables/Data/Functions (top right) - what is stored in your current RSession

# Really Advanced Calculators

At their core, R and RStudio are just calculators! Try the following

$$3 + 2 = ? \qquad \log(10) = ? \qquad \sqrt{32} = ?$$

```
> 3 + 2
```

```
> log(10)
```

```
> sqrt(32)
```

# Some terms we'll be using

- command: tell R to do something (e.g., add, subtract, print)
- variable: assign a value to an identifier

  ```
  > my_cool_integer <- 5
  > my_cool_string <- "hedgehog"
  ```

- object: often used interchangeably with variable
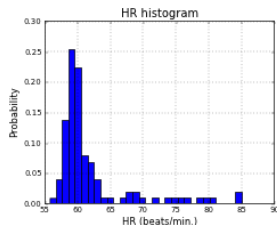- script: a file documenting many commands, can be rerun

- reproducibility

- share-ability

- automate the boring stuff

I compare 60 participants' heart rates before and after a stimulus

**Global analysis (time-domain parameters)**



HR histogram

**No. of beats:** *105.00*
**Mean HR:** *62.15 bps*
**STD HR:** *5.96 bps*
**Mean RR (AVNN):** *972.81 msec.*
**STD RR (SDNN):** *77.66 msec.*
**SDANN:** *--*
**SDNNIDX:** *--*

**pNN50:** *25.96%*
**rMSSD:** *74.99 msec.*
**IRRR:** *46.75 msec.*
**MADRR:** *24.00 msec.*
**TINN:** *48.25 msec.*
**HRV index:** *6.18*

| C | D | E | F | G | H |
|---|---|---|---|---|---|
| signal_length | lf_hf_ratio | percent_lf | percent_hf | stim | cond |
| 263.24 | 1.066 | 0.24372817 | 0.27757109 | none | rest |
| 290.18 | 1.133 | 0.33410858 | 0.29051973 | none | rest |
| 304.81 | 0.616 | 0.18962122 | 0.33726336 | none | rest |
| 302.69 | 0.752 | 0.2117973 | 0.28170506 | none | rest |

- function: often used interchangeably with "command"
  - though generally we would call "add(2, 3)" a function but wouldn't call "2 + 3" a function
  - you will eventually be able to write your own functions!
  - very flexible
- argument: a function takes arguments to perform it's task (e.g., in "add(2, 3)" the "2" and the "3" are arguments)
  - an argument can be a variable or an option (e.g., number_of_decimals = 2)
- working directory: where you are in the computers file structure

# Notes on R

- R is case-sensitive
- I require you to use the assignment operator `<-` instead of the equality operator `=` for all submitted code, even though both work, e.g.,

| Syntax | *Comments* |
|--------|-----------|
| `x <- 5` | standard syntax, required |
| `x = 5` | poor syntax, not permitted |
| `5 -> x` | awkward syntax, not permitted (but it works) |

# Basic R Help Functions

| Function | Action |
|---------:|--------|
| `?foo` | Help on the function `foo` |
| `??foo` | Search the help system for instances of the function `foo` |
| `data()` | List all available example datasets contained in currently loaded packages |
| `getwd()` | List the current working directory |
| `ls()` | List the objects in the current directory |

# Basic R Workspace Functions

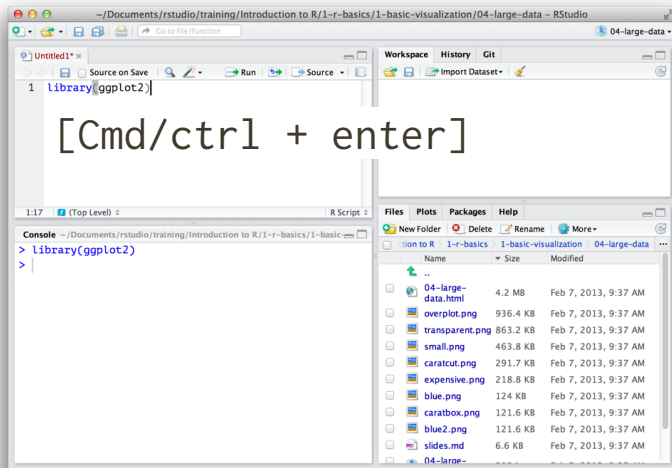| Function | Action |
|---|---|
| `getwd()` | List the current working directory. |
| `setwd("`*`mydirectory`*`")` | Change the current working directory to *mydirectory*. |
| `ls()` | List the objects in the current workspace. |
| `rm(`*`objectlist`*`)` | Remove (delete) one or more objects. |
| `help(options)` | Learn about available options. |
| `options()` | View or set current options. |
| `history(#)` | Display your last # commands (default = 25). |
| `savehistory("`*`myfile`*`")` | Save the commands history to *myfile* ( default = `.Rhistory`). |
| `loadhistory("`*`myfile`*`")` | Reload a command's history (default = `.Rhistory`). |
| `save.image("`*`myfile`*`")` | Save the workspace to myfile (default = `.RData`). |
| `save(`*`objectlist`*`, file="`*`myfile`*`")` | Save specific objects to a file. |
| `load("`*`myfile`*`")` | Load a workspace into the current session (default = `.RData`). |
| `q()` | Quit R. You'll be prompted to save the workspace. |

# Useful ʀ Keyboard Shortcuts: Execute Code
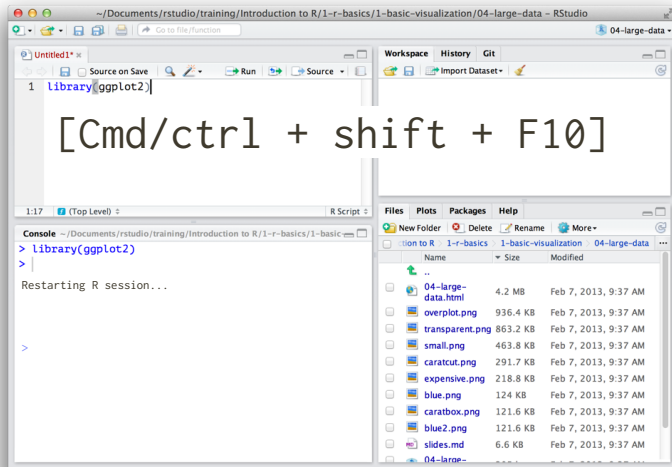
- When writing code, you want to clear your environment to ensure the fidelity of your results
- If your previous setting is correct, one way to do this is to close and re-open your `R` environment
- You can also run the following two lines of code

```
rm(list=ls())
cat("\014")
```

1. `rm(list=ls())` removes all objects in the current environment
2. On a Mac, `cat("\014")` clear the console windows (same as `ctrl + l`)

If you are sharing code you may not want to include this in your script

- R comes built in with multiple data sets you can play with

- Many (most?) packages also have data sets

- `data()` will bring up a list of all data sets available across all loaded packages

- `help(<nameOfDataSet>)` will provide you a detailed description of the data set in question

# How Big is *Big Data* in R?

- R holds data in memory, effectively limiting data to the amount of RAM a computer has access to

- It is not uncommon to work with a data set containing 100,000,000 elements (e.g., 100,000 observations of 1,000 variables or 1,000,000 observations of 100 variables) without difficulty

- Approximations depend on what type of data is contained in each variable, e.g., a data set with 2.2 million records and twenty variables, which takes approximately one minute to load into memory

# How Big is *Big Data* in R?

- Also depends on what techniques and/or functions will be applied to the data

- The more complex and memory intensive the task, the smaller the data will be required to be

- Basic plotting requires far less computational exertion than a complex statistical learning model

- **Common Definition**: *Big Data* refers to any data set that cannot be loaded into working memory on your personal computer