# An Introduction to R

**MacOdrum Library**

**February 23, 2021**

Carleton
UNIVERSITY

Canada's Capital University

**Go to:**

**http://madgic.library.carleton.ca/deposit/R Workshop/PhD%20Africa/**

**And download the four files**

**Save those files in
C:/Users/ *YOU*/Downloads**

- **About R**
- **What can R be used for**
- **The First Steps**
- **Activities**
- **Bonus Data (Time Permitting)**

- **Language and Environment originally designed for statistical computing and graphics**
- **Scripting Language**
- **Provides a lot of flexibility**
- **Can be used for multiple tasks beyond just statistical analysis**

- **Free**
- **Open Source**
- **Simple but powerful**
- **Lots of help online**
- **Can work with all different types of data/documents**
- **Lots of different "packages" that are used for specific analysis**
- **Replication is Easy**
- **So much that you can do with R**

- **Google**
- **Youtube**
- **Stack Overflow**
- **R Cheat sheets**
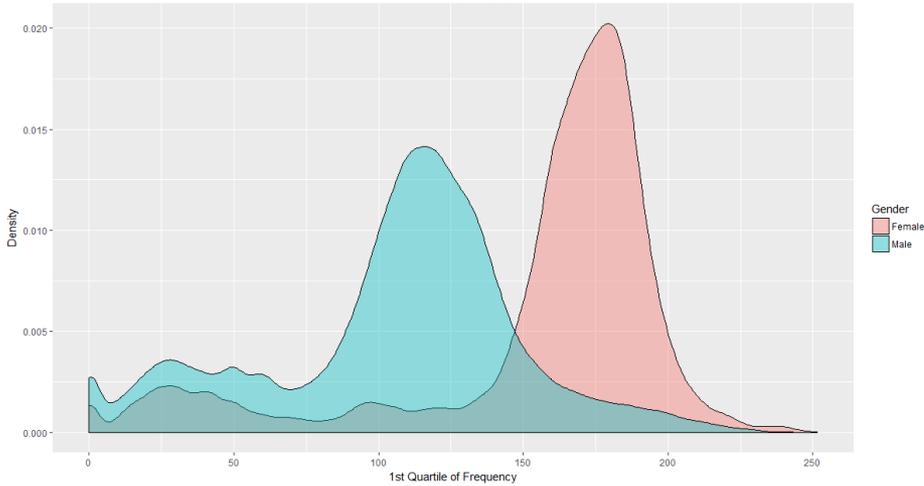- **Lots of eBooks found in your university library**

- **Statistical Analysis**
  - Basic Stats
  - Deep Learning, Machine Learning
  - Text Analysis
  - Geospatial Analysis
  - Audio Analysis

- **Cleaning Messy Data**
  - Cleaning up spreadsheets, Data files

- **Automation**
  - Web Scraping
  - Virtual Machine
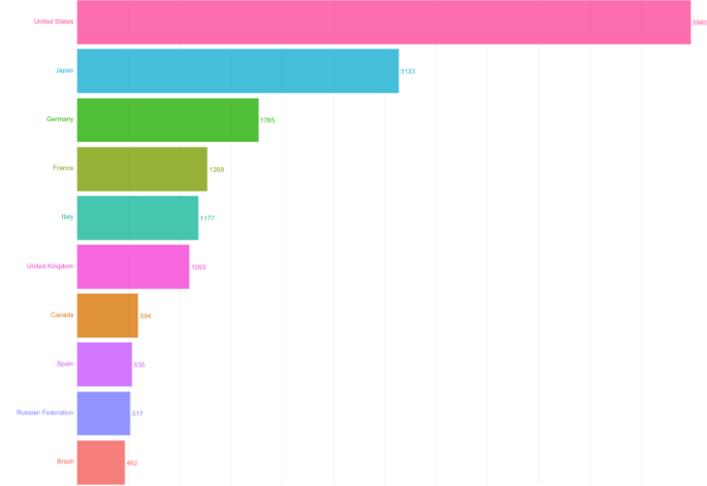  - Roll through multiple files in one go

- **Visualisations**

7

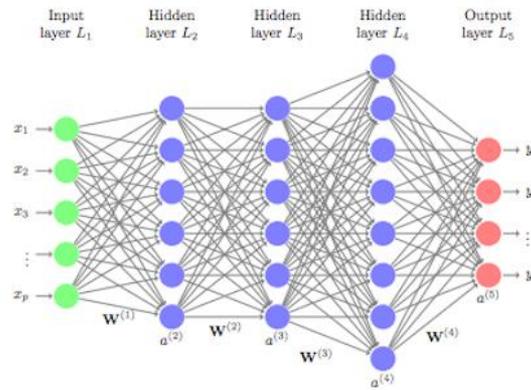1st Quartile of Frequency by Gender



GDP per Year : 1990
Top 10 Countries
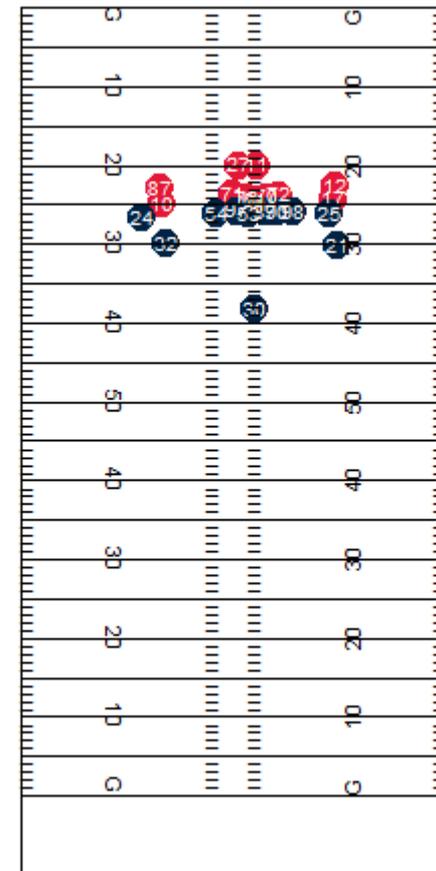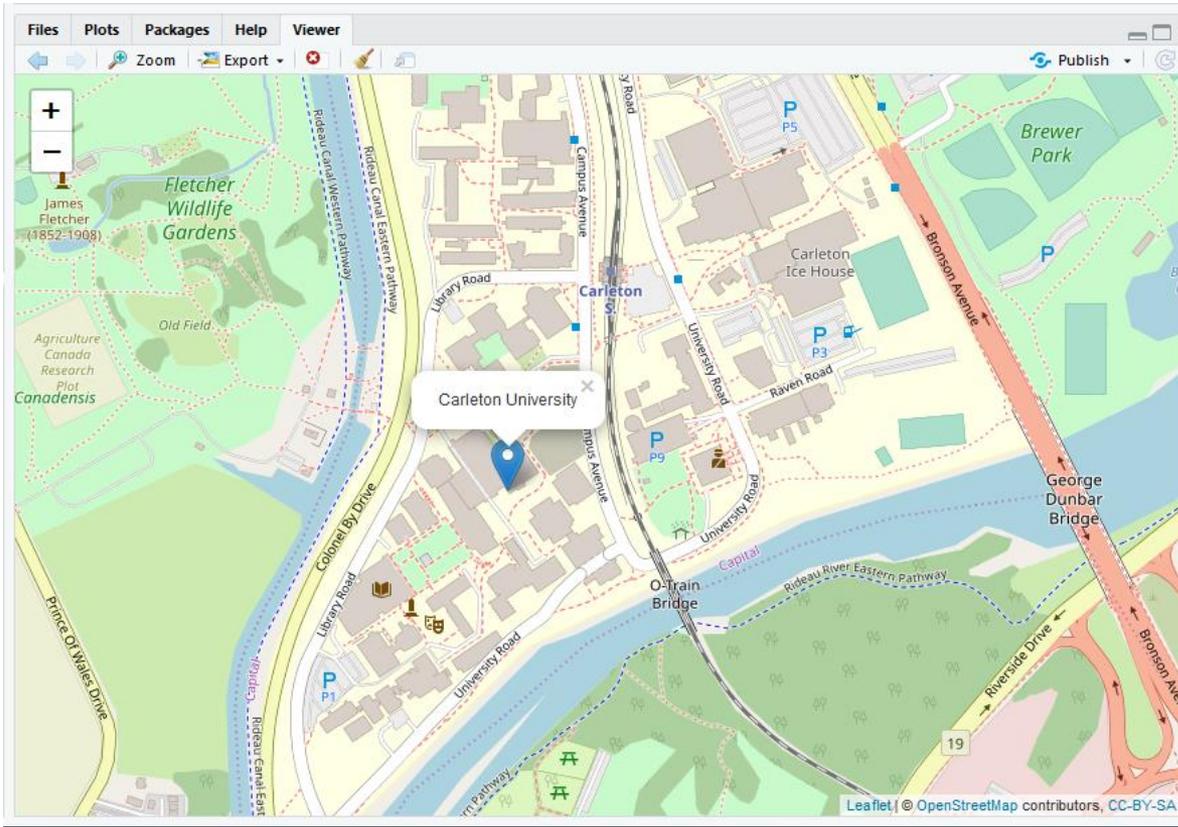
```
summary(mylogit)

## 
## Call:
## glm(formula = admit ~ gre + gpa + rank, family = "binomial", 
##     data = mydata)
## 
## Deviance Residuals:
##     Min       1Q   Median       3Q      Max  
## -1.627   -0.866   -0.639    1.149    2.079  
## 
## Coefficients:
##             Estimate Std. Error z value Pr(>|z|)    
## (Intercept) -3.98998    1.13995   -3.50  0.00047 ***
## gre          0.00226    0.00109    2.07  0.03847 *  
## gpa          0.80404    0.33182    2.42  0.01539 *  
## rank2       -0.67544    0.31649   -2.13  0.03283 *  
## rank3       -1.34020    0.34531   -3.88  0.00010 ***
## rank4       -1.55146    0.41783   -3.71  0.00020 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## 
## (Dispersion parameter for binomial family taken to be 1)
## 
##     Null deviance: 499.98  on 399  degrees of freedom
## Residual deviance: 458.52  on 394  degrees of freedom
## AIC: 470.5
## 
## Number of Fisher Scoring iterations: 4
```

8

**Carleton UNIVERSITY**

**Canada's Capital University**

## Negative Tweets



## Positive Tweets

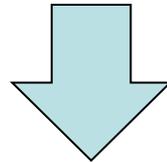- **Transform a Word Doc to Excel Spreadsheet**

Australia

Federal Institutions:
Administrative Appeals Tribunal (AAT) http://www.aat.gov.au/
AER Energy Made Easy http://www.energymadeeasy.gov.au/
www.energymadeeasy.gov.au/* Federal
Australian Antarctic Data Centre https://data.aad.gov.au/ data.aad.gov.au/* Federal
Australian Bureau of Statistics http://www.abs.gov.au/ www.abs.gov.au/* Federal

| Country | Name.Of.Agency | Agency.Website | Root.URL | Govt.Level | |
|---|---|---|---|---|---|
| Australia | | | | | |
| | Federal Institutions: | | | | |
| | Administrative Appeals Tribunal (AAT) | http://www.aat.gov.au/ | | | |
| | AER Energy Made Easy | http://www.energymadeeasy.gov.au/ | www.energymadeeasy.gov.au/* | Federal | |
| | Australian Antarctic Data Centre | https://data.aad.gov.au/ | data.aad.gov.au/* | Federal | |
| | Australian Bureau of Statistics | http://www.abs.gov.au/ | www.abs.gov.au/* | Federal | |

12

Canada's Capital University

**I removed the video that was here so that the slides could be sent through email.**

**If you are interested in learning more about web crawling, look up the RSelenium package**

# Scraping Data

**Grab Data from the Web and Bring it into R**

Canada's Capital University

- **Objects/Variables**
- **Data Types**
- **Functions**
- **Reading and Writing Data**
- **Operators**
- **Conditional Statements**
- **Loops**

# This is R



18

# This is R + RStudio

# Open RStudio

# 4 Main Areas (Each Corner)

| Script

| Console

| Environment/History

| Help/Plot

- **Object: Essentially a variable that can store a value**

```
1   # Different Variables:
2   x <- 89
3   Y <- "I'm a character variable"
4   A <- TRUE
5
6
7
```

- **Vector**
  - Collection of cells with a fixed size

- **Matrix**
  - Two-dimensional vector (row and column) "mathy"

- **Array**
  - A vector with one or more dimensions
    - *One dimensional array ~ A vector*
    - *Two dimensional array ~ A matrix*

- **List**
  - Can hold objects of different types

- **Data Frame**
  - A Table in which each column holds the same data type

**Function.Name(x=******, y = *******)**

**Function**          **Arguments**

**MyData <- read.csv("datafile.csv" , Header = TRUE)**

**Object**        **Function**                **Arguments**
**(Data Frame)**

**print(**"What we want to print"**)**

**mean(**A Numerical Object**)**

**strsplit(x=**A Character Vector, **split=**Where we want to split**)**

# For Help with functions:

- ## Within R
  - | Type ?read.csv
  - | help(read.csv)
  - | Check the Help box in bottom-right of RStudio

- ## Online
  - | Google!
  - | StackOverflow

## Different Functions:

- read.csv()
- readLines()

## Many Different Packages for Different Data

- sas7bdat
- foreign
- officer

- Flat Files
- Excel Files
- Statistical Software
- Databases
- Data from the Web

- **There are a few ways to grab a specific element or column/row of a data type**

- **We use square brackets [ ] to grab a specific element**

- **Elements are indexed in numerical order from 1 to n**

**We'll make a vector called My.Vector**

```
> My.Vector <- c(1:10)
> My.Vector
 [1]  1  2  3  4  5  6  7  8  9 10
>
```

**We can grab the first element by typing**

**My.Vector[1]**

**Or the fifth element by typing**

**My.Vector[5]**

**We'll make a DataFrame called DataFrame1**

```
> City <- c("Ottawa", "Montreal", "Calgary", "Edmonton")
> Province <- c("ON", "QC", "AB", "AB")
> Is.Capital <- c(TRUE, FALSE, FALSE, TRUE)
> Population <- c(994837, 1780000,1336000, 981280)
> DataFrame1 <- data.frame(City, Province, Is.Capital, Population)
>
> DataFrame1
      City Province Is.Capital Population
1   Ottawa       ON       TRUE     994837
2 Montreal       QC      FALSE    1780000
3  Calgary       AB      FALSE    1336000
4 Edmonton       AB       TRUE     981280
> |
```

**In a two-dimensional object, we have to index the row and column**

**To grab a specific element, we type DataFrame1[Row #, Column #]**

```
> DataFrame1
      City Province Is.Capital Population
1   Ottawa       ON       TRUE     994837
2 Montreal       QC      FALSE    1780000
3  Calgary       AB      FALSE    1336000
4 Edmonton       AB       TRUE     981280
```

**So to get "Calgary", we type DataFrame1[3,1]**

## To get a entire row:

```
> DataFrame1
      City Province Is.Capital Population
1   Ottawa       ON       TRUE     994837
2 Montreal       QC      FALSE    1780000
3  Calgary       AB      FALSE    1336000
4 Edmonton       AB       TRUE     981280
```

## To get the third row, type:

## DataFrame1[3,]

## (Just the row number, but make sure to add the comma)

## To get a entire column:

```
> DataFrame1
      City Province Is.Capital Population
1   Ottawa       ON      TRUE     994837
2 Montreal       QC     FALSE    1780000
3  Calgary       AB     FALSE    1336000
4 Edmonton       AB      TRUE     981280
```

## To get the third column, type:

## DataFrame1[,3]

## (Just the row number, but make sure to add the comma)

**OR**

**To get a entire row or column:**

```
> DataFrame1
        City Province Is.Capital Population
1   Ottawa       ON       TRUE     994837
2 Montreal       QC      FALSE    1780000
3  Calgary       AB      FALSE    1336000
4 Edmonton       AB       TRUE     981280
```

**Dataframes and Lists, have named objects like columns which can be grabbed by using $**

**Type**

**DataFrame$Column.Name**

**Type DataFrame1$Province will return:**

```
> DataFrame1$Province
[1] ON QC AB AB
Levels: AB ON QC
>
```

**If Statement is TRUE than do this**

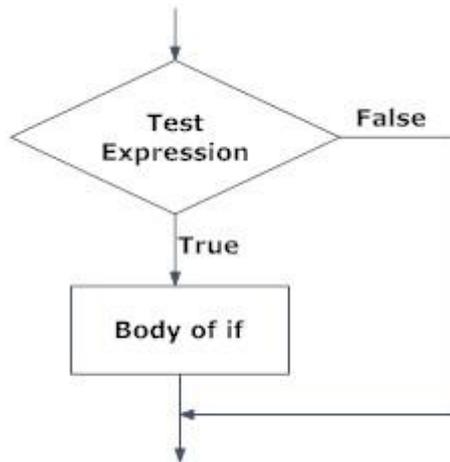**We use curly brackets {} to open and close our "if statement"**
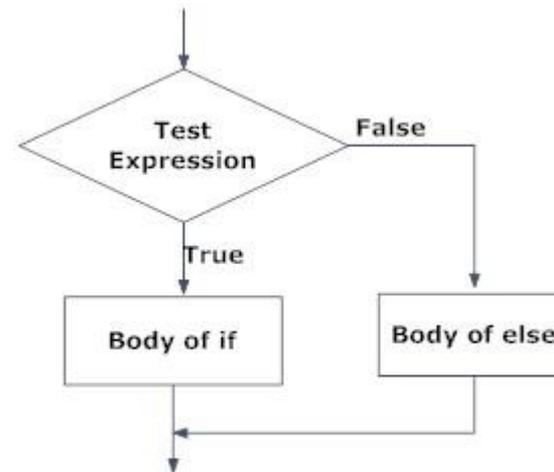


Fig: Operation of if statement

Fig: Operation of if...else statement

```
if(Conditional Statement) {
        thing to happen if Condition is TRUE
}




if(Conditional Statement){
        thing to happen if Condition is TRUE
} else {
        thing to happen if Condition is FALSE
}
```

- **Loops are used to repeat a specific task over a block of code**

- **For Loops**
  - Runs for a specific period of time (10 times, 20 times, 1000 times)

- **While Loops**
  - Runs until a specific condition is met (run until object is greater than 10, run until you encounter a specific object)

**Looks similar to our if statements**

```
for (i in 1:10){
        print(i)
}
```

**i is an object**

**1:10 is a range of numbers**

**print(i) will be done for each value of i**

- **Within R**
  - Type ?function
  - Check the Help box in bottom-right of RStudio
- **Online**
  - Google!
  - StackOverflow