# DS311 - R Lab Assignment

## Hanyang Xu

### 10/24/2022

## R Assignment 1

- In this assignment, we are going to apply some of the build in data set in R for descriptive statistics analysis.
- To earn full grade in this assignment, students need to complete the coding tasks for each question to get the result.
- After finished all the questions, knit the document into HTML format for submission.

**Question 1**

Using the **mtcars** data set in R, please answer the following questions.

```
# Loading the data

data(mtcars)
install.packages('plyr', repos = "http://cran.us.r-project.org")
```

```
## package 'plyr' successfully unpacked and MD5 sums checked
##
## The downloaded binary packages are in
##   C:\Users\Han\AppData\Local\Temp\Rtmpuq1fHB\downloaded_packages
```

```
install.packages("lifecycle", repos = "http://cran.us.r-project.org")
```

```
## package 'lifecycle' successfully unpacked and MD5 sums checked
##
## The downloaded binary packages are in
##   C:\Users\Han\AppData\Local\Temp\Rtmpuq1fHB\downloaded_packages
```

```
library("lifecycle")
```

```
## Warning: package 'lifecycle' was built under R version 4.2.2
```

```
library(rlang)
library(dplyr)
```

```
##
## Attaching package: 'dplyr'
```

```
## The following objects are masked from 'package:stats':
##
##     filter, lag


## The following objects are masked from 'package:base':
##
##     intersect, setdiff, setequal, union
```

```
install.packages("dplyr", repos = "http://cran.us.r-project.org")
```

```
## Warning: package 'dplyr' is in use and will not be installed
```

```
# Head of the data set
head(mtcars)
```

```
##                    mpg cyl disp  hp drat    wt  qsec vs am gear carb
## Mazda RX4         21.0   6  160 110 3.90 2.620 16.46  0  1    4    4
## Mazda RX4 Wag     21.0   6  160 110 3.90 2.875 17.02  0  1    4    4
## Datsun 710        22.8   4  108  93 3.85 2.320 18.61  1  1    4    1
## Hornet 4 Drive    21.4   6  258 110 3.08 3.215 19.44  1  0    3    1
## Hornet Sportabout 18.7   8  360 175 3.15 3.440 17.02  0  0    3    2
## Valiant           18.1   6  225 105 2.76 3.460 20.22  1  0    3    1
```

a. Report the number of variables and observations in the data set.

```
# Enter your code here!
dim(mtcars)
```

```
## [1] 32 11
```

```
# Answer:
print("There are total of 11 variables and 32 observations in this data set.")
```

```
## [1] "There are total of 11 variables and 32 observations in this data set."
```

b. Print the summary statistics of the data set and report how many discrete and continuous variables are in the data set.

```
# Enter your code here!
 summary(mtcars)
```

```
##       mpg             cyl             disp             hp
##  Min.   :10.40   Min.   :4.000   Min.   : 71.1   Min.   : 52.0
##  1st Qu.:15.43   1st Qu.:4.000   1st Qu.:120.8   1st Qu.: 96.5
##  Median :19.20   Median :6.000   Median :196.3   Median :123.0
##  Mean   :20.09   Mean   :6.188   Mean   :230.7   Mean   :146.7
##  3rd Qu.:22.80   3rd Qu.:8.000   3rd Qu.:326.0   3rd Qu.:180.0
##  Max.   :33.90   Max.   :8.000   Max.   :472.0   Max.   :335.0
##       drat             wt             qsec             vs
##  Min.   :2.760   Min.   :1.513   Min.   :14.50   Min.   :0.0000
```

```
##  1st Qu.:3.080    1st Qu.:2.581    1st Qu.:16.89    1st Qu.:0.0000
##  Median :3.695    Median :3.325    Median :17.71    Median :0.0000
##  Mean   :3.597    Mean   :3.217    Mean   :17.85    Mean   :0.4375
##  3rd Qu.:3.920    3rd Qu.:3.610    3rd Qu.:18.90    3rd Qu.:1.0000
##  Max.   :4.930    Max.   :5.424    Max.   :22.90    Max.   :1.0000
##        am              gear            carb
##  Min.   :0.0000   Min.   :3.000   Min.   :1.000
##  1st Qu.:0.0000   1st Qu.:3.000   1st Qu.:2.000
##  Median :0.0000   Median :4.000   Median :2.000
##  Mean   :0.4062   Mean   :3.688   Mean   :2.812
##  3rd Qu.:1.0000   3rd Qu.:4.000   3rd Qu.:4.000
##  Max.   :1.0000   Max.   :5.000   Max.   :8.000
```

```r
# Answer:
print("There are 2 discrete variables and 9 continuous variables in this data set.")
```

```
## [1] "There are 2 discrete variables and 9 continuous variables in this data set."
```

c. Calculate the mean, variance, and standard deviation for the variable **mpg** and assign them into variable names m, v, and s. Report the results in the print statement.

```r
# Enter your code here!
meann <- mean(mtcars$mpg)
vari <- var(mtcars$mpg)
sdev <- (mtcars$mpg)


print(paste("The average of Mile Per Gallon from this data set is ", meann , " with variance ", vari ,
```

```
##  [1] "The average of Mile Per Gallon from this data set is  20.090625  with variance  36.324102822580
##  [2] "The average of Mile Per Gallon from this data set is  20.090625  with variance  36.324102822580
##  [3] "The average of Mile Per Gallon from this data set is  20.090625  with variance  36.324102822580
##  [4] "The average of Mile Per Gallon from this data set is  20.090625  with variance  36.324102822580
##  [5] "The average of Mile Per Gallon from this data set is  20.090625  with variance  36.324102822580
##  [6] "The average of Mile Per Gallon from this data set is  20.090625  with variance  36.324102822580
##  [7] "The average of Mile Per Gallon from this data set is  20.090625  with variance  36.324102822580
##  [8] "The average of Mile Per Gallon from this data set is  20.090625  with variance  36.324102822580
##  [9] "The average of Mile Per Gallon from this data set is  20.090625  with variance  36.324102822580
## [10] "The average of Mile Per Gallon from this data set is  20.090625  with variance  36.324102822580
## [11] "The average of Mile Per Gallon from this data set is  20.090625  with variance  36.324102822580
## [12] "The average of Mile Per Gallon from this data set is  20.090625  with variance  36.324102822580
## [13] "The average of Mile Per Gallon from this data set is  20.090625  with variance  36.324102822580
## [14] "The average of Mile Per Gallon from this data set is  20.090625  with variance  36.324102822580
## [15] "The average of Mile Per Gallon from this data set is  20.090625  with variance  36.324102822580
## [16] "The average of Mile Per Gallon from this data set is  20.090625  with variance  36.324102822580
## [17] "The average of Mile Per Gallon from this data set is  20.090625  with variance  36.324102822580
## [18] "The average of Mile Per Gallon from this data set is  20.090625  with variance  36.324102822580
## [19] "The average of Mile Per Gallon from this data set is  20.090625  with variance  36.324102822580
## [20] "The average of Mile Per Gallon from this data set is  20.090625  with variance  36.324102822580
## [21] "The average of Mile Per Gallon from this data set is  20.090625  with variance  36.324102822580
## [22] "The average of Mile Per Gallon from this data set is  20.090625  with variance  36.324102822580
```

```
## [23] "The average of Mile Per Gallon from this data set is  20.090625  with variance  36.324102822580
## [24] "The average of Mile Per Gallon from this data set is  20.090625  with variance  36.324102822580
## [25] "The average of Mile Per Gallon from this data set is  20.090625  with variance  36.324102822580
## [26] "The average of Mile Per Gallon from this data set is  20.090625  with variance  36.324102822580
## [27] "The average of Mile Per Gallon from this data set is  20.090625  with variance  36.324102822580
## [28] "The average of Mile Per Gallon from this data set is  20.090625  with variance  36.324102822580
## [29] "The average of Mile Per Gallon from this data set is  20.090625  with variance  36.324102822580
## [30] "The average of Mile Per Gallon from this data set is  20.090625  with variance  36.324102822580
## [31] "The average of Mile Per Gallon from this data set is  20.090625  with variance  36.324102822580
## [32] "The average of Mile Per Gallon from this data set is  20.090625  with variance  36.324102822580
```

    d. Create two tables to summarize 1) average mpg for each cylinder class and 2) the standard deviation
       of mpg for each gear class.

```
# Enter your code here!
install.packages("qwraps2", repos = "http://cran.us.r-project.org")
```

```
## package 'qwraps2' successfully unpacked and MD5 sums checked
##
## The downloaded binary packages are in
##   C:\Users\Han\AppData\Local\Temp\Rtmpuq1fHB\downloaded_packages
```

```
library("qwraps2")
```

```
## Warning: package 'qwraps2' was built under R version 4.2.2
```

```
##
## Attaching package: 'qwraps2'
```

```
## The following object is masked from 'package:rlang':
##
##     ll
```

```
library("tidyverse")
```

```
## Warning: package 'tidyverse' was built under R version 4.2.2
```

```
## -- Attaching packages --------------------------------------- tidyverse 1.3.2 --
```

```
## v ggplot2 3.4.0      v purrr   0.3.5
## v tibble  3.1.8      v stringr 1.4.1
## v tidyr   1.2.1      v forcats 0.5.2
## v readr   2.1.3
```

```
## Warning: package 'ggplot2' was built under R version 4.2.2
```

```
## -- Conflicts ------------------------------------------ tidyverse_conflicts() --
## x purrr::%@%()         masks rlang::%@%()
## x purrr::as_function() masks rlang::as_function()
## x dplyr::filter()      masks stats::filter()
```

```
## x purrr::flatten()      masks rlang::flatten()
## x purrr::flatten_chr() masks rlang::flatten_chr()
## x purrr::flatten_dbl() masks rlang::flatten_dbl()
## x purrr::flatten_int() masks rlang::flatten_int()
## x purrr::flatten_lgl() masks rlang::flatten_lgl()
## x purrr::flatten_raw() masks rlang::flatten_raw()
## x purrr::invoke()       masks rlang::invoke()
## x dplyr::lag()          masks stats::lag()
## x qwraps2::ll()         masks rlang::ll()
## x purrr::splice()       masks rlang::splice()
```

```
mtcars %>% group_by(cyl) %>% summarize((Mean = mean(mpg)))
```

```
## # A tibble: 3 x 2
##     cyl '(Mean = mean(mpg))'
##   <dbl>              <dbl>
## 1     4              26.7
## 2     6              19.7
## 3     8              15.1
```

```
mtcars %>% group_by(gear) %>% summarize((SDD = sd(mpg)))
```

```
## # A tibble: 3 x 2
##    gear '(SDD = sd(mpg))'
##   <dbl>           <dbl>
## 1     3            3.37
## 2     4            5.28
## 3     5            6.66
```

e. Create a crosstab that shows the number of observations belong to each cylinder and gear class com-
   binations. The table should show how many observations given the car has 4 cylinders with 3 gears,
   4 cylinders with 4 gears, etc. Report which combination is recorded in this data set and how many
   observations for this type of car.

```
# Enter your code here!
library(tidyverse)
mtcars %>%
 select(cyl, gear) %>%
 table()
```

```
##      gear
## cyl  3  4  5
##   4  1  8  2
##   6  2  4  1
##   8 12  0  2
```

```
print("The most common car type in this data set is car with 3 cylinders and 8 gears. There are total o
```

```
## [1] "The most common car type in this data set is car with 3 cylinders and 8 gears. There are total
```

**Question 2**

Use different visualization tools to summarize the data sets in this question.

    a. Using the **PlantGrowth** data set, visualize and compare the weight of the plant in the three separated group. Give labels to the title, x-axis, and y-axis on the graph. Write a paragraph to summarize your findings.

```r
# Load the data set
data("PlantGrowth")

# Head of the data set
head(PlantGrowth)
```
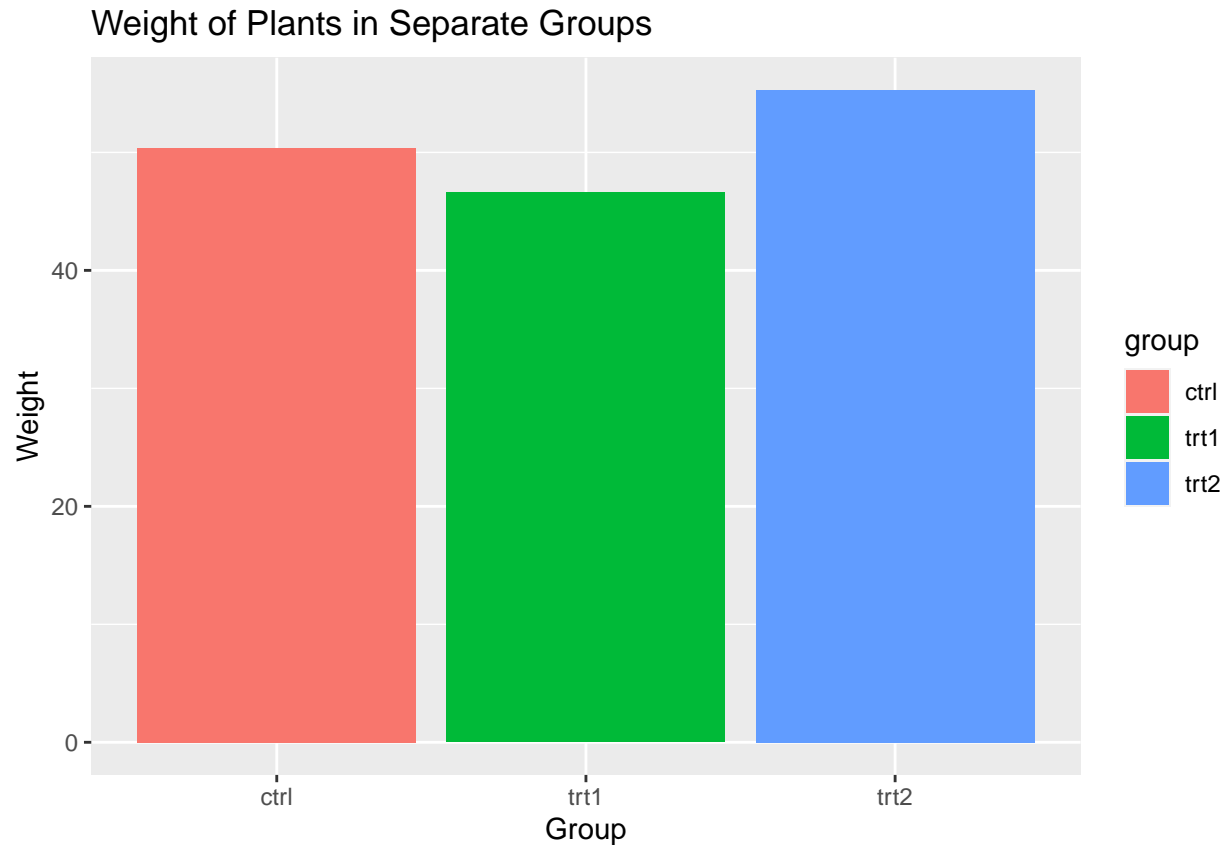
```
##   weight group
## 1   4.17  ctrl
## 2   5.58  ctrl
## 3   5.18  ctrl
## 4   6.11  ctrl
## 5   4.50  ctrl
## 6   4.61  ctrl
```

```r
# Enter your code here!
install.packages("tidyverse", repos = "http://cran.us.r-project.org")
```

```
## Warning: package 'tidyverse' is in use and will not be installed
```

```r
library(tidyverse)
PlantGrowth %>%
 ggplot(aes(x = group, y = weight, fill = group)) +
 geom_bar(stat = "identity") +
 labs(title = "Weight of Plants in Separate Groups",
      x = "Group",
      y = "Weight")
```

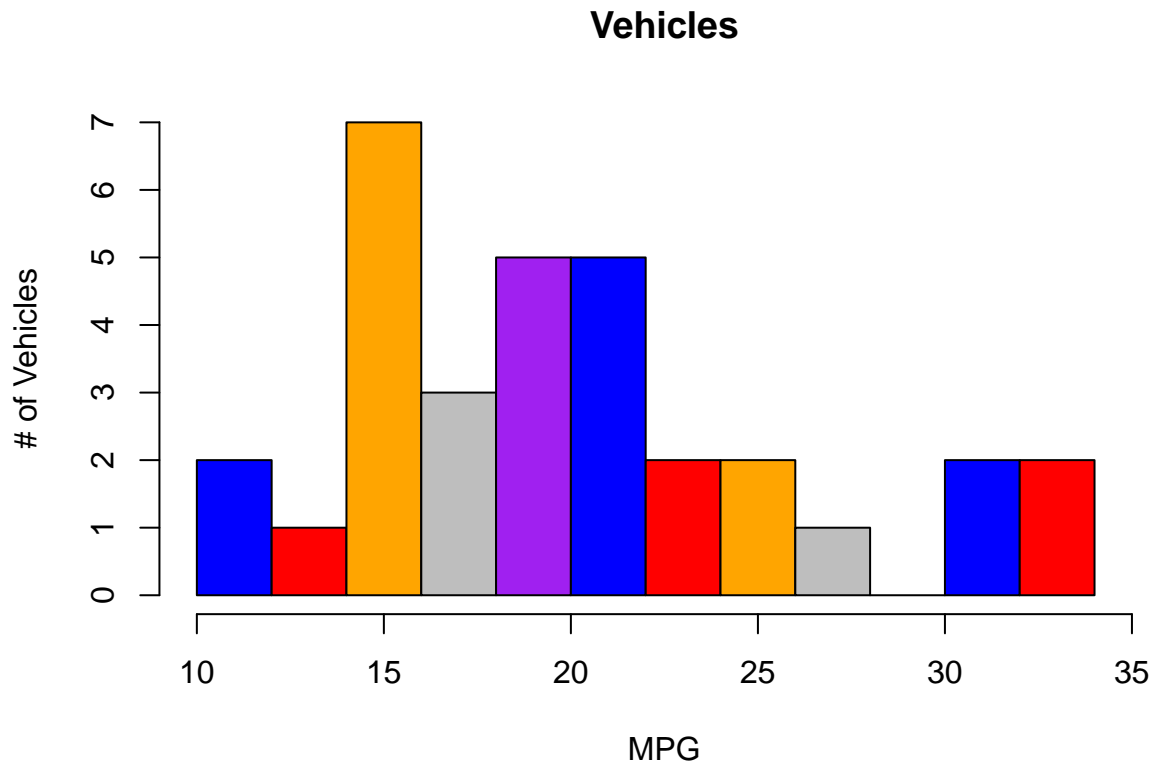# Weight of Plants in Separate Groups



Result:

=> Report a paragraph to summarize your findings from the plot! In these groups, trt2 has the heaviest plants and trt1 has the least weight. ctrl is the group that is in between of both groups.

b. Using the **mtcars** data set, plot the histogram for the column **mpg** with 10 breaks. Give labels to the title, x-axis, and y-axis on the graph. Report the most observed mpg class from the data set.

```
colors <- c("Blue", "red", "orange", "grey", "purple")

hist(mtcars$mpg,
    col=colors,
    main="Vehicles",
    breaks=10,
    xlim = range(10:35),
    xlab="MPG",
    ylab= "# of Vehicles")
```

# Vehicles



```
print("Most of the cars in this data set are in the class of 15 mile per gallon.")
```

```
## [1] "Most of the cars in this data set are in the class of 15 mile per gallon."
```

c. Using the **USArrests** data set, create a pairs plot to display the correlations between the variables in the data set. Plot the scatter plot with **Murder** and **Assault**. Give labels to the title, x-axis, and y-axis on the graph. Write a paragraph to summarize your results from both plots.

```r
# Load the data set
data("USArrests")

# Head of the data set
head(USArrests)
```
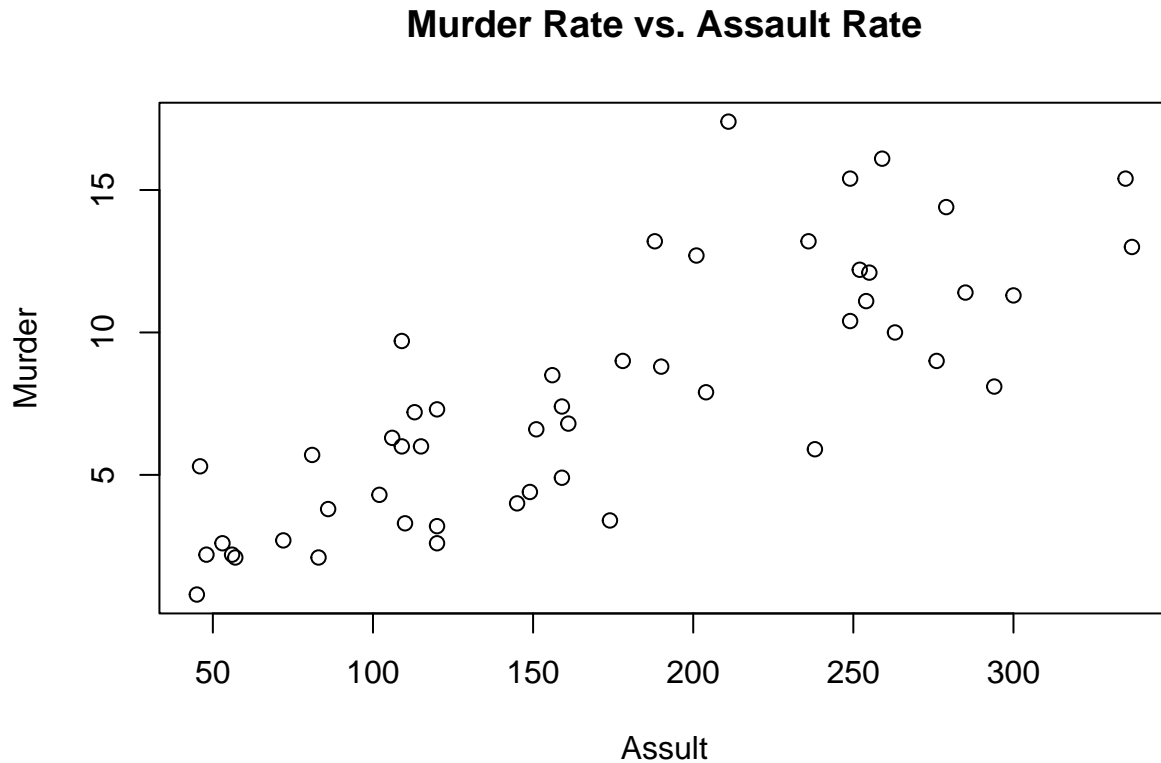
```
##            Murder Assault UrbanPop Rape
## Alabama      13.2     236       58 21.2
## Alaska       10.0     263       48 44.5
## Arizona       8.1     294       80 31.0
## Arkansas      8.8     190       50 19.5
## California    9.0     276       91 40.6
## Colorado      7.9     204       78 38.7
```

```
# Enter your code here!
plot(y = USArrests$Murder, x = USArrests$Assault, main = "Murder Rate vs. Assault Rate",xlab = "Assult"
```

## Murder Rate vs. Assault Rate



Result:

=> Report a paragraph to summarize your findings from the plot! The graph is uptrend and that means there is a positive relationship between assult and murder. If y goes up, then x will also gose up.

---

**Question 3**

Download the housing data set from www.jaredlander.com and find out what explains the housing prices in New York City.

Note: Check your working directory to make sure that you can download the data into the data folder.

   a. Create your own descriptive statistics and aggregation tables to summarize the data set and find any meaningful results between different variables in the data set.

```
# Head of the cleaned data set
head(housingData)
```

```
##    Neighborhood Market.Value.per.SqFt     Boro Year.Built
## 1    FINANCIAL               200.00 Manhattan       1920
```

```
## 2      FINANCIAL                  242.76 Manhattan         1985
## 4      FINANCIAL                  271.23 Manhattan         1930
## 5        TRIBECA                  247.48 Manhattan         1985
## 6        TRIBECA                  191.37 Manhattan         1986
## 7        TRIBECA                  211.53 Manhattan         1985
```

```
# Enter your code here!
summary(housingData)
```

```
##  Neighborhood        Market.Value.per.SqFt      Boro               Year.Built
##  Length:2530        Min.   : 10.66         Length:2530         Min.   :1825
##  Class :character   1st Qu.: 75.10         Class :character    1st Qu.:1926
##  Mode  :character   Median :114.89         Mode  :character    Median :1986
##                     Mean   :133.17                             Mean   :1967
##                     3rd Qu.:189.91                             3rd Qu.:2005
##                     Max.   :399.38                             Max.   :2010
```

Result:

=>The market minimum value of the market value per sqft is 10.66 and the maximum value is 399.38.The
oldest year built is 1825 and the newest year build is 2010. Both Boro and Neighborhood is character type.
***

   b. Create multiple plots to demonstrates the correlations between different variables. Remember to label
      all axes and give title to each graph.

```
# Enter your code here!

library(corrplot)
```

```
## Warning: package 'corrplot' was built under R version 4.2.2
```
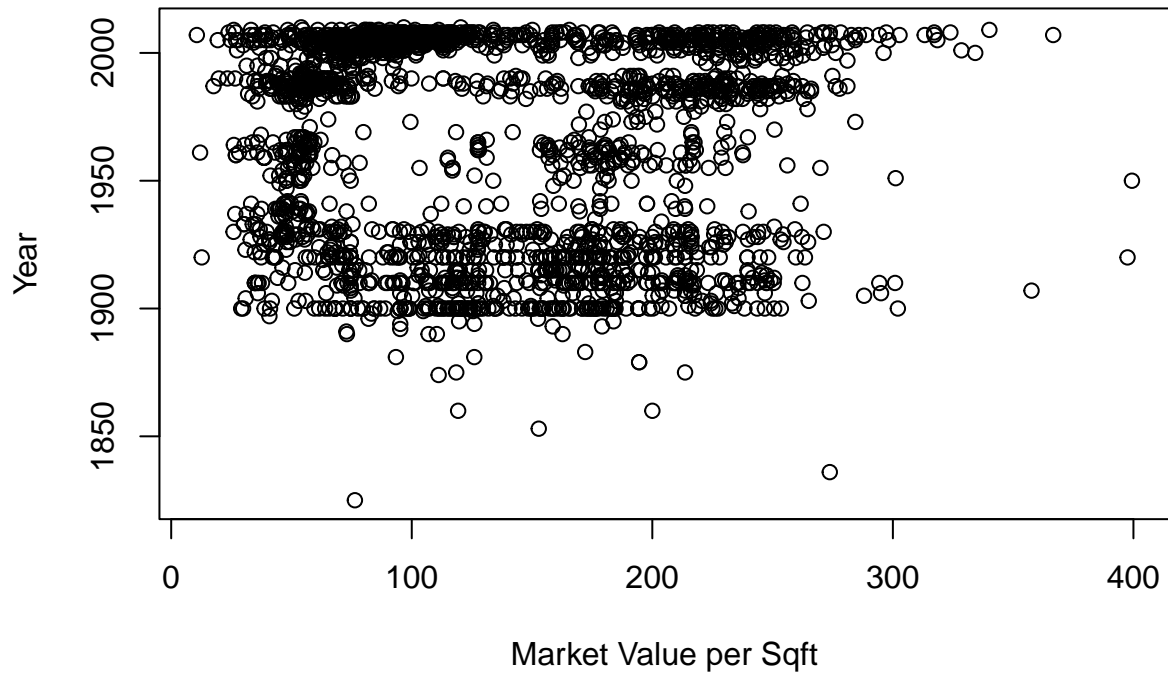
```
## corrplot 0.92 loaded
```

```
library(RColorBrewer)
library(ggplot2)
library(tidyverse)
plot(y = housingData$Year.Built, x = housingData$Market.Value.per.SqFt, main = "Market. Value vs. year.
```
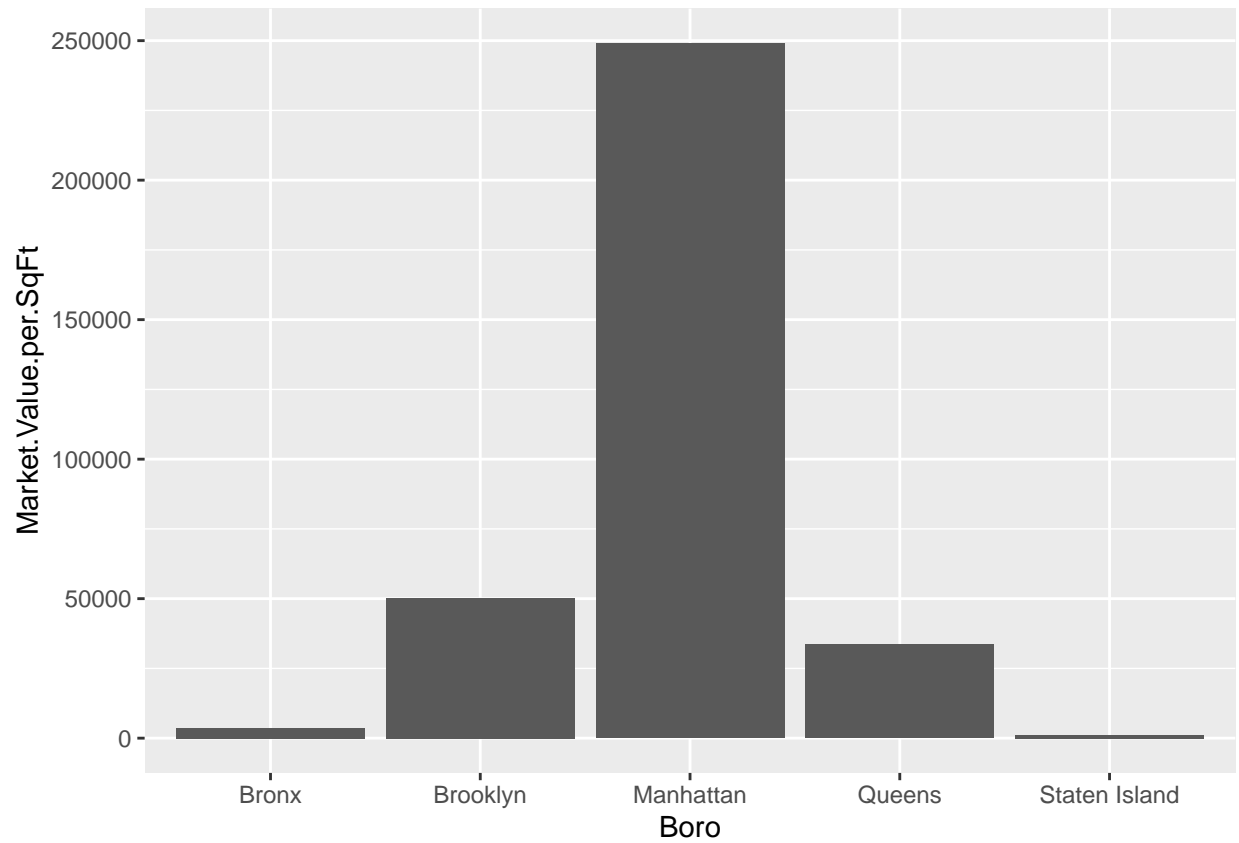
## Market. Value vs. year.Built



Market Value per Sqft

```
ggplot(housingData, aes(x=Boro, y=Market.Value.per.SqFt))+geom_bar(stat="identity")
```

c. Write a summary about your findings from this exercise.

=> Enter your answer here!

From the graphs, I found out that the there are more houses sell with year built after 2020, but the most expensive house is near 1950, and the cheapest house is the oldest. Most of the houses belong to Manhattan.