

# DS311 - R Lab Assignment

Tuan Pham

2023-04-09

## R Assignment 1

- In this assignment, we are going to apply some of the built-in data sets 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 finishing 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)
library(tidyverse)
```

```
## -- Attaching packages ----- tidyverse 1.3.2 --
## v ggplot2 3.4.0      v purrr   1.0.0
## v tibble  3.1.8      v dplyr  1.0.10
## v tidyr   1.2.1      v stringr 1.5.0
## v readr   2.1.3      v forcats 0.5.2
## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()    masks stats::lag()
```

```
# 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
```

```
str(mtcars)
```

```
## 'data.frame': 32 obs. of 11 variables:
## $ mpg : num 21 21 22.8 21.4 18.7 18.1 14.3 24.4 22.8 19.2 ...
## $ cyl : num 6 6 4 6 8 6 8 4 4 6 ...
## $ disp: num 160 160 108 258 360 ...
## $ hp : num 110 110 93 110 175 105 245 62 95 123 ...
## $ drat: num 3.9 3.9 3.85 3.08 3.15 2.76 3.21 3.69 3.92 3.92 ...
## $ wt : num 2.62 2.88 2.32 3.21 3.44 ...
## $ qsec: num 16.5 17 18.6 19.4 17 ...
## $ vs : num 0 0 1 1 0 1 0 1 1 1 ...
## $ am : num 1 1 1 0 0 0 0 0 0 0 ...
## $ gear: num 4 4 4 3 3 3 3 4 4 4 ...
## $ carb: num 4 4 1 1 2 1 4 2 2 4 ...
```

```
# Answer:
print("There are 5 discrete variables and 6 continuous variables in this data set.")
```

```
## [1] "There are 5 discrete variables and 6 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.

```
# Enter your code here!
m <- mtcars %>% summarize(mean(mpg))
v <- mtcars %>% summarize(var(mpg))
s <- mtcars %>% summarize(sd(mpg))

# print(paste("The average of Mile Per Gallon from this data set is ", 20.09062 , " with variance ", 3
```

- 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!
mtcars %>% group_by(cyl) %>% summarize(mean(mpg))
```

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

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

```
## # A tibble: 3 x 2
##   gear '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 combinations. 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!
mtcars %>% group_by(cyl,gear) %>% count()
```

```
## # A tibble: 8 x 3
## # Groups:   cyl, gear [8]
##   cyl gear     n
##   <dbl> <dbl> <int>
```

```
## 1      4      3      1
## 2      4      4      8
## 3      4      5      2
## 4      6      3      2
## 5      6      4      4
## 6      6      5      1
## 7      8      3     12
## 8      8      5      2
```

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

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

---

## Question 2

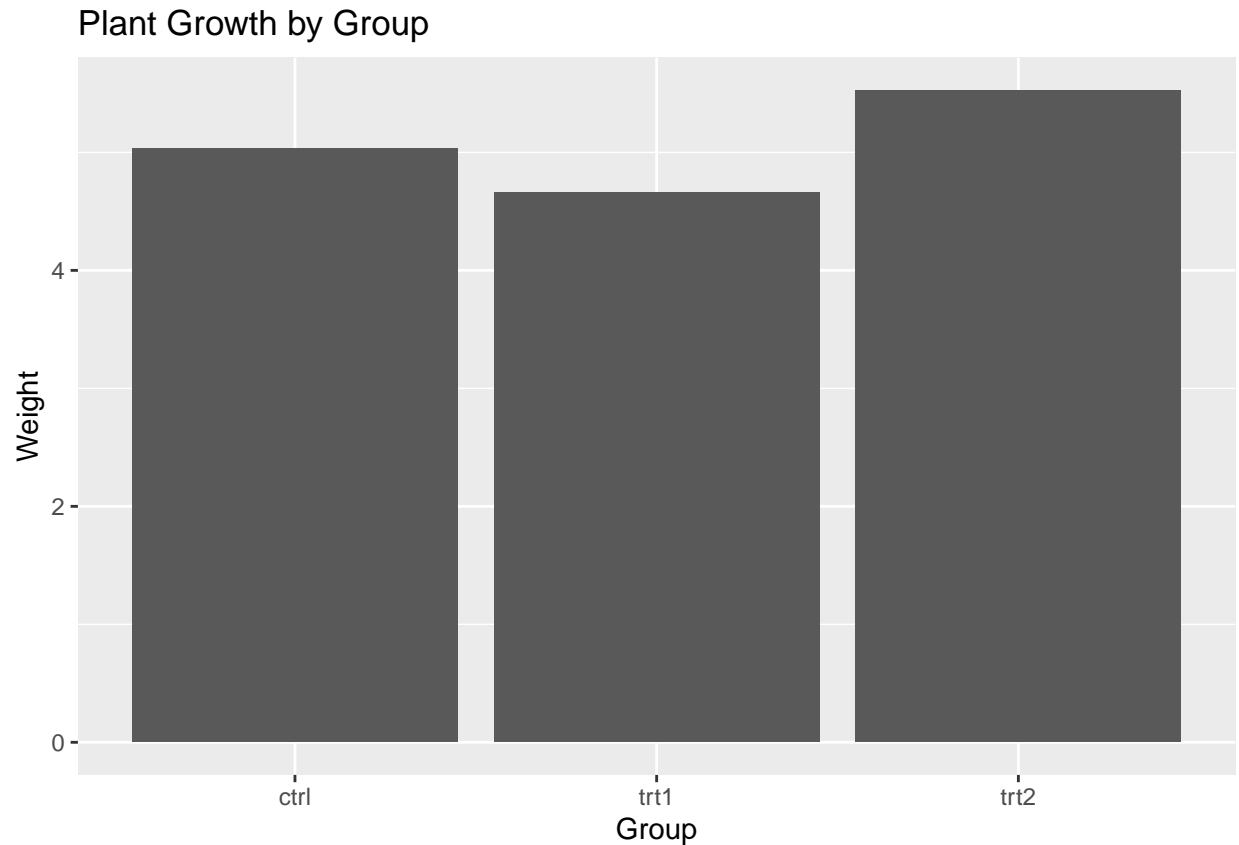
Use different visualization tools to summarize the data sets in this question.

- 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.

```
# Load the data set
data("PlantGrowth")
library(ggplot2)
# 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
```

```
# Enter your code here!
ggplot(PlantGrowth, aes(x = group, y = weight)) +
  geom_bar(stat = "summary", fun = "mean") +
  ggtitle("Plant Growth by Group") +
  xlab("Group") +
  ylab("Weight")
```

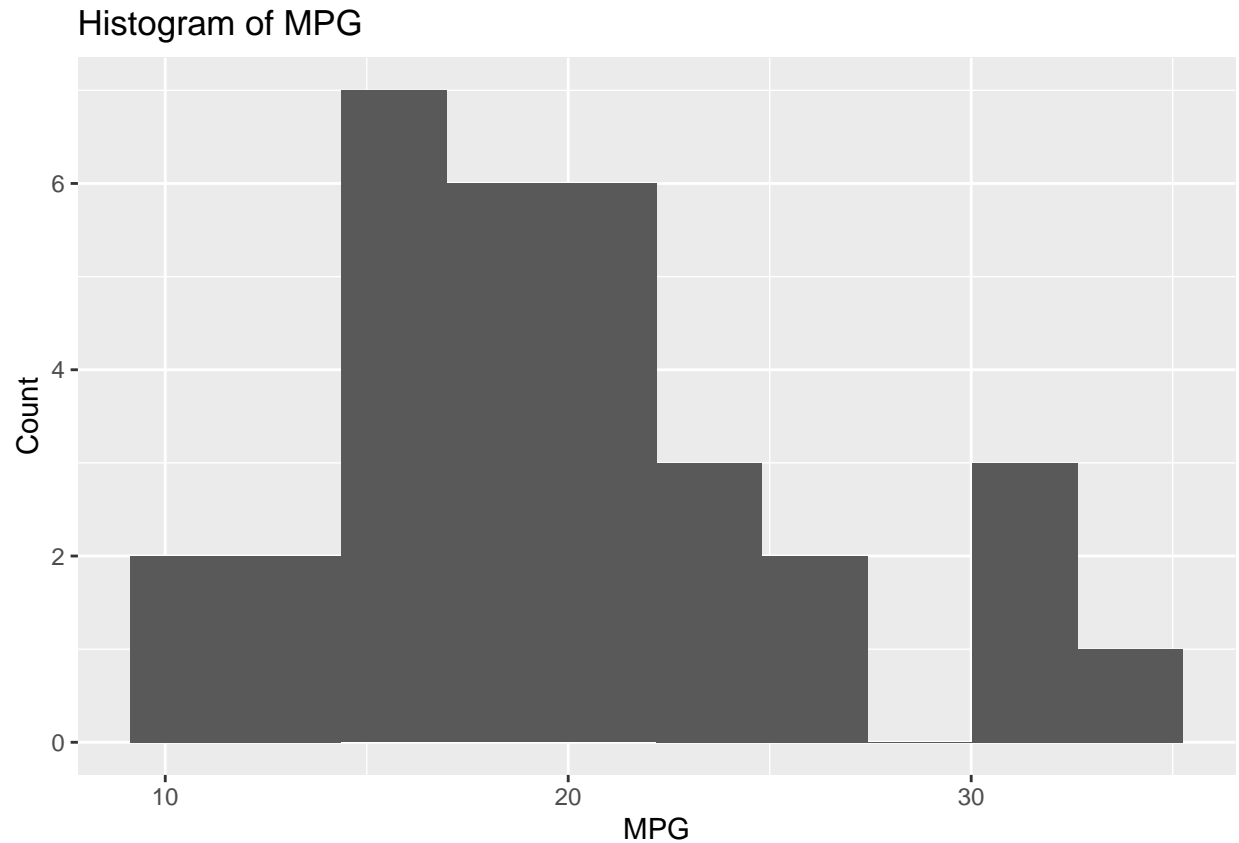


Result:

=> Report a paragraph to summarize your findings from the plot! Based on the bar graph, we can see that the heaviest average weight group is trt2. The group trt1 is the lightest average weight.

- 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.

```
ggplot(mtcars, aes(x = mpg)) +  
  geom_histogram(bins = 10) +  
  ggtitle("Histogram of MPG") +  
  xlab("MPG") +  
  ylab("Count")
```



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

```
## [1] "Most of the cars in this data set are in the class of around 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.

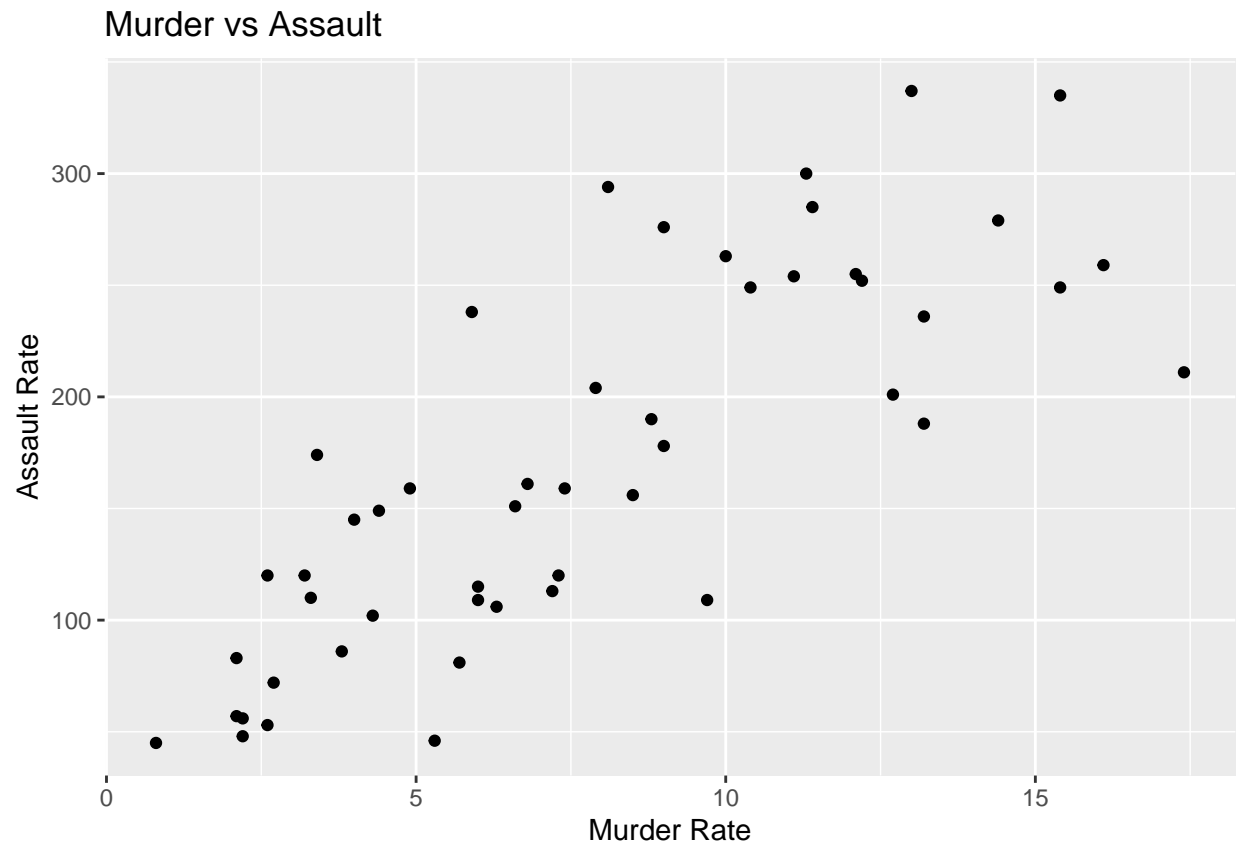
```
# Load the data set
data("USArrests")

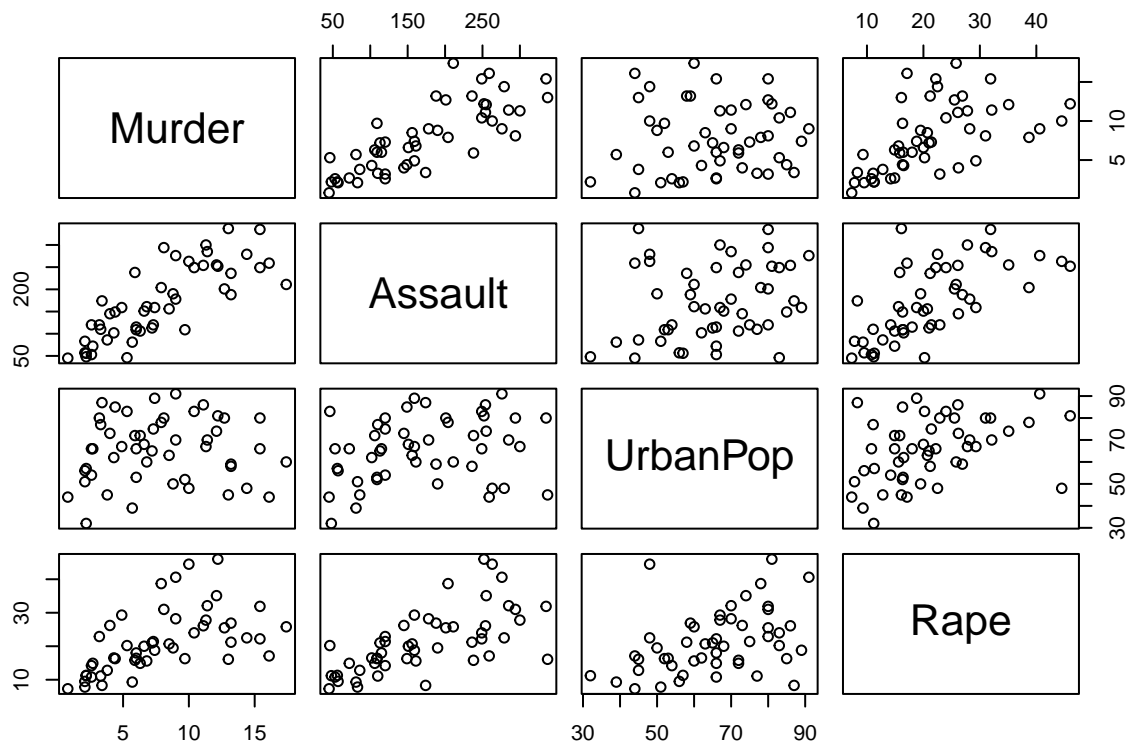
# Head of the data set
head(USArrests)
```

```
##           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!
```

```
ggplot(data = USArrests, aes(x = Murder, y = Assault)) +  
  geom_point() +  
  xlab("Murder Rate") +  
  ggtitle("Murder vs Assault") +  
  ylab("Assault Rate")
```





Result:

=> Report a paragraph to summarize your findings from the plot! Based on the scatter plot of Murder versus Assault, we can see that there is a positive correlation between these two variables, which suggests that states with higher murder rates tend to also have higher assault rates.

### Question 3

Download the housing data set from [www.jaredlander.com](http://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.

- 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
```

```
str(housingData)
```

```
## 'data.frame': 2530 obs. of 4 variables:
## $ Neighborhood : chr "FINANCIAL" "FINANCIAL" "FINANCIAL" "TRIBECA" ...
## $ Market.Value.per.SqFt: num 200 243 271 247 191 ...
## $ Boro : chr "Manhattan" "Manhattan" "Manhattan" "Manhattan" ...
## $ Year.Built : int 1920 1985 1930 1985 1986 1985 1986 1987 1985 1986 ...
## - attr(*, "na.action")= 'omit' Named int [1:96] 3 1395 1400 1412 1417 1425 1428 1429 1440 1445 ...
## ..- attr(*, "names")= chr [1:96] "3" "1395" "1400" "1412" ...
```

```
housingData %>% group_by(Neighborhood) %>% summarize(mean(Market.Value.per.SqFt))
```

```
## # A tibble: 148 x 2
## Neighborhood      'mean(Market.Value.per.SqFt)'
## <chr>              <dbl>
## 1 ALPHABET CITY      148.
## 2 ARROCHAR-SHORE ACRES 57.8
## 3 ASTORIA            91.5
## 4 BATH BEACH          70.3
## 5 BAY RIDGE           68.0
## 6 BAYSIDE             71.4
## 7 BEDFORD PARK/NORWOOD 38.2
## 8 BEDFORD STUYVESANT  83.2
## 9 BELMONT            56.4
## 10 BENSONHURST        71.7
## # ... with 138 more rows
```

```
housingData %>% group_by(Boro) %>% summarize(mean(Market.Value.per.SqFt))
```

```
## # A tibble: 5 x 2
## Boro      'mean(Market.Value.per.SqFt)'
## <chr>      <dbl>
## 1 Bronx      47.9
## 2 Brooklyn   80.1
## 3 Manhattan  181.
## 4 Queens     77.4
## 5 Staten Island 41.3
```

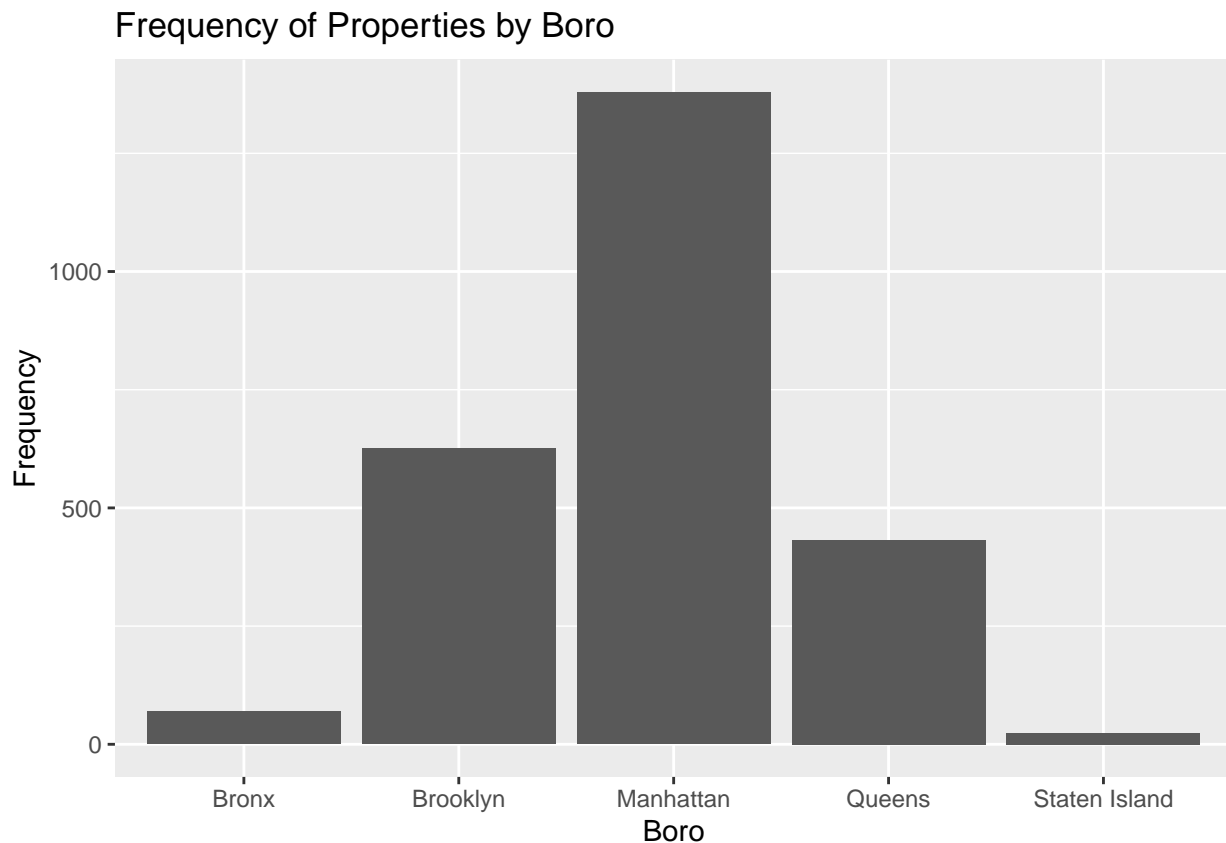
```
housingData %>% group_by(Boro) %>% count()
```

```
## # A tibble: 5 x 2
## # Groups:   Boro [5]
##   Boro      n
##   <chr>  <int>
## 1 Bronx      69
## 2 Brooklyn  626
## 3 Manhattan 1379
## 4 Queens    432
## 5 Staten Island 24
```

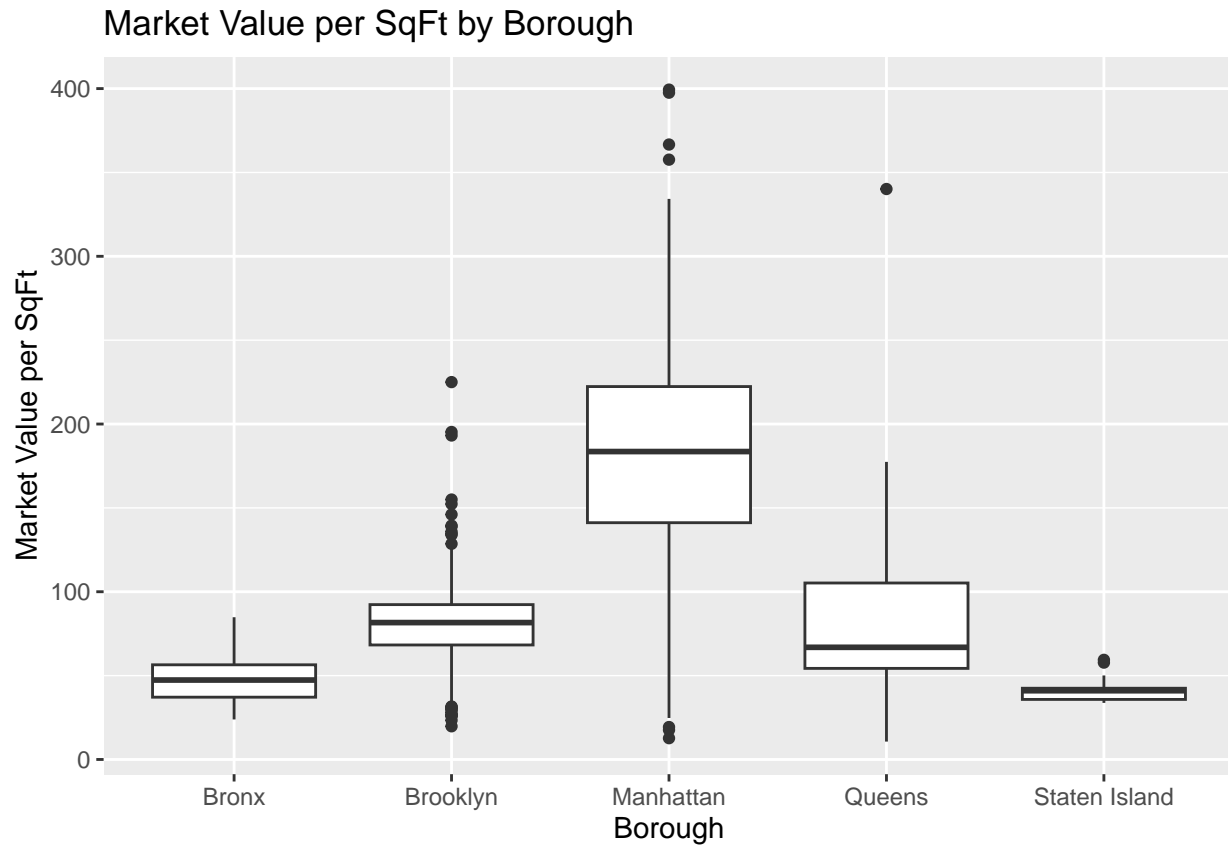
```
housing_summary <- housingData %>%
  group_by(Boro) %>%
  summarize(mean_market_value = mean(Market.Value.per.SqFt))
```

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

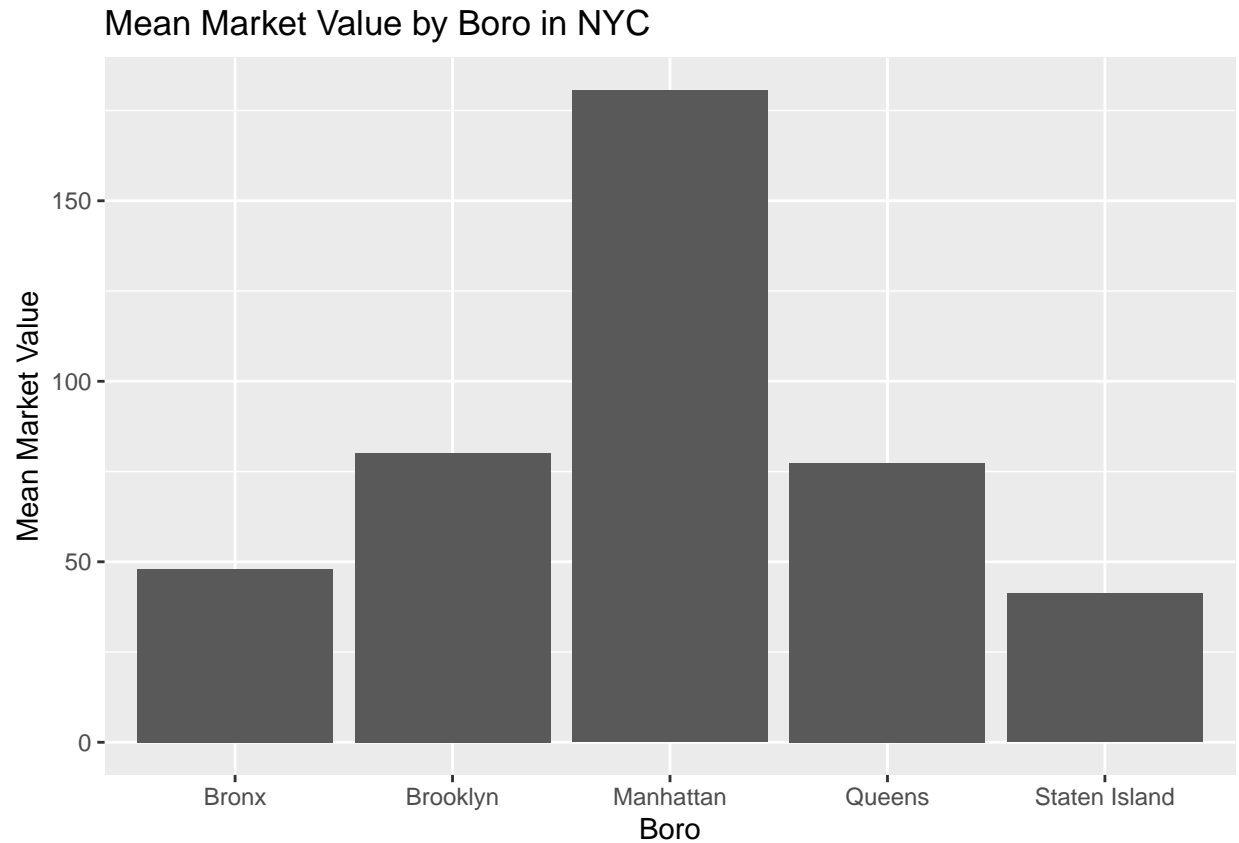
```
# Enter your code here!
ggplot(housingData, aes(x = Boro)) +
  geom_bar() +
  ggtitle("Frequency of Properties by Boro") +
  xlab("Boro") +
  ylab("Frequency")
```



```
ggplot(housingData, aes(x = Boro, y = Market.Value.per.SqFt)) +
  geom_boxplot() +
  ggtitle("Market Value per SqFt by Borough") +
  xlab("Borough") +
  ylab("Market Value per SqFt")
```



```
ggplot(housing_summary, aes(x = Boro, y = mean_market_value)) +
  geom_bar(stat="identity") +
  ggtitle("Mean Market Value by Boro in NYC") +
  xlab("Boro") +
  ylab("Mean Market Value")
```



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

According to the data, we can see that Manhattan has the most number of houses in New York City with 1379 houses among all the boroughs, followed by Brooklyn, Queens, Bronx, and Staten Island. The data is grouped by boro, and the mean market value is calculated for each boroughs. It shows that Manhattan is the boroughs having a mean market value per square feet of over 170. We can see that is a big gap when we compare to other boroughs.