

# DS311 - R Lab Assignment

Seyoung Kim

3/26/2023

## 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)
```

```
# 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!  
var<-dim(mtcars)[1]  
obser<-dim(mtcars)[2]
```

```
# Answer:  
print(paste("There are total of",var,"variables and ",obser,"observations in this data set."))
```

```
## [1] "There are total of 32 variables and  11 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
```

```
# Answer:
```

```
#discrete:cyl, vs,am,gear,carb
```

```
#continuous: mpg,disp,hp, drat,wt,qsec
```

```
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 <- mean(mtcars$mpg)
```

```
v <- var(mtcars$mpg)
```

```
s <- sd(mtcars$mpg)
```

```
print(paste("The average of Mile Per Gallon from this data set is ",m , " with variance ", v, " and s
```

```
## [1] "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("magrittr") # package installations are only needed the first time you use it
```

```
## Installing package into 'C:/Users/alluo/AppData/Local/R/win-library/4.2'  
## (as 'lib' is unspecified)
```

```
## package 'magrittr' successfully unpacked and MD5 sums checked
```

```
## Warning: cannot remove prior installation of package 'magrittr'
```

```
## Warning in file.copy(savedcopy, lib, recursive = TRUE): problem copying  
## C:\Users\alluo\AppData\Local\R\win-library\4.2\00LOCK\magrittr\libs\x64\magrittr.dll  
## to  
## C:\Users\alluo\AppData\Local\R\win-library\4.2\magrittr\libs\x64\magrittr.dll:  
## Permission denied
```

```
## Warning: restored 'magrittr'
```

```
##  
## The downloaded binary packages are in  
## C:\Users\alluo\AppData\Local\Temp\RtmpyKPD\downloaded_packages
```

```
install.packages("dplyr") # alternative installation of the %>%
```

```
## Installing package into 'C:/Users/alluo/AppData/Local/R/win-library/4.2'  
## (as 'lib' is unspecified)
```

```
## package 'dplyr' successfully unpacked and MD5 sums checked
```

```
## Warning: cannot remove prior installation of package 'dplyr'
```

```
## Warning in file.copy(savedcopy, lib, recursive = TRUE): problem copying  
## C:\Users\alluo\AppData\Local\R\win-library\4.2\00LOCK\dplyr\libs\x64\dplyr.dll  
## to C:\Users\alluo\AppData\Local\R\win-library\4.2\dplyr\libs\x64\dplyr.dll:  
## Permission denied
```

```
## Warning: restored 'dplyr'
```

```
##  
## The downloaded binary packages are in  
## C:\Users\alluo\AppData\Local\Temp\RtmpyKPD\downloaded_packages
```

```
library(magrittr) # needs to be run every time you start R and want to use %>%
```

```
## Warning: package 'magrittr' was built under R version 4.2.3
```

```
library(dplyr)
```

```
## Warning: package 'dplyr' was built under R version 4.2.3
```

```
##
```

```
## 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
```

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

```
## # A tibble: 3 x 2
```

```
##   cyl AVGMPG
```

```
##   <dbl> <dbl>
```

```
## 1     4  26.7
```

```
## 2     6  19.7
```

```
## 3     8  15.1
```

```
mtcars %>%  
  group_by(gear) %>%  
  summarize(AVGMPG=mean(mpg))
```

```
## # A tibble: 3 x 2
```

```
##   gear AVGMPG
```

```
##   <dbl> <dbl>
```

```
## 1     3  16.1
```

```
## 2     4  24.5
```

```
## 3     5  21.4
```

- 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.

```
Cyl<-mtcars$cyl  
Gear<-mtcars$gear  
  
table(Cyl,Gear)
```

```
##   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 8 cylinders and 3 gears. There are total of 14 cars with 8 cylinders and 3 gears.")
```

```
## [1] "The most common car type in this data set is car with 8 cylinders and 3 gears. There are total of 14 cars with 8 cylinders and 3 gears."
```

---

## 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.

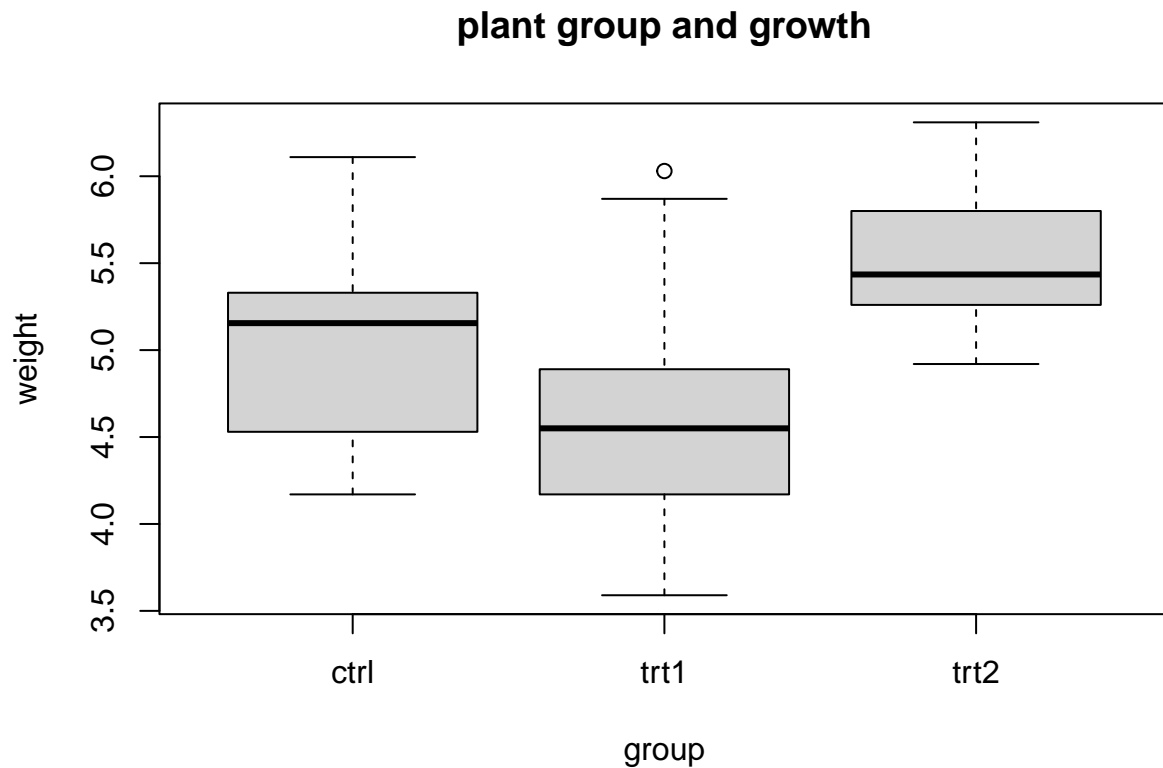
```
# Load the data set  
data("PlantGrowth")
```

```
# 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!
```

```
boxplot(weight ~ group, data = PlantGrowth, main="plant group and growth", xlab="group", ylab="weight")
```



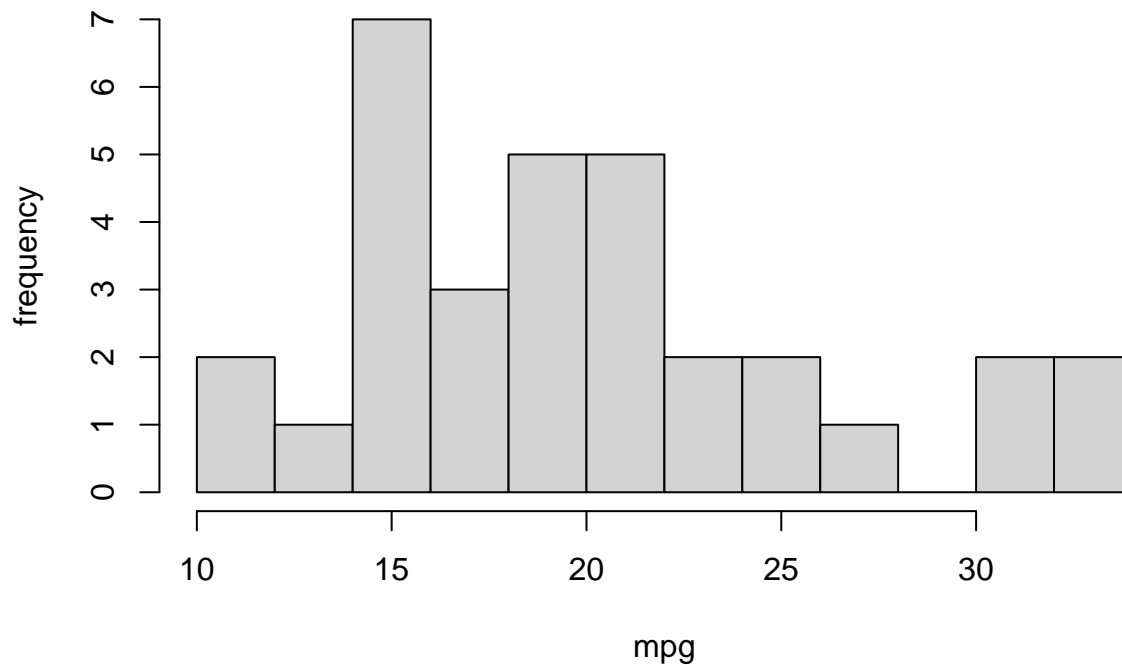
Result:

=> Report a paragraph to summarize your findings from the plot! The tallest group is trt2, and the shortest group is trt1. The trt1 group has the greatest variability in weight, while the trt2 group has the least variability in 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.

```
hist(mtcars$mpg,breaks=10,main="MPG of the cars",xlab="mpg",ylab="frequency")
```

## MPG of the cars



```
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.

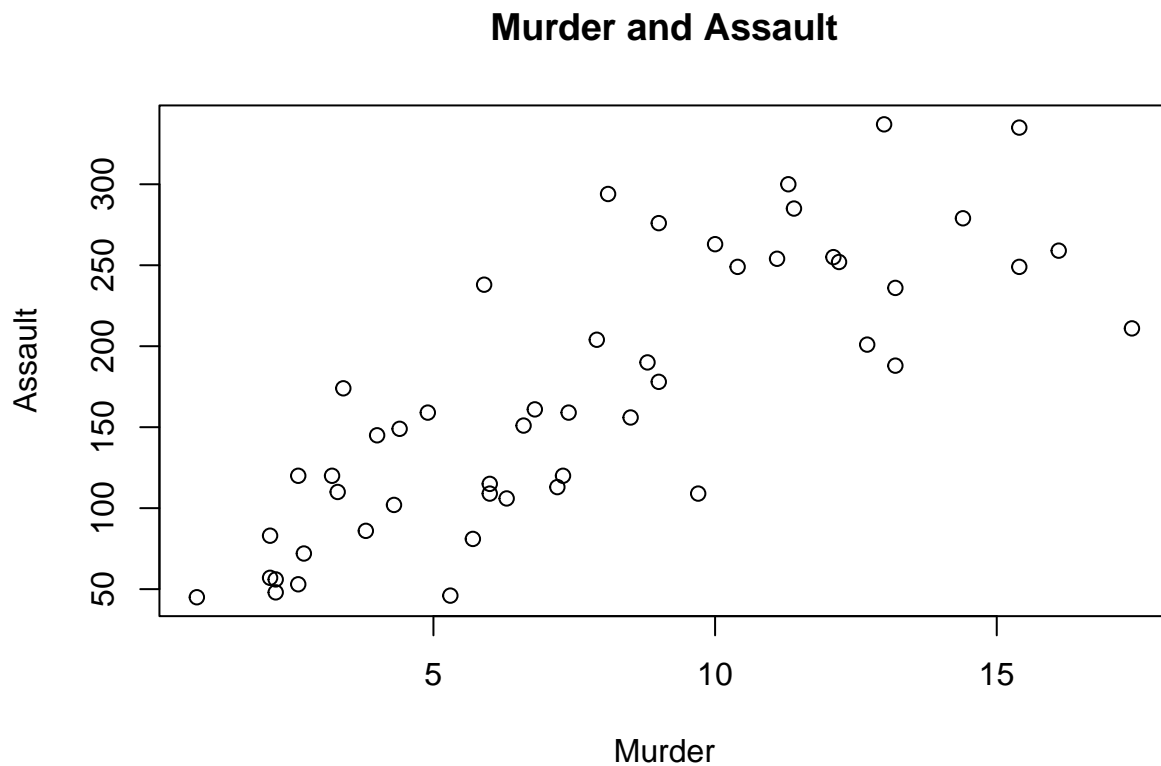
```
# Load the data set
data("USArrests")

# 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(USArrests$Murder,USArrests$Assault,main="Murder and Assault",xlab="Murder",ylab="Assault")
```



Result:

=> Report a paragraph to summarize your findings from the plot! As the value of Murder increases, there is a tendency for the value of Assault to also increase.

---

### 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.

```
# Load and clean the housing data set
wd1 <- getwd()
paste("Current Working Directory: ", wd1)
```

```
## [1] "Current Working Directory: C:/Users/alluo/Documents"
```



```
# Load and clean the housing data set
download.file(url='https://www.jaredlander.com/data/housing.csv',
             destfile='data/housing.csv', mode='wb')
housingData <- read.csv('./data/housing.csv')
housingData <- subset(housingData,
                     select = c("Neighborhood", "Market.Value.per.SqFt", "Boro", "Year.Built"))
housingData <- na.omit(housingData)
```

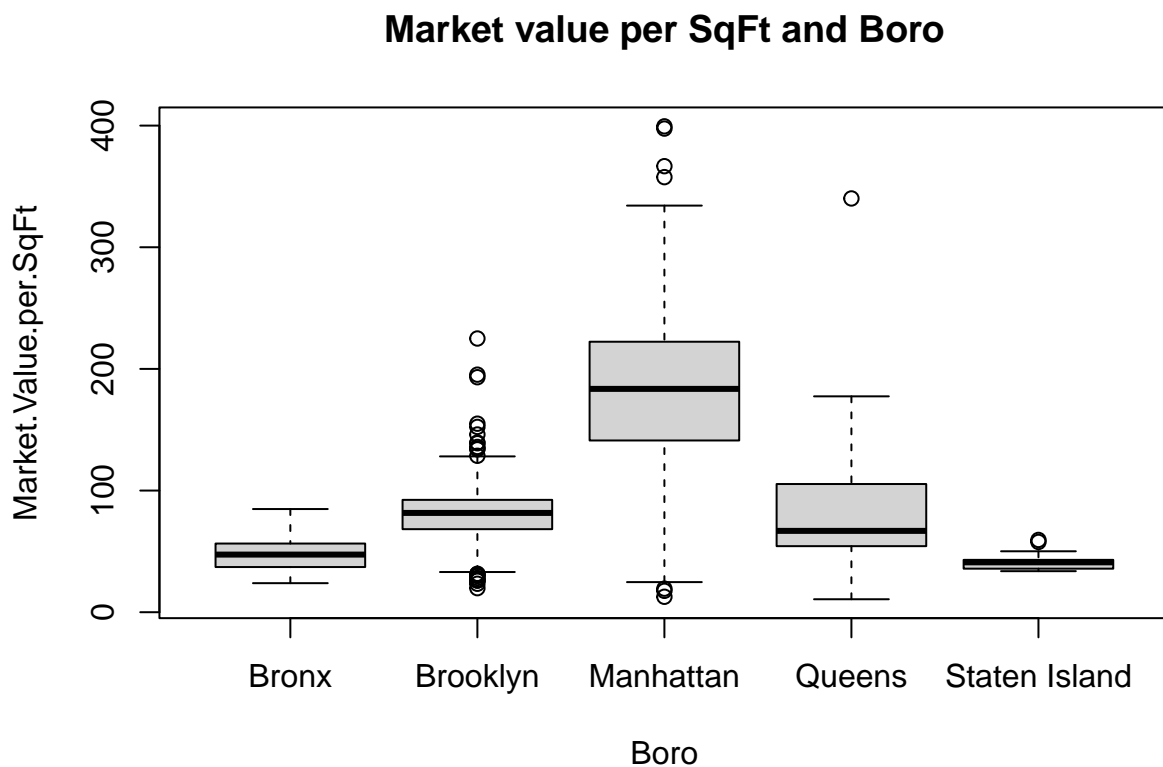
- a. Create your own descriptive statistics and aggregation tables to summarize the data set and find any meaningful results between different variables in the data set.

```
# 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!
```

```
boxplot(Market.Value.per.SqFt ~ Boro, main="Market value per SqFt and Boro", data = housingData, xlab="Boro")
```



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

```
# Enter your code here!  
install.packages("ggplot")
```

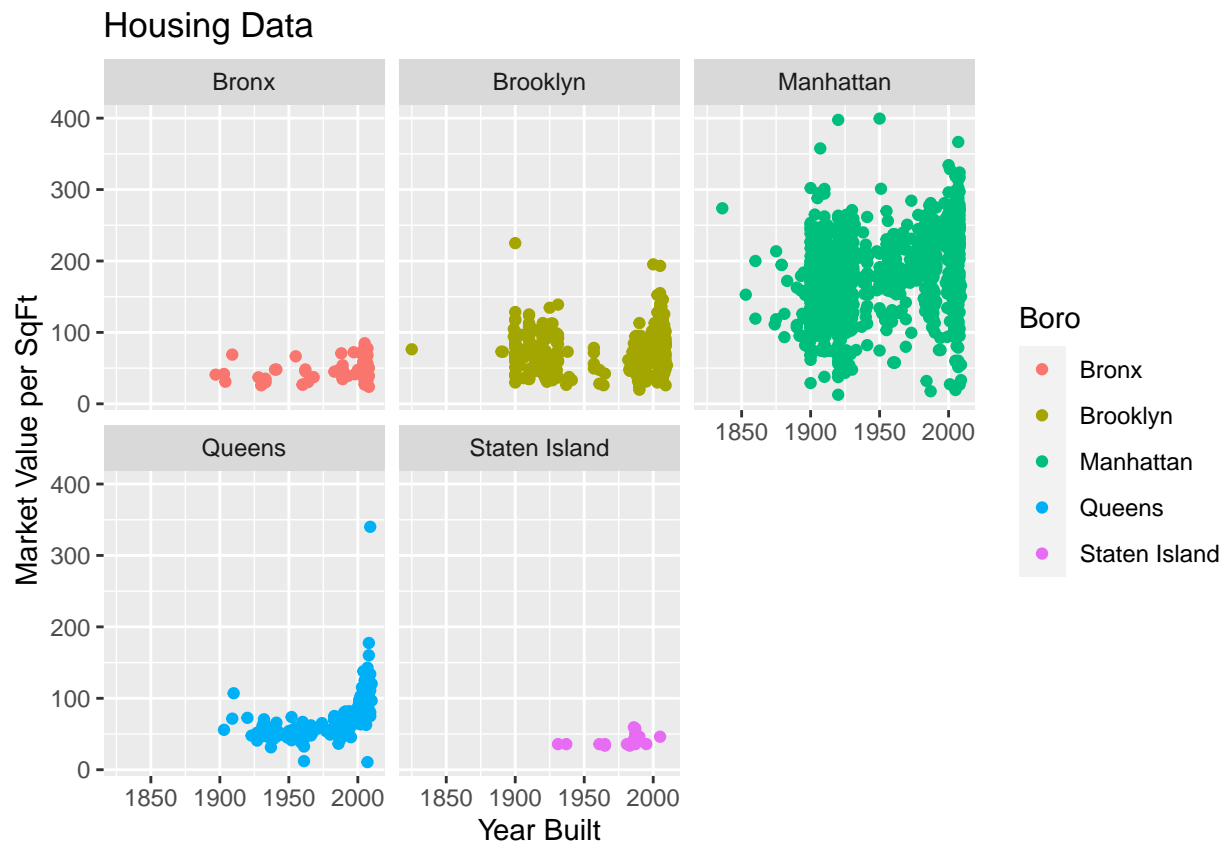
```
## Installing package into 'C:/Users/alluo/AppData/Local/R/win-library/4.2'  
## (as 'lib' is unspecified)
```

```
## Warning: package 'ggplot' is not available for this version of R  
##  
## A version of this package for your version of R might be available elsewhere,  
## see the ideas at  
## https://cran.r-project.org/doc/manuals/r-patched/R-admin.html#Installing-packages
```

```
library(ggplot2)
```

```
## Warning: package 'ggplot2' was built under R version 4.2.3
```

```
g<-ggplot(housingData,aes(x=Year.Built,y=Market.Value.per.SqFt))+ geom_point(aes(color=Boro))  
g+ facet_wrap(~Boro)+ labs(x = "Year Built", y = "Market Value per SqFt", title = "Housing Data")
```



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

=> I used plot,boxplot,hist functions to visualize data. Also, I used facet\_wrap function to categorize the data by Boro.