# DS311 - R Lab Assignment

Maximilian Leitschuh

8/22/2022

## 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)

# 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!
variables <- ncol(mtcars)
observations <- nrow(mtcars)
# Answer:
print(paste("There are total of", variables, "variables and", observations,"observations in this data s
```

```
## [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

# Enter your code here!
summ <- summary(mtcars)
summ
```

```
##       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
```

1

```
##  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 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.

```r
# Enter your code here!
m <- round(mean(mtcars$mpg), digits=2)
v <- round(var(mtcars$mpg), digits=2)
s <- round(sd(mtcars$mpg), digits=2)

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  20.09  with variance  36.32  and standard
```

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

```r
# Enter your code here!
tab1 <-mtcars%>%
        group_by(cyl)%>%
        summarize(average=mean(mpg))
tab1
```

```
## # A tibble: 3 x 2
##     cyl average
##   <dbl>   <dbl>
## 1     4    26.7
## 2     6    19.7
## 3     8    15.1
```

```r
tab2 <-mtcars%>%
        group_by(gear)%>%
        summarize(standardDeviation=sd(mpg))
tab2
```

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

2

```
##    <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!
cross <- crosstable(mtcars, cols=c(cyl), by=c(gear), total = "row")
cross
```

```
## # A tibble: 3 x 7
##    .id   label variable `3`          `4`         `5`         Total
##    <chr> <chr> <chr>    <chr>        <chr>       <chr>       <chr>
## 1 cyl    cyl   4        1 (9.09%)    8 (72.73%) 2 (18.18%) 11 (34.38%)
## 2 cyl    cyl   6        2 (28.57%)   4 (57.14%) 1 (14.29%) 7 (21.88%)
## 3 cyl    cyl   8        12 (85.71%)  0 (0%)     2 (14.29%) 14 (43.75%)
```

```
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
```

---

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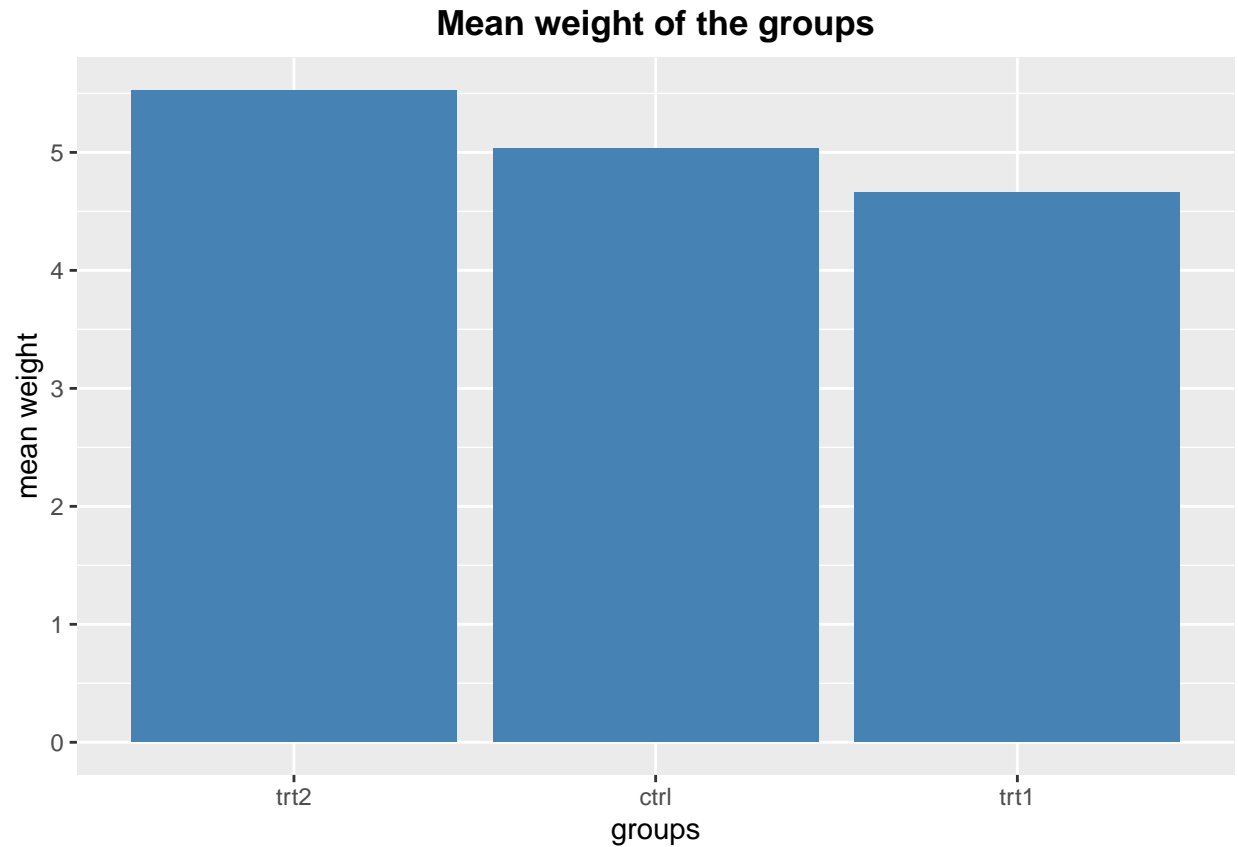
```
# Load the data set
data("PlantGrowth")
```

```
# 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!
ggplot(PlantGrowth, aes(x=reorder(group,-weight), y=weight),stat = "summary") +
  geom_bar(stat="summary", fill="steelblue", fun="mean")+
  labs(y="mean weight", x="groups")+
  ggtitle("Mean weight of the groups")+
  scale_y_continuous(breaks = seq(0,6,1))+
  theme(plot.title = element_text(hjust=0.5, face="bold"))
```
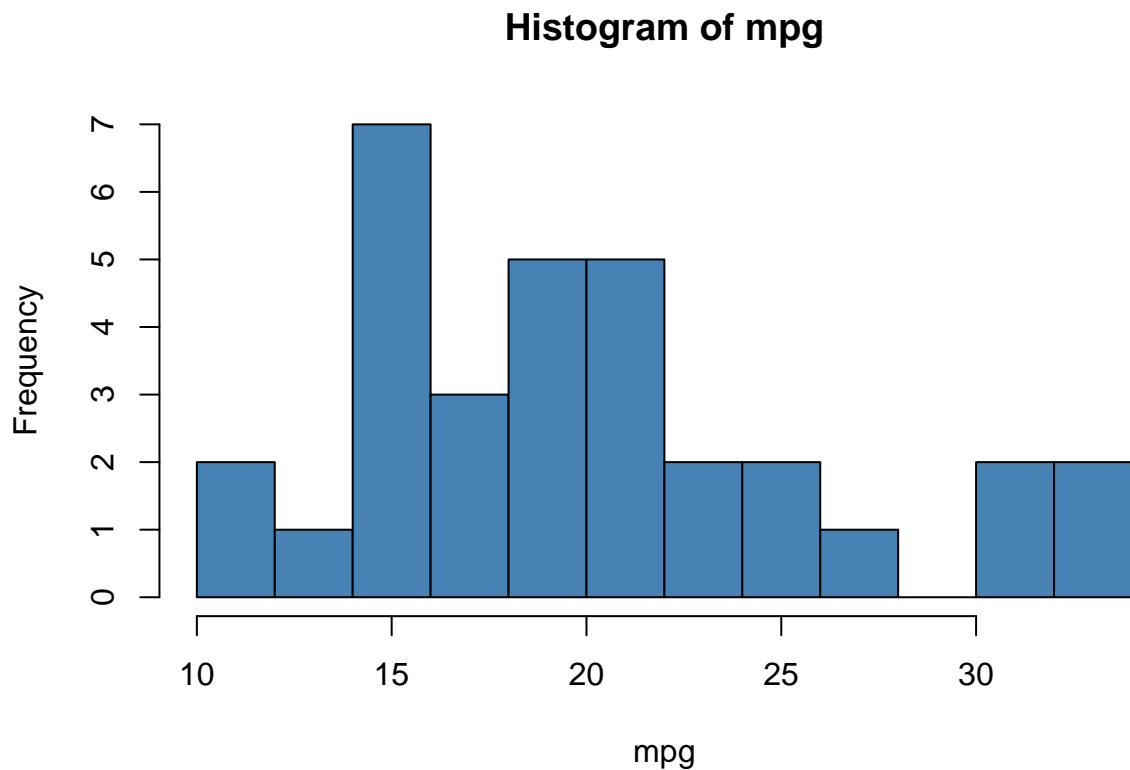
**Mean weight of the groups**

Result:

=> On average, plants in group "trt2" weigh the most with around 5.5 pound. This is followed by the plants of group "ctrl" with about 5 pounds and the lightest are the plants of group "trt1" with an average of around 4.6 pounds.

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,
     col='steelblue',
     main='Histogram of mpg',
     xlab='mpg',
     breaks=10,
     ylab='Frequency')
```

# Histogram of mpg



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

```
## [1] "Most of the cars in this data set are in the class of 15 miles 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)
```
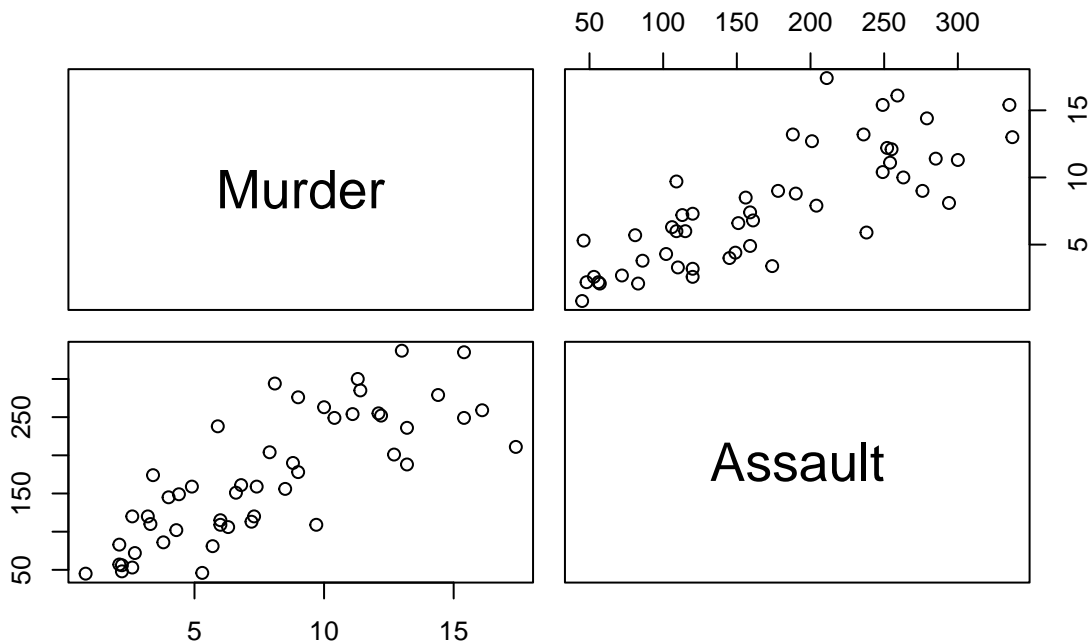
```
##            Murder Assault UrbanPop Rape
## Alabama      13.2     236       58 21.2
## Alaska       10.0     263       48 44.5
## Arizona       8.1     294       80 31.0
## Arkansas      8.8     190       50 19.5
## California    9.0     276       91 40.6
## Colorado      7.9     204       78 38.7
```

```
# Enter your code here!
pairs(~ Murder + Assault, data = USArrests, main = "Correlation between Murder and Assault")
```

# Correlation between Murder and Assault



Result:

=> The correlations within the pairplot indicate that the more people arrested for assault, the more people arrested for murder, and vice versa. In my opinion, this is logical, because with an increase in arrests for assault, people's propensity to violence continues to increase.

---

**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
```

```r
# Enter your code here!
#descriptive statistics
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
```

```r
#Mean market value per sqft in the different boroughs of new york
value_boroughs <- aggregate(housingData$Market.Value.per.SqFt, list(housingData$Boro), FUN=mean)
value_boroughs[order(value_boroughs$x, decreasing = TRUE),]
```

```
##          Group.1          x
## 3      Manhattan 180.59265
## 2       Brooklyn  80.13439
## 4         Queens  77.38137
## 1          Bronx  47.93232
## 5 Staten Island  41.26958
```

```r
#Mean market value per sqft in the 5 most expensive neighborhoods of new york
value_neighborhood <- aggregate(housingData$Market.Value.per.SqFt, list(housingData$Neighborhood), FUN=r
top5_neighborhoods <- value_neighborhood[order(value_neighborhood$x, decreasing = TRUE)[1:5],]
colnames(top5_neighborhoods) <- c("Neighborhood","mean")
top5_neighborhoods
```

```
##                   Neighborhood     mean
## 92                 MIDTOWN CBD 234.3615
## 49                    FLATIRON 223.3031
## 94                MIDTOWN WEST 222.0649
## 130 UPPER EAST SIDE (59-79) 216.8372
## 23                     CHELSEA 215.9493
```

```r
#information of the most expensive building
housingData[which.max(housingData$Market.Value.per.SqFt),]
```

```
##          Neighborhood Market.Value.per.SqFt      Boro Year.Built
## 191 LOWER EAST SIDE                 399.38 Manhattan       1950
```

```r
#information of the most cheapest building
housingData[which.min(housingData$Market.Value.per.SqFt),]
```

```
##           Neighborhood Market.Value.per.SqFt   Boro Year.Built
## 2126 LONG ISLAND CITY                 10.66 Queens       2007
```

b. Create multiple plots to demonstrates the correlations between different variables. Remember to label all axes and give title to each graph.
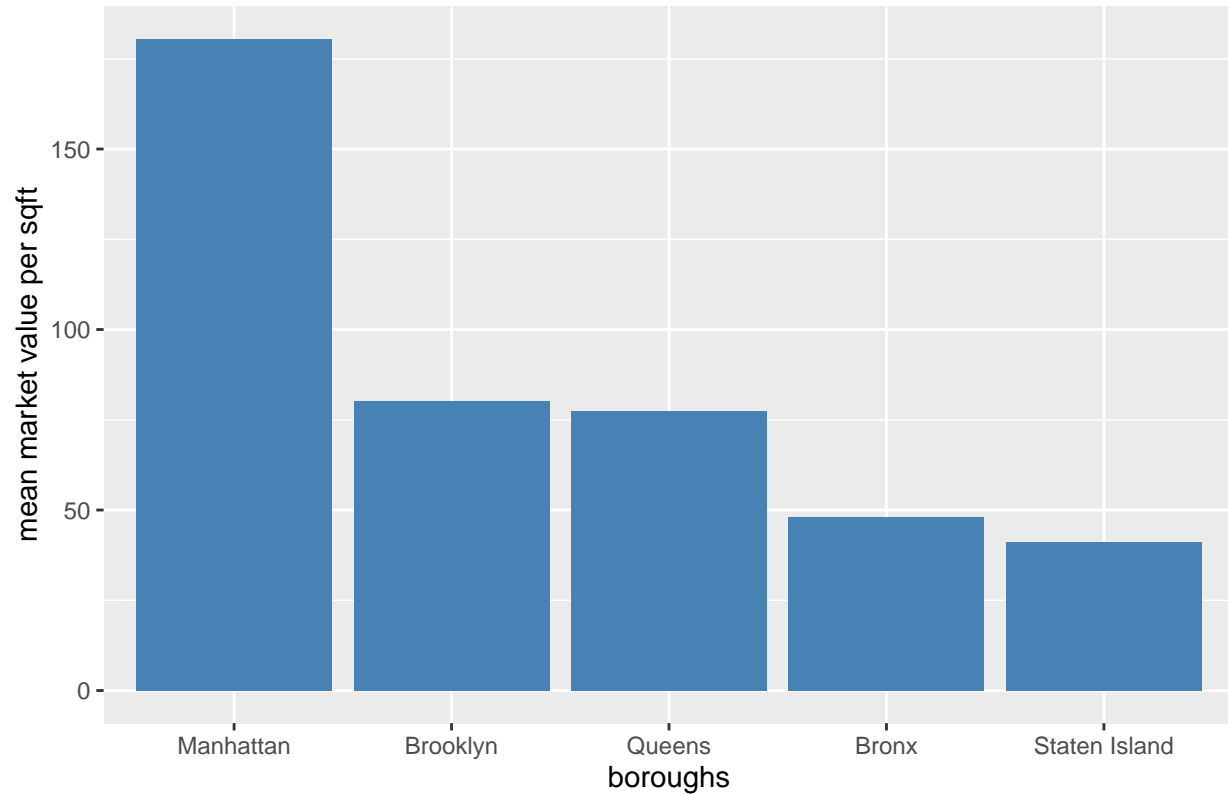
```r
# Enter your code here!

#plot 1 - Mean market value per sqft in the different boroughs of new york
ggplot(housingData, aes(x=reorder(Boro,-Market.Value.per.SqFt), y=Market.Value.per.SqFt)) +
  geom_bar(stat = "summary", fill="steelblue")+
  labs(y="mean market value per sqft", x="boroughs")+
  ggtitle("Mean market value per sqft in the different boroughs of new york")+
```

```
        theme(plot.title = element_text(hjust=0.5, face="bold"))
```
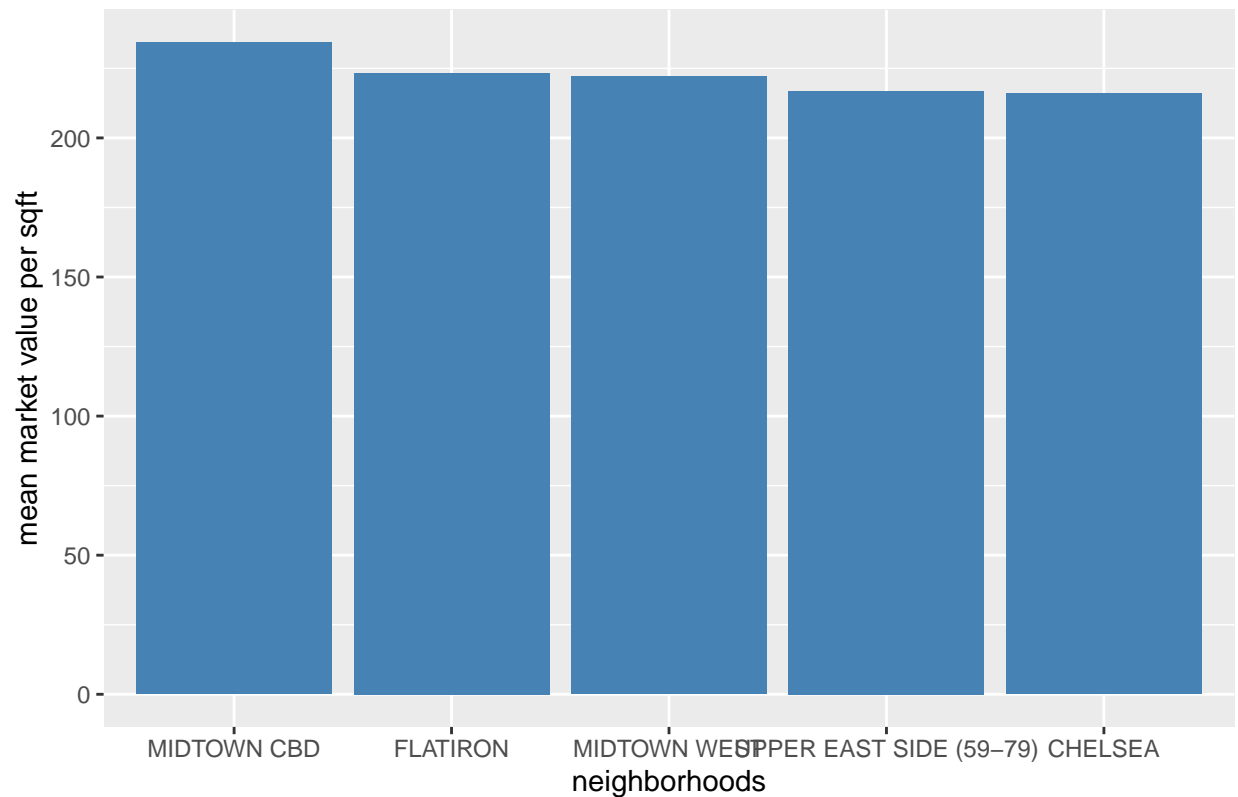
## No summary function supplied, defaulting to `mean_se()`

**Mean market value per sqft in the different boroughs of new york**
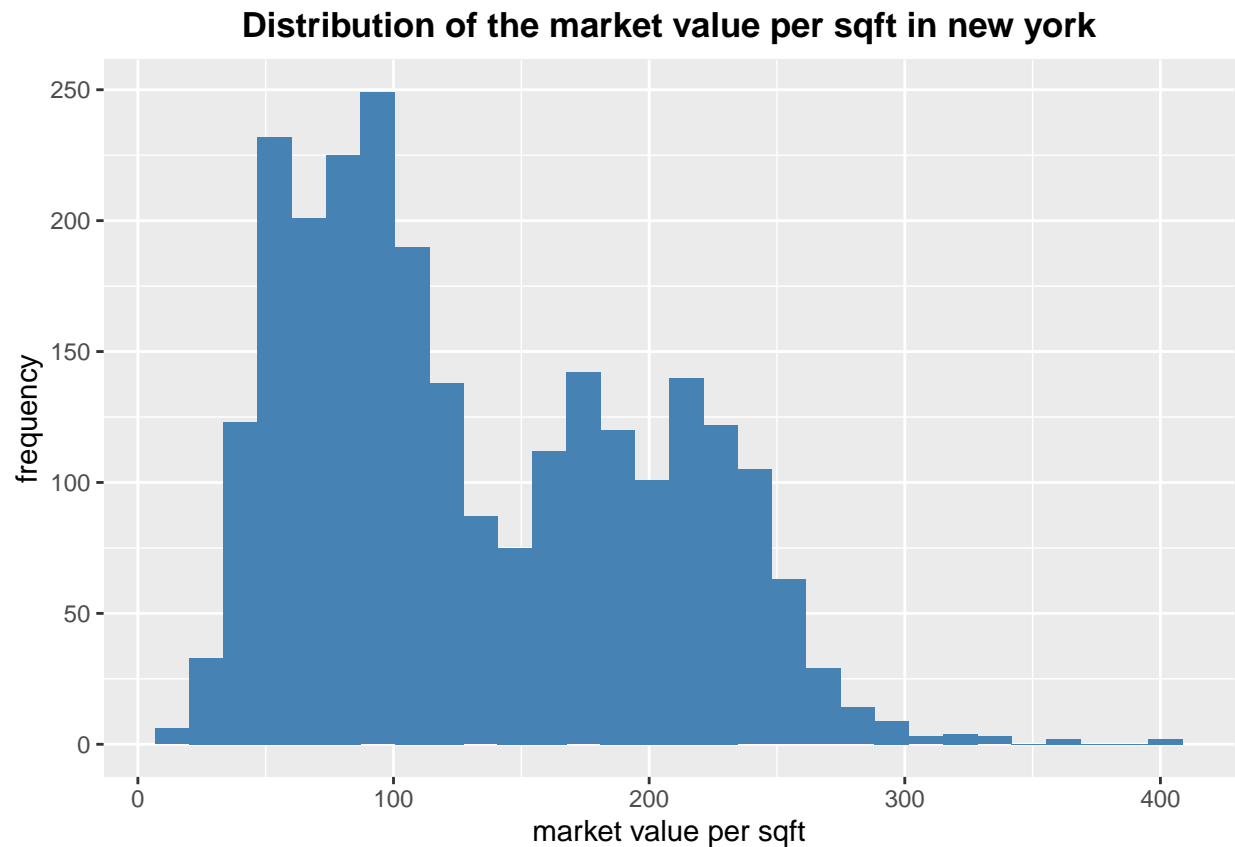


```
#plot 2 - Mean market value per sqft in the top 5 most expensive neighborhoods of new york
ggplot(top5_neighborhoods, aes(x=reorder(Neighborhood,-mean), y=mean)) +
  geom_bar(stat = "identity", fill="steelblue")+
  labs(y="mean market value per sqft", x="neighborhoods")+
  ggtitle("Mean market value per sqft in the different neighborhoods of new york")+
  theme(plot.title = element_text(hjust=0.5, face="bold"))
```

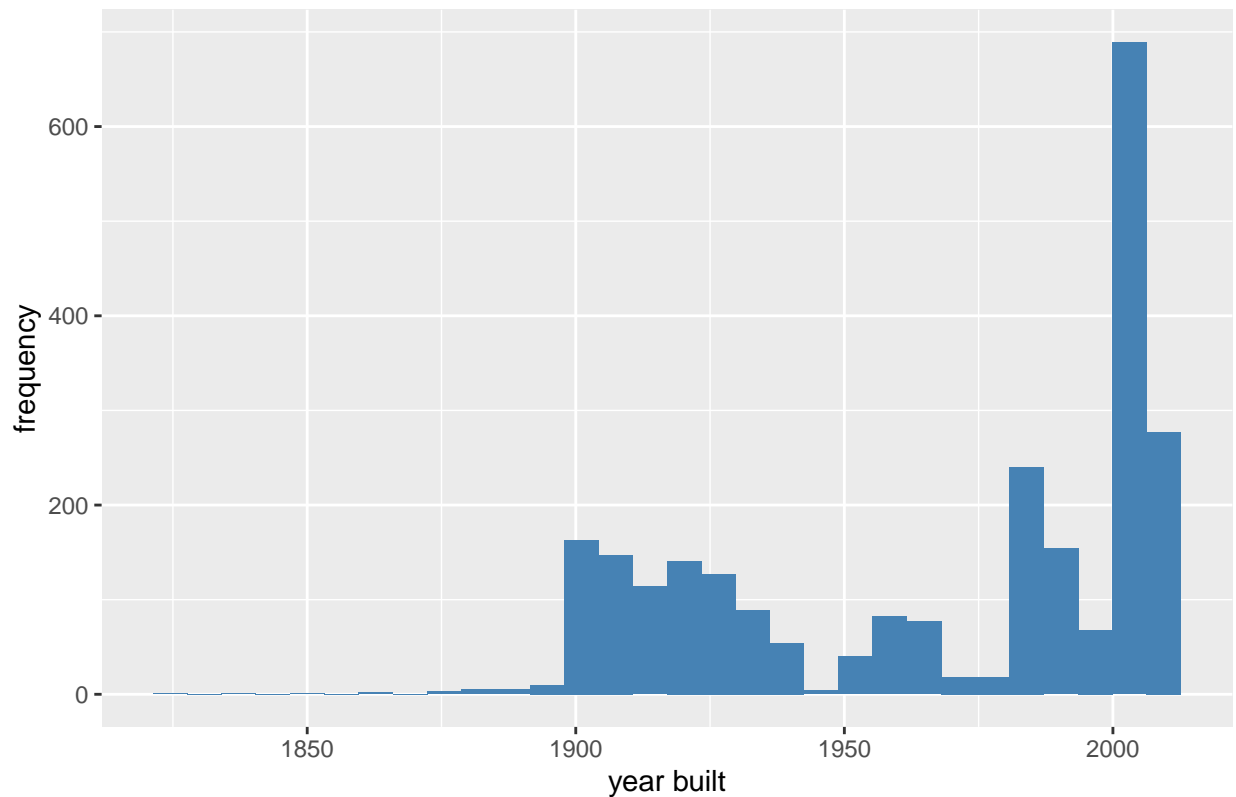**Mean market value per sqft in the different neighborhoods of new york**



```
#plot 3 - histogram distribution of the market value per sqft
ggplot(data=housingData, aes(x=Market.Value.per.SqFt))+
  geom_histogram(fill="steelblue", bins=30)+
  labs(y="frequency", x="market value per sqft")+
  ggtitle("Distribution of the market value per sqft in new york")+
  theme(plot.title = element_text(hjust=0.5, face="bold"))
```

**Distribution of the market value per sqft in new york**



```
#plot 4 - histogram distribution of the years the buildings were built
ggplot(data=housingData, aes(x=Year.Built))+
  geom_histogram(fill="steelblue", bins=30)+
  labs(y="frequency", x="year built")+
  ggtitle("Distribution of the years the buildings were built in new york")+
  theme(plot.title = element_text(hjust=0.5, face="bold"))
```

**Distribution of the years the buildings were built in new york**



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

=> Manhattan is by far the most expensive borough in New York with around 189 dollar market value per square feet. On the second place is Brooklyn with around 80 dollar market value per square feet and on the third place is Queens with around 77 dollar market value per square feet. Midtown, Flatiron, Midtown West, Upper East Side and Chelsea are the most expensive neighborhoods in New York and they are all in Manhattan. So, it makes sense that Manhattan is the most expensive borough in New York. Their average market value per square feet is with around 230 dollar way higher than the average in Manhattan. So those are the best neighborhoods in Manhattan. In general, the most buildings have a market value per square feet from around 100$ and were build around 2000.