

DS311 - R Lab Assignment

Paul Hartung

11/10/2022

R Assignment 1

- In this assignment, we are going to apply some of the built-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!
```

```
rowcount = ncol(mtcars)
```

```
columnscout = nrow(mtcars)
```

```
# Answer:
```

```
print(paste("There are total of ",rowcount," variables and "
            ,columnscout," 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!
```

```
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
## 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:
```

```
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!
```

```
average = mean(mtcars$mpg)
```

```
variance = var(mtcars$mpg)
```

```
standard_deviation = sd(mtcars$mpg)
```

```
print(paste("The average of Mile Per Gallon from this data set is ",round(average,4), " with variance
```

```
## [1] "The average of Mile Per Gallon from this data set is 20.0906 with variance 36.3241 and stan
```

- 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$gear=as.character(mtcars$gear)
```

```
mtcars$cyl=as.character(mtcars$cyl)
```

```
mtcars %>%
```

```
  group_by(cyl)%>%
```

```
  summarise(average=mean(mpg))
```

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

```
##   cyl   average
##   <chr>   <dbl>
## 1 4       26.7
## 2 6       19.7
## 3 8       15.1
```

```
mtcars %>%
  group_by(gear)%>%
  summarise(standard_deviation=sd(mpg))
```

```
## # A tibble: 3 x 2
##   gear standard_deviation
##   <chr>             <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)%>%
  summarise(number_of_cylinders=n(), .groups = 'drop')%>%
  arrange(desc(number_of_cylinders))
```

```
## # A tibble: 8 x 3
##   cyl   gear number_of_cylinders
##   <chr> <chr>             <int>
## 1 8     3                 12
## 2 4     4                 8
## 3 6     4                 4
## 4 4     5                 2
## 5 6     3                 2
## 6 8     5                 2
## 7 4     3                 1
## 8 6     5                 1
```

```
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.

- 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.

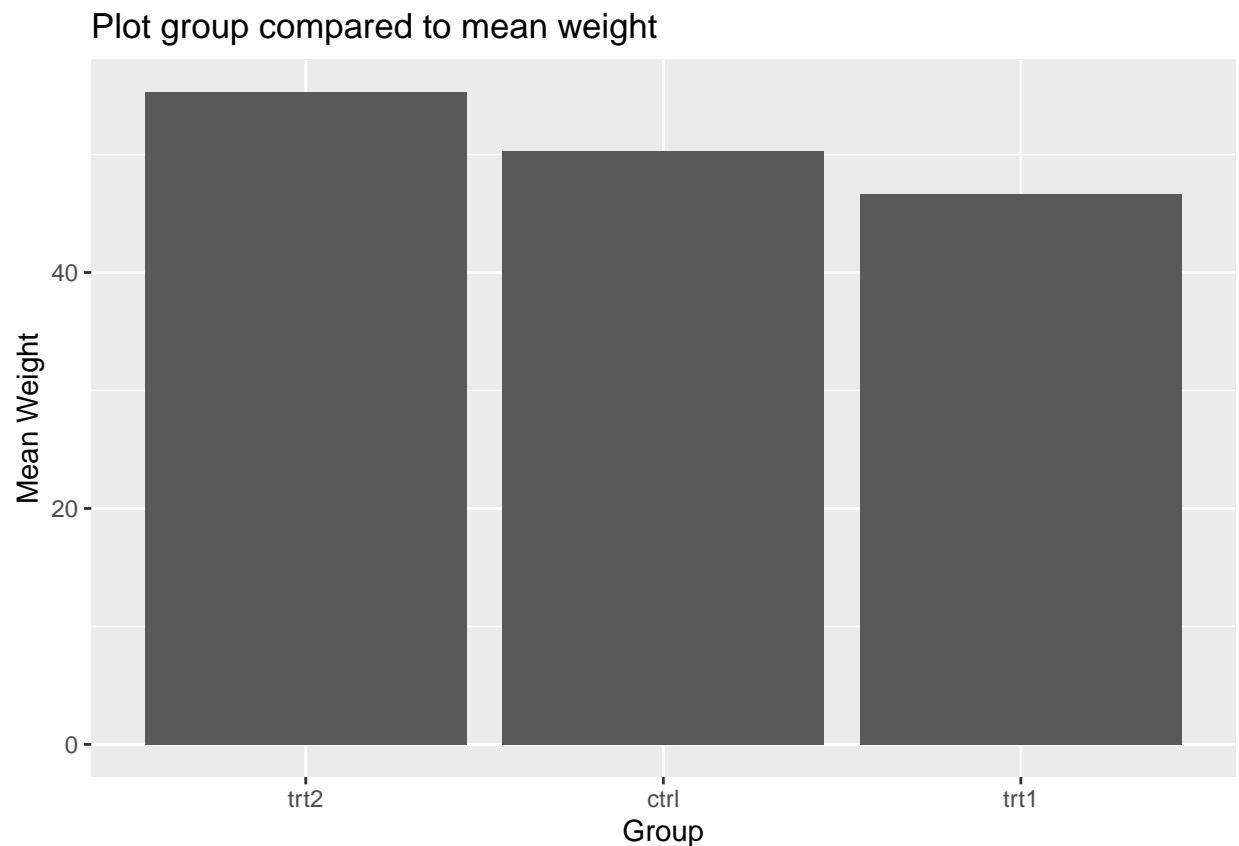
```
# Load the data set
data("PlantGrowth")

# 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!

```
PlantGrowth%>%
  ggplot(aes(x=reorder(group,-weight), y=weight))+
  geom_bar(stat="identity")+
  labs(title="Plot group compared to mean weight")+
  xlab("Group")+
  ylab("Mean Weight")
```



Result:

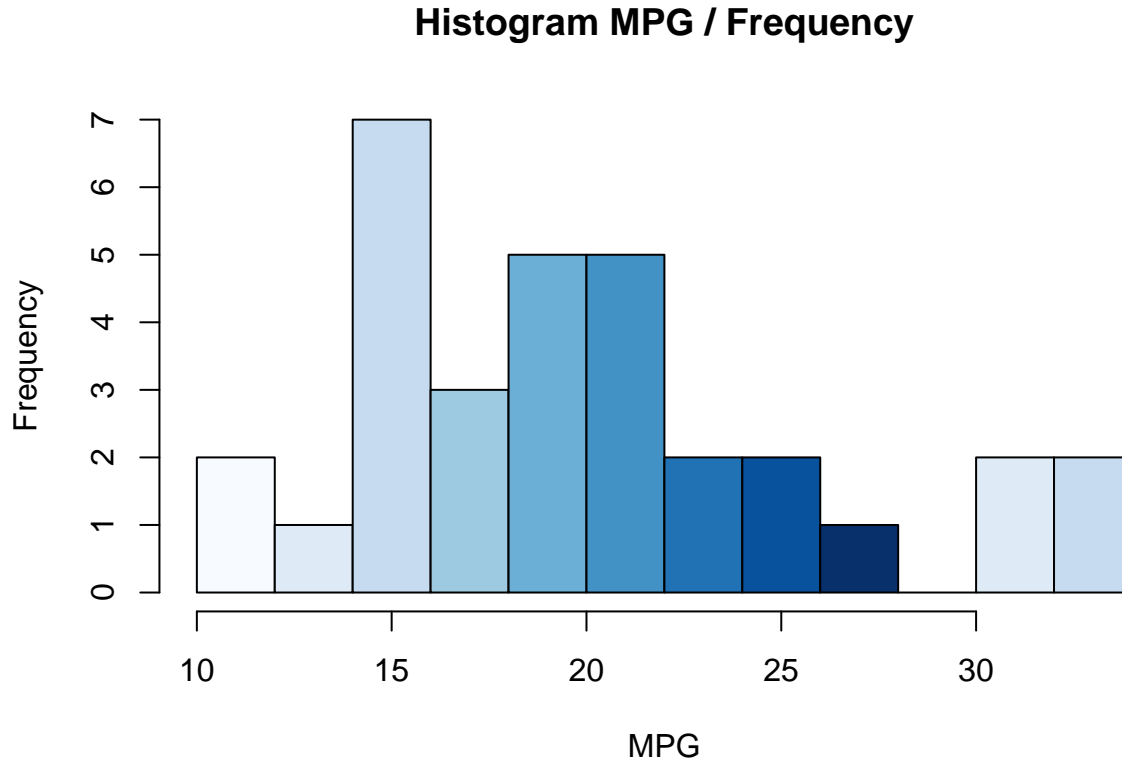
=> Report a paragraph to summarize your findings from the plot!

Looking at my chart, I can see the following: Plants in group “trt2” are the heaviest at around 55. Plants in group “ctrl” are the medium heavy champions with about 50. Plants in group “trt1” are the lightest at just over 45.

All measured in average 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,main='Histogram MPG / Frequency',xlab='MPG',breaks=10,ylab='Frequency', col = brewer.pa
```



```
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."
```

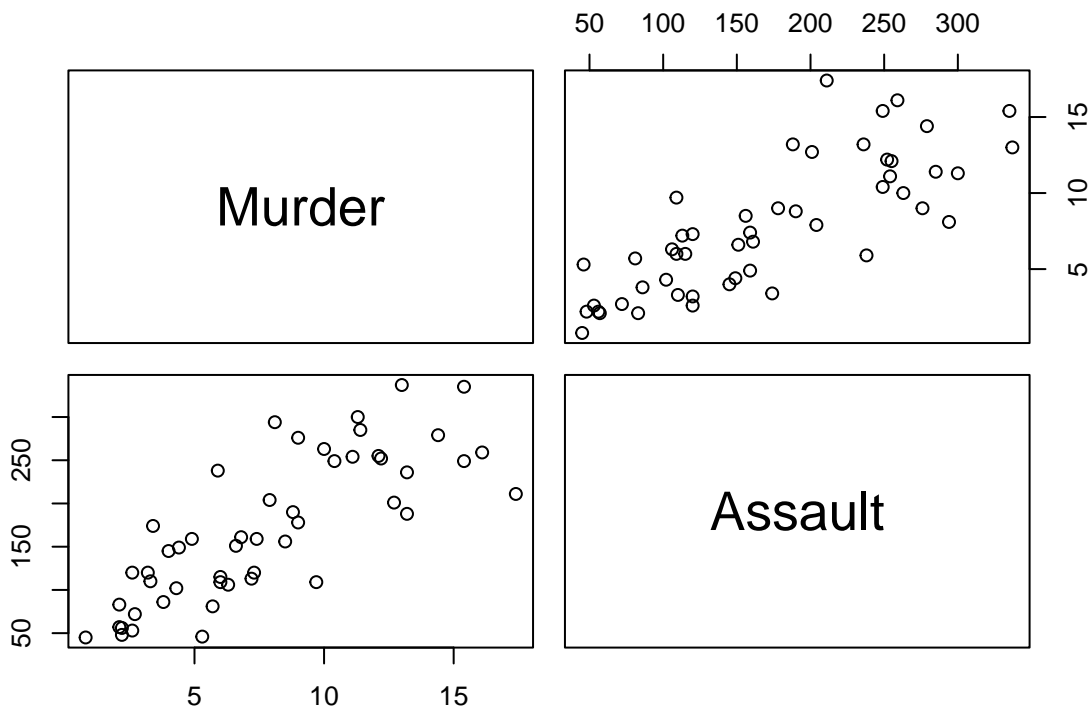
- c. Using the **USArrests** data set, create a pairs plot to display the correlations between the variables in the data set. Plot the scatter plot with **Murder** and **Assault**. Give labels to the title, x-axis, and y-axis on the graph. Write a paragraph to summarize your results from both plots.

```
# 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!
pairs(~ Murder + Assault, data = USArrests)
```



Result:

=> Report a paragraph to summarize your findings from the plot!

Looking at my spreadsheet, I can see the following: The more people arrested for assault, the more people arrested for murder. The same is true the other way around. This shows that assault and murder cases are positive correlated.

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.

- 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!
housingData%>%
  group_by(Neighborhood)%>%
  summarise(avg_price=mean(Market.Value.per.SqFt))%>%
  arrange(desc(avg_price))
```

```
## # A tibble: 148 x 2
##   Neighborhood      avg_price
##   <chr>            <dbl>
## 1 MIDTOWN CBD      234.
## 2 FLATIRON        223.
## 3 MIDTOWN WEST    222.
## 4 UPPER EAST SIDE (59-79) 217.
## 5 CHELSEA         216.
## 6 MIDTOWN EAST    211.
## 7 EAST VILLAGE    207.
## 8 MURRAY HILL     206.
## 9 UPPER EAST SIDE (79-96) 202.
## 10 GREENWICH VILLAGE-WEST 202.
## # ... with 138 more rows
```

```
housingData%>%
  group_by(Boro)%>%
  summarise(avg_price=mean(Market.Value.per.SqFt))%>%
  arrange(desc(avg_price))
```

```
## # A tibble: 5 x 2
##   Boro      avg_price
##   <chr>      <dbl>
## 1 Manhattan  181.
## 2 Brooklyn  80.1
## 3 Queens    77.4
## 4 Bronx     47.9
## 5 Staten Island 41.3
```

```
housingData%>%
  group_by(Year.Built)%>%
  summarise(avg_price=mean(Market.Value.per.SqFt))%>%
  arrange(desc(avg_price))
```

```
## # A tibble: 124 x 2
##   Year.Built avg_price
##   <int>      <dbl>
## 1 1836      274.
## 2 1978      255.
## 3 1970      215.
## 4 1934      204.
## 5 1975      201.
## 6 1973      197.
## 7 1879      195.
## 8 1905      188.
## 9 1948      186.
## 10 1972      186.
## # ... with 114 more rows
```

b. Create multiple plots to demonstrates the correlations between different variables. Remember to label

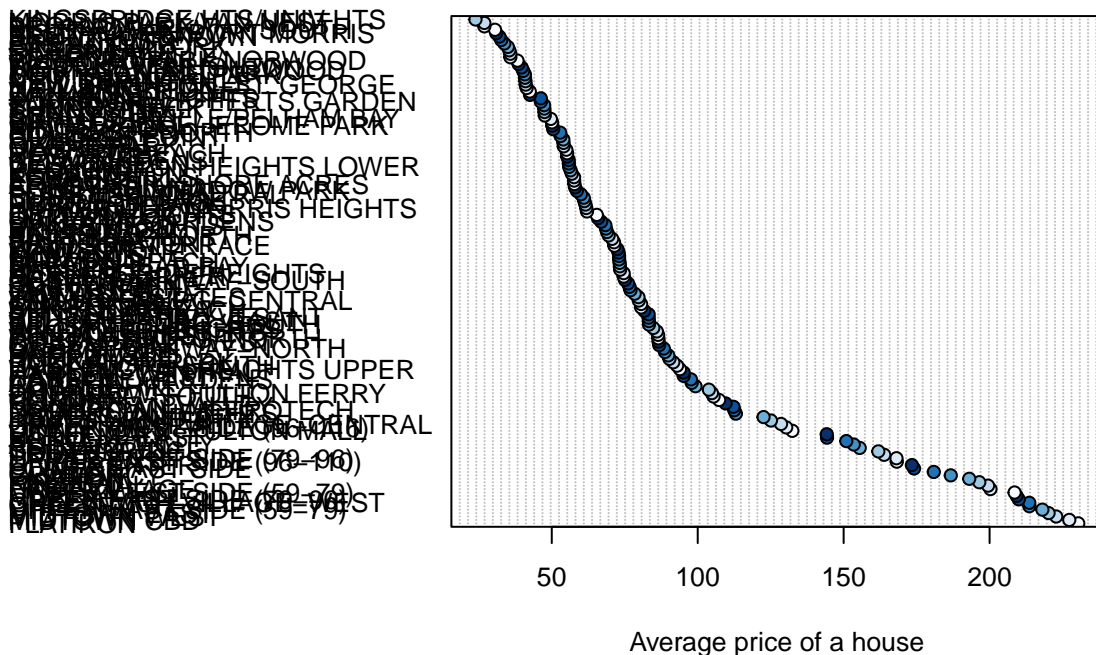
all axes and give title to each graph.

```
# Enter your code here!
neighbourhoods = tapply(housingData$Market.Value.per.SqFt, housingData$Neighborhood, median)
neighbourhoods = sort(neighbourhoods, decreasing = TRUE)

dotchart(neighbourhoods, pch = 21, bg = brewer.pal(9, "Blues"),
         cex = 0.85,
         xlab="Average price of a house",
         main = "Which neighborhood is the most expensive to buy a house in?")

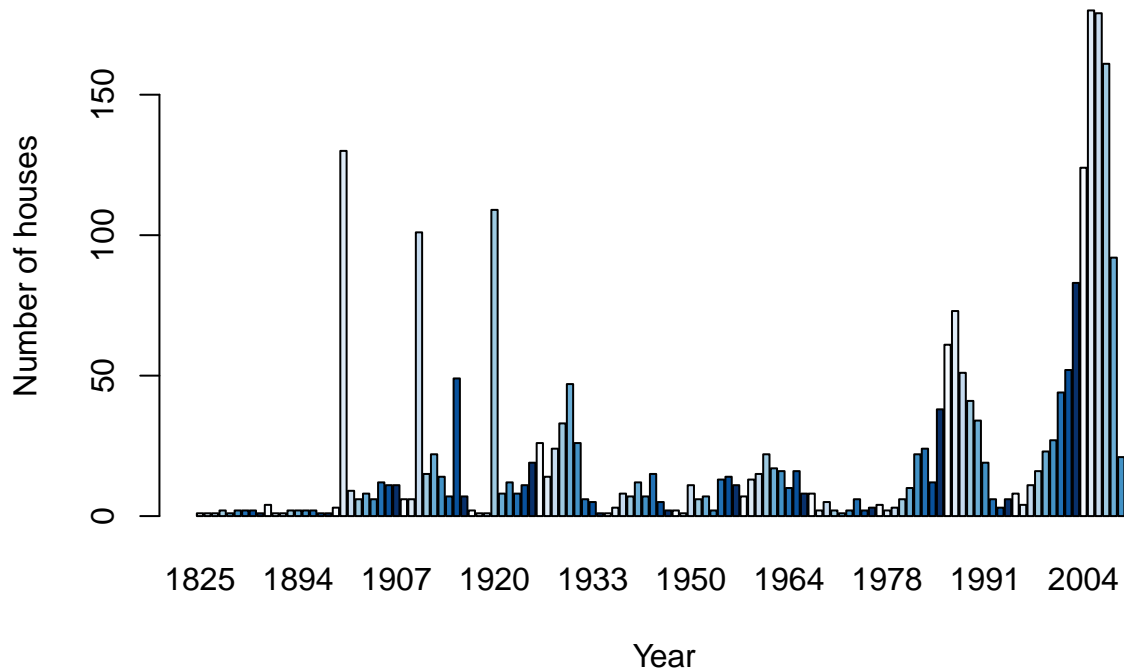
## Warning in dotchart(neighbourhoods, pch = 21, bg = brewer.pal(9, "Blues"), : 'x'
## is neither a vector nor a matrix: using as.numeric(x)
```

Which neighborhood is the most expensive to buy a house



```
barplot(table(housingData$Year.Built),
       main = "When were the most houses built?",
       xlab = "Year",
       ylab = "Number of houses",
       col = brewer.pal(9, "Blues"))
```


When were the most houses built?



=> Enter your answer here!

Looking at my spreadsheet, I can see the following: MIDTOWN CBD has highest avg price per sqft with 234.36. *Manhattan highest avg price per sqft with 180.59.* 1978 & 1970 & 1973 & 1975 are in the Top 6 years with highest avg price per sqft. Plus minus 2 years around 2007, the most houses were built.