

# DS311 - R Lab Assignment

Norman Lo

2/20/2023

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

```
# Answer:  
print(paste("There are total of ", obs, " variables and ", var, " 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:
print("There are 0 discrete variables and 11 continuous variables in this data set.")
```

```
## [1] "There are 0 discrete variables and 11 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 ", round(m,2) , " with variance ", round(v,2) , " and standard deviation ", round(s,2)))
```

```
## [1] "The average of Mile Per Gallon from this data set is 20.09 with variance 36.32 and standard deviation 6.024"
```

- 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!
d1 <- setNames(aggregate(mpg~cyl, data=mtcars, FUN=mean), c('Cylinder', 'MPG_Avg'))
d2 <- setNames(aggregate(mpg~gear, data=mtcars, FUN=sd), c('Gear', 'MPG_Std'))

print(d1)
```

```
##   Cylinder  MPG_Avg
## 1         4 26.66364
## 2         6 19.74286
## 3         8 15.10000
```

```
print(d2)
```

```
##   Gear  MPG_Std
## 1    3 3.371618
## 2    4 5.276764
## 3    5 6.658979
```

- 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!
e <- xtabs(~gear+cyl, data=mtcars)
print(e)
```

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

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

```
## [1] "The most common car type in this data set is a 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 in this graph.

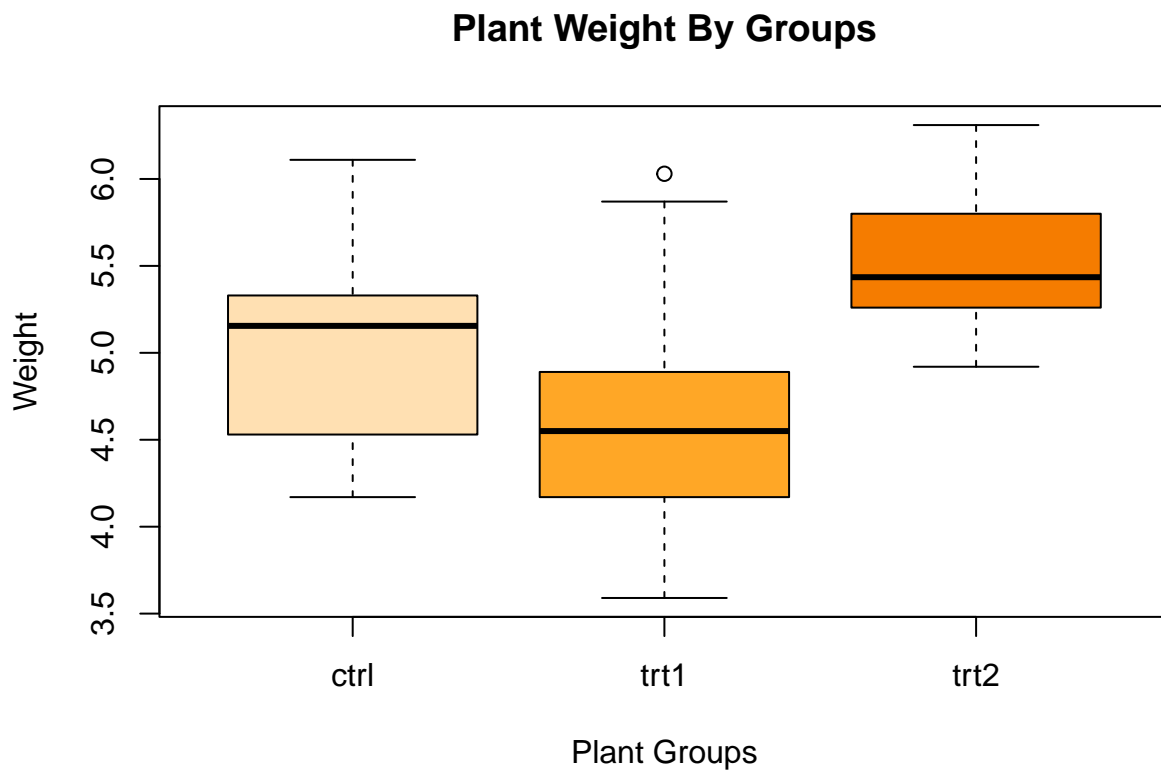
```
# Load the data set
data("PlantGrowth")

# 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(PlantGrowth$weight ~ PlantGrowth$group,
        col = c("#FFE0B2", "#FFA726", "#F57C00"),
        main = "Plant Weight By Groups",
        xlab = "Plant Groups",
        ylab = "Weight")
```



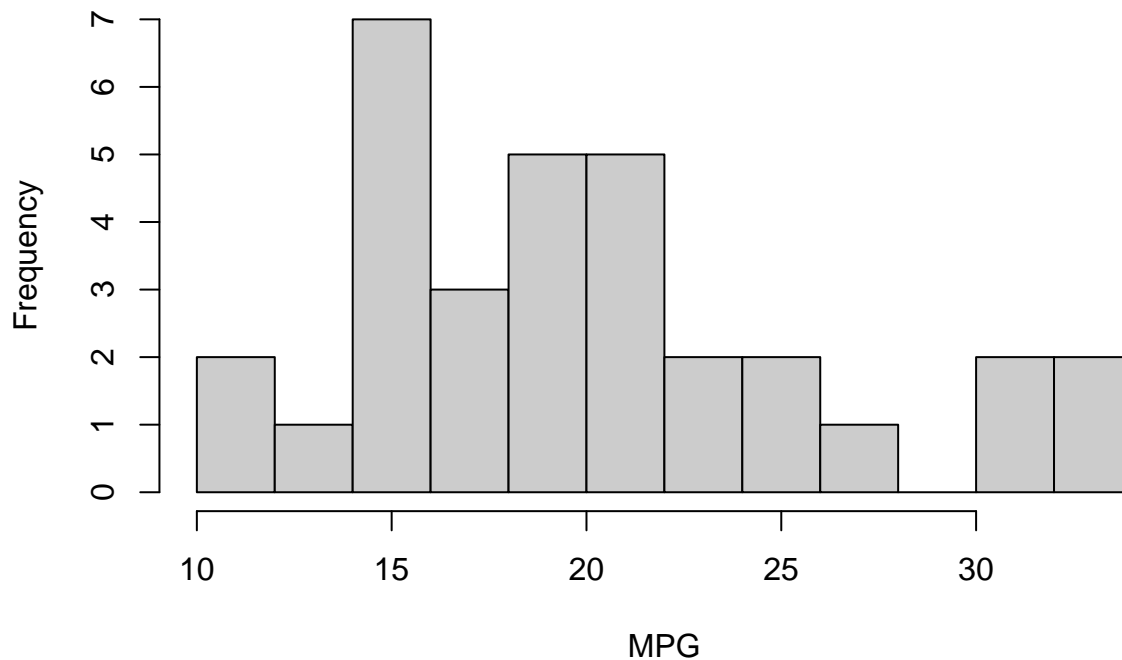
Result:

=> Enter your results here!

- 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,
     col="grey80",
     main="MPG Histogram",
     xlab="MPG",
     ylab="Frequency")
```

## MPG Histogram



```
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."
```

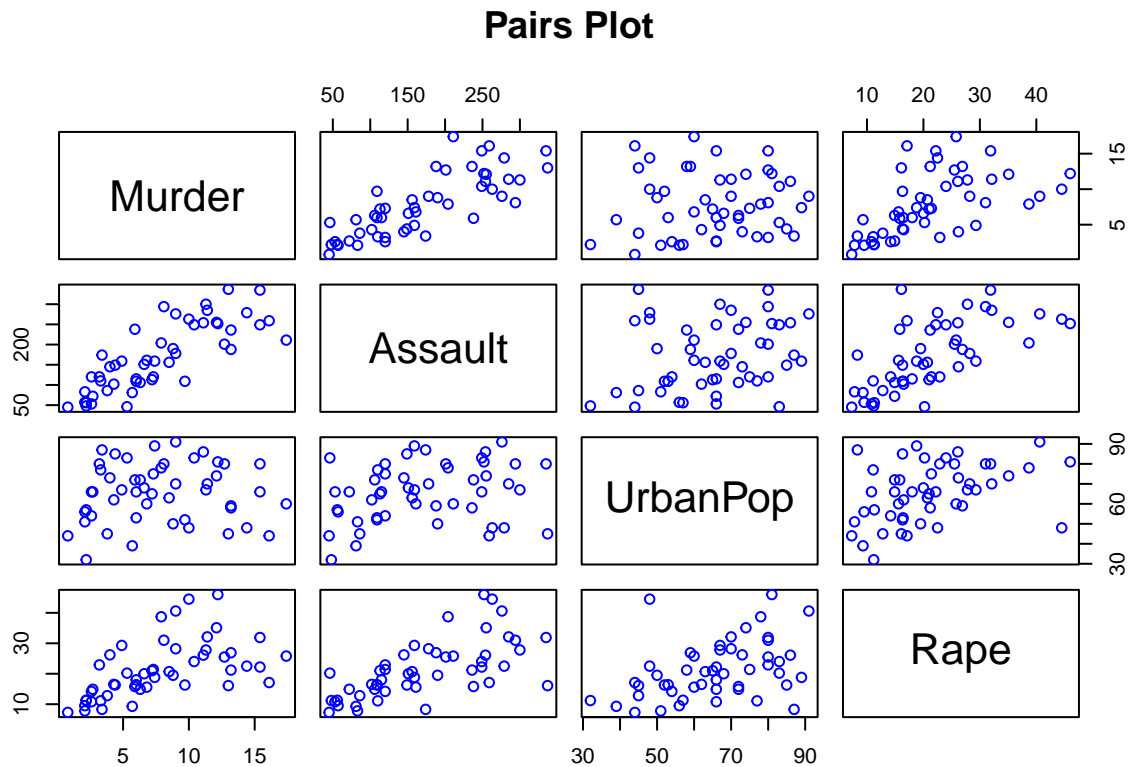
- c. Using the **USArrests** data set, create a pairs plot to display the correlations between the variables in the data set. Plot the scatter plot graph of **Murder** and **Assault**. Give labels to the title, x-axis, and y-axis on the graph. Write a paragraph to summarize your results from both plots.

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

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

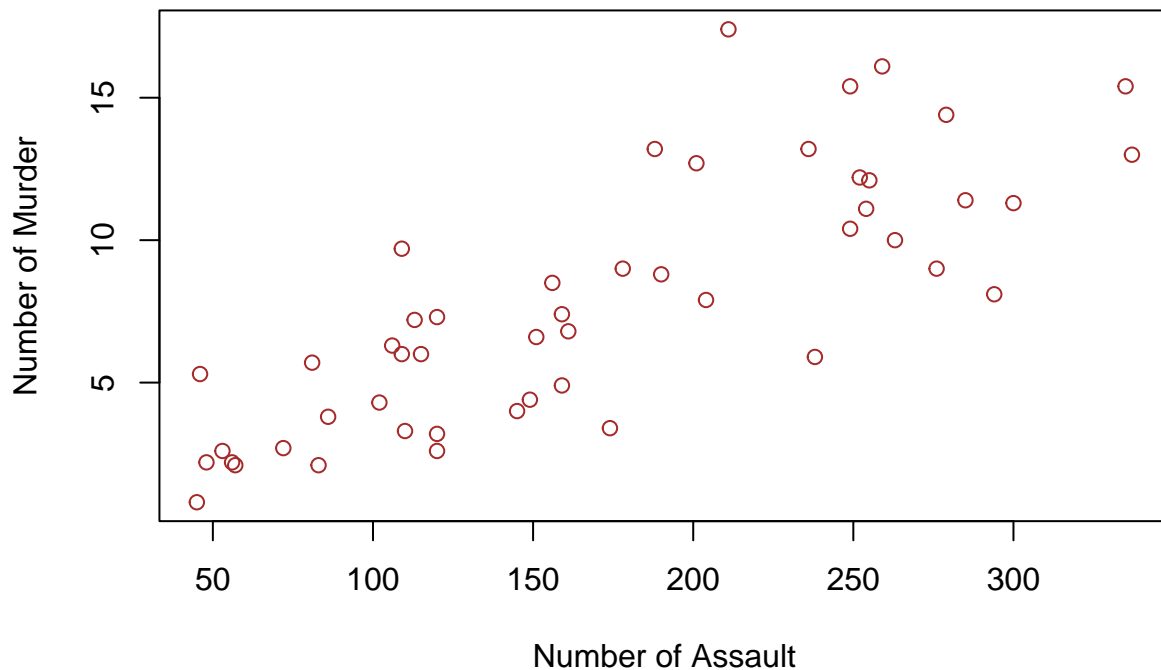
```
##           Murder Assault UrbanPop Rape
## Alabama      13.2    236      58 21.2
## Alaska       10.0    263      48 44.5
## Arizona       8.1    294      80 31.0
## Arkansas      8.8    190      50 19.5
## California    9.0    276      91 40.6
## Colorado     7.9    204      78 38.7
```

```
# Enter your code here!
pairs(USArrests,
      col="blue",
      main="Pairs Plot")
```



```
plot(USArrests$Assault, USArrests$Murder,
      col="brown",
      main="Correlation between Assault and Murder",
      xlab="Number of Assault",
      ylab="Number of Murder")
```

## Correlation between Assault and Murder



Result:

=> Enter your result here!

---

### 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.

- 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!
```

- 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!
```

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

Enter your answer here!