

DS311 - R Lab Assignment

Barbara Wallen

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!
```

```
ncol(mtcars)
```

```
## [1] 11
```

```
nrow(mtcars)
```

```
## [1] 32
```

```
# Answer:
```

```
print("There are total of __32__ variables and __32__ observations in this data set.")
```

```
## [1] "There are total of __32__ 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
```

Answer:

```
print("There are __6__ discrete variables and __5__ continuous variables in
this data set.")
```

```
## [1] "There are __6__ discrete variables and __5__ 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 <- sqrt(v)
```

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

```
## [1] "The average of Mile Per Gallon from this data set is 20.090625 with
variance 36.3241028225806 and standard deviation 6.0269480520891 ."
```

- 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!

```
library(dplyr)

##
## 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

library(tidyr)

mtcars %>%
  group_by(cyl, gear) %>%
  tally() %>%
  spread(cyl, n)

## # A tibble: 3 × 4
##   gear   `4`   `6`   `8`
##   <dbl> <int> <int> <int>
## 1     3     1     2    12
## 2     4     8     4     NA
## 3     5     2     1     2

avgclass = mtcars %>% group_by(mtcars$cyl) %>% summarise(mean_mpg =
mean(mtcars$mpg))
print(avgclass)

## # A tibble: 3 × 2
##   `mtcars$cyl` mean_mpg
##   <dbl>      <dbl>
## 1         4      20.1
## 2         6      20.1
## 3         8      20.1
```

- 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!

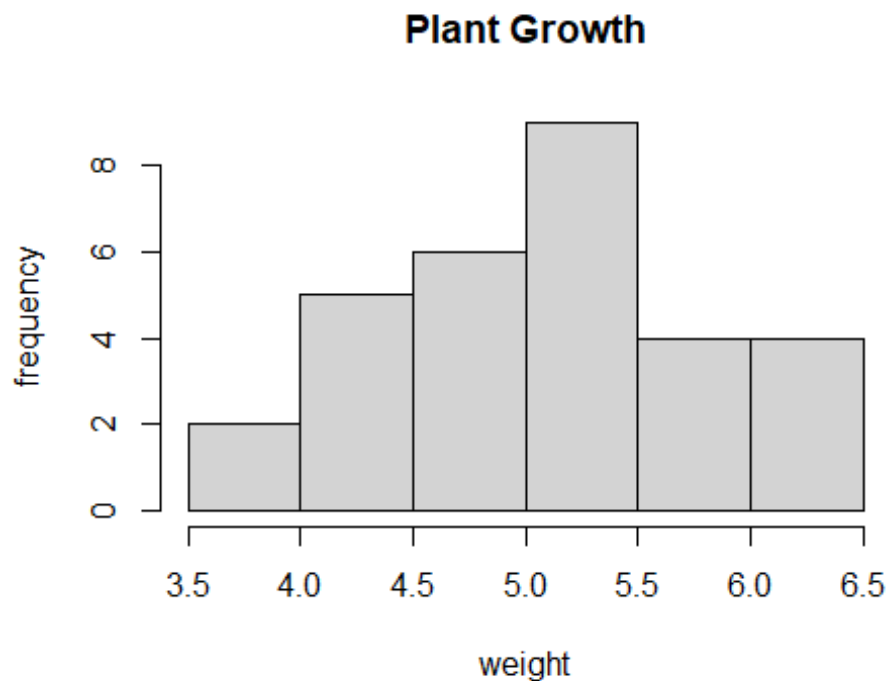
```
print("The most common car type in this data set is car with ____ cylinders  
and ____ gears. There are total of ____ cars belong to this specification in  
the data set.")  
  
## [1] "The most common car type in this data set is car with ____ cylinders  
and ____ gears. There are total of ____ cars belong to this specification in  
the data set."
```

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!  
hist <- PlantGrowth$weight  
hist(hist,  
      main = 'Plant Growth',  
      xlab = 'weight',  
      ylab = 'frequency' )
```



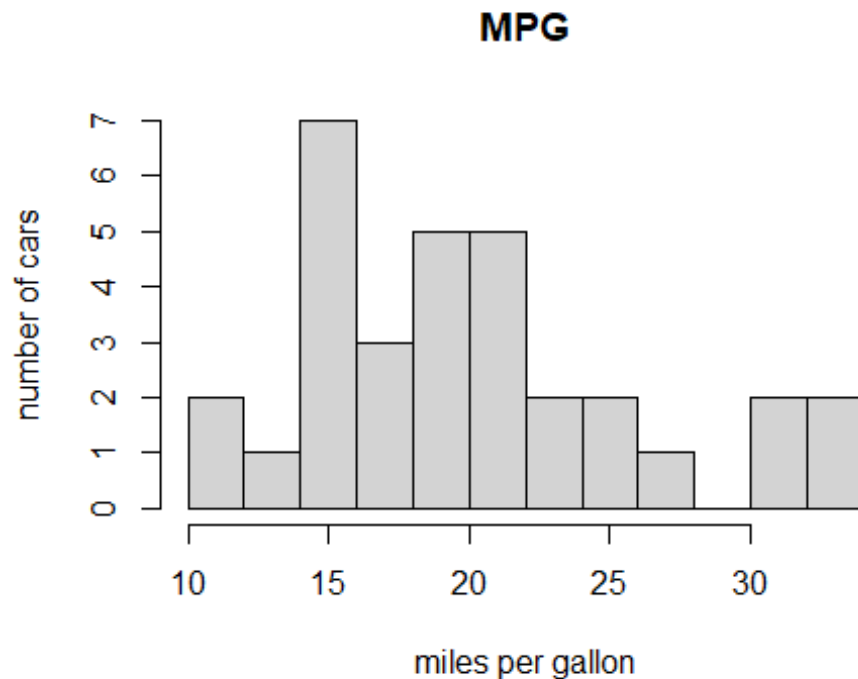
Result:

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

There is 8 plants with the weight between 5 to 5.5 and 2 plants with weight between 3.5 to 4.0

- 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.

```
miles <- mtcars$mpg
hist(miles,
     breaks = 10,
     main = 'MPG',
     xlab = 'miles per gallon',
     ylab = 'number of cars' )
```



```
print("Most of the cars in this data set are in the class of __7__ mile per gallon.")
```

```
## [1] "Most of the cars in this data set are in the class of __7__ 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.

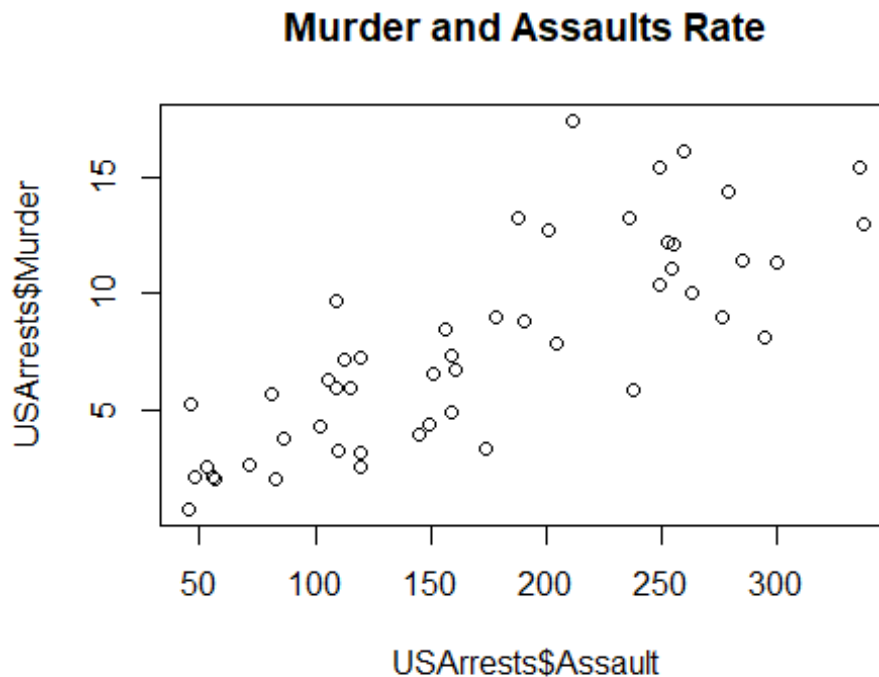
```
# Load the data set
data("USArrests")
```

```
# 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!
library(ggplot2)
```

```
plot(y= USArrests$Murder, x= USArrests$Assault, main = "Murder and Assaults Rate")
```



Result:

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

The ratio between murders is lower than assaults, meaning that assaults has a high rate

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!

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" ...

min(housingData$Market.Value.per.SqFt)

## [1] 10.66

max(housingData$Market.Value.per.SqFt)

## [1] 399.38

mean(housingData$Market.Value.per.SqFt)

## [1] 133.1731

min(housingData$Year.Built)

## [1] 1825

max(housingData$Year.Built)

## [1] 2010

mean(housingData$Year.Built)

## [1] 1967.46
```

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

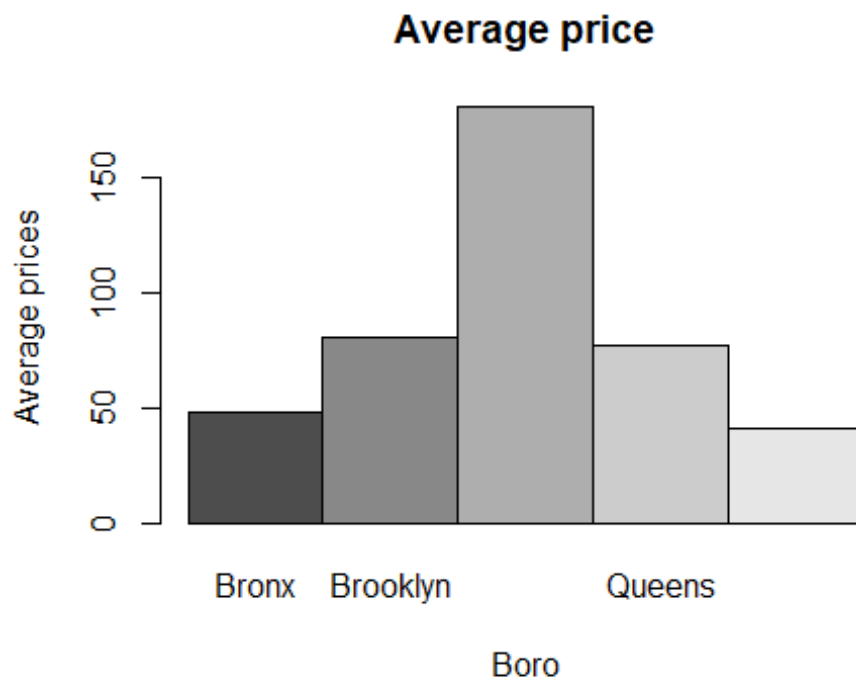
```
# Enter your code here!
avg <- aggregate(housingData$Market.Value.per.SqFt, list(housingData$Boro),
FUN=mean)
avg

##      Group.1      x
## 1      Bronx 47.93232
```



```
## 2      Brooklyn  80.13439
## 3      Manhattan 180.59265
## 4       Queens   77.38137
## 5  Staten Island  41.26958

barplot(matrix(avg$x), beside=T, main='Average price', names.arg =
avg$Group.1, ylab="Average prices", xlab="Boro")
```



- c. Write a summary about your findings from this exercise. The houses from Manhattan is the most expensive (180k) from this data set, and the cheapest house (\$41k) is found on State Island neighborhood.

The oldest house is from 1825.

=> Enter your answer here!