

# Aradhyा Mathur Homework0

Aradhyा Mathur

2022-09-05

```
library(readr)
library(ggplot2)
library(moments)
```

Question 1. Getting familiar with the dataset via exploratory data analysis  
Read the data into RStudio and summarize the data with the `summary()` function.

```
data1 <- read_csv("car_sales.csv")

## # Rows: 152 Columns: 11
## -- Column specification -----
## Delimiter: ","
## chr (2): Manufacturer, Model
## dbl (9): price, Engine_size, Horsepower, Wheelbase, Width, Length, Curb_weig...
##
## i Use `spec()` to retrieve the full column specification for this data.
## i Specify the column types or set `show_col_types = FALSE` to quiet this message.

summary(data1)
```

```
##   Manufacturer      Model       price     Engine_size
##   Length:152      Length:152    Min.   : 9235    Min.   :1.000
##   Class :character Class :character  1st Qu.:17889   1st Qu.:2.300
##   Mode  :character Mode  :character  Median :22747    Median :3.000
##                               Mean   :27332    Mean   :3.049
##                               3rd Qu.:31939   3rd Qu.:3.575
##                               Max.  :85500    Max.  :8.000
##   Horsepower      Wheelbase      Width       Length
##   Min.   : 55.0    Min.   : 92.6    Min.   :62.60    Min.   :149.4
##   1st Qu.:147.5   1st Qu.:102.9   1st Qu.:68.38   1st Qu.:177.5
##   Median :175.0   Median :107.0   Median :70.40   Median :186.7
##   Mean   :184.8   Mean   :107.4   Mean   :71.09   Mean   :187.1
##   3rd Qu.:211.2   3rd Qu.:112.2   3rd Qu.:73.10   3rd Qu.:195.1
##   Max.  :450.0   Max.  :138.7   Max.  :79.90   Max.  :224.5
##   Curb_weight     Fuel_capacity  Fuel_efficiency
##   Min.   :1.895   Min.   :10.30   Min.   :15.00
##   1st Qu.:2.965   1st Qu.:15.78   1st Qu.:21.00
##   Median :3.336   Median :17.20   Median :24.00
##   Mean   :3.376   Mean   :17.96   Mean   :23.84
##   3rd Qu.:3.821   3rd Qu.:19.80   3rd Qu.:26.00
##   Max.  :5.572   Max.  :32.00   Max.  :45.00
```

```
data2 <- data1[c('price')]
Price <- as.numeric(data2$price)
```

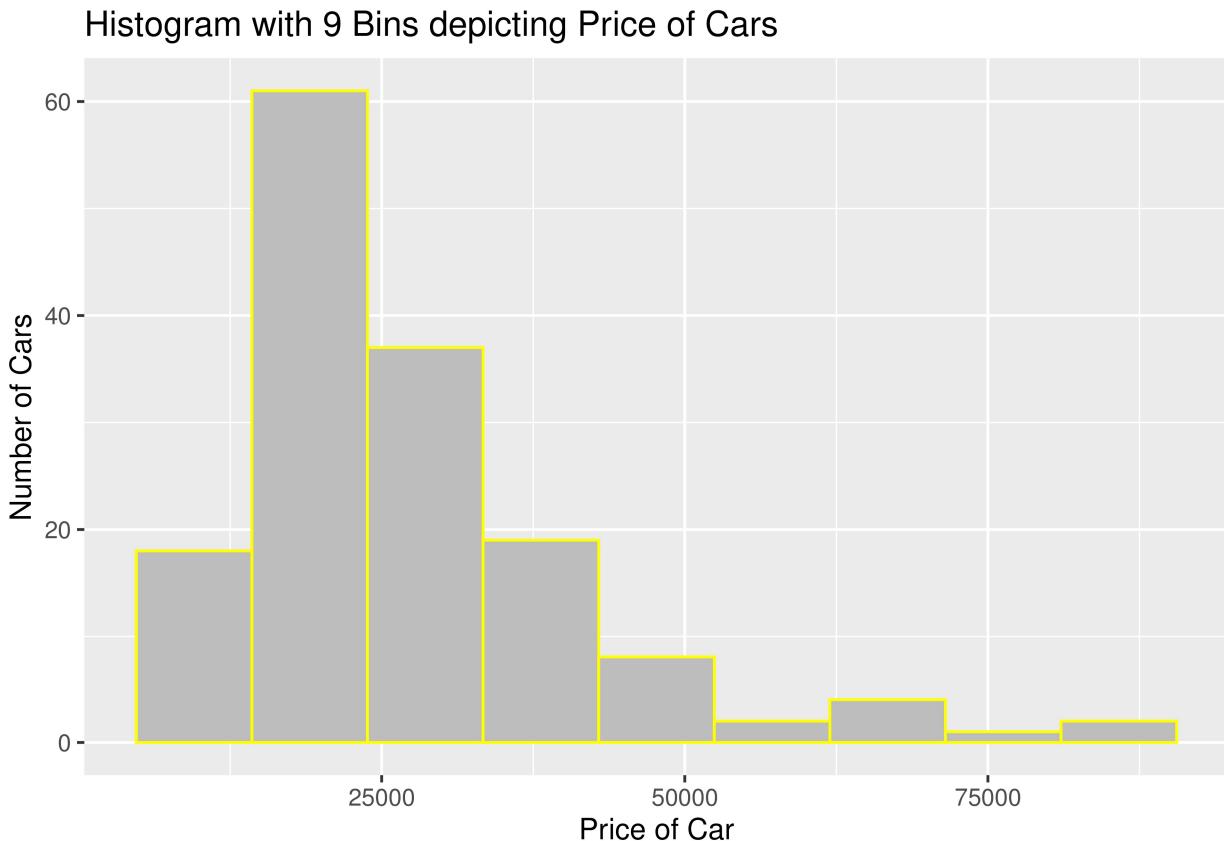
b. How many bins does Sturges' formula suggest we use for a histogram of `price`? Show your work

```
n <- 152 # 152 observations are seen in data2
bins <- ceiling(log2(n)) + 1 # Using Sturges' Formula
bins <- ceiling(log2(152)) + 1 #Substituting value of n
bins <- ceiling(7.24792751344) + 1 # log2(152) = 7.24792751344
bins <- 8 + 1 # ceiling(7.24792751344) = 8
bins
```

```
## [1] 9
```

c. Create a histogram of `price` using the number of bins suggested by Sturges' formula in 1b. Make sure to appropriately title the histogram and label the axes. Comment on the center, shape, and spread.

```
histo <- ggplot(data1,aes(x=price)) + geom_histogram(bins = bins,color="yellow",
                                                       fill="gray") +
  ggtitle("Histogram with 9 Bins depicting Price of Cars") +
  labs(y= "Number of Cars", x = "Price of Car")
histo
```



This histogram is positive skewed (right skew), unimodal, asymmetric and wide spread. Median should be used to find the center, approximate median would be around 22000.

2. Measures of center and spread for the selling price of cars.

- a. Calculate the mean, median, and 10% trimmed mean of the selling price. Report the mean, median, and 10% trimmed mean on the histogram. In particular, create a red vertical line on the histogram at the mean, and report the value of the mean in red next to the line using the form “ $\bar{x} =$ ”. Create a blue vertical line on the histogram at the median, and report the value of the median in blue next to the line using the form “ $\tilde{x} =$ ”. Create a green vertical line on the histogram at the 10% trimmed mean, and report the value of the 10% trimmed mean in green next to the line using the form “ $\bar{x}_{10} =$ ” (to get  $\bar{x}_{10}$  to print on the plot, use `bar(x)[10]` within the `paste()` function).

```

mean <- mean(data2$price)
mean

## [1] 27331.82

#is the mean
median <- median(data2$price)
median

## [1] 22747

#is the median
trim <- mean(data2$price, trim=0.1)
trim

## [1] 25077.88

#is the trimmed mean

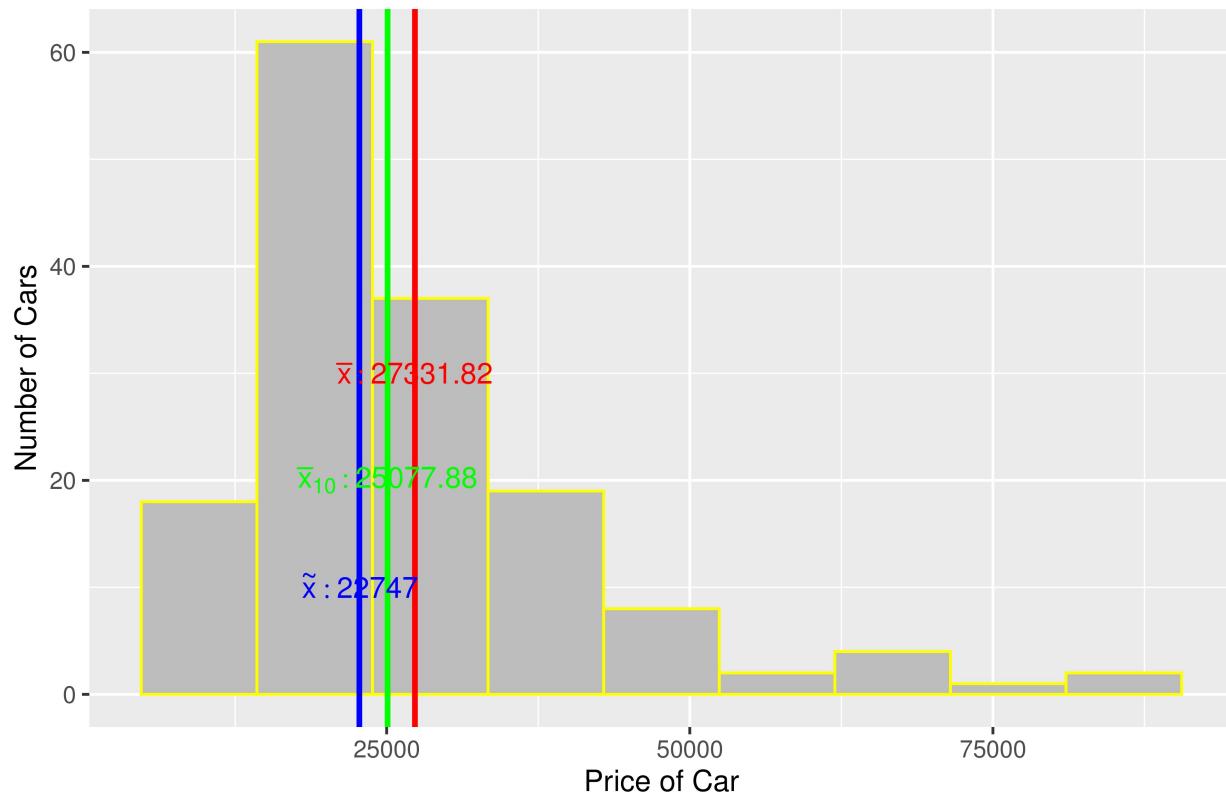
data1 <- read_csv("car_sales.csv")

## Rows: 152 Columns: 11
## -- Column specification -----
## Delimiter: ","
## chr (2): Manufacturer, Model
## dbl (9): price, Engine_size, Horsepower, Wheelbase, Width, Length, Curb_weig...
##
## i Use `spec()` to retrieve the full column specification for this data.
## i Specify the column types or set `show_col_types = FALSE` to quiet this message.

histo + geom_vline(aes(xintercept=mean(price)), color="red", size=1) +
  geom_vline(aes(xintercept=median(price)), color="blue", size=1) +
  geom_vline(aes(xintercept=mean(price,trim=0.1)), color="green", size=1) +
  annotate(geom = "text", x = mean, y = 30, parse = TRUE, label = paste("bar(x) :", mean), size = 4, col =

```

## Histogram with 9 Bins depicting Price of Cars



b. Calculate and report the 25th and 75th percentiles.

```
quantile <- quantile(data2$price, probs = c(.25, .75))
quantile
```

```
##      25%      75%
## 17888.75 31938.75
```

```
#are the 25th and 75th Percentile
```

c. Calculate and report the interquartile range.

```
IQR <- IQR(data2$price)
IQR
```

```
## [1] 14050
```

```
#is the IQR
```

d. Calculate and report the standard span, the lower fence, and the upper fence.

```

span <- IQR*1.5
span

## [1] 21075

#is the Standard Span
percentile1 <- quantile(data2$price, probs = c(.25))
percentile2 <- quantile(data2$price, probs = c(.75))
lfence <- percentile1 - (IQR*1.5)
lfence

##      25%
## -3186.25

#is the Lower Fence
ufence <- percentile2 + (IQR*1.5)
ufence

##      75%
## 53013.75

#is the Upper Fence

```

- e. Are there any outliers? Subset the outlying points. Use code based on the following:

```

#
outlier1 <- data2[data2$price >= 53013.75, ]
outlier2 <- data2[data2$price <= -3186.25, ]
outlier1

## # A tibble: 9 x 1
##   price
##   <dbl>
## 1 71020
## 2 74970
## 3 69725
## 4 54005
## 5 62000
## 6 85500
## 7 82600
## 8 69700
## 9 60105

outlier2

## # A tibble: 0 x 1
## # ... with 1 variable: price <dbl>

```

```
#Are the outliers
```

- f. Calculate and report the variance, standard deviation, and coefficient of variation of car prices

```
var <- var(data2$price)
var
```

```
## [1] 207898012
```

```
#is the Variance
sd <- sd(data2$price)
sd
```

```
## [1] 14418.67
```

```
#is the Standard Deviation
cv <- sd/mean*100
cv
```

```
## [1] 52.75414
```

```
#is the is the coefficient of variation.
```

- g. We have seen from the histogram that the data are skewed. Calculate and report the skewness.  
Comment on this value and how it matches with what you visually see in the histogram.

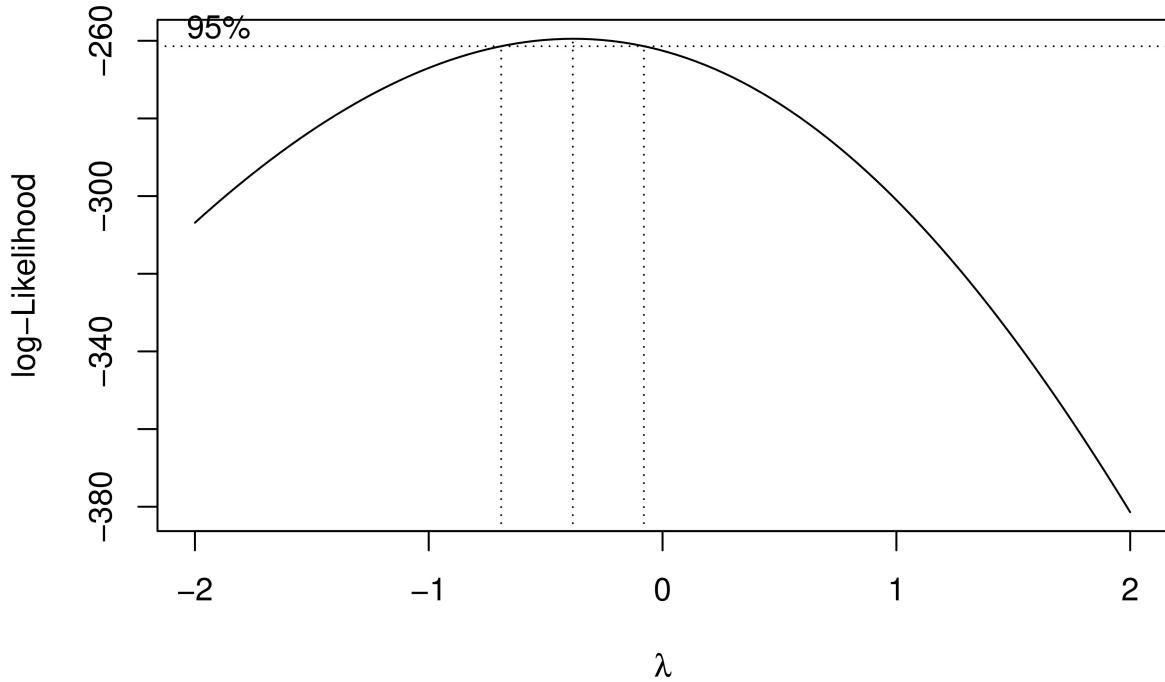
```
skew <- skewness(Price)
skew #is the Skewness
```

```
## [1] 1.760286
```

As initially observed the histogram was right skewed (positive skew) and the value generated (1.76) depicts the same

Question 3: Transforming the data. a. Use a Box-Cox power transformation to appropriately transform the data. In particular, use the `boxcox()` function in the `MASS` library. Report the recommended transformation. Do not apply this transformation to the data yet. (Note: the `boxcox` function automatically produces a plot. You do NOT need to make this in `ggplot2`.)

```
library(MASS)
bxco <- boxcox(data1$price ~ 1)
```



```
lambda <- bxcx$x[bxcx$y==max(bxcx$y)]
lambda
```

```
## [1] -0.3838384
```

- b. Apply the exact Box-Cox recommended transformation (rounded to four decimal places) to the data (this transformation is hereon referred to as the Box-Cox transformed data). Use the `summary()` function to summarize the results of this transformation.

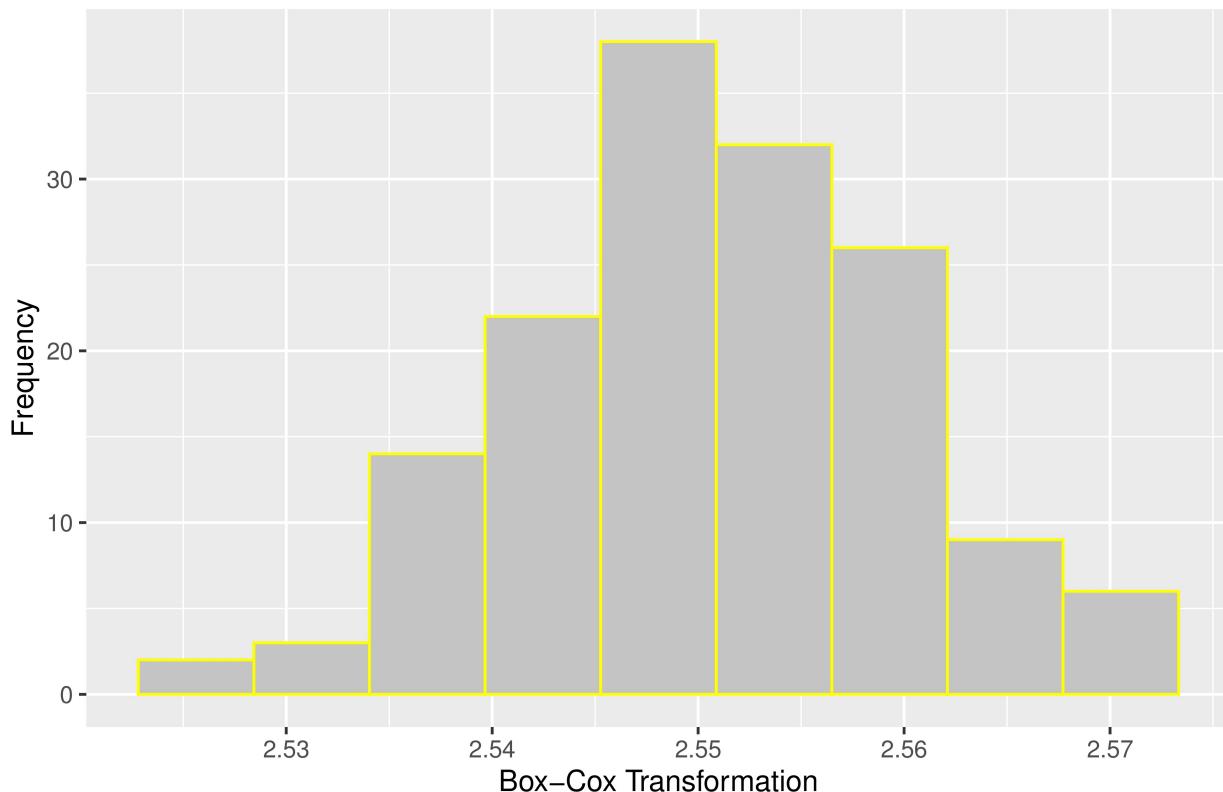
```
boxdata <- round (( data1$price^lambda -1 )/lambda ,4)
summary(boxdata)
```

```
##      Min. 1st Qu. Median    Mean 3rd Qu.    Max.
##    2.527   2.545   2.550   2.551   2.557   2.572
```

- c. Create a histogram of the Box-Cox transformed data using the number of bins suggested by Sturges' formula. On this histogram, report the mean, median, and 10% trimmed mean using the same formatting options as in part 2a above. Comment on the center, shape, and spread.

```
boxhist <- ggplot(data2,aes(x=boxdata)) +
  geom_histogram(bins = bins,color="yellow",fill="grey77") +
  ggtitle(" Boxcox Transformed Histogram with 9 Bins") +
  labs(y= "Frequency", x = "Box-Cox Transformation")
boxhist
```

## Boxcox Transformed Histogram with 9 Bins



```
tmean <- mean(boxdata)
tmean # Mean of BoxCox transformed data

## [1] 2.550573

tmedian <- median(boxdata)
tmedian # Median of BoxCox transformed data

## [1] 2.54985

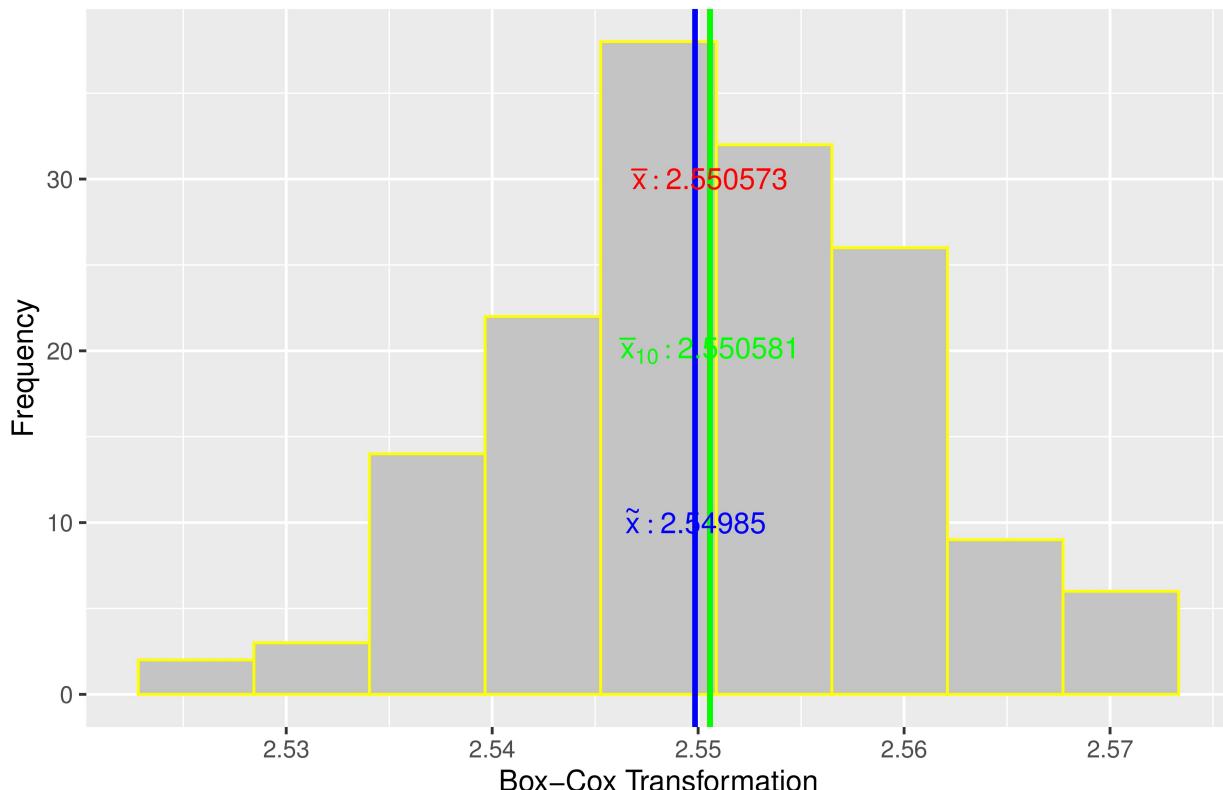
ttrim <- mean(boxdata, trim=0.1)
ttrim # Trimmed mean of BoxCox transformed data

## [1] 2.550581

boxhist + geom_vline(aes(xintercept=mean(boxdata)), color="red", size=1) + geom_vline(aes(xintercept=median(boxdata)), color="blue", size=1) + geom_vline(aes(xintercept=mean(boxdata, trim=0.1)), color="green", size=1)

# Add annotations for the mean, median, and trimmed mean
annotate(geom = "text", x = tmean, y = 30, parse = TRUE,
         label= paste("bar(x) :", tmean), size = 4, col ="red") +
  annotate(geom = "text", x = tmedian, y = 10, parse = TRUE, label=paste("tilde(x) :", tmedian),
           size = 4, color= "blue") +
  annotate(geom = "text", x = ttrim, y = 20, parse = TRUE, label=paste("bar(x)[10] :", ttrim),
           size = 4, col="green")
```

## Boxcox Transformed Histogram with 9 Bins



Histogram is Unimodal, Symmetric, Uniformly spread. Mean can be used to find the center.

d. As an alternative to the Box-Cox transformation, let's also use a log transformation. Apply the log transformation to the original 'price' data (this transformation is hereon referred to as the log transformed data). Use the 'summary()' function to summarize the results of this transformation.

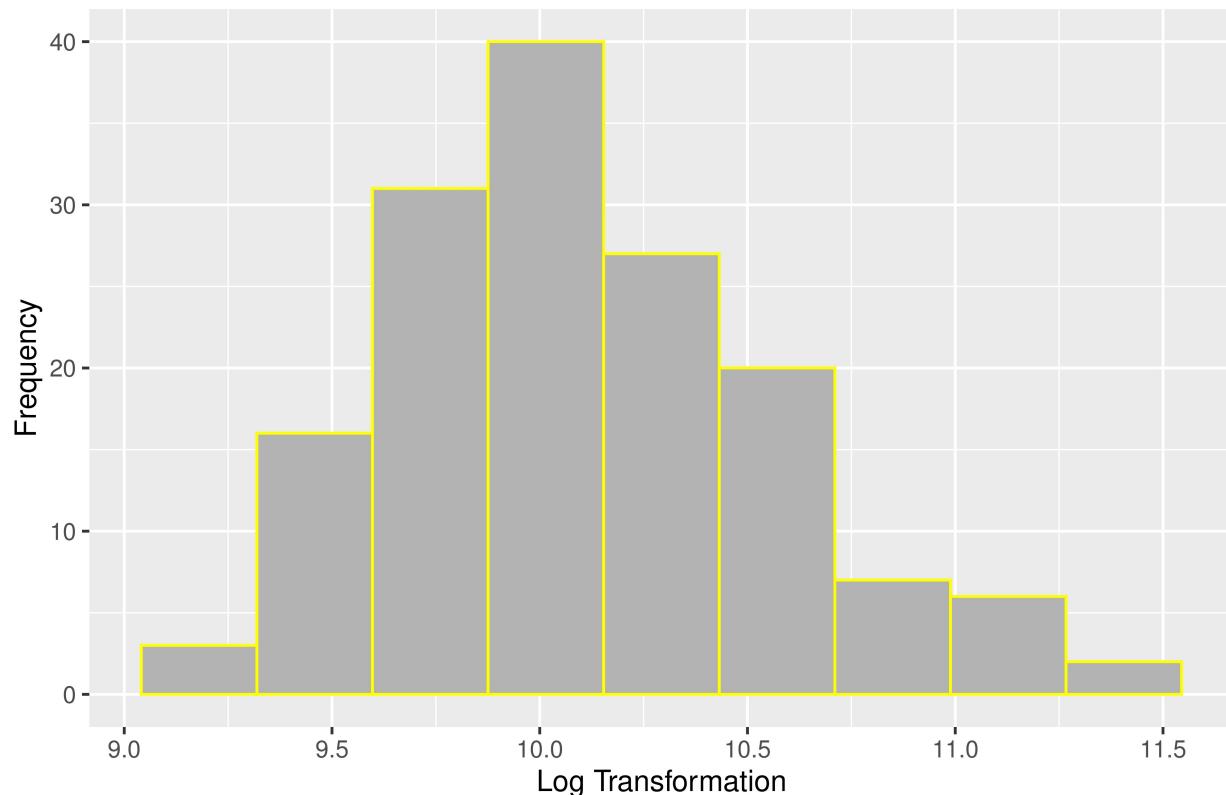
```
```r
logdata <- log(data1$price)
summary(logdata)

##      Min. 1st Qu. Median     Mean 3rd Qu.    Max.
##    9.131   9.792 10.032 10.105 10.372 11.356
```

e. Create a histogram of the log transformed data using the number of bins suggested by Sturges' formula. On this histogram, report the mean, median, and 10% trimmed mean using the same formatting options as in part 2a and 3c above. Comment on the center shape and spread.

```
loghist <- ggplot(data1,aes(x=logdata)) + geom_histogram(bins = bins,
color="yellow",fill="grey70") + ggtitle(" Log Transformed Histogram with 9 Bins") +
labs(y= "Frequency", x = "Log Transformation")
loghist
```

## Log Transformed Histogram with 9 Bins



```
lmean <- mean(logdata)
lmean # Mean of Log transformed data

## [1] 10.10457

lmedian <- median(logdata)
lmedian # Median of Log transformed data

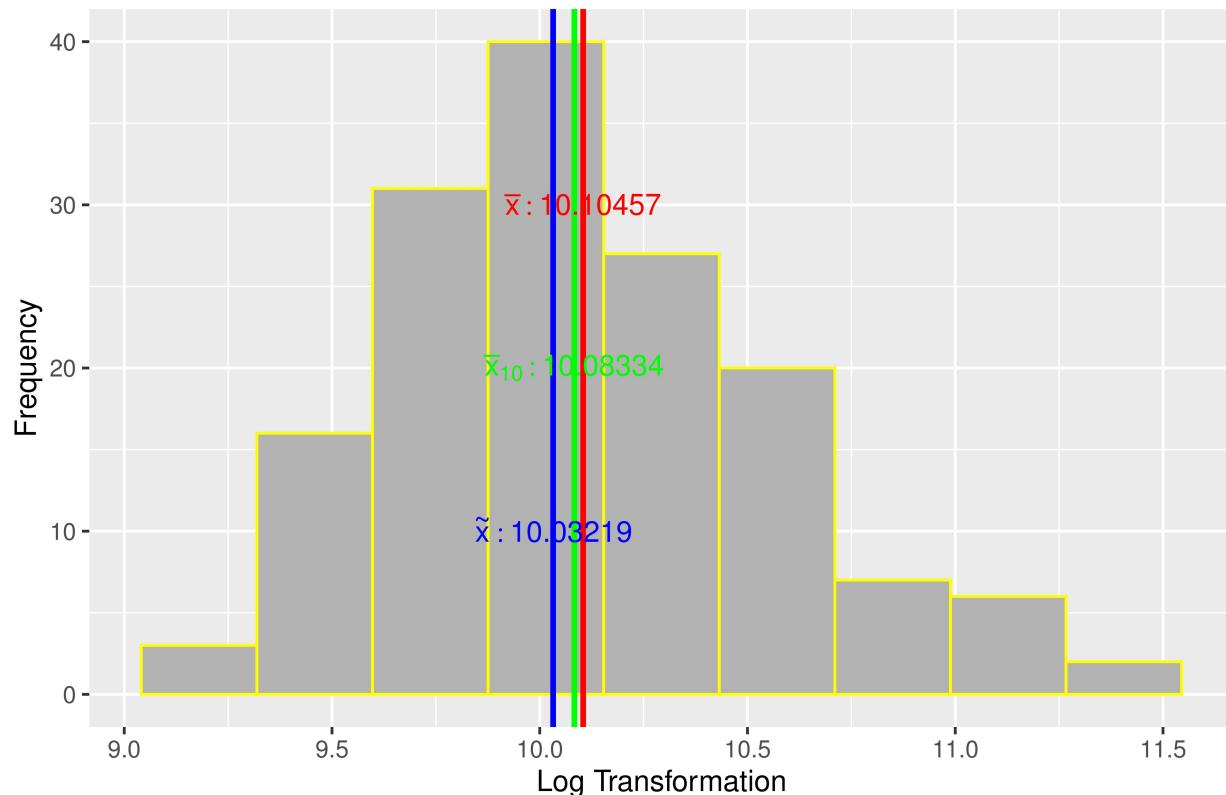
## [1] 10.03219

ltrim <- mean(logdata, trim=0.1)
ltrim # Trimmed Mean of Log transformed data

## [1] 10.08334

loghist + geom_vline(aes(xintercept=mean(logdata)), color="red", size=1) + geom_vline(aes(xintercept=me
```

## Log Transformed Histogram with 9 Bins



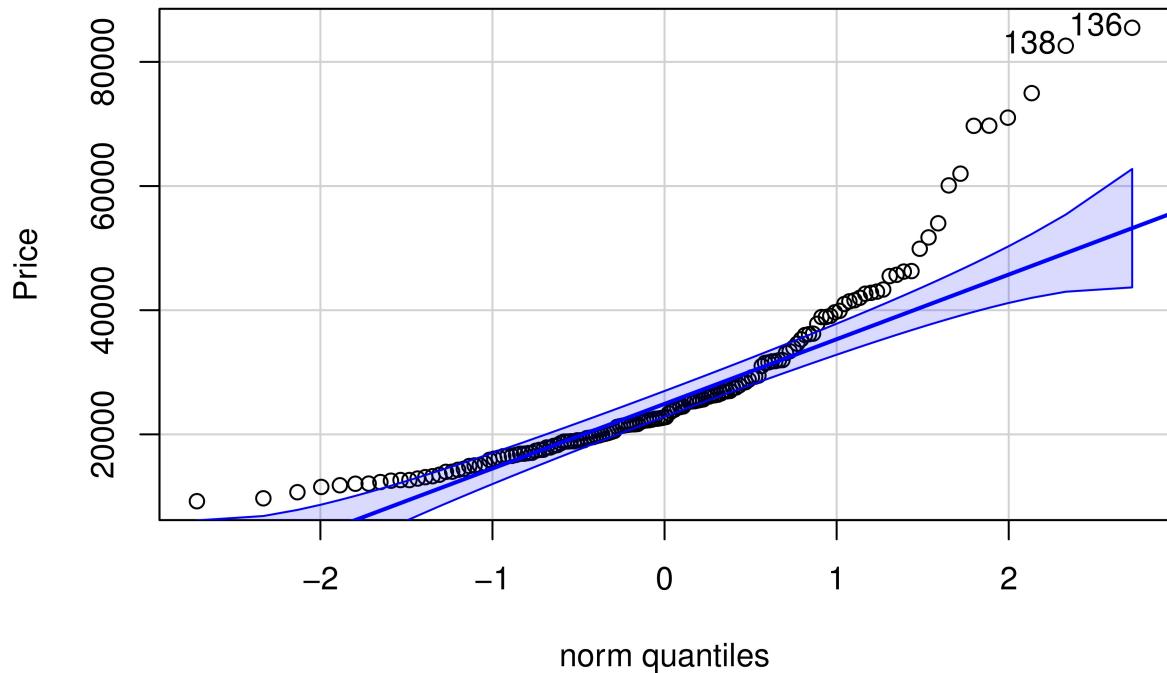
Histogram is unimodal, evenly spread, almost symmetric.

- f. Create a qqplot for the original data, a qqplot for the Box-Cox transformed data, and a qqplot of the log transformed data. Comment on the results.

```
library(car) #importing library
```

```
## Loading required package: carData
```

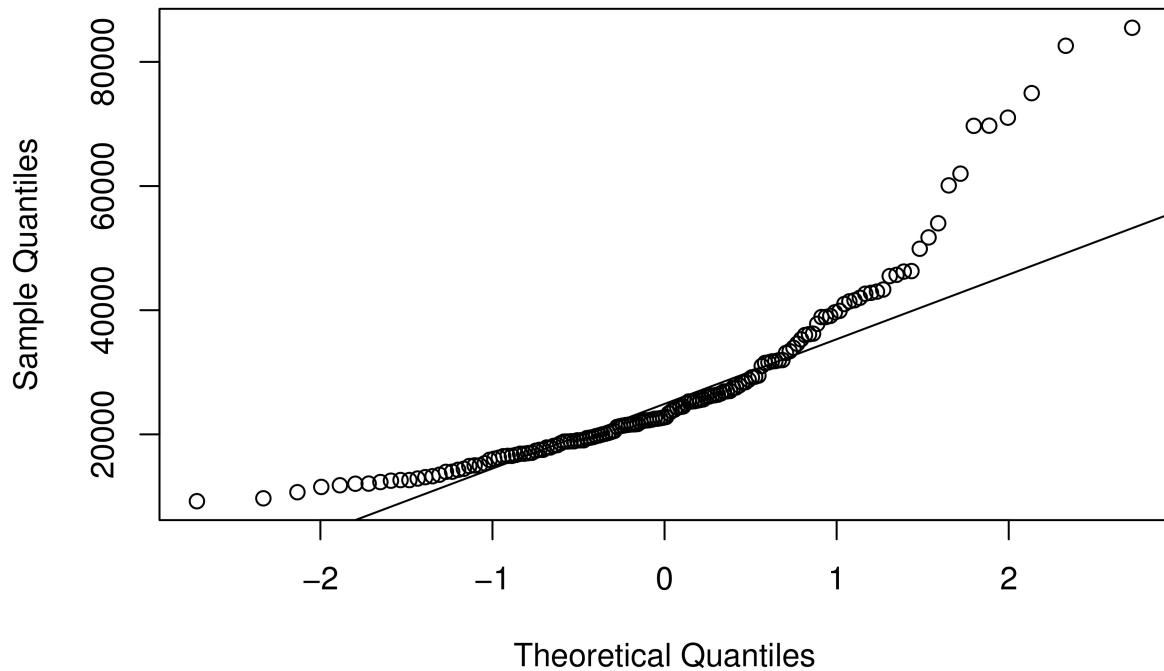
```
qqPlot(Price) #qqPlot of the original data
```



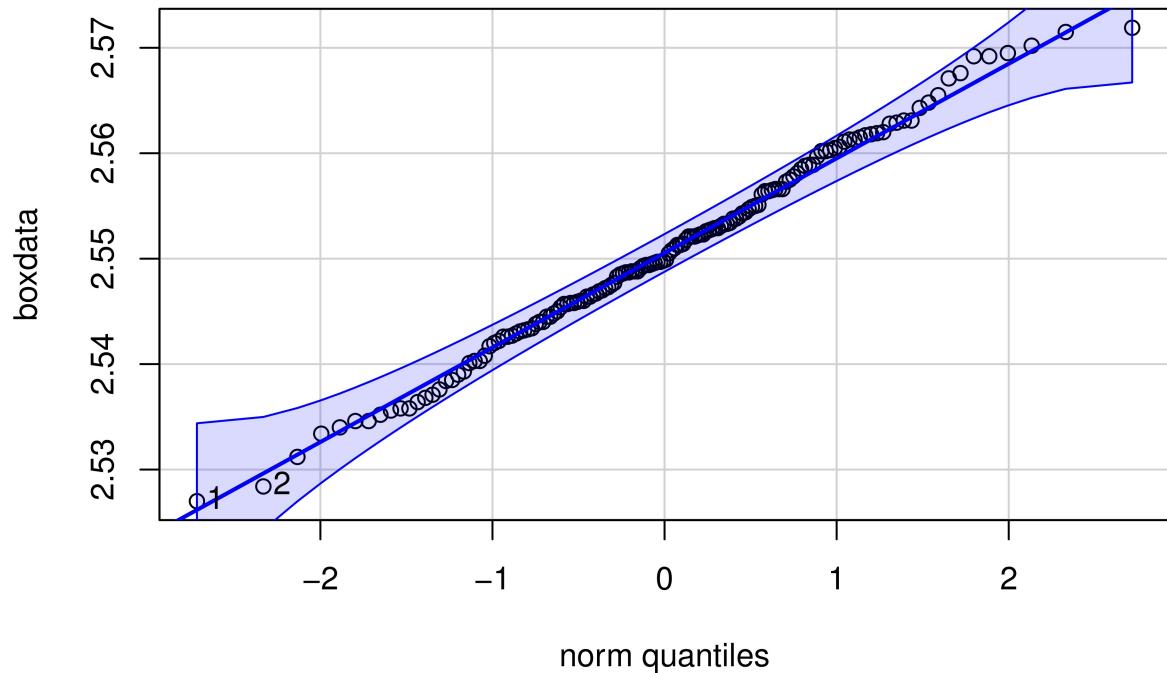
```
## [1] 136 138
```

```
qqnorm(Price); qqline(Price)
```

## Normal Q-Q Plot

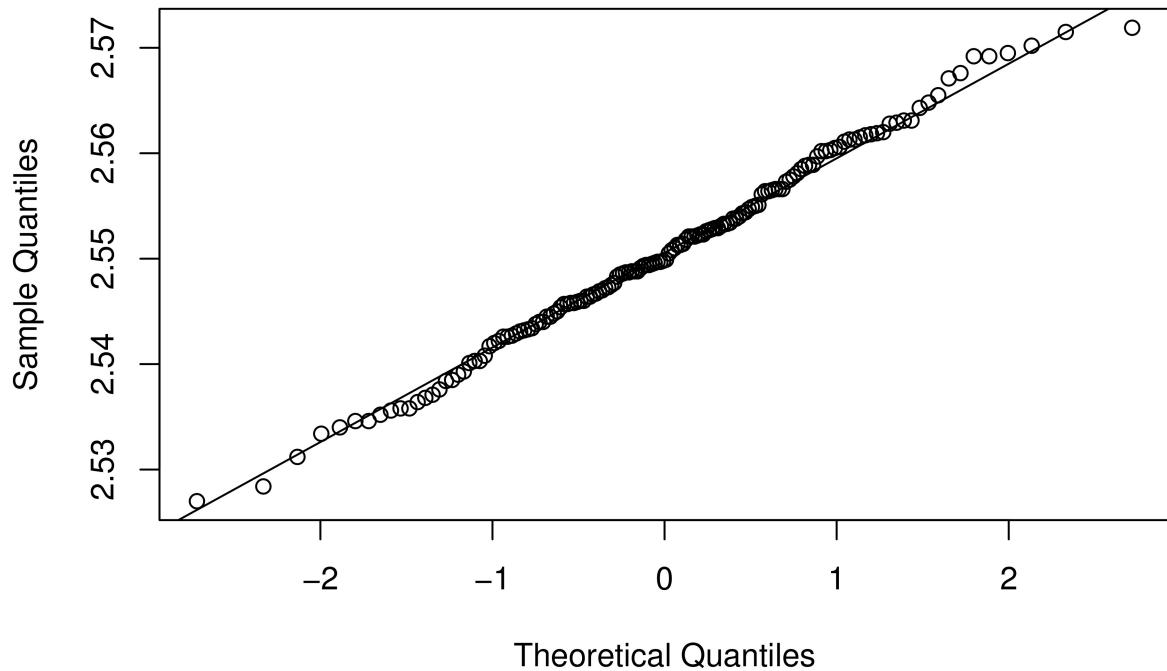


```
qqPlot(boxdata) #qqPlot of BoxCox transformed data
```

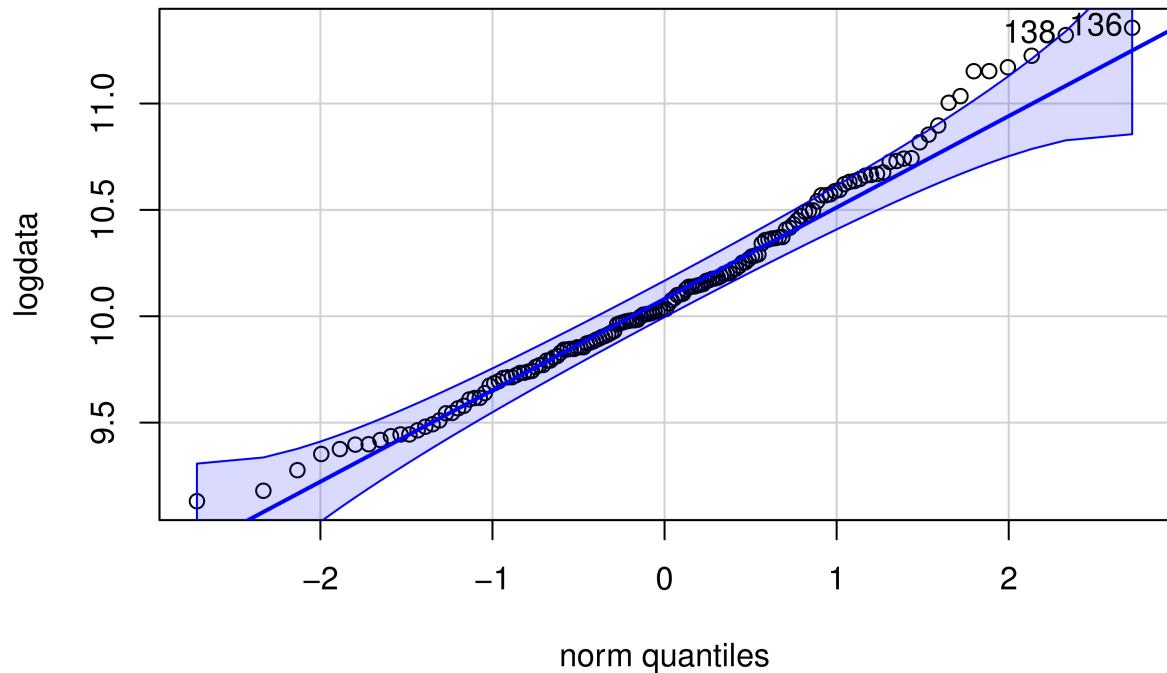


```
## [1] 1 2  
qqnorm(boxdata); qqline(boxdata)
```

## Normal Q-Q Plot



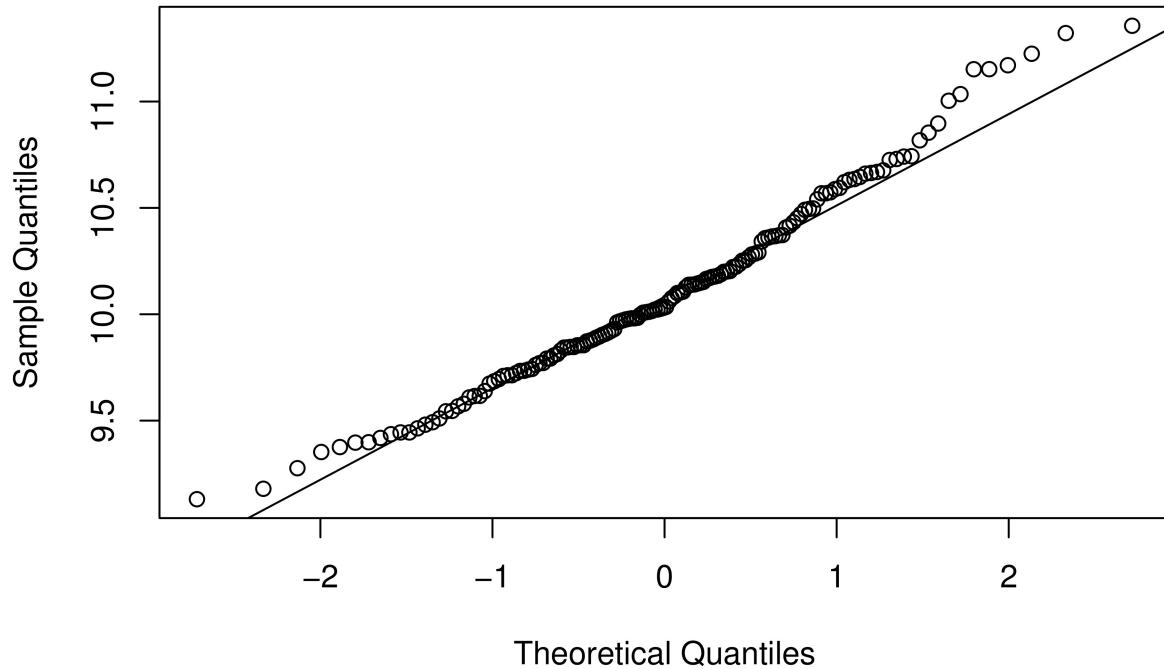
```
qqPlot(logdata) #qqPlot of Log Transformed data
```



```
## [1] 136 138
```

```
qqnorm(logdata); qqline(logdata)
```

## Normal Q-Q Plot



- g. Evaluate the empirical rule for the original data, the Box-Cox transformed data, and the log transformed data. In particular, make a table similar to that on slide 71 of the Chapter 2 notes. Comment on the results. Do either of the transformed data seem to be “better” to work with? Note, you can use code similar to the following to answer this question: