

Lab 2 (100 points + 10 points BONUS) - Describing Distributions with Graphs and Numbers

Objectives: Numerical Summaries, Histograms, and Boxplots

Remember to use the cleaned data set that you generated in Lab 1.

A. (10 points) Online Prelab

B (45 pts) Average Test Scores (Data Set: USData cleaned) We are interested in the graphical and numeric summaries for the adjusted average test score for a college admission exam (TestScore).

1. (10 points) Code. The code is the script (commands) either in R or in the Editor in SAS. Remember that you need to include the code or procedure for inputting the data into your software package.

Solution:

```
library(ggplot2)
# Read in data using the interface:
# Import Dataset --> From CSV --> browse to the file
# --> set delimiter to tab --> Change name to USData --> Import
# It is acceptable to read in the data via code
#
# 2. Five-number summary
fivenum(USData$TestScore)
#
# 3. Outliers
#
# What are the fences?
LowerFence <- fivenum(USData$TestScore)[2] -
              1.5*(fivenum(USData$TestScore)[4] -
                  fivenum(USData$TestScore)[2])
UpperFence <- fivenum(USData$TestScore)[4] +
              1.5*(fivenum(USData$TestScore)[4] -
                  fivenum(USData$TestScore)[2])

LowerFence
UpperFence

# what are the outliers?
USData[which(USData$TestScore < LowerFence |
             USData$TestScore > UpperFence),
       c("State", "CountyIndex", "TestScore")]
#
# 4. Boxplot
windows()
ggplot(USData, aes(x = "", y = TestScore)) +
  stat_boxplot(geom = "errorbar") +
  geom_boxplot() +
  ggtitle("Boxplot of TestScore") +
  stat_summary(fun.y = mean, color = "black", geom = "point", size = 3)
#
```

```
# 5. Histogram
#   for bins: Because this is a large data set, using the sqrt method
#   is usually not valid. You may use any where from 20 to the default
#   value of 30
windows()
xbar <- mean(USData$TestScore)
s <- sd(USData$TestScore)
ggplot(USData, aes(TestScore)) +
  geom_histogram(aes(y=..density..),
                 bins=20, fill="grey", col="black") +
  geom_density(col="red", lwd=1) +
  stat_function(fun=dnorm, args=list(mean=xbar, sd=s),
               col="blue", lwd=1) +
  ggtitle("Histogram of TestScore")
#
# Mean and SD
#
xbar
s
```

2. (2 points) Find the five-number summary for these data.

Solution:

```
> fivenum(USData$TestScore)
[1] 1377.151 1563.421 1604.037 1643.348 1754.276
```

Min: 1377.151 Q_1 : 1563.421 Median: 1604.037 Q_3 : 1643.348 Max: 1754.276

3. (5 points) Calculate the 1.5 IQR upper and lower limits for the outliers. Are there any outliers according to the 1.5 IQR rule (just answer yes or no, and explain why you know)? This part may be done by hand. If done by hand, all work needs to be provided. If done via computer code, then the code must be listed.

Solution:

$$IQR = Q_3 - Q_1 = 1643.348 - 1563.421 = 79.927$$

$$LowerFence = Q_1 - 1.5(IQR) = 1563.421 - 1.5(79.927) = 1443.531$$

$$UpperFence = Q_3 + 1.5(IQR) = 1643.348 + 1.5(79.927) = 1763.238$$

Since the minimum 1377.151 is less than the lower limit of 1443.531, we know there is at least one outlier. We can get more specific information with R (and do all of the calculations):

```

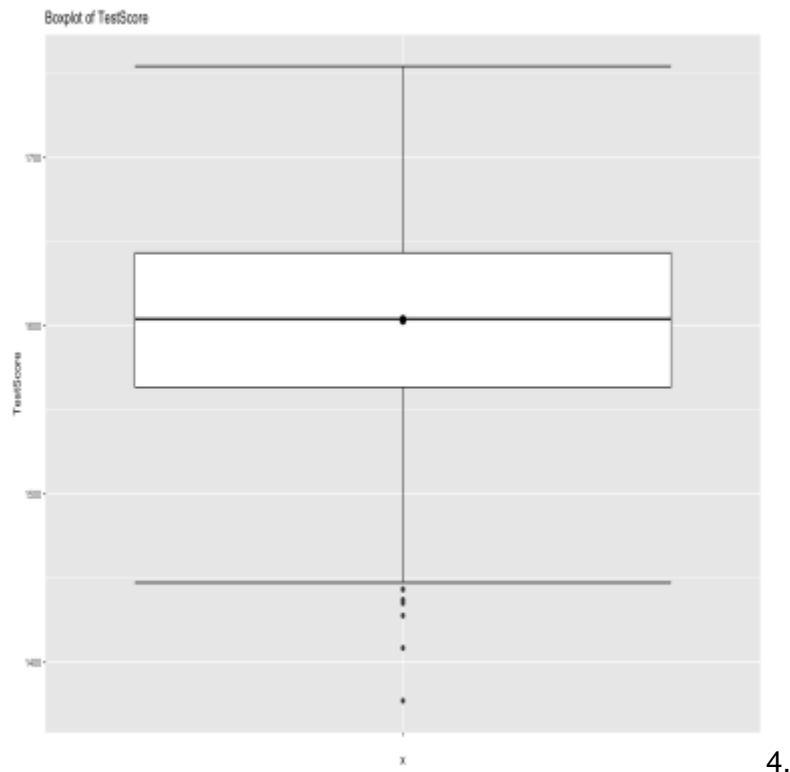
> LowerFence <- fivenum(USData$TestScore)[2] -
+               1.5*(fivenum(USData$TestScore)[4] -
+                   fivenum(USData$TestScore)[2])
> UpperFence <- fivenum(USData$TestScore)[4] +
+               1.5*(fivenum(USData$TestScore)[4] -
+                   fivenum(USData$TestScore)[2])
> LowerFence
[1] 1443.53
> UpperFence
[1] 1763.239

> USData[which(USData$TestScore < LowerFence |
+             USData$TestScore > UpperFence),
+        c("State", "CountyIndex", "TestScore")]
# A tibble: 8 x 3
   State CountyIndex TestScore
  <chr>      <int>      <dbl>
1 Alabama         3  1436.830
2 Texas          12  1435.146
3 West Virginia   2  1437.173
4 West Virginia   4  1443.278
5 West Virginia   5  1377.151
6 West Virginia   6  1443.493
7 West Virginia   7  1427.778
8 West Virginia   9  1408.523

```

It looks like 6 counties in West Virginia are outliers on the low end. Plus one from Alabama and one from Texas.

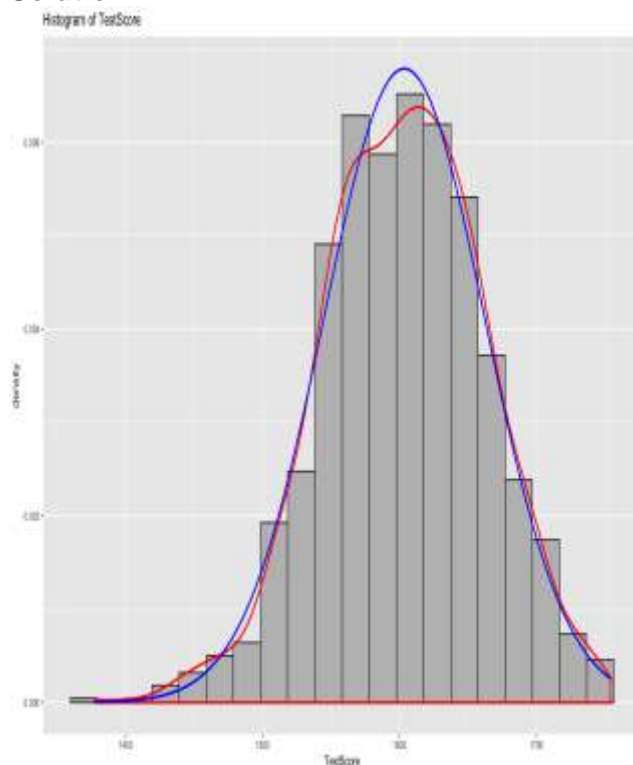
4. (5 points) Make a modified boxplot. Describe the distribution by stating whether it is symmetrical, left, or right skewed, and if there are outliers. Please indicate which features of the plot you used to classify the type of skewness and determine whether outliers are present.

Solution:

The distribution is approximately symmetric; the boxplot does not stretch out in one direction more than another direction.

It looks like there is no outlier at a higher value and at least six outliers at lower values. It is hard to tell the number because one of the circles looks bold which means that there is more than one point there. This is consistent with the optional work shown in Question 3.

5. (5 points) Make a histogram of the data. Describe the distribution by stating whether it is symmetrical, left, or right skewed, and if there are outliers. Please indicate which features of the plot you used to classify the type of skewness and determine whether outliers are present.

Solution:

The distribution looks symmetric though it has a slight tail at lower values. The distribution exhibits a very similar shape on both sides of its mode, suggesting symmetry. In addition, the two curves are close. I do not see any points that are not close to the rest of the data; therefore, there are no outliers.

6. (5 points) Are the data points you considered outliers in the histogram and the boxplot the same or different? If they are different, please provide a possible explanation for the difference.

Solution:

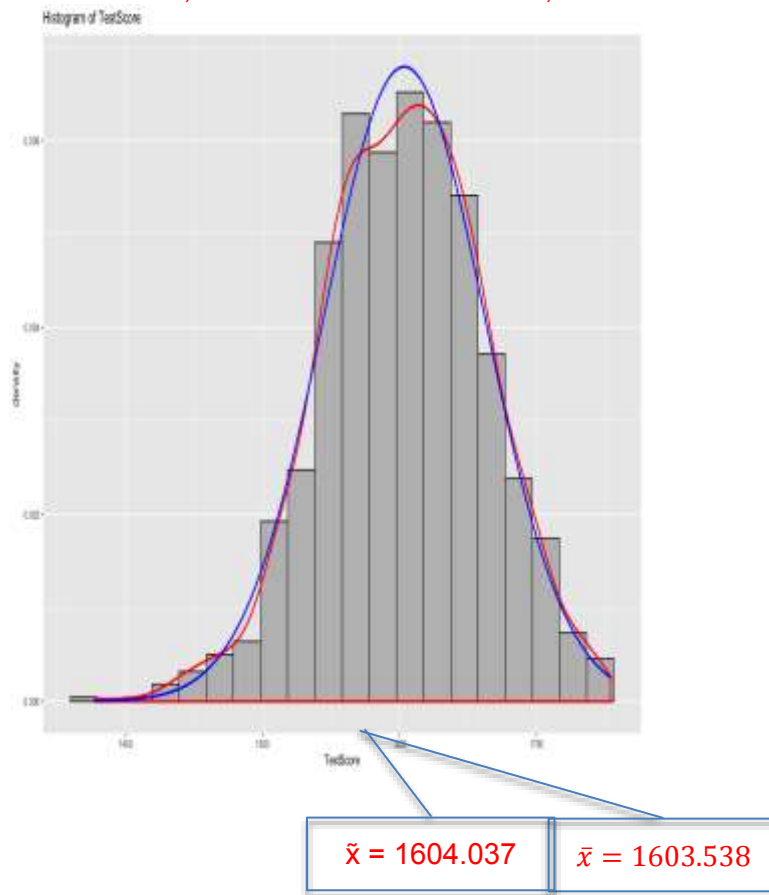
They are different. I identified more than six outliers using the boxplot and no outlier(s) using the histogram. One possible reason for the difference is that the boxplot employs a “rule of thumb” based on the IQR and its criterion for identifying outliers is not guaranteed to align with one’s impression obtained from visually inspecting a histogram. In addition, the outliers are not far above or below the fence. It is not surprising that different methods will suggest different categorizations of these “boundary” points.

7. (5 points) Obtain the sample mean, \bar{x} , and the sample standard deviation, s . Indicate the location of the mean and median on the histogram in Question 5. This may be done by hand or computer software. \bar{x} is the average of the average values and the standard deviation of the average values for the counties in the sample.

Solution:

```
> xbar
[1] 1603.538
> s
[1] 58.74789
```

$\bar{x} = 1603.538$, $s = 58.74789$. From above, the median is 1604.037.



8. (3 points) Do you think the median is close to the mean? Please explain your rationale.

Solution:

Yes, the median is very close to the mean both by just looking at the numbers, where they are placed on the histogram and looking at the boxplot.

$$\frac{\text{mean} - \text{median}}{\text{maximum} - \text{minimum}} = \frac{1603.538 - 1604.037}{1754.276 - 1377.151} = -0.001323$$

$$\frac{\text{mean} - \text{median}}{\text{standard deviation}} = \frac{1603.538 - 1604.037}{58.74789} = -0.00849$$

9. (5 points) If you only had one measure to describe the central location of the distribution of TestScore, which would you choose? Please explain your answer.

Solution:

I would choose the mean. Both the mean and median are appropriate when the distribution is approximately symmetric. However, the mean is easier to calculate so is usually used in symmetric situations.

C (45 pts). Larcenies (Data Set: USData Cleaned). We are interested in the graphical and numeric summaries of number of larcenies (thefts of personal property) out of 100,000 people (LarceniesPerPopulation).

1. (10 points) Code

Solution:

```
# 2. Five number summary
fivenum(USData$LarceniesPerPopulation)
# What are the fences?
LowerFence <- fivenum(USData$LarceniesPerPopulation)[2] -
  1.5*(fivenum(USData$LarceniesPerPopulation)[4] -
    fivenum(USData$LarceniesPerPopulation)[2])
UpperFence <- fivenum(USData$LarceniesPerPopulation)[4] +
  1.5*(fivenum(USData$LarceniesPerPopulation)[4] -
    fivenum(USData$LarceniesPerPopulation)[2])

LowerFence
UpperFence

# what are the outliers?
nrow(USData[which(USData$LarceniesPerPopulation < LowerFence |
  USData$LarceniesPerPopulation > UpperFence),
  c("State", "CountyIndex", "TestScore")])

# 4. Boxplot
windows()
ggplot(USData, aes(x = "", y = LarceniesPerPopulation)) +
  stat_boxplot(geom = "errorbar") +
  geom_boxplot() +
  ggtitle("Boxplot of Larcenies Per Population") +
  stat_summary(fun.y=mean, colour="black", geom="point", size = 3)
#
# 5. Histogram
#   for bins: see comment for part B.
windows()
xbar <- mean(USData$LarceniesPerPopulation)
s <- sd(USData$LarceniesPerPopulation)
ggplot(USData, aes(LarceniesPerPopulation)) +
  geom_histogram(aes(y=..density..),
    bins=20, fill="grey", col="black") +
  geom_density(col="red", lwd=1) +
  stat_function(fun=dnorm, args=list(mean=xbar, sd=s),
    col="blue", lwd=1) +
  ggtitle("Histogram of Larcenies Per Population")
#
# 6. Mean and SD
#
```

xbar
s

2. (2 points) Find the five-number summary.

Solution:

```
> fivenum(USData$LarceniesPerPopulation)
[1] 170.16 2028.42 3071.26 4370.17 12274.59
```

Min: 170.16 Q₁: 2028.42 Median: 3071.26 Q₃: 4370.17 Max: 12274.59

3. (5 points) Calculate the 1.5 IQR upper and lower limits for the outliers. Are there any outliers according to the 1.5 IQR rule (just answer yes or no, and explain why you know)? This part may be done by hand. If done by hand, all work needs to be provided. If done via computer code, then the code must be listed.

Solution:

$$IQR = Q_3 - Q_1 = 4370.17 - 2028.42 = 2341.75$$

$$LowerFence = Q_1 - 1.5(IQR) = 2028.42 - 1.5(2341.75) = -1484.205$$

$$UpperFence = Q_3 + 1.5(IQR) = 4370.17 + 1.5(2341.75) = 7882.795$$

Since the Lower Fence is negative and the data have to be positive, there cannot be lower outliers according to this rule. However, since the maximum of 12274.59 is greater than the 7882.795, there is at least one outlier with a higher value.

You can use R to determine what the fences are and the number of outliers.

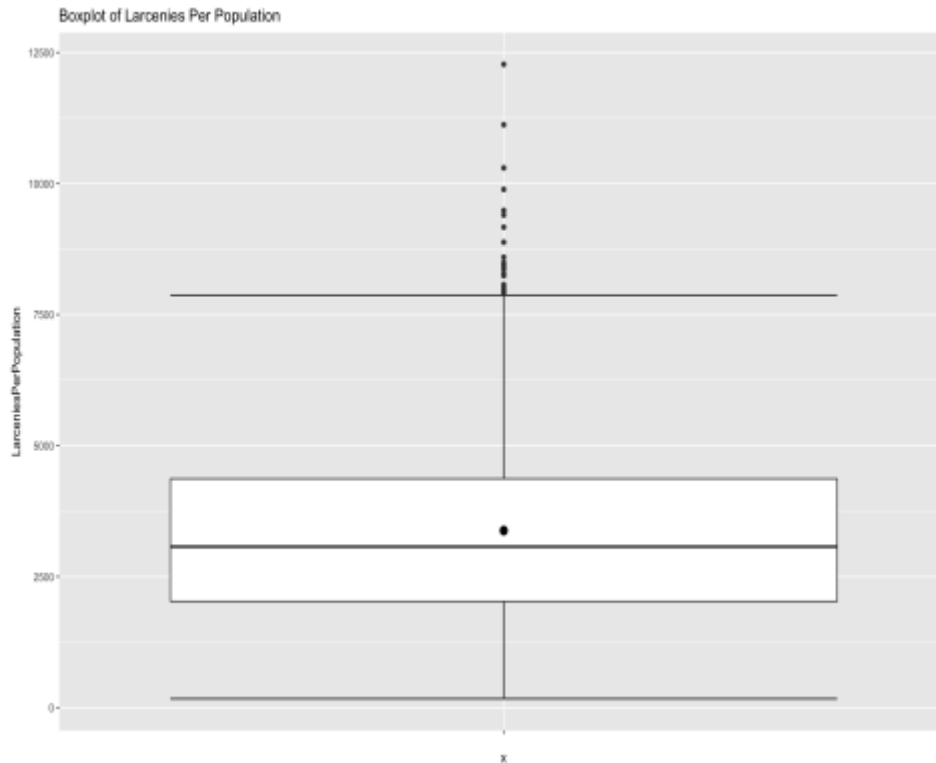
```
> LowerFence <- fivenum(USData$LarceniesPerPopulation)[2] -
+ 1.5*(fivenum(USData$LarceniesPerPopulation)[4] -
+ fivenum(USData$LarceniesPerPopulation)[2])
> UpperFence <- fivenum(USData$LarceniesPerPopulation)[4] +
+ 1.5*(fivenum(USData$LarceniesPerPopulation)[4] -
+ fivenum(USData$LarceniesPerPopulation)[2])
> LowerFence
[1] -1484.205
> UpperFence
[1] 7882.795

> nrow(USData[which(USData$LarceniesPerPopulation < LowerFence |
+ USData$LarceniesPerPopulation > UpperFence),
+ c("State", "CountyIndex", "TestScore")])
[1] 25
```

There are 25 outliers for this variable.

4. (5 points) Make a modified boxplot. Describe the distribution by stating whether it is symmetrical, left, or right skewed, and if there are outliers. Please indicate which features of the plot you used to classify the type of skewness and determine whether outliers are present.

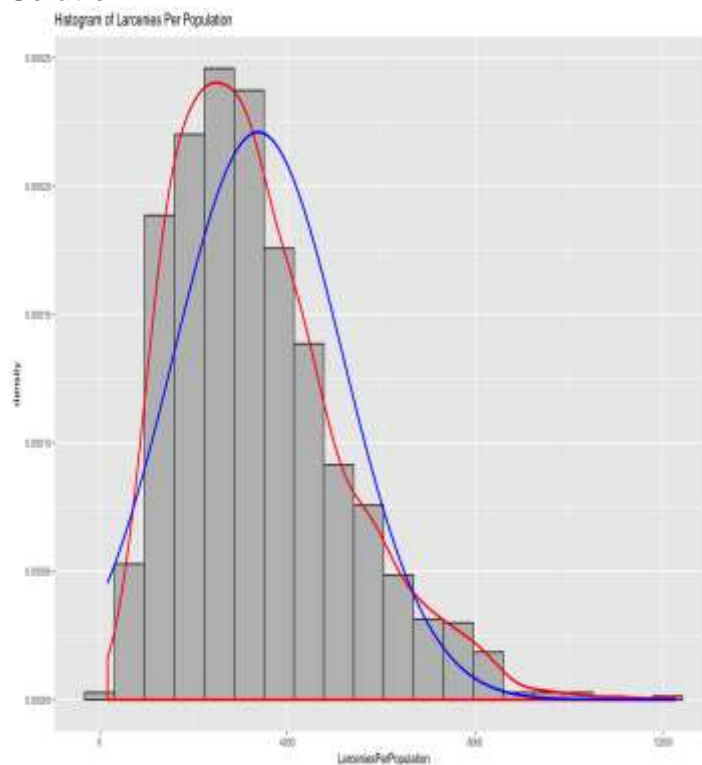
Solution:



The distribution is right skewed. The boxplot seems more “stretched out” above the median than below. Moreover, the mean is slightly larger than the median.

There appear to be many outliers, as evidenced by the many dots above the upper whisker.

5. (5 points) Make a histogram of the data. Describe the distribution by stating whether it is symmetrical, left, or right skewed, and if there are outliers. Please indicate which features of the plot you used to classify the type of skewness and determine whether outliers are present.

Solution:

The distribution is right skewed, as can be seen by the long right tail. Remember that R automatically determines the scale; therefore, if the range is larger than you see points, there have to be data there.

There is one bump on the far right which would definitely indicate that there are outliers.

There is one bump on the far right which would definitely indicate that there are outliers.

6. (5 points) Are the data points you considered outliers in the histogram and the boxplot the same or different? If they are different, please provide a possible explanation for the difference.

Solution:

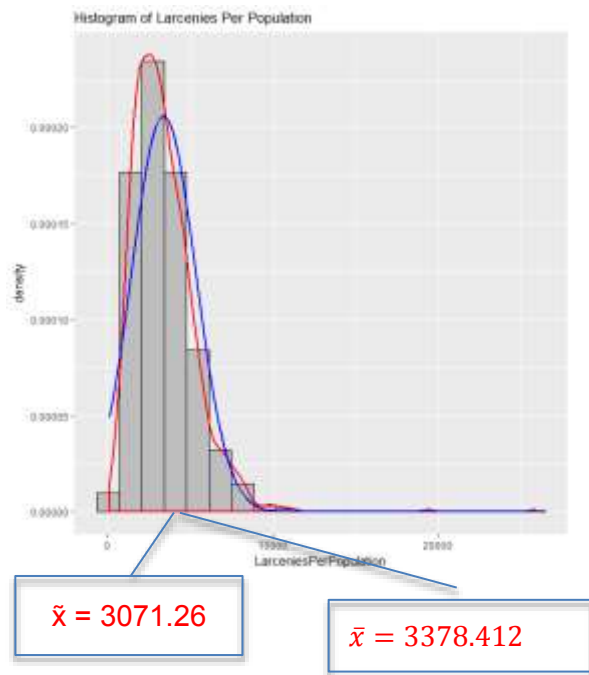
They are different. The histogram doesn't seem to have so many outliers that can be identified as in the boxplot. One possible reason for the difference is that the boxplot employs a "rule of thumb" based on the IQR and its criterion for identifying outliers is not guaranteed to align with one's impression obtained from visually inspecting a histogram. In addition, some outliers above the upper inner fence are not far above the fence. It is not surprising that different methods will suggest different categorizations of these "boundary" points.

7. (5 points) Obtain the sample mean, \bar{x} , and the sample standard deviation, s . Indicate the location of the mean and median on the histogram generated in Question 5. This may be done by hand or computer software.

Solution:

```
> xbar
[1] 3378.412
> s
[1] 1804.218
```

$\bar{x} = 3378.412$, $s = 1804.218$. From above, the median is 3071.26.



8. (3 points) Do you think the median is close to the mean? Please explain your rationale.

Solution:

Yes, the median looks close to the mean but not as close as in part B.

$$\frac{\text{mean} - \text{median}}{\text{maximum} - \text{minimum}} = \frac{3378.412 - 3071.26}{12274.59 - 170.16} = 0.0254$$

$$\frac{\text{mean} - \text{median}}{\text{standard deviation}} = \frac{3378.412 - 3071.26}{1804.218} = 0.17$$

Though they look close and the percentage difference with respect to the range is small, the percentage with respect to the standard deviation is fairly large.

This emphasizes the problems of using the comparison of mean and median to determine the skewedness of a distribution. It is much better just to look at the histogram or boxplot to determine skewedness.

9. (5 points) If you only had one measure to describe the central location of the distribution of `LarceniesPerPopulation`, which would you choose? Please explain your answer.

Solution:

I would use the median since the data is right skewed. The mean is only a good measure if the distribution is close to being symmetric.

D. (10 points) BONUS. We do not discuss how to make graphs of categorical variables in this class; however, this is very important. Make a pie chart for `Region`. Note: You will not get coding help for bonus questions.

1. (5 points) Code

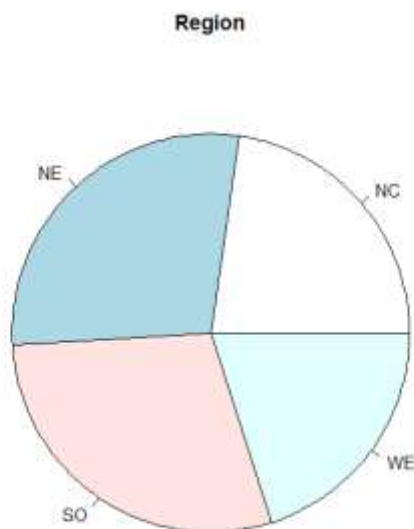
Solution:

```
windows()  
pie(table(USData$Region), main = "Region")
```

2. (5 points) If the labels are abbreviations which are not clear, state what each label represents. What information did you learn from the pie chart? Please explain your answer.

Solution:

If you provide all of the information required above, you will receive full credit.



Below is the correspondence between the abbreviations and the region names:

NC: North Central

NE: North Eastern

SO: South

WE: Western

Approximately one-fourth of the counties are in each region (which is good ☺). However, the North East and South regions have slightly more counties and the North Central and West have slightly fewer counties. This make sense because we are sampling by state and the states are smaller in the North East and South regions and larger in the North Central and West regions.