# STATISTICAL METHODS FOR DATA SCIENCE
# MINI PROJECT #2

SPRING 2018

# NAMES OF GROUP MEMBERS:

ADRITA DUTTA(axd172930)
NEETHU ANTONY(nxa171330)

## Contribution of team members:

ADRITA:

Learned R coding
Tried the R Codes
Documented the Codes and its Output
Tried out different plots and Boxplots
Explored the Data using plots

NEETHU:

Learned R coding
Tried the R Codes
Explained the Plots and Graphs
Wrote Code explanation for section-1
Explored the Data using plots

8.

a)read file college.csv

b)

->see the look of it



->add row names from the 1st column of the data set that stores the university names.

->look of the edited table, now each row has a name corresponding to the appropriate university marked by the column row.names.

->delete the first column that has college names in it from the data set, since this is treated as part of the data. We have already stored this data as row.names.

   ->look at edited table



## c.i) summary of all data sets:-

Here, since we asked for the summary of the entire data set, summary of each column is returned in the output. The column private has only 2 values- Yes or No, hence the summary will be the number of data having 'Yes' value for column Private and number of data having 'No' value for column Private. For all other columns, we have the 5 number summary which includes Minimum, 1st Quartile, Median(2nd Quartile), 3rd Quartile and Maximum.

## ii) scatterplot matrix of first 10 columns:

Considering the first 10 columns we see that it forms a matrix that is symmetric along its diagonal. Each scatterplot represents the correlations between the 2 elements.

On consideration of each scatterplot separately we can observe a few different kinds of relations among the variables:-

1)Horizontal lines:no correlation. The vertical axis is a value axis and the horizontal axis is a category axis. Instead of displaying values, a category axis shows evenly spaced groupings (categories) of data. Because my data has only values and no categories.

2)Vertical lines:no correlation.The horizontal axis is a value axis and the vertical axis is a category axis. Instead of displaying values, a category axis shows evenly spaced groupings (categories) of data. Because my data has only values and no categories.

3)diagonal line(/):a line with a positive slope indicates that the 2 variables have a positive correlation. A line with a negative slope would represent a negative correlation, but there is no such relation in this matrix.

4)approximately diagonal line: These show low positive correlation. The plot is almost diagonal but is not exactly diagonal.

# iii) side-by-side boxplots of Outstate vs Private:

       Boxplots are only able to deal with one quantitative variable. However, side-by-side box plots can be applied to data sets with one quantitative and one categorical variable, which makes them especially useful for many real-world statistical problems.

       Here Private is a categorical variable in the data set and Outstate is a quantitative variable. Private variable can take up two values or Categories(Yes or No). Hence there will be 2 boxplots drawn side-by-side as shown below. Side-by-side boxplots present all of the information that box plots do for each instance of a categorical variable. Box plots corresponding to each instance of the categorical variable (in this case two instances namely, Yes and No) summarize the data in five different numbers:-

- The Median
- The Lower Quartile
- The Upper Quartile
- The Minimum
- The Maximum

From the below side-by-side boxplot, we can infer that the Private Universities has a higher Out-of-state tuition compared to the public universities.

## iv) create Elite, by binning the Top10perc variable

    ->see how many Elite universities are there using summary: This gives the number of universities that are marked as Yes in Elite column and number of universities that are marked as No in the Elite Column.
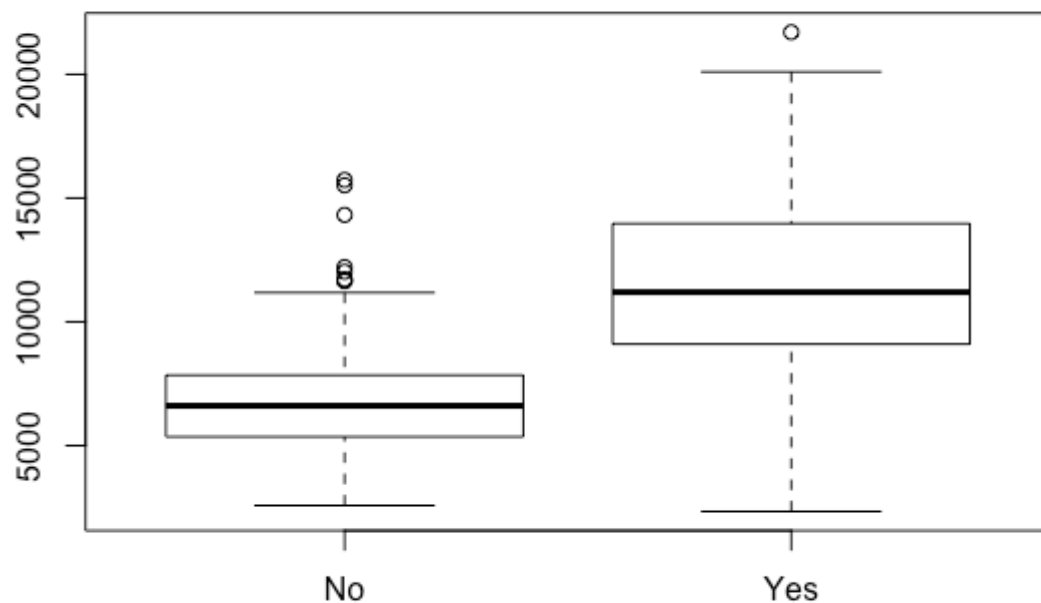
| NO | YES |
|---|---|
| 699 | 78 |

## ->side-by-side boxplots of Outstate vs Elite

Boxplots are only able to deal with one quantitative variable. However, side-by-side box plots can be applied to data sets with one quantitative and one categorical variable, which makes them especially useful for many real-world statistical problems.

    Here Private is a categorical variable in the data set and Outstate is a quantitative variable. Private variable can take up two values or Categories(Yes or No). Hence there will be 2 boxplots drawn side-by-side as shown below. Side-by-side boxplots present all of the information that box plots do for each instance of a categorical variable. Box plots corresponding to each instance of the categorical variable (in this case two instances namely, Yes and No) summarize the data in five different numbers:-

- The Median
- The Lower Quartile
- The Upper Quartile
- The Minimum
- The Maximum

From the below side-by-side boxplot, we can infer that the Elite Universities has a higher Out-of-state tuition compared to the other universities. Also, the box plot corresponding to Elite universities is left skewed, since the median is closer to 3rd Quartile, this again presses the fact that most of the Elite universities has a higher out-of-state tuition.

v) histograms of a few variables with different number of bins.

(a) Histogram of the quantitative variable Room.Board with differing number of bins.

From the below plots, it is clear that the distribution is right skewed. We can see that in most of the universities the Room and Board costs lies between 2000 and 7000. There are a few universities to the right that have a Room and Board costs higher than 7000 and these are the bars that make the data have a shape that is skewed right.

Histogram of college$Room.Board

Histogram of college$Room.Board

Histogram of college$Room.Board

Histogram of college$Room.Board

(b) Histogram of the quantitative variable PhD with differing number of bins.

From the below plots, it is clear that the distribution is left skewed. We can see that in most of the universities the percent of faculties with PhD's is between 40 and 100. There are a few universities to the left that have a percent of faculties with PhD's less than 40 and these are the bars that make the data have a shape that is skewed left.
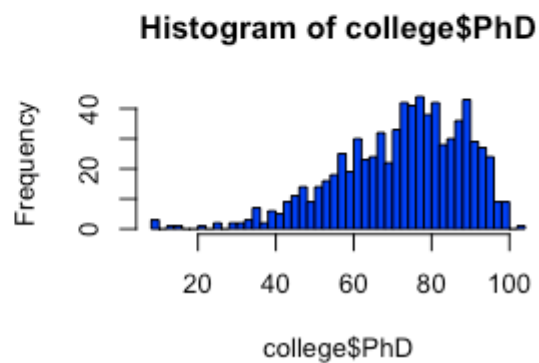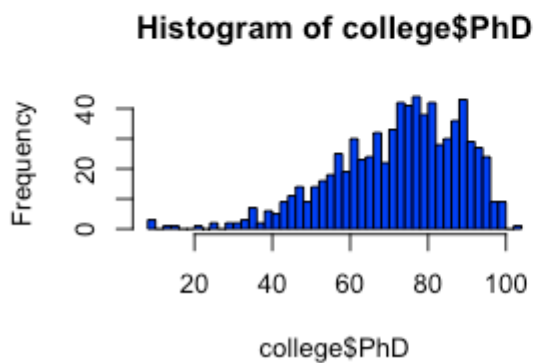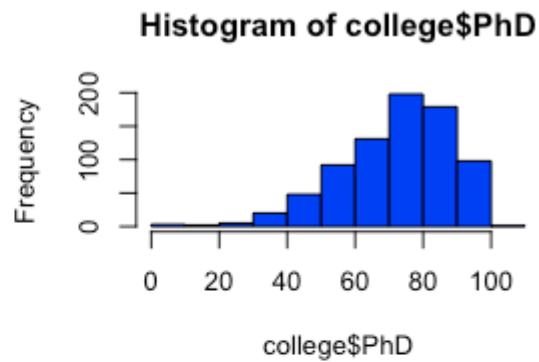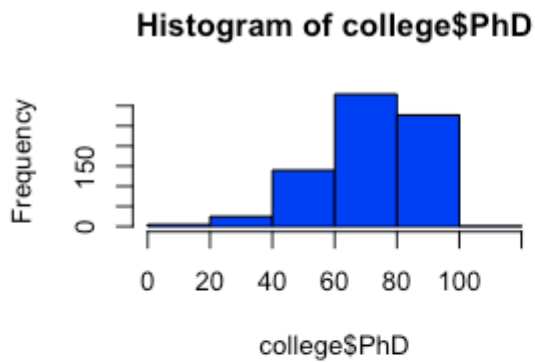
Histogram of college$PhD

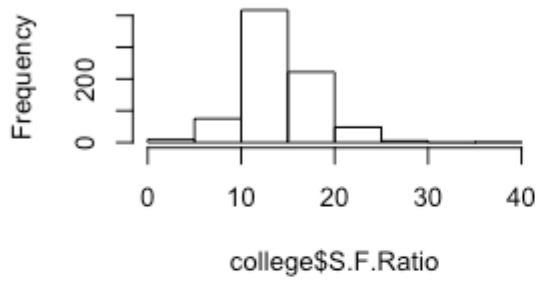(c) Histogram of the quantitative variable S.F.Ratio with differing number of bins.

From the below plots, it is clear that the distribution is left skewed. We can see that in most of the universities the percent of faculties with PhD's is between 40 and 100. There are a few universities to the left that have a percent of faculties with PhD's less than 40 and these are the bars that make the data have a shape that is skewed left.

**Histogram of college$S.F.Ratio**

Frequency

200

0

0   10   20   30   40

college$S.F.Ratio

**Histogram of college$S.F.Ratio**

Frequency

100

0

10   20   30   40

college$S.F.Ratio

**Histogram of college$S.F.Ratio**

Frequency

80

40

0

10   20   30   40

college$S.F.Ratio

**Histogram of college$S.F.Ratio**

Frequency

50

20

0

10   20   30   40

college$S.F.Ratio

**Histogram of college$Room.Board**

Frequency

college$Room.Board

**Histogram of college$PhD**

Frequency

college$PhD

**Histogram of college$S.F.Ratio**

Frequency

college$S.F.Ratio

**Histogram of college$Top25perc**

Frequency

college$Top25perc

vi) On further experimenting with the given data it is observed that:
   ->the higher the acceptance, graduation rate has lesser variance and is thus lower in average



->more applicants if more number of faculty have phd

->more variance in costs of books when lesser students enrolled



->as personal expenses increase donations have lower variance and thus lower average values

## Section 2:

## R-Code:

```r
#a)read file college.csv
college<-read.csv("C:/Users/adrit/Desktop/utd/sem2/stats for
ds/project/PROJ2/College.csv");

#b)
#fix file, see the look of it
fix(college);

#now we set names to each row(R will not perform any operations on row
names:)
#extra column for row name added for each row
rownames(college)<-college[,1];

#see the look of the table again(this time with an extra column for names)
fix(college);

#now that we have 2 columns with uni name, delete 1
#delete the original 1 as calculations can be done on that.
college<-college[,-1];

#see table again:
fix(college);

#c)
#i)summary of data set(MIN,Q1,Q2,MEAN,Q3,MAX)
summary(college);


#ii)scatterplot matrix of first 10 columns:
pairs(college[,1:10]);
```
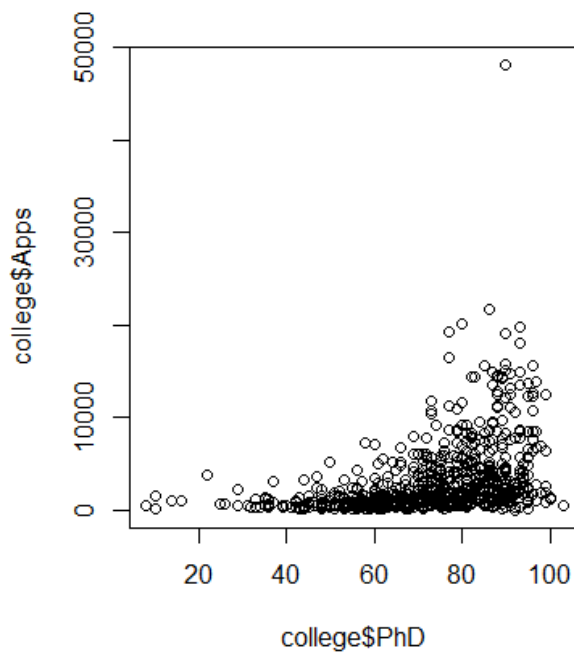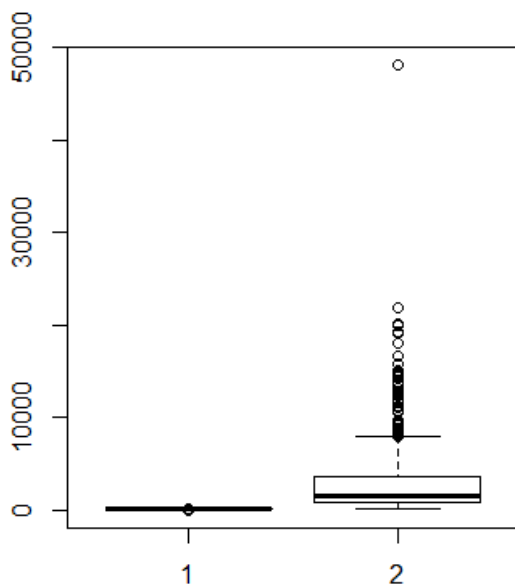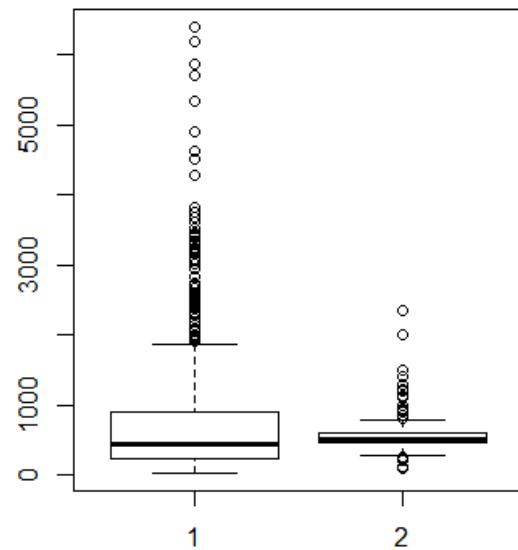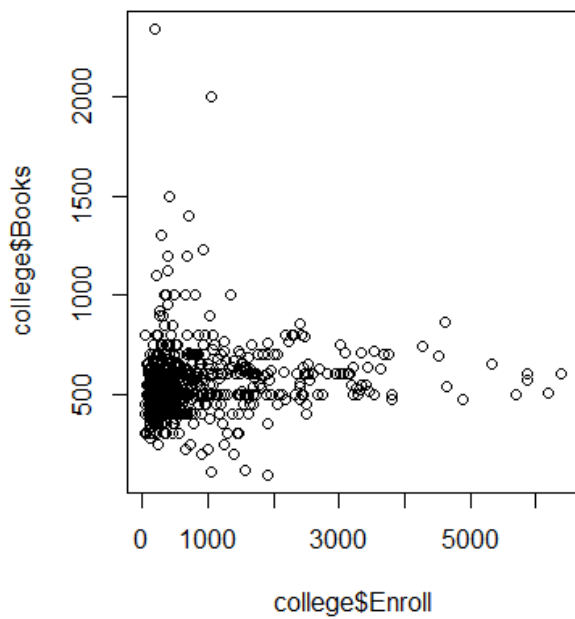
```r
#iii)side by side boxplots of Outstate vs Private
plot(college$Private, college$Outstate)

#iv)creating Elite(top 10 percentile)
Elite=rep("No",nrow(college))
Elite[college$Top10perc>50]="Yes"
Elite=as.factor(Elite)
college=data.frame(college,Elite)


#summary of Elite universities
summary(college$Elite)

#boxplots of outstate vs Elite
plot(college$Elite, college$Outstate)

#v)histograms of a few variables with different bin size
#dividing display page
par(mfrow=c(2,2))

#histogram of room and board costs
hist(college$Room.Board)

#histogram of PhD students
hist(college$PhD, col=4, breaks = 2)

#histogram of student faculty ratio
hist(college$S.F.Ratio, breaks =10)
hist(college$S.F.Ratio, breaks =25)
hist(college$S.F.Ratio, breaks =50)
hist(college$S.F.Ratio, breaks =100)


#histogram of top 25 percentile
hist(college$Top25perc, col=2)
```

#vi)doing other experiments on data:
#boxplot of no of acceptances and graduation rate
#we see here the higher the acceptance, graduation rate has lesser variance)
par(mfrow=c(1,2))
plot(college$Accept, college$Grad.Rate)
boxplot(college$Accept, college$Grad.Rate)


#boxplot of applicant no vs percentage of faculty with PHD
#more applicants if more number of faculty have phd
plot(college$PhD, college$Apps)
boxplot(college$PhD, college$Apps)


#boxplot of Enroll and Books
#more variance in costs of books when lesser students enrolled
plot(college$Enroll, college$Books)
boxplot(college$Enroll, college$Books)


#personal exp vs percentage of alumni who donate
#as personal expenses increase donations have lower variance and thus lower average values.
#if they spent less in college they are more likely to donate
plot(college$Personal, college$perc.alumni)