

CDA 532 Homework 1

UB ID: 50413349

Name: Shreeram G S

Q1) First we shall load all essential packages such as ISLR which holds the College data set, and caret, tidyverse, ggplot (for histograms). To load the dataset of college, we use the data() function. The head() function below gives us from 0 - 5 tuples of our dataset. We split the data set into 2 parts for our pairwise plots further, where we compare the two segments.

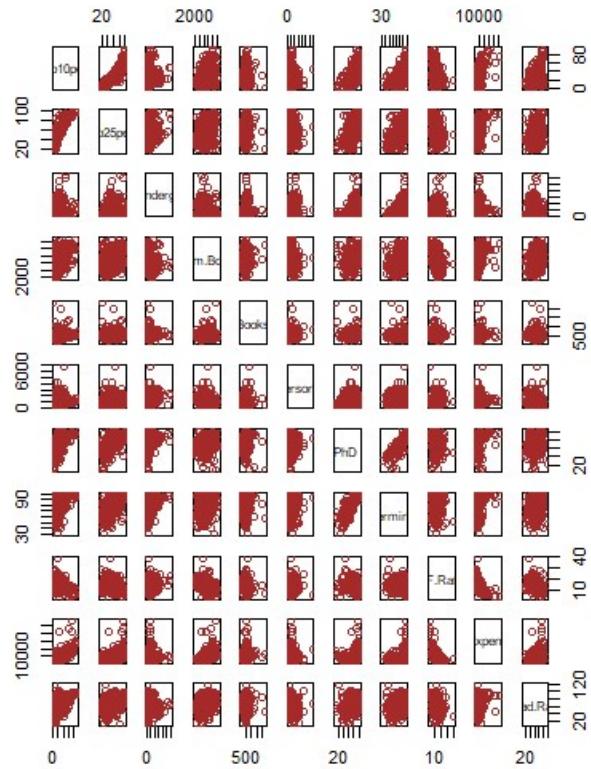
```
> require(ISLR)
> require(caret)
> require(tidyverse)
> library(ggplot2)
> data(College)
> head(College)

Private APPS Accept Enroll Top10perc Top25perc F.Undergrad P.Undergrad Outstate Room.Board Books Personal PhD Terminal S.F.Ratio perc.alumni
Abilene Christian University Yes 1660 1232 721 25 52 2885 537 7440 3200 450 3200 70 78 18.1 12
Adelphi University Yes 2186 1924 512 16 29 2683 1227 12280 6450 750 1500 29 30 12.2 16
Adrian College Yes 1428 1097 336 22 50 1036 99 11250 3750 400 1165 53 66 12.9 30
Agnes Scott College Yes 417 349 137 60 89 510 63 12960 6450 450 875 92 97 7.7 37
Alaska Pacific University Yes 193 146 55 16 44 249 869 7560 4120 800 1500 76 72 11.9 2
Albertson College Yes 587 479 158 38 62 678 41 13500 3335 500 675 67 73 9.4 11
Albertson College Expnd Grad.Rate
Abilene Christian University 7041 60
Adelphi University 7077 36
Adrian College 8735 54
Agnes Scott College 19016 59
Alaska Pacific University 10922 15
Albertson College 9727 55
> set.seed(1)
> var <- createDataPartition(College$Accept , p = 0.5, list = FALSE)
> |
```

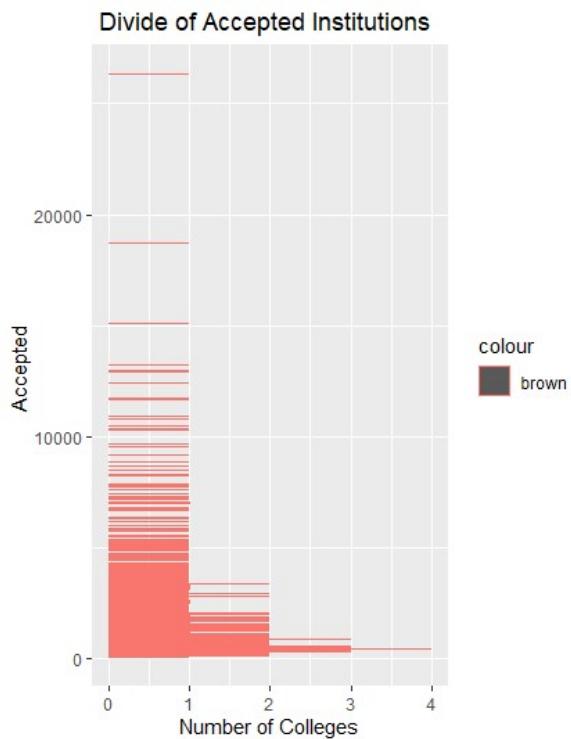
- a) After splitting we shall use the pairs() function to perform pairwise plots and use ggplot() to make histograms. Here we have taken Accept column and phD column from the college data set to create histograms.

```
> #Q1.a)
> #pairwise plots
> pairs(College[, -var], col="brown")
> #histograms: Accepted vs number of college
> var1 <- ggplot(College, aes(y=Accept, col="brown"))+
+   geom_bar()+
+   ggtitle(" Divide of Accepted Institutions")+
+   labs(x="Number of Colleges",
+        y="Accepted ")
> var1
>
> #histograms of colleges phd candidates
> var1 <- ggplot(College, aes(x=PhD))+
+   geom_bar()+
+   ggtitle("Count of Institutions per phd candidates")+
+   labs(x="PhD",
+        y=" Count of Colleges")
> var1
```

***PAIRWISE PLOTS:**

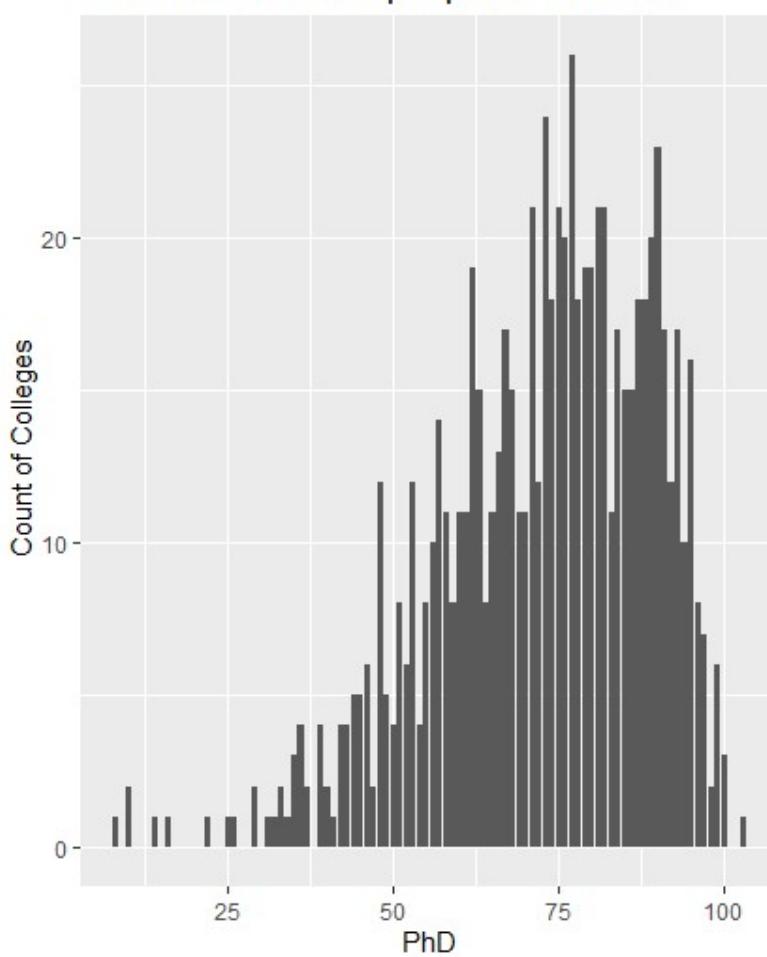


***Histogram of Accepted vs number of college:**



Histograms of colleges PhD candidates:

Count of Institutions per phd candidates



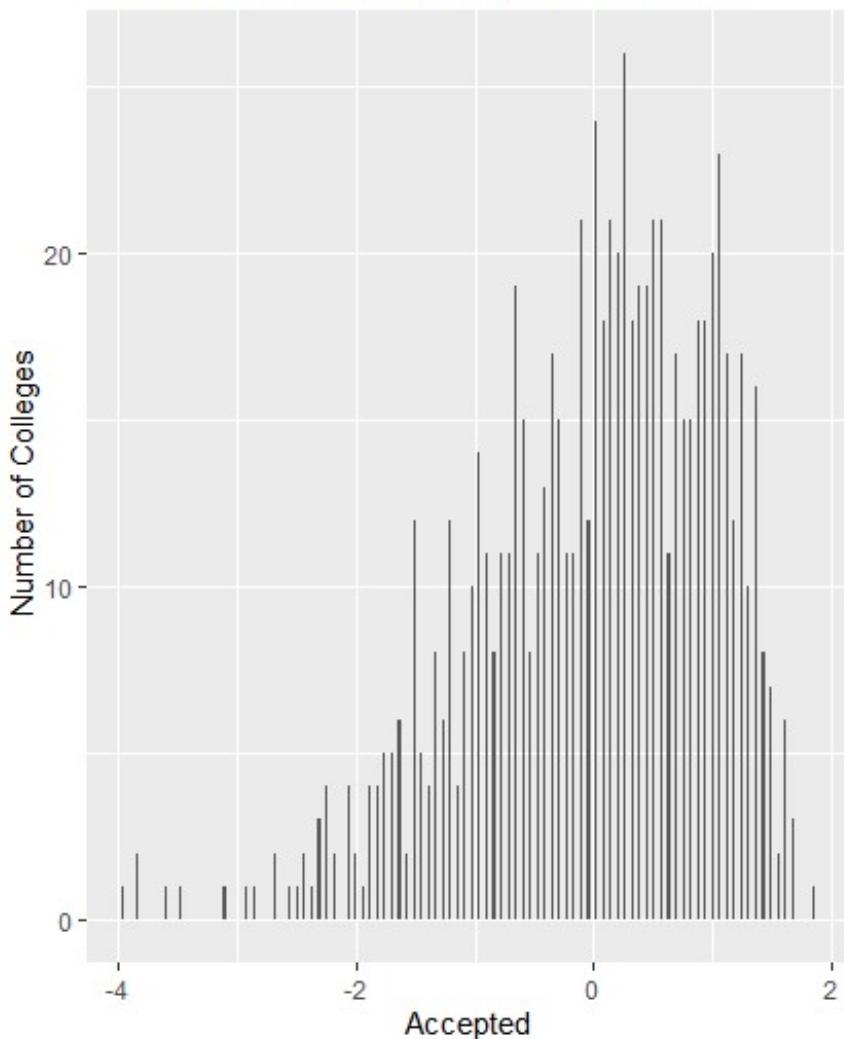
Yes, clearly even after viewing both the columns of the dataset, we see it is difficult to come to conclusions as the graph has too much noise. It would be cleaner to perform scaling to get an evenly distributed graph, which would make it easier to read and understand hence consequently allowing to make decisions based on it.

b) Scaling Accepted number of the colleges and plotting histogram:

```
> #Scaling PhD candidates of the colleges and plotting histogram
> var1<-scale(college$PhD, center = TRUE, scale = TRUE)
> var2 <- ggplot(college, aes(x=var1))+
+   geom_bar()+
+   ggtitle("Count of Institutions per phd candidates")+
+   labs(x="Accepted",
+        y=" Number of colleges")
> var2
```

Now after scaling we see each of the accepted applications are more sparsely seen, and we can differentiate and make much more accurate decisions. Now we see the range between [-4,2] and making it simpler to read the graph.

Count of Institutions per phd candidates



- c) Here we are splitting the private column of yes/no into private and public and converting the main data set into 2 parts. We use the save() function to save it as .Rdata in the file = “ ” part inside the save function. Now these data set can further separately be used to be loaded or to perform further operations on it.

- e) Now we eliminate the rest of the records which are not of median of Top25perc column of the college data set. Then we shall store it in a new variable called “private/public”

```
#1.e) Eliminate all colleges of median less than Top25%
private_eliminated<-College[median(private$Top25perc),]
summary(private)
public_eliminated<-College[median(public$Top25perc),]
summary(public)
```

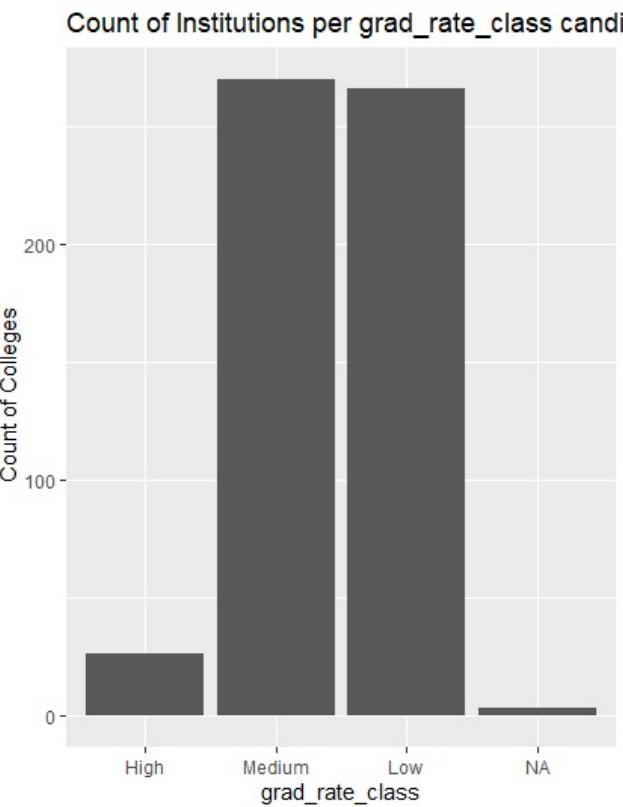
- f) Here we use the ordered to order function after the data is been cut. The cut() function allows us to cut data into bins and specify 'cut labels'. There should n+1 bins if there are n labels. So here as you can see, the bins are of the size 15,40,70,100 and the labels are high low and medium. Further we have append it to public and private datasets.

Later we plot histogram and we see that now we can make even easier and simpler decisions as there are only 3 labels.

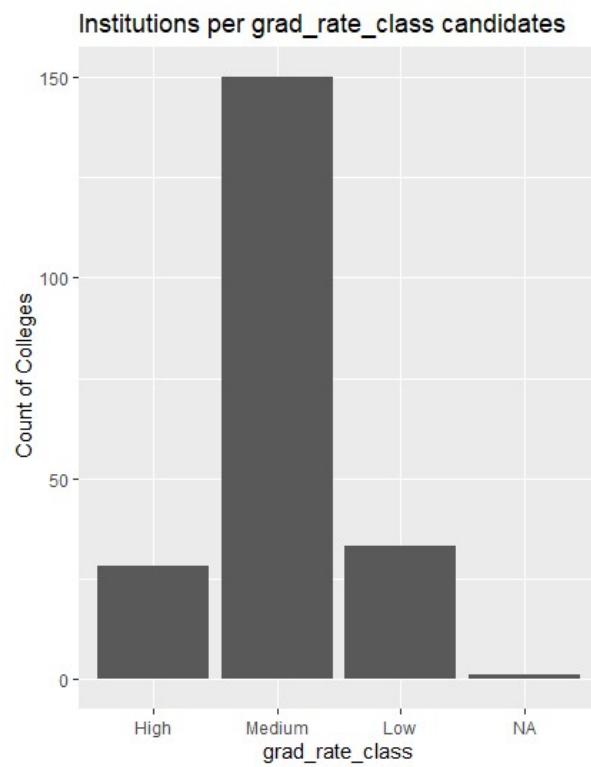
```
#1.f) graduation rate into high or low
#grad_rate_class of private
grad_rate_class <- ordered(cut(private$Grad.Rate, c(15,40, 70, 100)), labels = c("High", "Medium", "Low"))
summary(grad_rate_class)
#append to private
private$grad_rate_class <- grad_rate_class
head(private)
var1 <- ggplot(private, aes(x=grad_rate_class))+
  geom_bar()+
  ggtitle("Count of Institutions per grad_rate_class candidates")+
  labs(x="grad_rate_class",
       y="Count of Colleges")
var1
#grad_rate_class of public
grad_rate_class <- ordered(cut(public$Grad.Rate, c(15,40, 70, 100)), labels = c("High", "Medium", "Low"))
summary(grad_rate_class)
#append to public
public$grad_rate_class <- grad_rate_class
head(public)
var1 <- ggplot(public, aes(x=grad_rate_class))+
  geom_bar()+
  ggtitle("Count of Institutions per grad_rate_class candidates")+
  labs(x="grad_rate_class",
       y="Count of Colleges")
var1
```

Below there are plots of private and public data sets after being appended, respectively.

Private:



Public:



Q2)

i) We use the load() function to load the marketing.RData dataset, and use the dim() function to view the shape of the dataset. The summary function below gives us the statistical overview of the dataset.

Next we shall use the na.omit() function that removes all incomplete cases of the marketing data frame. Next when we use the dim() we see that the shape is now changed with 6876 records and 14 variables.

```
> library(randomForest)
> library("rpart") #install.packages("rpart")
> load("marketing.RData")
>
> dim(marketing) # 8993 x 14
[1] 8993 14
> summary(marketing)
   Income      Sex      Marital      Age      Edu      Occupation      Lived      Dual_Income      Household      Household18 
Min. :1.000  Min. :1.000 
1st Qu.:2.000 1st Qu.:1.000 1st Qu.:1.000 1st Qu.:2.000 1st Qu.:3.000 1st Qu.:1.000 1st Qu.:4.000 1st Qu.:1.000 1st Qu.:2.000 1st Qu.:0.0000 
Median :5.000  Median :2.000  Median :3.000  Median :3.000  Median :4.000  Median :4.000  Median :5.000  Median :1.000  Median :3.000  Median :0.0000 
Mean   :4.895  Mean   :1.547  Mean   :3.031  Mean   :3.415  Mean   :3.835  Mean   :3.788  Mean   :4.198  Mean   :1.545  Mean   :2.852  Mean   :0.6669 
3rd Qu.:7.000 3rd Qu.:2.000 3rd Qu.:5.000 3rd Qu.:4.000 3rd Qu.:5.000 3rd Qu.:6.000 3rd Qu.:5.000 3rd Qu.:2.000 3rd Qu.:4.000 3rd Qu.:1.0000 
Max.   :9.000  Max.   :2.000  Max.   :5.000  Max.   :7.000  Max.   :6.000  Max.   :9.000  Max.   :5.000  Max.   :3.000  Max.   :9.000  Max.   :9.0000 
NA's    :160     NA's    :86     NA's    :136    NA's    :913     NA's    :375    
   Status      Home_Type      Ethnic      Language
Min. :1.000  Min. :1.000  Min. :1.000  Min. :1.000 
1st Qu.:1.000 1st Qu.:1.000 1st Qu.:5.000 1st Qu.:1.000 
Median :2.000  Median :1.000  Median :7.000  Median :1.000 
Mean   :1.837  Mean   :1.856  Mean   :5.956  Mean   :1.127 
3rd Qu.:2.000 3rd Qu.:3.000 3rd Qu.:7.000 3rd Qu.:1.000 
Max.   :3.000  Max.   :5.000  Max.   :8.000  Max.   :3.000 
NA's    :240     NA's    :357    NA's    :68     NA's    :359    
>
> marketing <- na.omit(marketing)
> dim(marketing) # 6876 x 14
[1] 6876 14
>
> randomized_data <- marketing # for question i)
> permuted_data <- marketing # for question ii)
```

Here we use the sample() function to generate randomised values.

```
for (i in 1:14) {
  min_val = min(marketing[,i])
  max_val = max(marketing[,i])
  randomized_data[,i] = floor(runif(nrow(marketing), min=min_val, max=max_val + 1))
  permuted_data[,i] = sample(min_val: max_val, nrow(marketing), replace = TRUE)
}

marketing$result <- 1
randomized_data$result <- 0
permuted_data$result <- 0

new_data_1 <- rbind(marketing, randomized_data)
new_data_2 <- rbind(marketing, permuted_data)
```

Here we keep the minsplit as 350 and we get the following classification tree when plotted using plot() function.

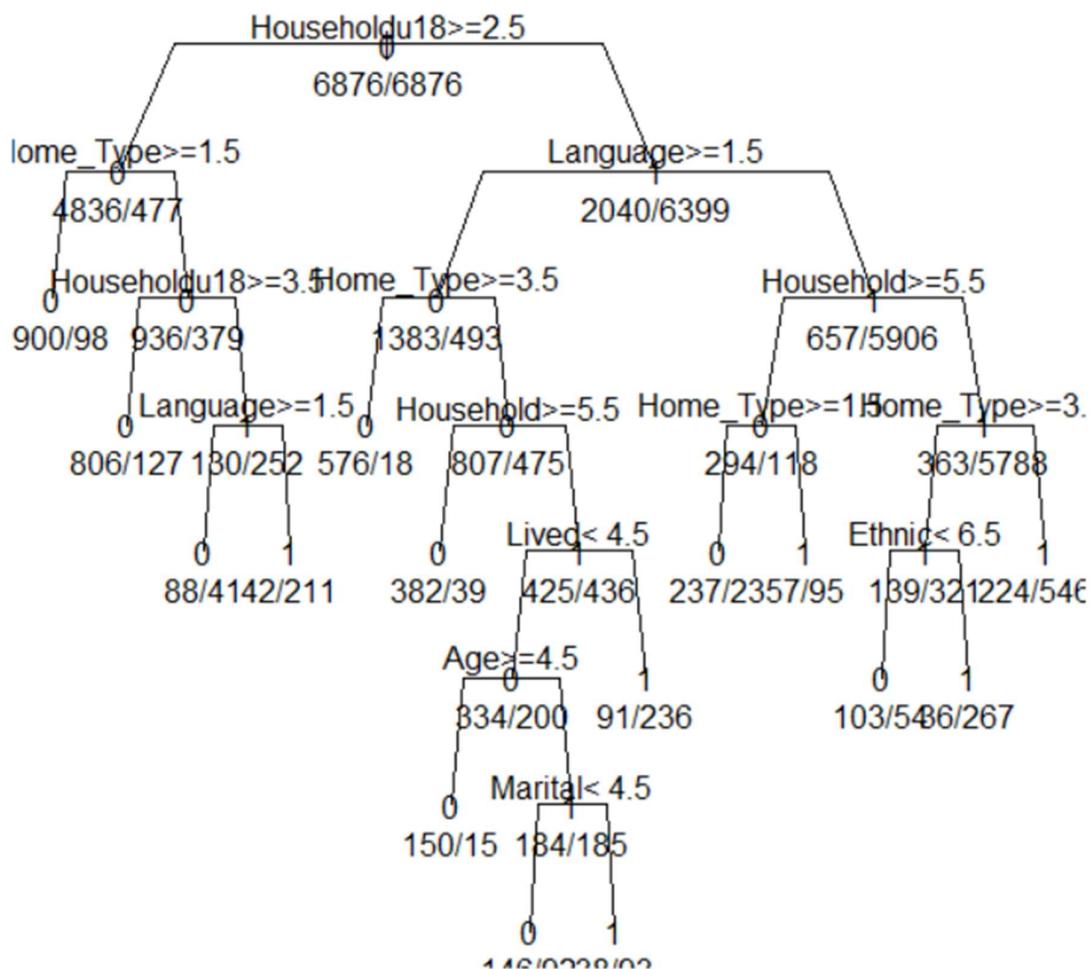
```

model <- rpart.control(minsplit = 350, xval = 40, cp = 0)
classif <- rpart(result~, data = new_data_1, method = "class", control = model)

x11()
plot(classif, branch = .8, uniform = T, compress = T, main="Classification Tree")
text(classif, use.n = T, all = T, cex = 1)
summary(classif )

```

Classification Tree

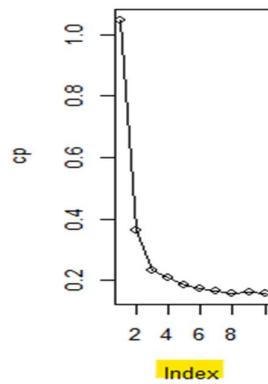


```

plot(classif$cptable[,4], type = "o", main = "Cp for Model selection", ylab="cp")

```

Cp for Model Selection

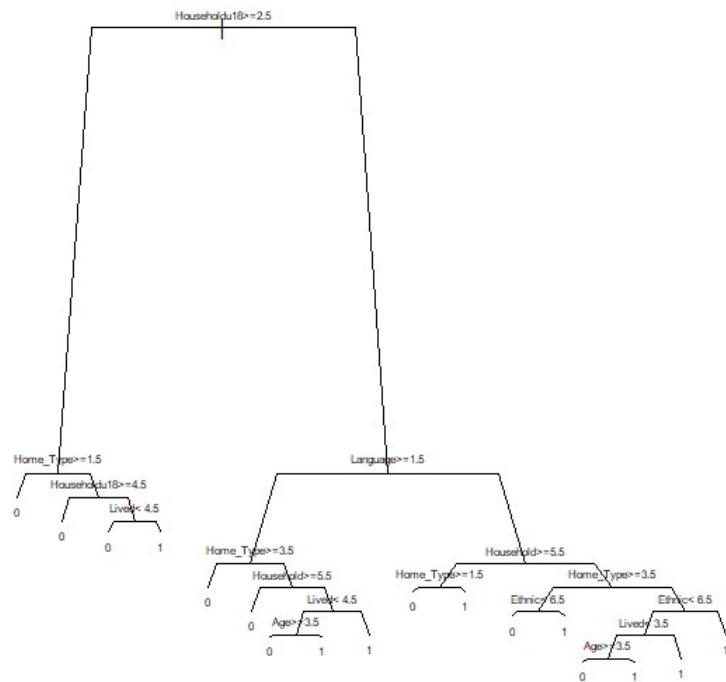


Next we generate a pruned tree

```
# Pruning
min_cp = which.min(classif$cptable[,5])
pruned_fit <- prune(classif, cp = classif$cptable[min_cp,1])

x11()
plot(pruned_fit, branch = 0.8, compress = T, main = "Pruned Tree")
text(pruned_fit, cex = 0.5)
```

Pruned Tree



ii)

```
model.control <- rpart.control(minsplit = 5, xval = 10, cp = 0)
classif <- rpart(result~, data = new_data_2, method = "class", control = model.control)

x11()
plot(classif, branch = .4, uniform = T, compress = T, main="Classification Tree")
text(classif, use.n = T, all = T, cex = 1)

# Pruning
min_cp = which.min(classif$cptable[,4])
pruned_fit <- prune(classif, cp = classif$cptable[min_cp,1])

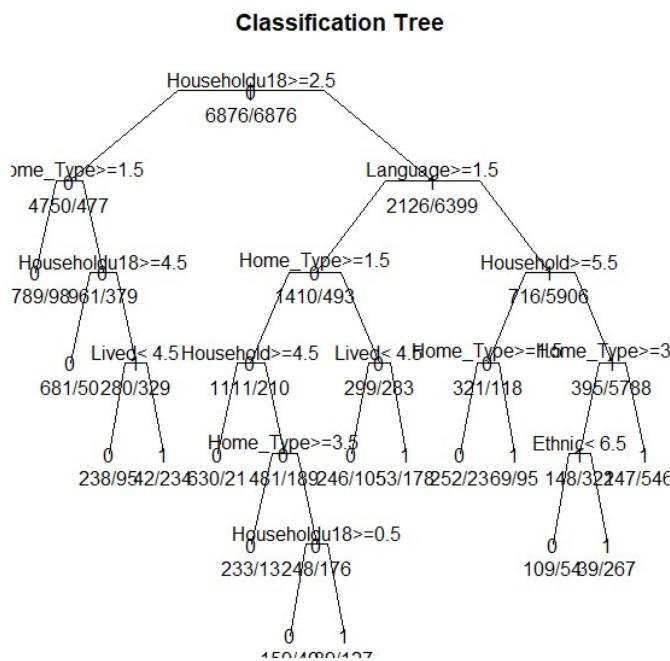
x11()
plot(pruned_fit, branch = 0.3, compress = T, main = "Pruned Tree")
text(pruned_fit, cex = 0.5)



---

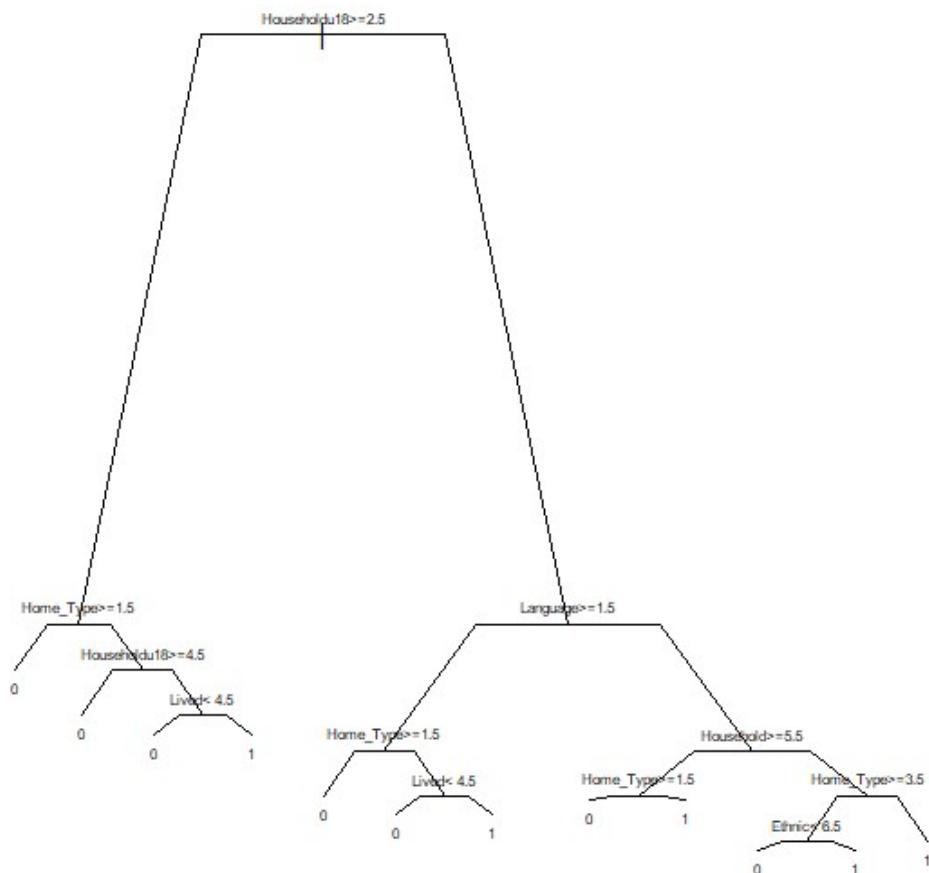

model.control <- rpart.control(minsplit = 350, xval = 10, cp = 0)
classif <- rpart(result~, data = new_data_2, method = "class", control = model.control)

x11()
plot(classif, branch = .4, uniform = T, compress = T, main="Classification Tree")
text(classif, use.n = T, all = T, cex = 1)
```



```
# Pruning  
min_cp = which.min(classif$cptable[,4])  
pruned_fit <- prune(classif, cp = classif$cptable[min_cp,1])  
  
x11()  
plot(pruned_fit, branch = 0.5, compress = T, main = "Pruned Tree")  
text(pruned_fit, cex = 0.5)
```

Pruned Tree



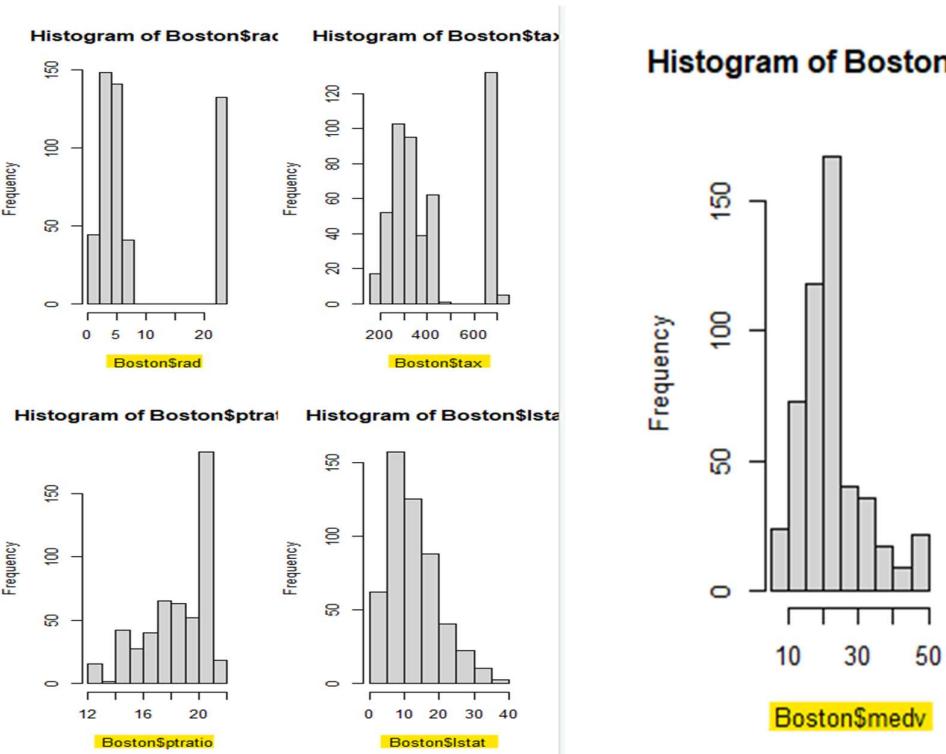
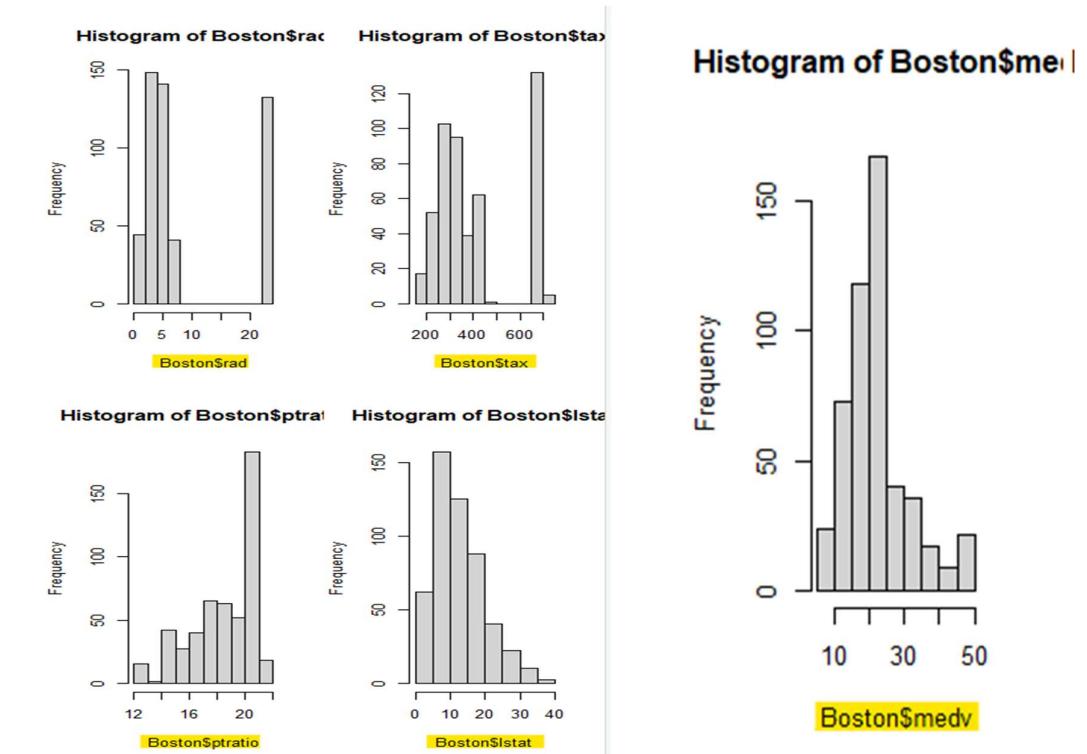
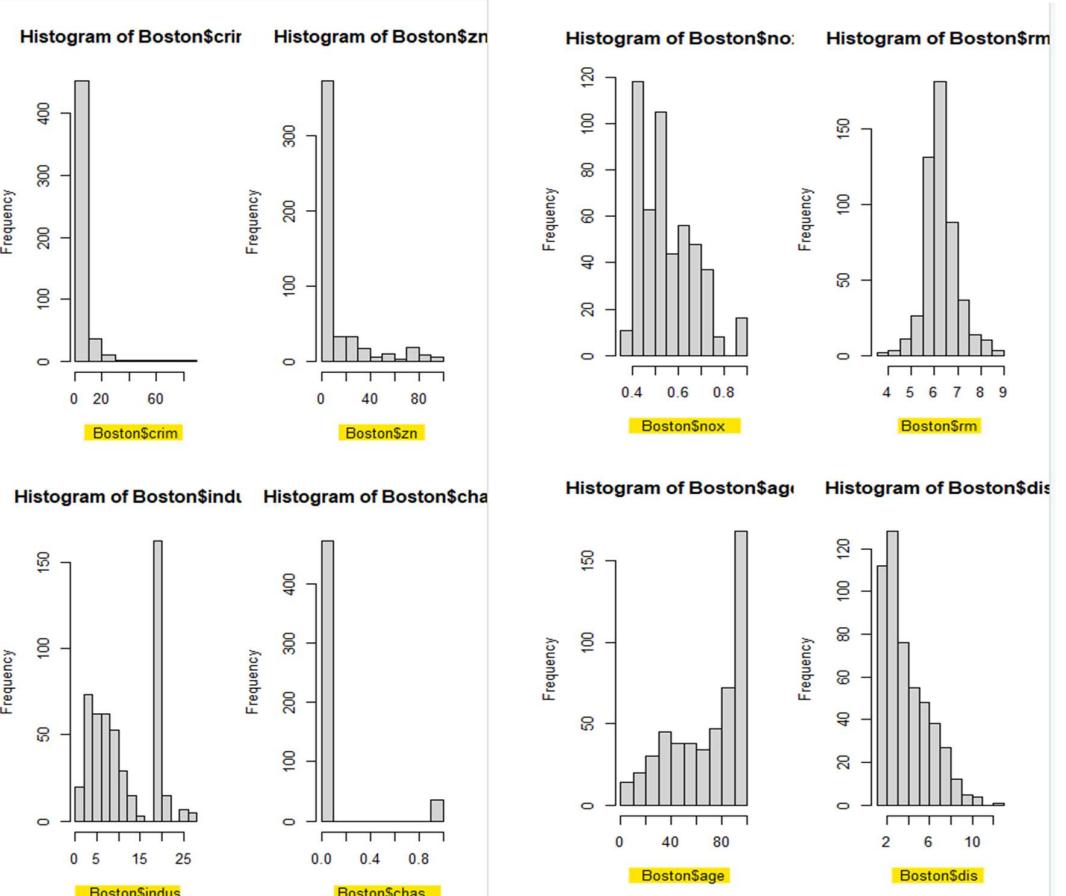
Q3

) a) Here we use setwd() to set the working directory and use the head function to view till tuple number 5, i.e, 6 rows in the dataset. We check the shape of the dataset using the dim() function and we see that the data set is comprised of 506 tuples and 13 variables. We further move forward to store the data set into a dataframe variable.

```
> rm(list=ls())
> library(arules)
> setwd("/Users/sriram/Desktop/R projects/Sem 2/Hw1/")
> library(ISLR2)
> head(Boston)
  crim zn indus chas nox rm age dis rad tax ptratio lstat medv
1 0.00632 18 2.31 0 0.538 6.575 65.2 4.0900 1 296 15.3 4.98 24.0
2 0.02731 0 7.07 0 0.469 6.421 78.9 4.9671 2 242 17.8 9.14 21.6
3 0.02729 0 7.07 0 0.469 7.185 61.1 4.9671 2 242 17.8 4.03 34.7
4 0.03237 0 2.18 0 0.458 6.998 45.8 6.0622 3 222 18.7 2.94 33.4
5 0.06905 0 2.18 0 0.458 7.147 54.2 6.0622 3 222 18.7 5.33 36.2
6 0.02985 0 2.18 0 0.458 6.430 58.7 6.0622 3 222 18.7 5.21 28.7
> dim(Boston)
[1] 506 13
> ?Boston
> Boston_data<-Boston
```

Next we shall view all the histograms of the variables in the dataframe, using the hist() function.

```
par(mfrow = c(2,2))
hist(Boston$crim)
hist(Boston$zn)
hist(Boston$indus)
hist(Boston$chas)
hist(Boston$nox)
hist(Boston$rm)
hist(Boston$age)
hist(Boston$dis)
hist(Boston$rad)
hist(Boston$tax)
hist(Boston$ptratio)
hist(Boston$lstat)
hist(Boston$medv)
|
```



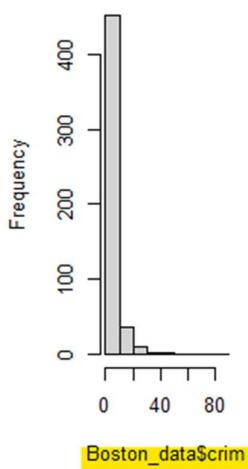
Next we plot box plots of each variables and add labels to it. Below you can see the boxplot's and hist() of the following variable.

```

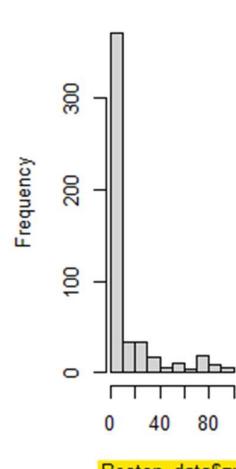
> hist(Boston_data$crim)
>
> boxplot.stats(Boston_data$crim)$stats
[1] 0.00632 0.08199 0.25651 3.67822 8.98296
> Boston_data$crim = ordered(cut(Boston_data[["crim"]],c(10,20,40,60,100)), labels = c("Low-Threat", "Medium-Threat", "High-alert", "CodeRed"))
>
>
> hist(Boston_data$zn)
> Boston_data[["zn"]]<- ordered (cut(Boston_data[["zn"]], c (0,15,50,100)), labels= c("no-zone", "less-zone", "High-Zone"))
>
> hist(Boston_data$indus)
> boxplot.stats(Boston_data$indus)$stats
[1] 0.46 5.19 9.69 18.10 27.74
> Boston_data$indus = ordered(cut(Boston_data[["indus"]],c(0,10,20,40)), labels = c('Low', 'Moderate', 'High'))
>
> hist(Boston_data$chas)
> Boston_data[["chas"]]<-ordered(Boston_data$chas,labels=c(" Not bound to river", "bound to river"))
>
>
> hist(Boston_data$nox)
> boxplot.stats(Boston_data$nox)$stats
[1] 0.385 0.449 0.538 0.624 0.871
> Boston_data$nox = ordered(cut(Boston_data[["nox"]],c(0.1,0.3,0.4,0.8)), labels = c('Low-Concentration', 'Moderate-Concentration', 'High-Concentration'))
Error in factor(x, ..., ordered = TRUE) :
  invalid 'labels'; length 3 should be 1 or 2
>
>
> hist(Boston_data$rm)
> boxplot.stats(Boston_data$rm)$stats
[1] 4.8800 5.8850 6.2085 6.6250 7.6910
> Boston_data$rm = ordered(cut(Boston_data[["rm"]],c(0,4,7,9)), labels = c('small houses', 'Medium Houses', 'Big Houses'))
>
>
> hist(Boston_data$age)
> boxplot.stats(Boston_data$age)$stats
[1] 2.9 45.0 77.5 94.1 100.0
> Boston_data$age = ordered(cut(Boston_data[["age"]],c(0,25,50,70,100)), labels = c('silver', 'gold', 'diamond', 'too old'))

```

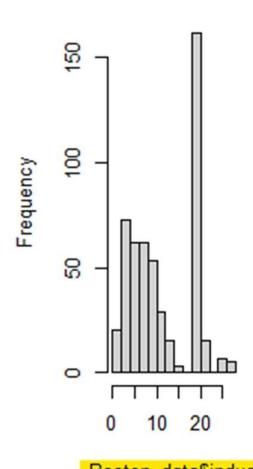
Histogram of Boston_data\$crim



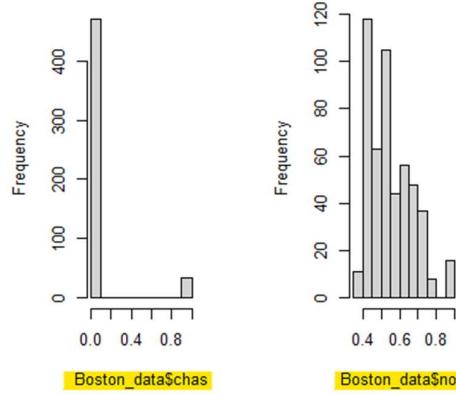
Histogram of Boston_data\$zn



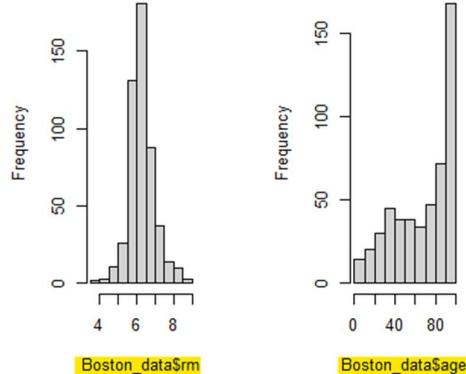
Histogram of Boston_data\$indus



Histogram of Boston_data\$ dis Histogram of Boston_data\$ rad



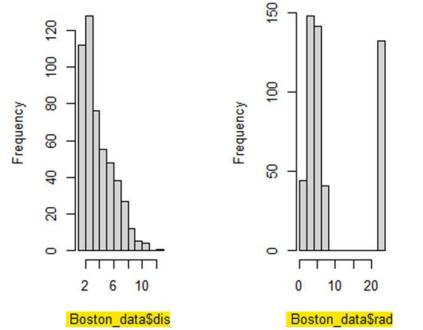
Histogram of Boston_data\$ tax Histogram of Boston_data\$ ptratio



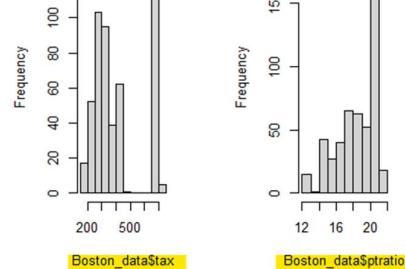
```
Boston_data$age = ordered(cut(Boston_data[["age"]],c(0,25,50,70,100)), labels = c("silver","gold","diamond","too old"))
levels: silver < gold < diamond < too old

hist(Boston_data$dis)
boxplot.stats(Boston_data$dis)$stats
1] 1.12960 2.10000 3.20745 5.21190 9.22290
Boston_data$dis = ordered(cut(Boston_data[["dis"]],c(0,4,7,14)), labels = c('Near','Far','Very Far'))
hist(Boston_data$rad)
boxplot.stats(Boston_data$rad)$stats
1] 1 4 5 24 24
Boston_data$rad = ordered(cut(Boston_data[["rad"]],c(0,2,10,26)), labels = c('Less Accessible','Moderately Accessible','Highly Accessible'))
hist(Boston_data$tax)
boxplot.stats(Boston_data$tax)$stats
1] 187 279 330 666 711
Boston_data$tax = ordered(cut(Boston_data[["tax"]],c(187,279,660,711)), labels = c('Low','Moderate','High'))
hist(Boston_data$ptratio)
boxplot.stats(Boston_data$ptratio)$stats
1] 13.60 17.40 19.05 20.20 22.00
Boston_data$ptratio = ordered(cut(Boston_data[["ptratio"]],c(10,19,26)), labels = c('Low','High'))
```

Histogram of Boston_data\$ dis Histogram of Boston_data\$ rad



Histogram of Boston_data\$ tax Histogram of Boston_data\$ ptratio

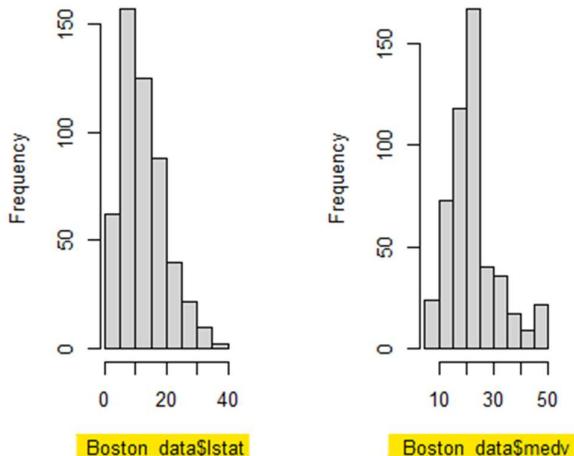


```

hist(Boston_data$lstat)
boxplot.stats(Boston_data$lstat)$stats
[1] 1.73 6.93 11.36 16.96 31.99
Boston_data$lstat = ordered(cut(Boston_data[["lstat"]],c(0,12,40)), labels = c('Low','High'))
hist(Boston_data$medv)
boxplot.stats(Boston_data$medv)$stats
[1] 5.0 17.0 21.2 25.0 37.0
Boston_data$medv = ordered(cut(Boston_data[["medv"]],c(0,22,51)), labels = c('Low','High'))

```

Histogram of Boston_data\$lstat Histogram of Boston_data\$medv



- b) Next we use the apriori() in order to calculate the apriori rules. We see that here we obtain 50 support counts.

```

> incidence_matrix = as(Boston_data,"transactions")
warning message:
Column(s) 5 not logical or factor. Applying default discretization (see '? discretizedDF').
> incidence_matrix
transactions in sparse format with
 506 transactions (rows) and
 37 items (columns)
> itemFrequencyPlot(incidence_matrix,support=0.001,cex.name= 0.8)
> rules<-apriori(incidence_matrix,parameter= list(support = 0.1,confidence = 0.3))
Apriori

Parameter specification:
confidence minval smax arem  aval originalSupport maxtime support minlen maxlen target ext
      0.3     0.1     1 none FALSE           TRUE      5     0.1     1    10   rules TRUE

Algorithmic control:
filter tree heap memopt load sort verbose
  0.1 TRUE TRUE FALSE TRUE     2     TRUE

Absolute minimum support count: 50

set item appearances ...[0 item(s)] done [0.00s].
set transactions ...[37 item(s), 506 transaction(s)] done [0.00s].
sorting and recoding items ... [25 item(s)] done [0.00s].
creating transaction tree ... done [0.00s].
checking subsets of size 1 2 3 4 5 6 7 8 9 10 done [0.00s].
writing ... [17755 rule(s)] done [0.01s].
creating s4 object ... done [0.01s].

```

The summary() function is used to view the statistical description of the apriori rules.

```

> summary(rules)
set of 17755 rules

rule length distribution (lhs + rhs):sizes
  1   2   3   4   5   6   7   8   9   10
17 311 1616 3564 4283 3588 2443 1328 495 110

  Min. 1st Qu. Median  Mean 3rd Qu. Max.
1.000 4.000 5.000 5.405 6.000 10.000

summary of quality measures:
  support confidence coverage lift count
Min. :0.1008 Min. :0.3000 Min. :0.1008 Min. :0.4699 Min. : 51.00
1st Qu.:0.1522 1st Qu.:0.8104 1st Qu.:0.1838 1st Qu.:1.1500 1st Qu.: 77.00
Median :0.1877 Median :0.9277 Median :0.2095 Median :1.6382 Median : 95.00
Mean :0.1871 Mean :0.8597 Mean :0.2266 Mean :1.7839 Mean : 94.69
3rd Qu.:0.2115 3rd Qu.:0.9778 3rd Qu.:0.2431 3rd Qu.:2.0220 3rd Qu.:107.00
Max. :0.9308 Max. :1.0000 Max. :1.0000 Max. :3.8333 Max. :471.00

mining info:
  data ntransactions support confidence
incidence_matrix      506       0.1          0.3 apriori(data = incidence_matrix, parameter = list(support = 0.1, confidence = 0.3)) call

```

c) Next we see the rules of medv being low with respect to dis being near, where lift is greater than 1.2 . We also use the inspect() function to view the support confidence and lift values, and head() function with n =3, i.e., to view the first 3 records .

```

> #3c
> rules_medv = subset(rules, subset = lhs %in% "dis=Near" & rhs %in% "medv=Low" & lift>1.2)
> inspect(head(sort(rules_medv, by = "confidence"), n = 3))
  lhs
[1] {rm=Medium Houses, age=Too old, dis=Near, ptratio=High, lstat=High}
[2] {chas= Not bound to river, rm=Medium Houses, age=Too old, dis=Near, ptratio=High, lstat=High}
[3] {nox=[0.597,0.871], rm=Medium Houses, dis=Near, ptratio=High, lstat=High}
>

  rhs      support  confidence coverage  lift  count
=> {medv=Low} 0.2786561 0.9657534 0.2885375 1.720673 141
=> {medv=Low} 0.2727273 0.9650350 0.2826087 1.719393 138
=> {medv=Low} 0.2173913 0.9649123 0.2252964 1.719175 110

```

d)

```

> #3d
> rules_ptratio_low<-subset(rules,subset= rhs %in% "ptratio=Low" & lift>1.2)
> inspect(head(sort(rules_ptratio_low, by="confidence"),n=3))
  lhs      rhs      support  confidence coverage  lift  count
[1] {indus=Moderate, dis=Near, rad=Moderately Accesible} => {ptratio=Low} 0.1067194 1.0000000 0.1067194 2.000000 54
[2] {indus=Low, rm=Big Houses}                         => {ptratio=Low} 0.1047431 0.9636364 0.1086957 1.927273 53
[3] {indus=Low, rm=Big Houses, medv=High}              => {ptratio=Low} 0.1047431 0.9636364 0.1086957 1.927273 53
>

```