

Experiment - 09

AIM :- Analysing Boston Data.

Commands :-

- > install.packages ("tree")
- > library(MASS)
- > options (super = 4)
- > view (Boston)

O/p :-

It shows 1 to 11 of 506 entries, 14 total columns in Boston char on age---

-
- > set.seed(1)
 - > train <- sample(1:nrow(Boston), nrow(Boston)/2)

O/P

Values:-

- train int [1:253] 505 234 167 129 ...
- > tree.boston <- tree (medv ~., Boston, subset = train)
- > install.packages ("tree")
- > tree.boston <- tree (medv ~ . Boston, subset = train)

O/p

- > tree.boston data.frame 130 obs of 5 variables
- > summary (tree.boston)

O/P

Regression tree.

tree (formula = medv ~ ., Boston, subset = train)

Variables actually used in tree construction
[1] "rm" "lstat" "tax" "dis" "age"

Number of terminal nodes : 7

Residual mean deviance: $10.38 = 2555/246$

Distribution of Residuals

Min	1 st Q	Median	Mean	3 rd Q	Max
-10.180	-1.77	-0.1775	0.000	1.9230	16.580

→ Run cv to find best level at which to prune

> cv.boston <- cv.tree(tree.Boston)

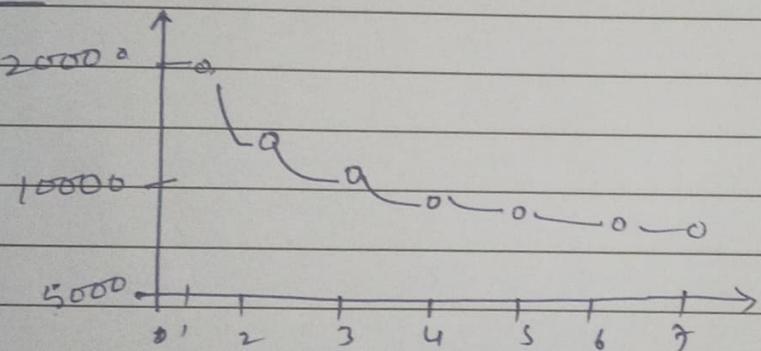
O/P

cv.boston list of 4

→ Construct a plot (dev = MSE on y-axis)

plot(cv.boston \$ size, cv.boston \$ dev, type='b')

O/P :-



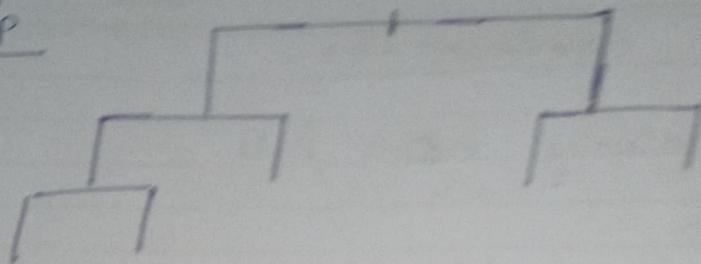
> prune.boston <- prune.tree(tree.boston, best=5)

O/P

prune.boston list of 6

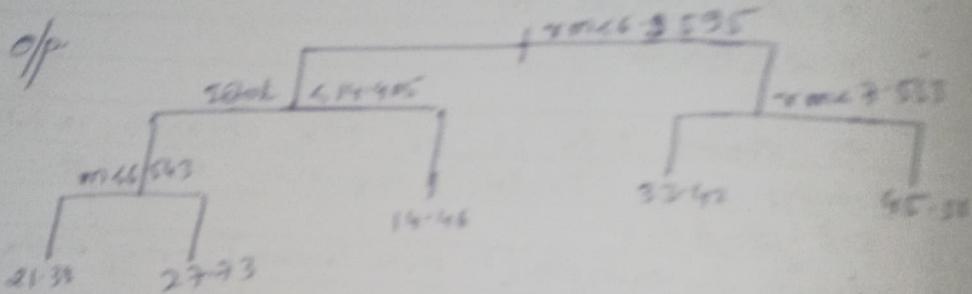
> plot(prune.boston)

o/p



> text (pruvia.boston, pretty=0)

o/p



Experiment - 10

AIM :- To given dataset (electrical consumption) of an organization it contains the monthly electrical consumption & annual average for various years

Program :-

Plot the various across the years . Is the consumption pattern comparable across the years ?

```
> df <- read.csv ("C:/Users/IRahul/Desktop")
print(df)
> install.packages("ggplot2")
> x = df[,1]
> y = df[,2]
plot(x)
plot(y)
> x = df[1:5,]
> y = df[1,4]
plot(x)
> plot(x, y, col='blue', xlab="years", ylab="Jan")
> df = (c("Jan", "Feb", "March", "April", "May", "June", "July", "August",
  September, October, November, December"))
> plot(x, df, col="red")
> cv.boston <- cv.tree(tree.Boston)
> plot(cv.boston <- prune.tree(tree.boston, best=5))
> plot(prune.boston)
> text(prune.boston, pretty=0)
```

Experiment-11

AIM :- Classification on Boston Dataset

Program :-

```

> library(MASS)
> installed.packages("tree")
> library(tree)
> options(scipar=4)
> view(Boston)
> set.seed(1)
> train <- sample(1:nrow(Boston), nrow(Boston)/2)
> tree.Boston <- tree(medv ~ ., Boston, subset=train)
> summary(tree.Boston)

```

O/P

Regression tree

tree(formula = medv ~ ., data = Boston, subset = train)

variables actually used in tree construction

[1] "rm" "Istat" "lrim" "age"

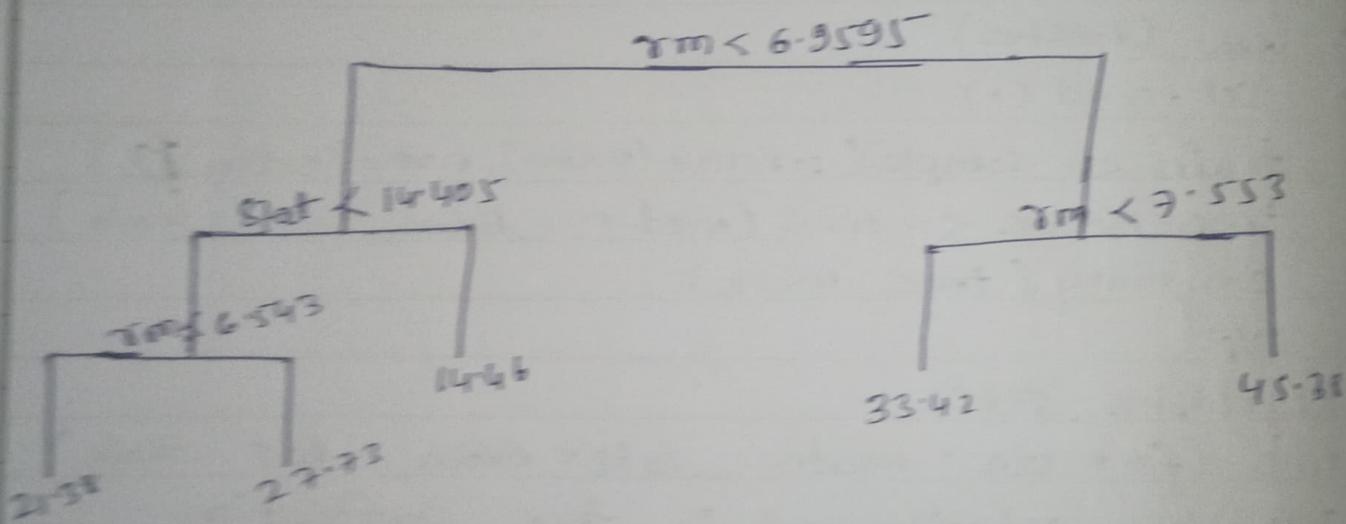
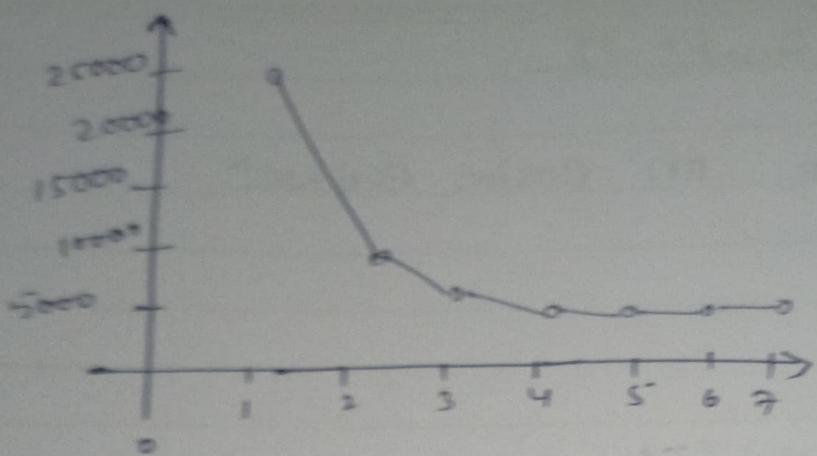
No. of terminal nodes: 7

Residual mean deviance: 10.38 - 2555/246

Distribution of residuals.

Min	1st Qu	Median	Mean	3rd Qu	Max
-10.080	-1.770	-0.1775	0.000	1.9830	16.580

→ Run CV (cross validation) to find best level at which to prune.



Km out list of 9

> table (Km.out of cluster, class, 16)

-	1	1	2
1	2	2	
2	4	23	

Experimental - 12

AIM :- To Learn & practice more ~~about~~ about K-means clustering.

Programs :-

```
> library ("ggplot2")
> set.seed(2)
> x <- matrix(rnorm(50*2), ncol=2)
> gplot (x=x[,1], y=x[,2])
> x[1:25,1] <- x[1:25,1]+1
> x[1:25,2] <- x[1:25,2]-1
> gplot (x=x[,1], y=x[,2])
```

Shift the 1st 25 pts to have mean (1,-1) instead of (0,0)

```
> x[1:25,1] <- x[1:25,1]+4
> x[1:25,2] <- x[1:25,2]-4
> gplot (x=x[,1], y=x[,2])
> class.1b1 <- as.factor(c(rep(1,25), rep(2,25)))
```

class. 1b1 factor w/ 2 levels "1", "2" : 11111...

```
> gplot(x=x[,1], y=x[,2], color = class.1b1,
       size = I(3)) + theme_bw()
> km.out <- kmeans(x[,2], nstart=20)
```

Experiment - 13

AIM :- Hadoop Implementation on Single Node, on Windows

Program :-

In Qall Java, Hadoop :-

Step 1:-

→ Verify the Java installed :-\$ javac -version

→ Step 2 :- Extract Hadoop at C:\Hadoop

→ Step 3:- setting up the Hadoop_Home Variable by using Windows environment variable. setting for Hadoop Path setting.

Step 4:- Set Java_Home Variable.

Use Windows enviornment variable. setting for Java Path setting.

Step 5:- Set Hadoop & Java bin directory Path.
to copy the path & of the file & create a new in edit environment variable.

Step 6 :- Hadoop configuration

for Hadoop configuration we need to modify six files that are below:-

- (1) core-site.xml (2) Mapred-site.xml
- (3) HDFS-site.xml (4) YARN-site.xml.

Hadoop-env.cmd

Create two folder datanode & namenode.

Step 7:- Format the namenode folder :-

Open command window & give the code
\$ hdfs namenode -format

→ This will show that we have successfully
install Hadoop.

Step 8:- Testing The setup :-(Running the Hadoop)-

\$ start-all.cmd

& start-dfs.cmd and \$ start-yarn.cmd.

Step 9:-

Output :- localhost active:-

Started: Tue Jan 12 14:15:15 ICT 2023

Version: 2.7.1

compiled: 2023-01-17 by Rahul Yadav

cluster ID: CFD 9.2636

Block Pool ID: BP-27..

Step 8 :- Stop Hadoop :-

\$ stop-dfs.cmd.

Experiment-14

AIM :- Visualization & modification on Diabetes dataset .

Program :-

① determine the structure of the data (

~~diabetes =~~ `read.csv ("diabetes.csv")`

~~summary (diabetes)~~

`str (diabetes)` & `head (diabetes)`

`Mean_Glucose ← mean (diabetes $ Glucose)`

`median_Glucose ← median (diabetes $ Glucose)`

② ~~sd_Glucose ← sd (diabetes \$ Glucose)~~

~~var_Glucose ← var (diabetes \$ Glucose)~~

~~numeric.var ← supply~~

③ Perform correlation against each other (install corrplot)

→ `install.packages ('corrplot')`

→ `library ('corrplot')`

→ `numeric.var <- supply (diabetes , is.numeric)`

→ `Corr.matrix <- cor (diabetes [numeric.var])`

→ `corrplot (corr.matrix , main = "Incorrelation Plot for"`

→ `"Numerical variables" , order = "h.cust" , tl.col = "black"`

→ `, tl.cex = 0.8 , cl.cex = 0.8)`

→ `box (which = "outer" , lty = "solid")`

Decision Tree:- `library (diabetes)`

`model.1 ← spt (outcome ~ Pressure + Glucose + BMI`

`data = tree , method = "jtree")`

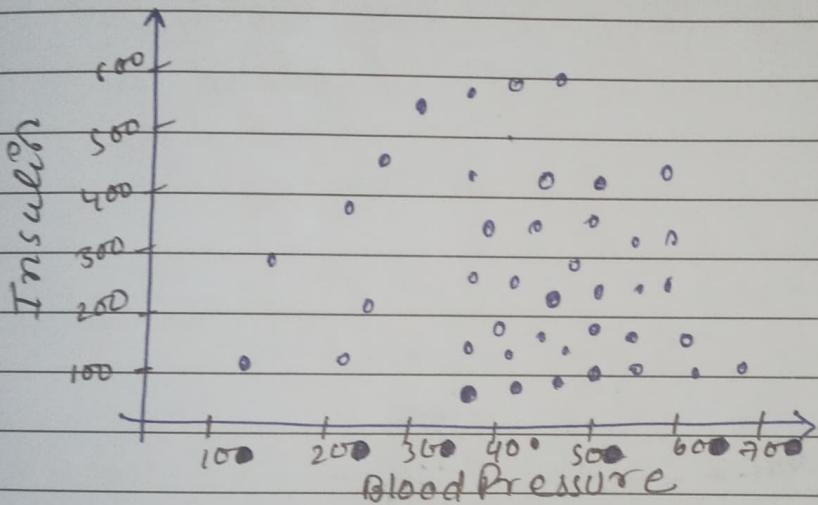
~~Decision Tree:-~~ plot (model1, $\text{vm} \rightarrow \text{true}$)
 box (which = 'outer', $\text{Ity} = \text{'solid'}$).

Scatterplot :-

`sns.set_style('darkgrid')`

`grid = sns.JointGrid(x = "Blood Pressure", y = "Insulin",
 data = diabetes);`

`grid.plot(sns.scatterplot)` ~~sns.jointplot~~



OUTPUT:-

74.56% of accuracy in ~~set~~ decision tree

The data has various parameters & decision tree is used because it can consider all these parameters as target & its contribution tell us which comes literally off diabetes in set class.

Experiment-15

AIM :- Visualize & Analyze on cat data set in R.

Program :-

```
library(MASS)
data(cats, package = 'MASS')
str(cats)
head(cats)
library(lattice)
```

```
myplot(Hwt ~ Bwt | sex, data = cats, style = c("p", "g", "l"),
xlab = "Baby cat (kg)", ylab = "Hect by (J)")
```

```
x.y Plot(Hwt ~ Bwt, q.b = sex, data = cats, style = c("p", "g", "l"),
xlab = "using(keyword)", plot = "Bwt (group), as = True).)
```

→ install.packages("e1071") . → library("e1071")

```
m <- svm(sex ~, data = cats)
```

```
:fcat(m, data = cats)
```

```
summary(m)
```

```
data.frame(plot(m), cats & sex)
```

```
mt <- svm(sex ~, data = cats, gamma = 0.5, cost = 678)
plot(sex ~, data = cats)
```

```
table(true = cats & sex, future = predict(x2))
```

OUTPUT:-

- The initial unpruned model gave a function accuracy of about 53%.
- The pruned model with Gamma = 0.5 & cost = 678 had a particular of 98.9% accuracy.
- No SVM is not suitable. for these costs, sex departs the weights of height weight & size. Therefore we need to use a multiple regression model which ~~will~~ ^{which} assist ~~to~~ less outliers.

Output Prediction

TRUE	F	M
F	33	14
M	11	86.

