

Data points are plotted in the main area of the window. At the top are two drop-down list buttons for selecting the axes to plot. The one on the left shows which attribute is used for the x-axis; the one on the right shows which is used for the y-axis.

Beneath the x-axis selector is a drop-down list for choosing the colour scheme. This allows you to colour the points based on the attribute selected. Below the plot area, a legend describes what values the colours correspond to. If the values are discrete, you can modify the colour used for each one by clicking on them and making an appropriate selection in the window that pops up.

To the right of the plot area is a series of horizontal strips. Each strip represents an attribute, and the dots within it show the distribution of values of the attribute. These values are randomly scattered vertically to help you see concentrations of points. You can choose what axes are used in the main graph by clicking on these strips. Left-clicking an attribute strip changes the x-axis to that attribute, whereas right-clicking changes the y-axis. The ‘ X’ and ‘ Y’ written beside the strips shows what the current axes are (‘ B’ is used for ‘ both X and Y’).

Above the attribute strips is a slider labelled Jitter, which is a random displacement given to all points in the plot. Dragging it to the right increases the amount of jitter, which is useful for spotting concentrations of points. Without jitter, a million instances at the same point would look no different to just a single lonely instance.

9.3 Selecting Instances

There may be situations where it is helpful to select a subset of the data using the visualization tool. (A special case of this is the User Classifier in the Classify panel, which lets you build your own classifier by interactively selecting instances.)

Below the y-axis selector button is a drop-down list button for choosing a selection method. A group of data points can be selected in four ways:

- 1. Select Instance.** Clicking on an individual data point brings up a window listing its attributes. If more than one point appears at the same location, more than one set of attributes is shown.
- 2. Rectangle.** You can create a rectangle, by dragging, that selects the points inside it.
- 3. Polygon.** You can build a free-form polygon that selects the points inside it. Left-click to add vertices to the polygon, right-click to complete it. The polygon will always be closed off by connecting the first point to the last.
- 4. Polyline.** You can build a polyline that distinguishes the points on one side from those on the other. Left-click to add vertices to the polyline, right-click to finish. The resulting shape is open (as opposed to a polygon, which is always closed).

Once an area of the plot has been selected using Rectangle, Polygon or Polyline, it turns grey. At this point, clicking the Submit button removes all instances from the plot except those within the grey selection area. Clicking on the Clear button erases the selected area without affecting the graph.

Once any points have been removed from the graph, the Submit button changes to a Reset button. This button undoes all previous removals and returns you to the original graph with all points included. Finally, clicking the Save button allows you to save the currently visible instances to a new ARFF file.

Part-A: Data Mining

1. Demonstration of preprocessing on dataset student.arff

Aim: This experiment illustrates some of the basic data preprocessing operations that can be performed using WEKA-Explorer. The sample dataset used for this example is the student data available in arff format.

Weka - Preprocessing the Data. The data that is collected from the field contains many unwanted things that leads to wrong analysis. For example, the data may contain null fields, it may contain columns that are irrelevant to the current analysis, and so on.

Step1: Loading the data. We can load the dataset into weka by clicking on open button in preprocessing interface and selecting the appropriate file.

Step2: Once the data is loaded, weka will recognize the attributes and during the scan of the data weka will compute some basic strategies on each attribute. The left panel in the above figure shows the list of recognized attributes while the top panel indicates the names of the base relation or table and the current working relation (which are same initially).

Step3: Clicking on an attribute in the left panel will show the basic statistics on the attributes for the categorical attributes the frequency of each attribute value is shown, while for continuous attributes we can obtain min, max, mean, standard deviation and deviation etc.,

Step4: The visualization in the right button panel in the form of cross-tabulation across two attributes.

Note:we can select another attribute using the dropdown list.

Step5:Selecting or filtering attributes Removing an attribute-When we need to remove an attribute,we can do this by using the attribute filters in weka.In the filter model panel,click on choose button,This will show a popup window with a list of available filters. Scroll down the list and select the “ weka.filters.unsupervised.attribute.remove” filters.

Step 6:a)Next click the textbox immediately to the right of the choose button.In the resulting dialog box enter the index of the attribute to be filtered out.

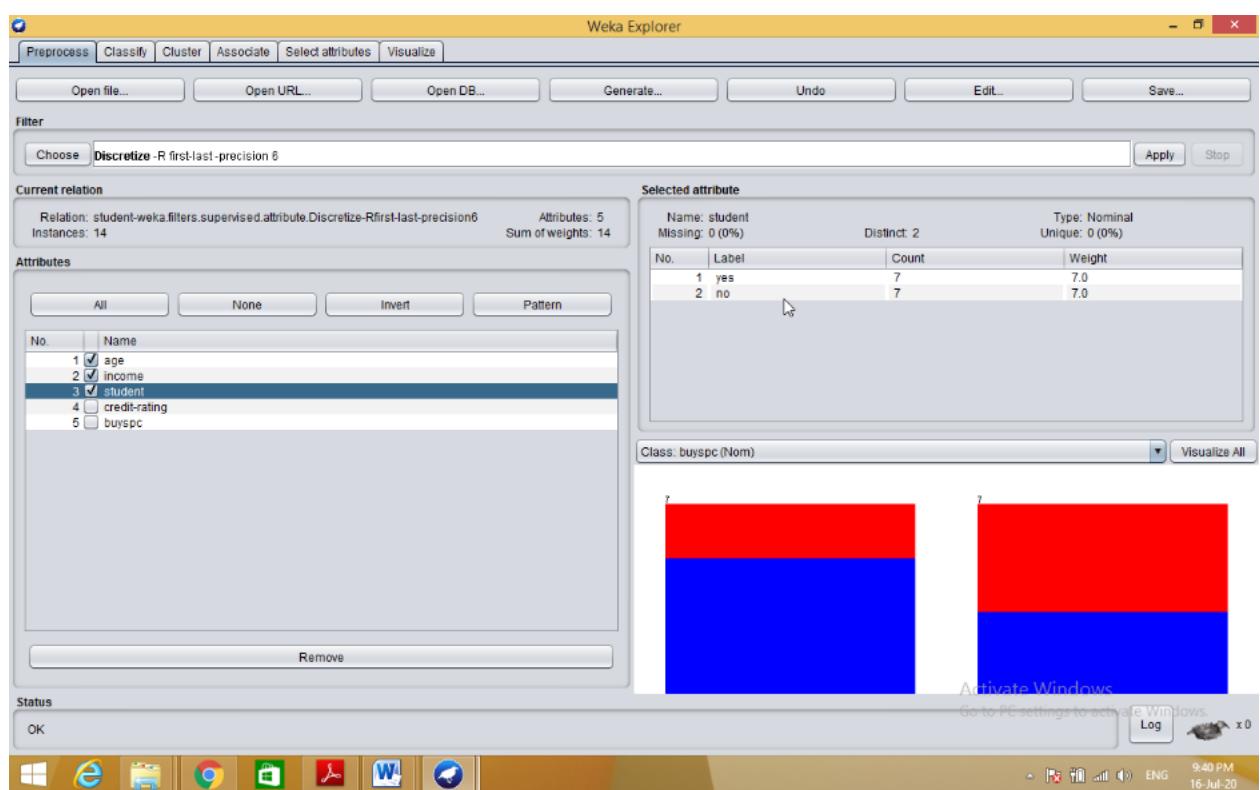
b)Make sure that invert selection option is set to false.The click OK now in the filter box,you will see “ Remove-R-7” .

c)Click the apply button to apply filter to this data.This will remove the attribute and create new working relation.

d)Save the new working relation as an arff file by clicking save button on the top(button)panel.(student.arff)

Procedure:

Choose Tab \square Weka \square filters \square Supervised \square attribute \square Discretize



Dataset student.arff

```
@relation student
@attribute age {<30,30-40,>40}
@attribute income {low, medium, high}
@attribute student {yes, no}
@attribute credit-rating {fair, excellent}
@attribute buyspc {yes, no}
@data
%
<30, high, no, fair, no
<30, high, no, excellent, no
30-40, high, no, fair, yes
>40, medium, no, fair, yes
>40, low, yes, fair, yes
>40, low, yes, excellent, no
30-40, low, yes, excellent, yes
<30, medium, no, fair, no
<30, low, yes, fair, no
>40, medium, yes, fair, yes
<30, medium, yes, excellent, yes
30-40, medium, no, excellent, yes
30-40, high, yes, fair, yes
>40, medium, no, excellent, no
%
```

2. Demonstration of Association rule process on dataset contactlenses.arff using apriori algorithm.

Aim: This experiment illustrates some of the basic elements of association rule mining using WEKA. The sample dataset used for this example is contactlenses.arff

Step1: Open the data file in Weka Explorer. It is presumed that the required data fields have been discretized. In this example it is age attribute.

Step2: Clicking on the associate tab will bring up the interface for association rule algorithm.

Step3: We will use apriori algorithm. This is the default algorithm.

Step4: Inorder to change the parameters for the run (example support, confidence etc) we click on the text box immediately to the right of the choose button.

Description:

The Apriori Algorithm: Finding Frequent Itemsets Using Candidate Generation

Apriori is a seminal algorithm proposed by R. Agrawal and R. Srikant in 1994 for mining frequent itemsets for Boolean association rules. The name of the algorithm is based on the fact that the algorithm uses *prior knowledge* of frequent itemset properties. Apriori employs an iterative approach known as a *level-wise* search, where k -itemsets are used to explore $(k+1)$ -itemsets.

Apriori property: All nonempty subsets of a frequent itemset must also be frequent.

It has two steps.

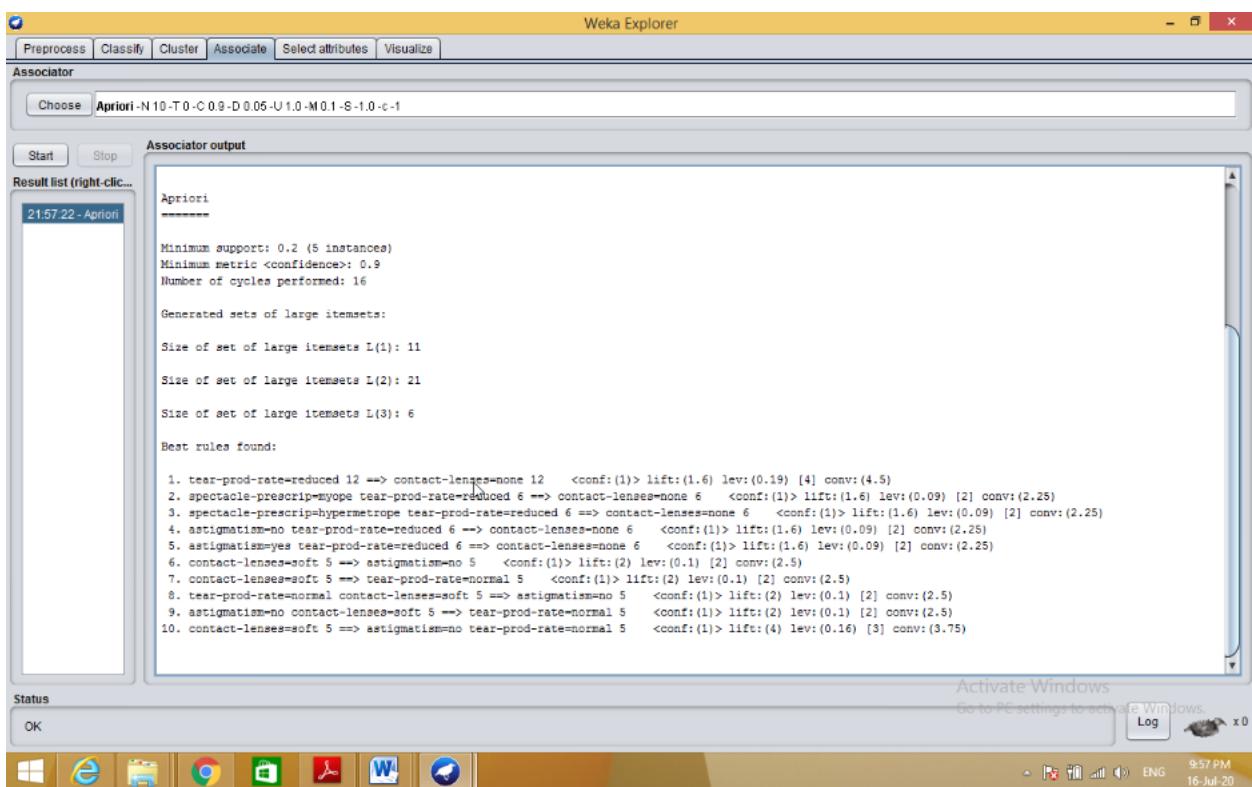
1) Join step

2) Prune step

Procedure:

Step 1: Load appropriate dataset into weka

Step 2: Select associate tab and select apriori algorithm



3. Demonstration of classification rule process on dataset employee.arff using j48 algorithm.

Aim: This experiment illustrates the use of j-48 classifier in weka. The sample data set used in this experiment is “ student” data available at arff format. This document assumes that appropriate data pre processing has been performed.

Steps involved in this experiment:

Step-1: We begin the experiment by loading the data (student.arff)into weka.

Step2: Next we select the “ classify” tab and click “ choose” button to select the “ j48” classifier.

Step3: Now we specify the various parameters. These can be specified by clicking in the text box to the right of the chose button. In this example, we accept the default values. The default version does perform some pruning but does not perform error pruning.

Step4: Under the “ text” options in the main panel. We select the 10-fold cross validation as our evaluation approach. Since we don’ t have separate evaluation data set, this is necessary to get a reasonable idea of accuracy of generated model.

Step-5: We now click " start" to generate the model .the Ascii version of the tree as well as evaluation statistic will appear in the right panel when the model construction is complete.

Step-6: Note that the classification accuracy of model is about 69%.this indicates that we may find more work. (Either in preprocessing or in selecting current parameters for the classification)

Step-7: Now weka also lets us a view a graphical version of the classification tree. This can be done by right clicking the last result set and selecting " visualize tree" from the pop-up menu.

Step-8: We will use our model to classify the new instances.

Step-9: In the main panel under " text" options click the " supplied test set" radio button and then click the " set" button. This will pop-up a window which will allow you to open the file containing test instances.

Description:

Classification is a **data mining** function that assigns items in a collection to target categories or classes. The goal of **classification** is to accurately predict the target class for each case in the **data**. For example, a **classification** model could be used to identify loan applicants as low, medium, or high credit risks.

Entropy characterizes the (im) purity of an arbitrary collection of examples

Information Gain is the expected reduction in entropy caused by partitioning the examples according to a given attribute

If we have a set with k different values in it, we can calculate the entropy as follows:

Where $P(\text{value}_i)$ is the probability of getting the i^{th} value when randomly selecting one from the set. So, for the set $R = \{a,a,a,b,b,b,b\}$

Kappa statistic is a measure of how closely the instances classified by the machine learning classifier matched the **data** labeled as ground truth, controlling for the accuracy of a random classifier as measured by the expected accuracy.

Decision tree builds classification or regression models in the form of a tree structure. It breaks down a dataset into smaller and smaller subsets while at the same time an associated decision tree is incrementally developed. The final result is a tree with **decision nodes** and **leaf nodes**.

Classifier Accuracy Measures

The **confusion matrix** is a useful tool for analyzing how well your classifier can recognize tuples of different classes.

F measure (F1 score or F score) is a **measure** of a test's accuracy and is defined as the weighted harmonic mean of the precision and recall of the test.

TP Rate: rate of true positives (instances correctly classified as a given class)

FPRate: rate of false positives (instances falsely classified as a given class)

Precision is the fraction of retrieved data that are [relevant](#) to the query:

Recall is the fraction of the relevant data that are successfully retrieved.

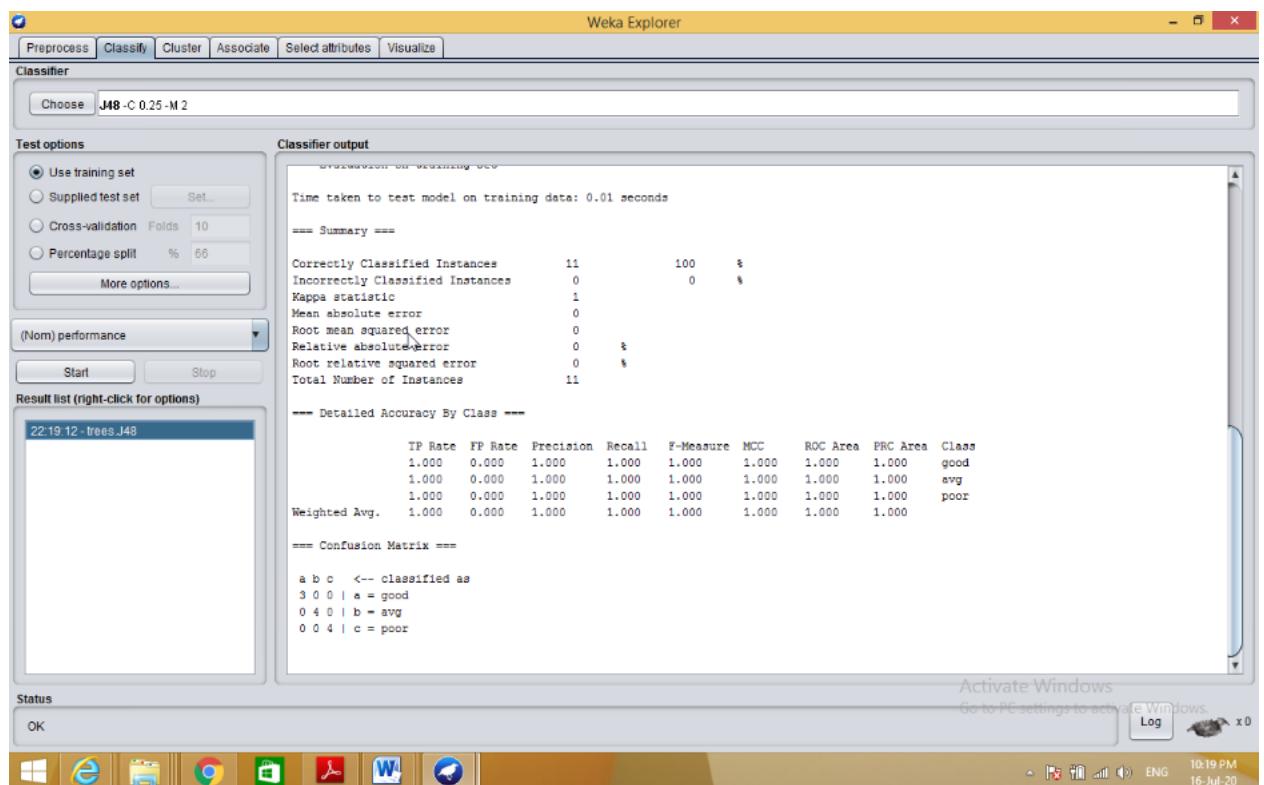
Cross-validation: In k -fold **cross-validation**, the initial data are randomly partitioned into k mutually exclusive subsets or “ folds,” $D_1, D_2: D_k$, each of approximately equal size.

Predictor Error Measures

Procedure:

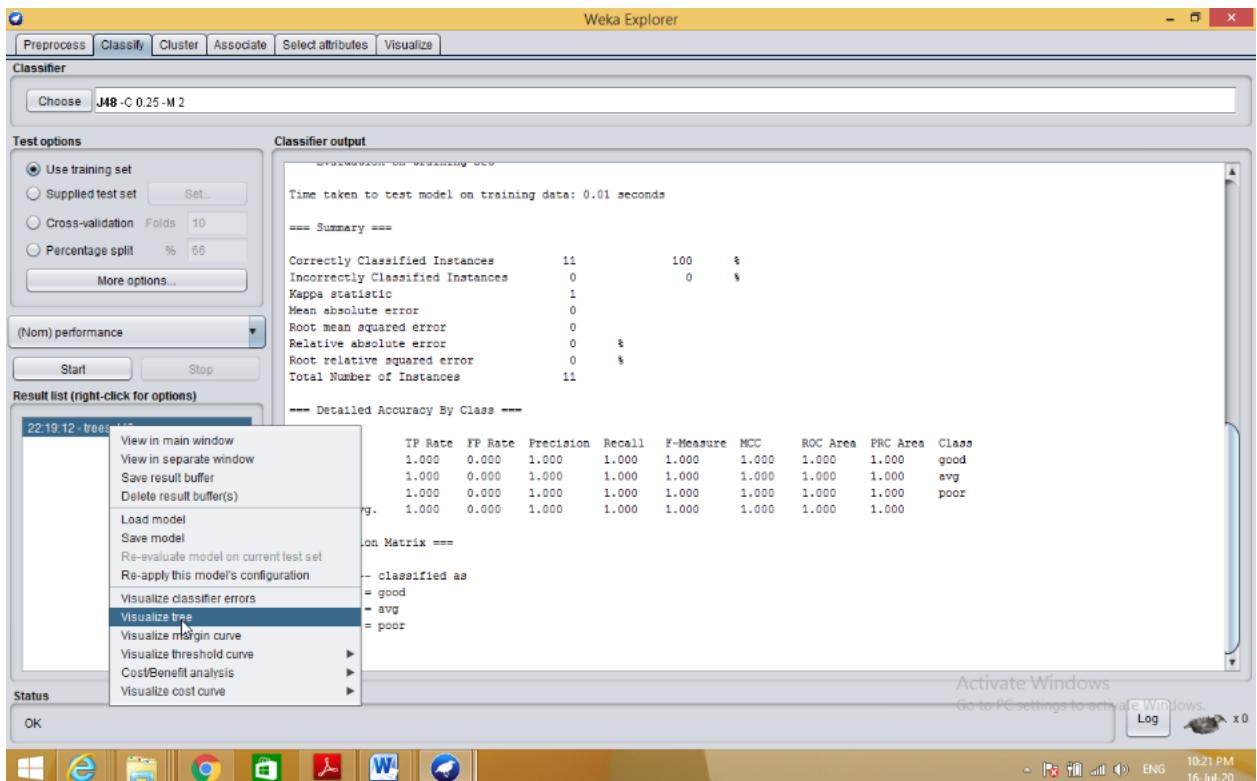
Step 1: Load appropriate dataset into WEKA.

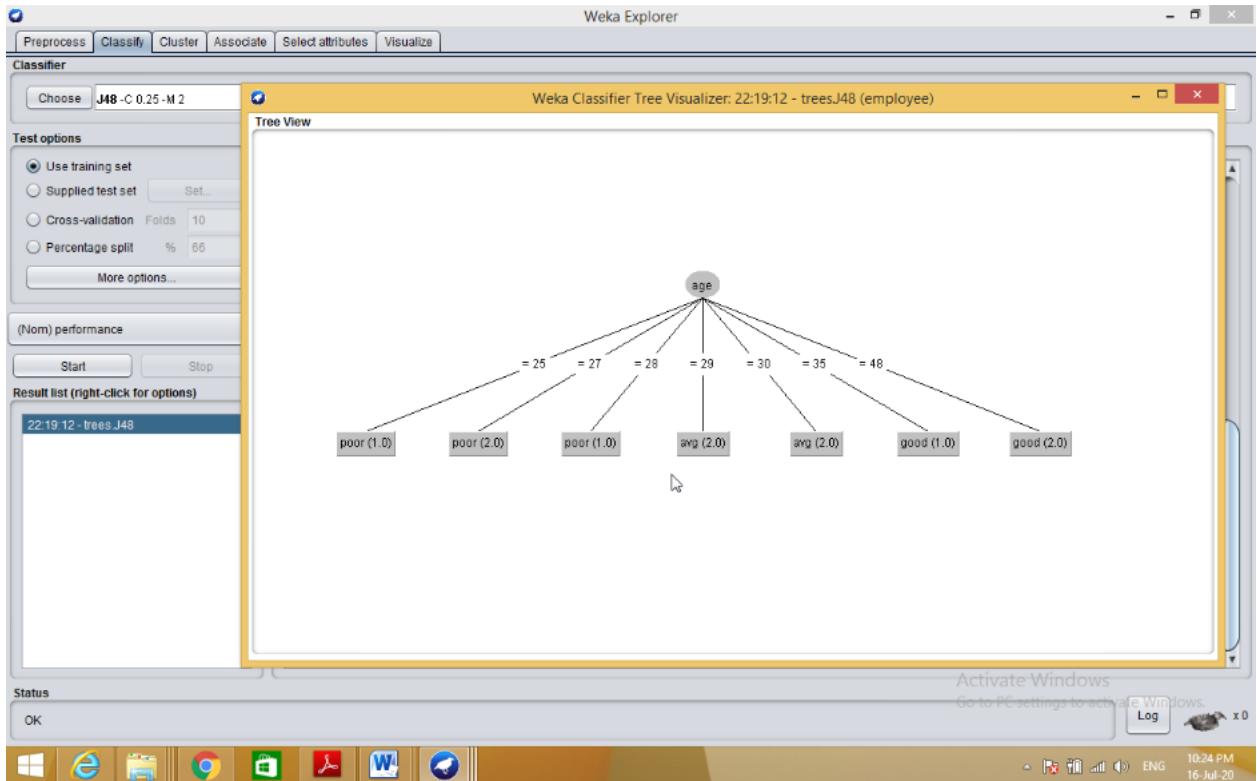
Step 2: Select classify TAB and select J48 algorithm.



Generating Decision Tree:-

Right click on result list & select visualize tree option





4. Demonstration of classification rule process on dataset employee.arff using id3 algorithm.

Aim: This experiment illustrates the use of j-48 classifier in weka. the sample data set used in this experiment is “ employee” data available at arff format. This document assumes that appropriate data pre processing has been performed.

Steps involved in this experiment:

Step 1: We begin the experiment by loading the data (employee.arff) into weka.

Step2: Next we select the “ classify” tab and click “ choose” button to select the “ j48” classifier.

Step3: Now we specify the various parameters. These can be specified by clicking in the text box to the right of the chose button. In this example, we accept the default values the default version does perform some pruning but does not perform error pruning.

Step4: Under the “ text “ options in the main panel. We select the 10-fold cross validation as our evaluation approach. Since we don’ t have separate evaluation data set, this is necessary to get a reasonable idea of accuracy of generated model.

Step-5: We now click " start" to generate the model .the ASCII version of the tree as well as evaluation statistic will appear in the right panel when the model construction is complete.

Step-6: Note that the classification accuracy of model is about 69%.this indicates that we may find more work. (Either in preprocessing or in selecting current parameters for the classification)

Step-7: Now weka also lets us a view a graphical version of the classification tree. This can be done by right clicking the last result set and selecting " visualize tree" from the pop-up menu.

Step-8: We will use our model to classify the new instances.

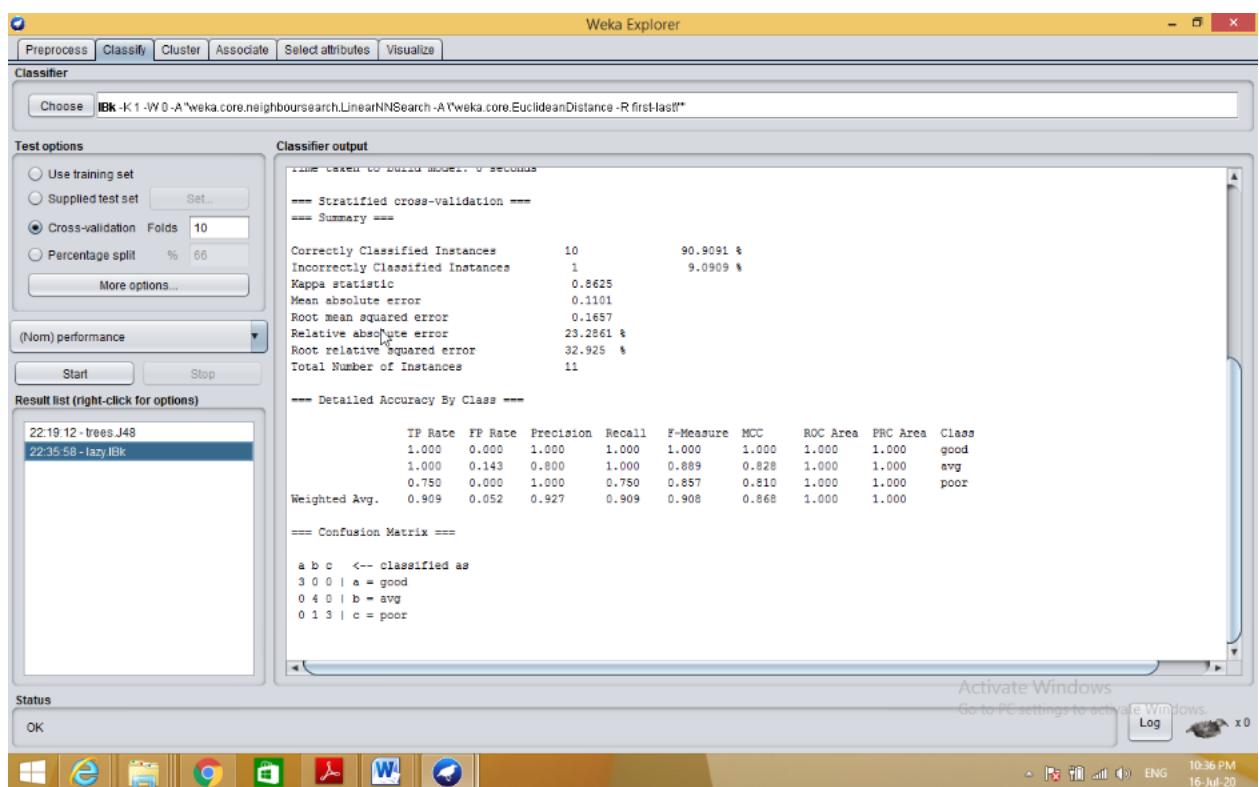
Step-9: In the main panel under " text " options click the " supplied test set" radio button and then click the " set" button. This will pop-up a window which will allow you to open the file containing test instances.

Description: Classification is a data mining function that assigns items in a collection to target categories or classes. The goal of classification is to accurately predict the target class for each case in the data. For example, a classification model could be used to identify loan applicants as low, medium, or high credit risks.

Data set employee.arff:

```
@relation employee  
@attribute age {25, 27, 28, 29, 30, 35, 48}  
@attribute salary{10k,15k,17k,20k,25k,30k,35k,32k}  
@attribute performance {good, avg, poor}  
@data  
%  
25, 10k, poor  
27, 15k, poor
```

27, 17k, poor
 28, 17k, poor
 29, 20k, avg
 30, 25k, avg
 29, 25k, avg
 30, 20k, avg
 35, 32k, good
 48, 35k, good
 48, 32k,good
 %



5. Demonstration of classification rule process on dataset employee.arff using naïve bayes algorithm.

Aim: This experiment illustrates the use of naïve bayes classifier in weka. The sample data set used in this experiment is “ employee” data available at arff format. This document assumes that appropriate data pre processing has been performed.

Steps involved in this experiment:

1. We begin the experiment by loading the data (employee.arff) into weka.

Step2: next we select the “ classify” tab and click “ choose” button to select the “ id3” classifier.

Step3: now we specify the various parameters. These can be specified by clicking in the text box to the right of the chose button. In this example, we accept the default values his default version does perform some pruning but does not perform error pruning.

Step4: under the “ text “ options in the main panel. We select the 10-fold cross validation as our evaluation approach. Since we don’ t have separate evaluation data set, this is necessary to get a reasonable idea of accuracy of generated model.

Step-5: we now click” start” to generate the model .the ASCII version of the tree as well as evaluation statistic will appear in the right panel when the model construction is complete.

Step-6: note that the classification accuracy of model is about 69%.this indicates that we may find more work. (Either in preprocessing or in selecting current parameters for the classification)

Step-7: now weka also lets us a view a graphical version of the classification tree. This can be done by right clicking the last result set and selecting “ visualize tree” from the pop-up menu.

Step-8: we will use our model to classify the new instances.

Step-9: In the main panel under “ text “ options click the “ supplied test set” radio button and then click the “ set” button. This will show pop-up window which will allow you to open the file containing test instances.

Description:

Bayesian classification is based on Baye’ s Theorem. Bayesian classifiers are the statistical classifiers. Bayesian classifiers can predict class membership probabilities such as the probability that a given tuple belongs to a particular class.

Baye's Theorem

Bayes' Theorem is named after Thomas Bayes. There are two types of probabilities

Posterior Probability [P(H/X)]

Prior Probability [P(H)]

where X is data tuple and H is some hypothesis.

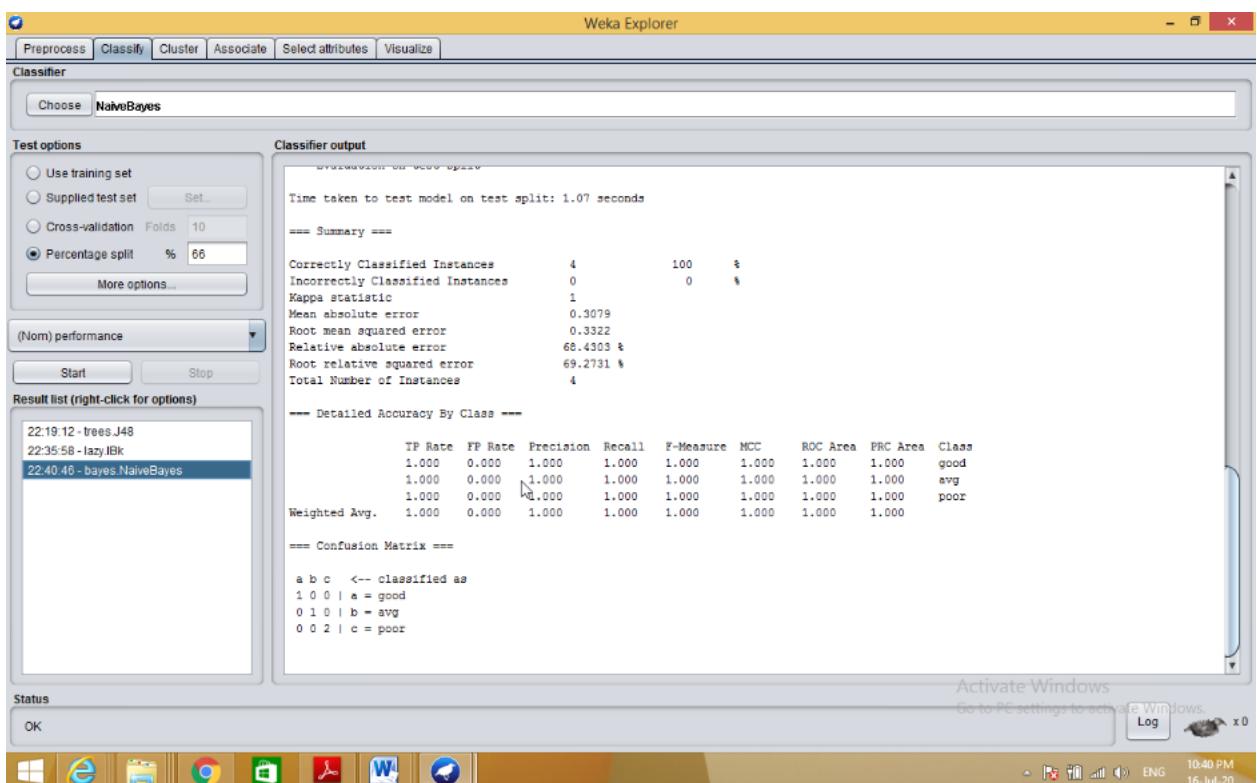
According to Bayes' Theorem,

$$P(H/X) = P(X/H)P(H) / P(X)$$

Procedure:

Step 1: Load appropriate dataset into WEKA

Step 2: Go to Classify tab → Bayes → NaiveBayes



6. Demonstration of clustering rule process on dataset iris.arff using simple k-means.

Aim: This experiment illustrates the use of simple k-mean clustering with Weka explorer. The sample data set used for this example is based on the iris data available in ARFF format. This document assumes that appropriate preprocessing has been performed. This iris dataset includes 150 instances.

Steps involved in this Experiment

Step 1: Run the Weka explorer and load the data file iris.arff in preprocessing interface.

Step 2: Inorder to perform clustering select the ‘ cluster’ tab in the explorer and click on the choose button. This step results in a dropdown list of available clustering algorithms.

Step 3 : In this case we select ‘ simple k-means’ .

Step 4: Next click in text button to the right of the choose button to get popup window shown in the screenshots. In this window we enter six on the number of clusters and we leave the value of the seed on as it is. The seed value is used in generating a random number which is used for making the internal assignments of instances of clusters.

Step 5 : Once of the option have been specified. We run the clustering algorithm there we must make sure that they are in the ‘ cluster mode’ panel. The use of training set option is selected and then we click ‘ start’ button. This process and resulting window are shown in the following screenshots.

Step 6 : The result window shows the centroid of each cluster as well as statistics on the number and the percent of instances assigned to different clusters. Here clusters centroid are means vectors for each clusters. This clusters can be used to characterized the cluster. For eg, the centroid of cluster1 shows the class iris.versicolor mean value of the sepal length is 5.4706, sepal width 2.4765, petal width 1.1294, petal length 3.7941.

7. Demonstration of clustering rule process on dataset student.arff using hierarchical clustering.

Aim: This experiment illustrates the use of hierarchical clustering with Weka explorer. The sample data set used for this example is based on the student data available in ARFF format. This document assumes that appropriate preprocessing has been performed. This student dataset includes 14 instances.

Steps involved in this Experiment

Step 1: Run the Weka explorer and load the data file student.arff in preprocessing interface.

Step 2: Inorder to perform clustering select the ‘ cluster’ tab in the explorer and click on the choose button. This step results in a dropdown list of available clustering algorithms.

Step 3 : In this case we select ‘ hierarchical’ .

Step 4: Next click in text button to the right of the choose button to get popup window shown in the screenshots. In this window we enter six on the number of clusters and we leave the value of the seed on as it is. The seed value is used in generating a random number which is used for making the internal assignments of instances of clusters.

Step 5 : Once of the option have been specified. We run the clustering algorithm there we must make sure that they are in the ‘ cluster mode’ panel. The use of training set option is selected and then we click ‘ start’ button. This process and resulting window are shown in the following screenshots.

Step 6 : The result window shows the centroid of each cluster as well as statistics on the number and the percent of instances assigned to different clusters. Here clusters centroid are means vectors for each clusters. This clusters can be used to characterized the cluster.

Step 7: Another way of understanding characteristics of each cluster through visualization ,we can do this, try right clicking the result set on the result. List panel and selecting the visualize cluster assignments. Interpretation of the above visualization From the above visualization, we can understand the distribution of age and instance number in each cluster. For instance, for each cluster is dominated by age. In this case by changing the color dimension to other attributes we can see their distribution with in each of the cluster.

Step 8: We can assure that resulting dataset which included each instance along with its assign cluster. To do so we click the save button in the visualization window and save the result student hierarchical .The top portion of this file is shown in the following figure.

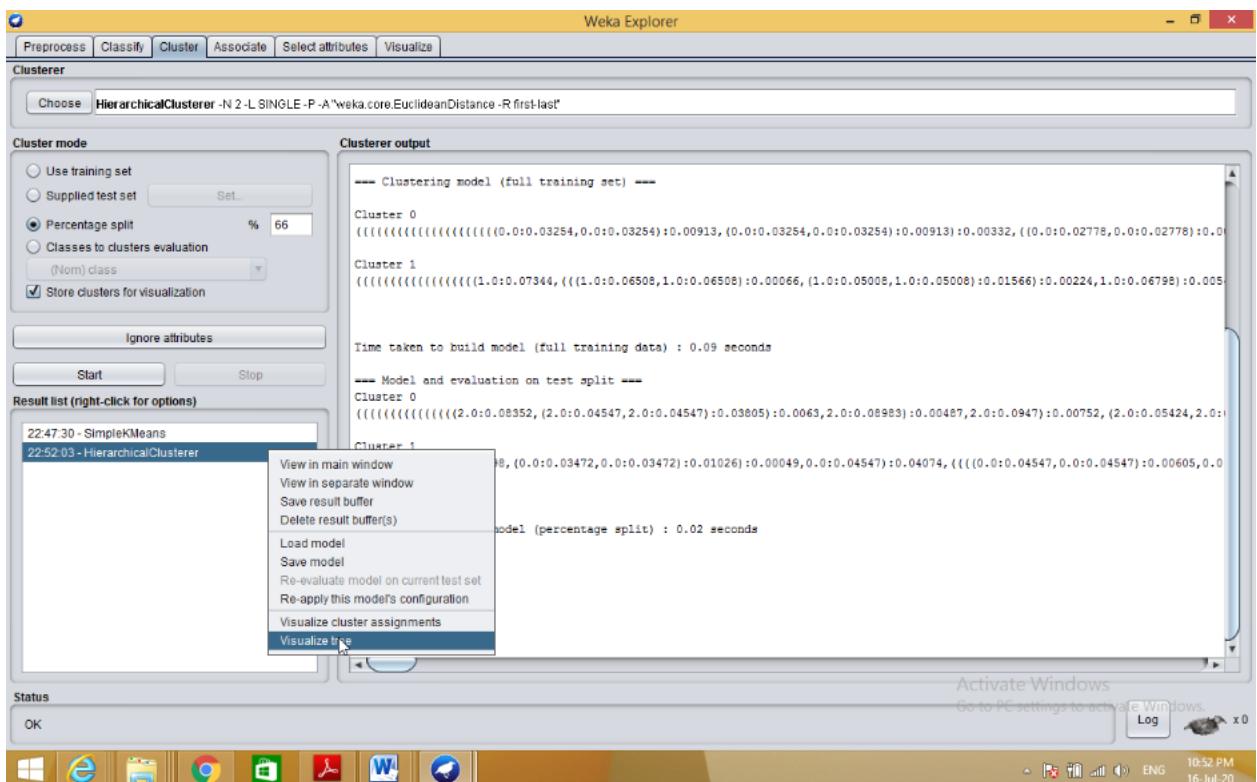
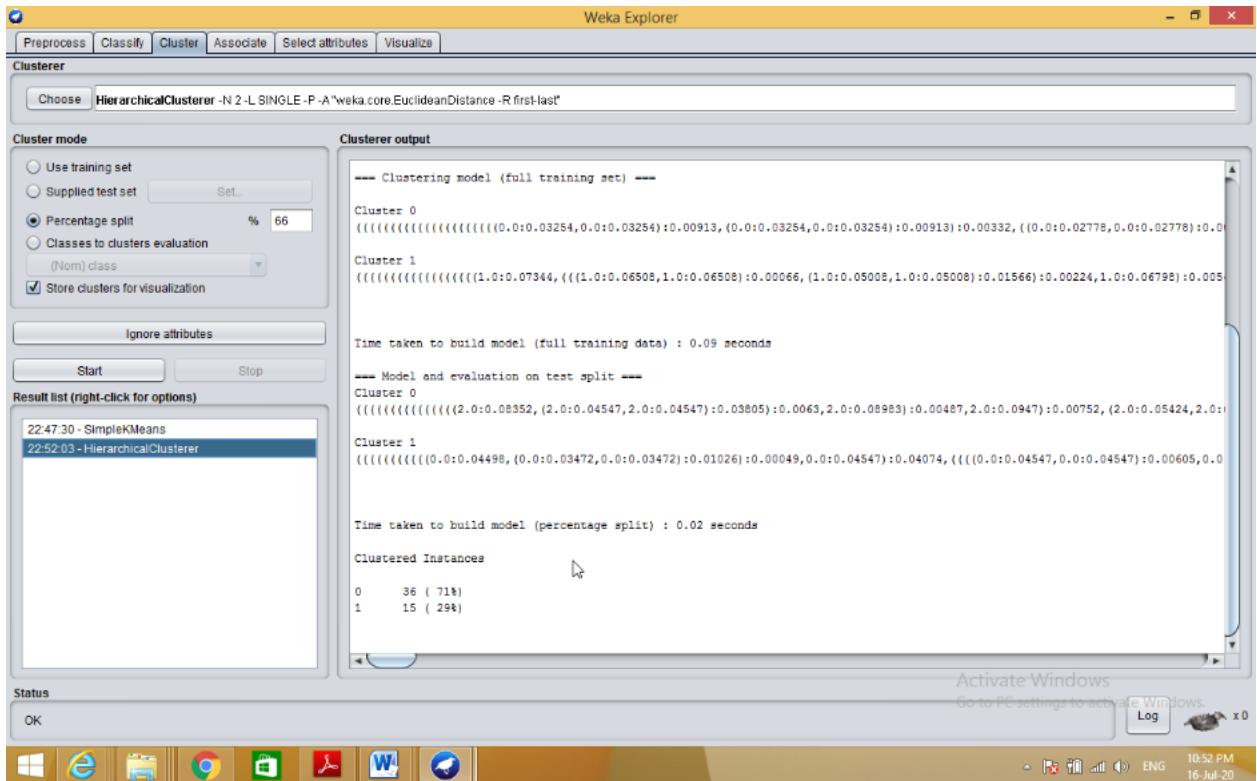
Description:

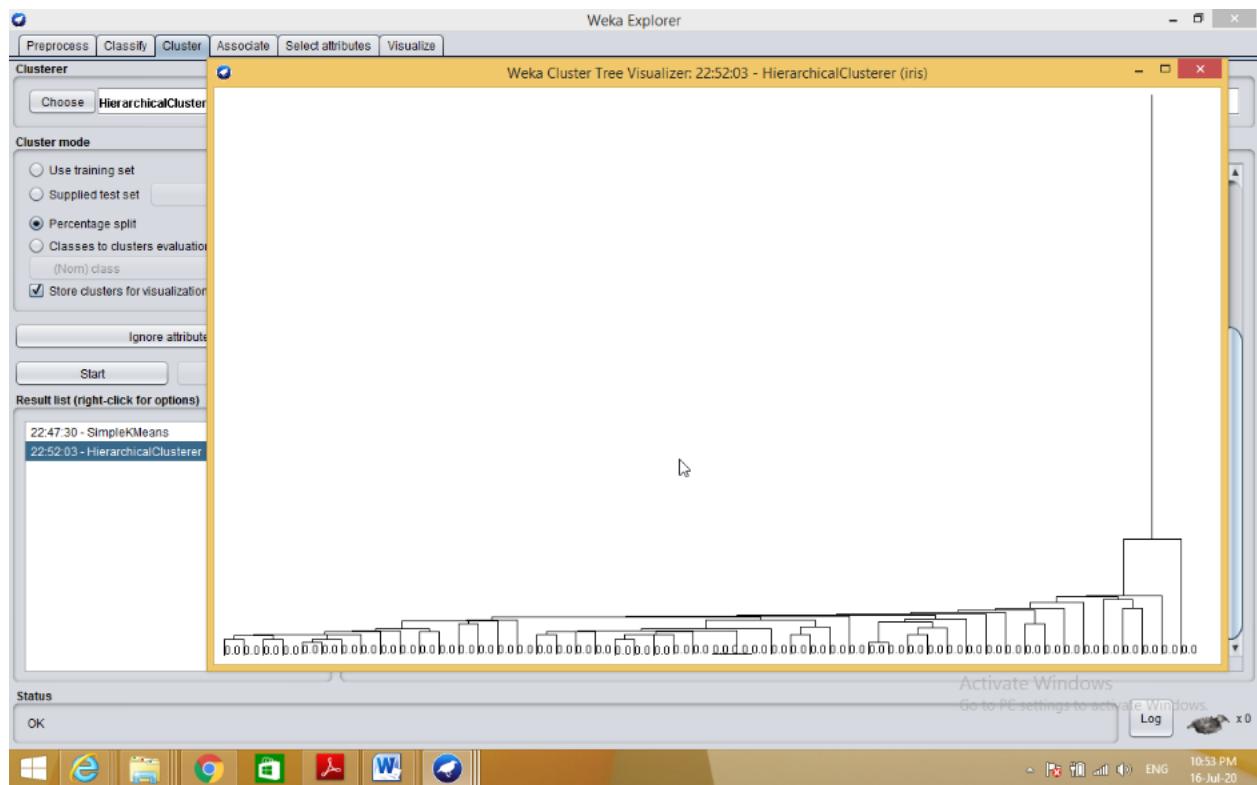
Hierarchical clustering involves creating **clusters** that have a predetermined ordering from top to bottom. For example, all files and folders on the hard disk are organized in a **hierarchy**. There are two types of **hierarchical clustering**, Divisive and Agglomerative.

Procedure:

Step 1: Load appropriate dataset into WEKA

Step 2: Go to Cluster tab ☰choose ☰Hierarchical clustering





Part-B: Data Analytics

1. a) Write R program to find R-Mean, Median & Mode with the sample data.

Source code:

```
install.packages("datasets")
library(datasets)

airquality <- datasets::airquality
summary(airquality)

summary(airquality)

var(airquality$Solar.R,na.rm = T)

summary(airquality$Temp)

summary(airquality$Wind)

plot(airquality$Wind)
plot(airquality$Temp,airquality$Wind)
plot(airquality)

plot(airquality$Wind, type= "b") # p: points, l: lines,b: both
plot(airquality$Wind, xlab = 'ozone Concentration',
     ylab = 'No of Instances', main = 'Ozone levels in NY city',
     col = 'blue')

# Horizontal bar plot
barplot(airquality$Ozone, main = 'Ozone Concentration in air',
         ylab = 'ozone levels', col= 'blue',horiz = F)

#Histogram
hist(airquality$Temp)
hist(airquality$Temp,
     main = 'Solar Radiation values in air',
     xlab = 'Solar rad.', col='blue')

#Single box plot
boxplot(airquality$Temp,main="Boxplot")
```

```
# Multiple box plots
boxplot(airquality[1:4],main='Multiple')

output:

summary(airquality)
   Ozone      Solar.R      Wind       Temp
Min. : 1.00  Min. : 7.0  Min. :1.700  Min. :56.00
1st Qu.:18.00 1st Qu.:115.8 1st Qu.:7.400 1st Qu.:72.00
Median:31.50 Median:205.0 Median:9.700 Median:79.00
Mean :42.13 Mean :185.9 Mean :9.958 Mean :77.88
3rd Qu.:63.25 3rd Qu.:258.8 3rd Qu.:11.500 3rd Qu.:85.00
Max. :168.00 Max. :334.0 Max. :20.700 Max. :97.00
NA's :37    NA's :7

   Month      Day
Min. :5.000  Min. :1.0
1st Qu.:6.000 1st Qu.:8.0
Median:7.000 Median:16.0
Mean :6.993 Mean :15.8
3rd Qu.:8.000 3rd Qu.:23.0
Max. :9.000 Max. :31.0

var(airquality$Solar.R,na.rm=T)
[1] 8110.519

summary(airquality$Temp)
   Min. 1st Qu. Median  Mean 3rd Qu.  Max.
56.00  72.00  79.00  77.88  85.00  97.00
```

b) Write R program to find Analysis and Covariance with the sample data and visualize the regression graphically.

Source code:

```
input <- mtcars[,c("am","mpg","hp")]
print(head(input))
```

```
# Get the dataset.
input <- mtcars
```

```
# Create the regression model.
result <- aov(mpg~hp*am,data = input)
print(summary(result))
```

```
# Create the regression model.
result <- aov(mpg~hp+am,data = input)
print(summary(result))
```

```
# Create the regression models.
result1 <- aov(mpg~hp*am,data = input)
result2 <- aov(mpg~hp+am,data = input)
```

```
# Compare the two models.
print(anova(result1,result2))
```

Output:

```
result <- aov(mpg~hp*am,data = input)
> print(summary(result))
   Df Sum Sq Mean Sq F value    Pr(>F)
hp      1  678.4  678.4 77.391150e-09 ***
am      1 202.2  202.2 23.072475e-05 ***
hp:am   1  0.0   0.0  0.001  0.981
Residuals 28 245.4   8.8
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
result <- aov(mpg~hp+am,data = input)
> print(summary(result))
   Df Sum Sq Mean Sq F value    Pr(>F)
hp      1  678.4  678.4 80.15763e-10 ***
am      1 202.2  202.2 23.89346e-05 ***
Residuals 29 245.4   8.5
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
result1 <- aov(mpg~hp*am,data = input)
> result2 <- aov(mpg~hp+am,data = input)
> # Compare the two models.
```

```
>print(anova(result1,result2))
Analysis of Variance Table

Model1: mpg ~ hp * am
Model2: mpg ~ hp + am
Res.Df  RSS Df Sum of Sq   F Pr(>F)
1    28 245.43
2    29 245.44 -1.00525156e-04 0.9806
```

>

2. Write R program to find the following Regressions with the sample data and visualize the regressions graphically.

Source code:

a) Linear Regression

```
#Load data
Nd<-read.csv("C:/Users/User/Desktop/EXCELR/simple
linear/material/NewspaperData.csv")

# Visualization
install.packages("lattice")
library(lattice)
dotplot(Nd$sunday, main="Dot Plot of Sunday Circulations",col="dodgerblue4")
dotplot(Nd$daily, main="Dot Plot of Daily Circulations", col="dodgerblue4")
boxplot(Nd$sunday,col="dodgerblue4")
boxplot(Nd$daily,col="dodgerblue4")

#Regression equation
#Syntax model<-lm(y~x,data set name)

#column names
colnames(Nd)
#Model building
model<- lm(sunday~daily,data =Nd)
summary(model)

#Preparing new data frame with new data
new_daily=data.frame(daily=c(300,200))

#Predict for the new data
sun1=predict(model,new_daily)

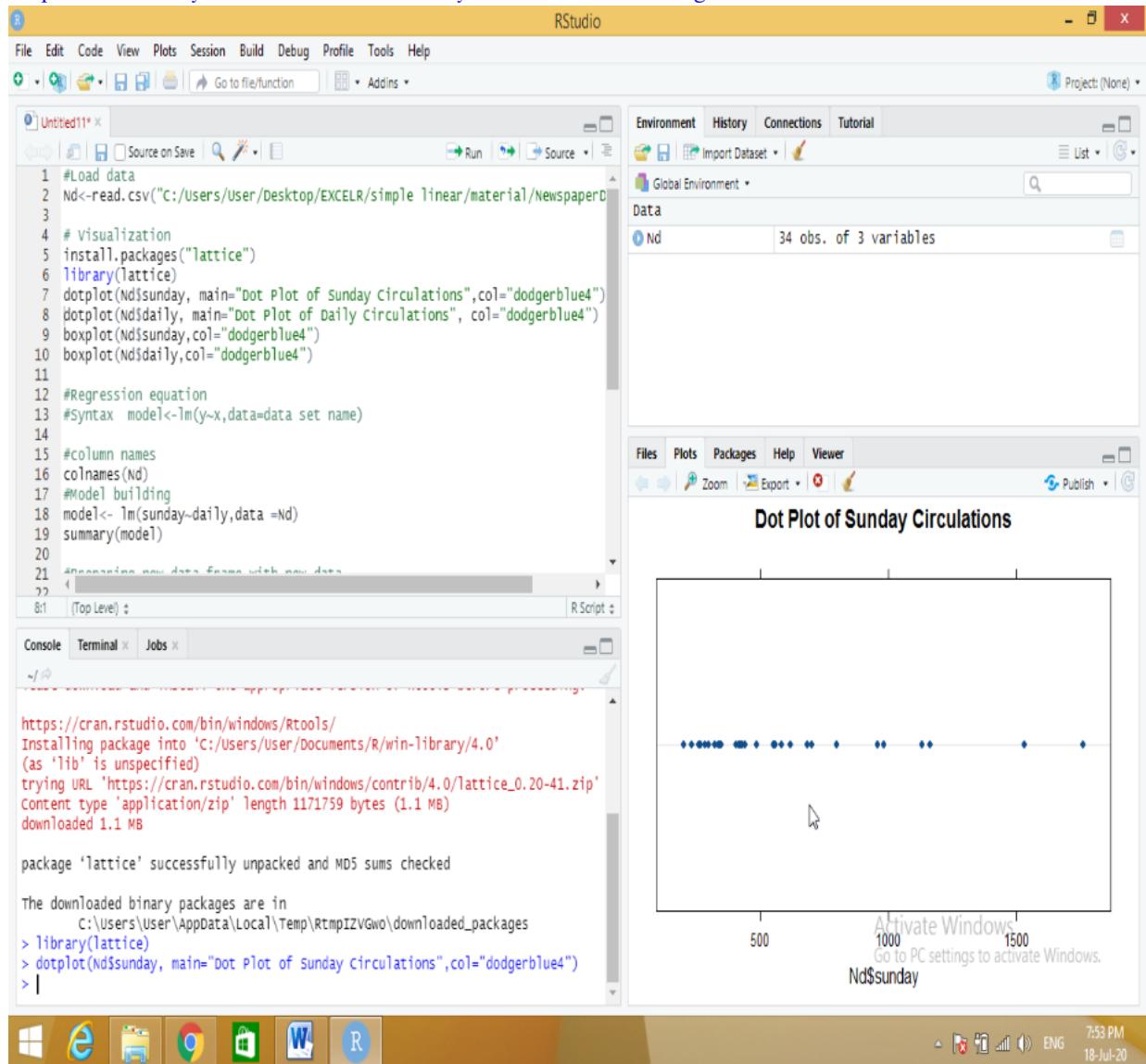
#Print predicted value
sun1

#Predict for daily variable from the historical data
pred<-predict(model)
#Print predicted values
pred
#Prepare a new data frame with pred and error
```

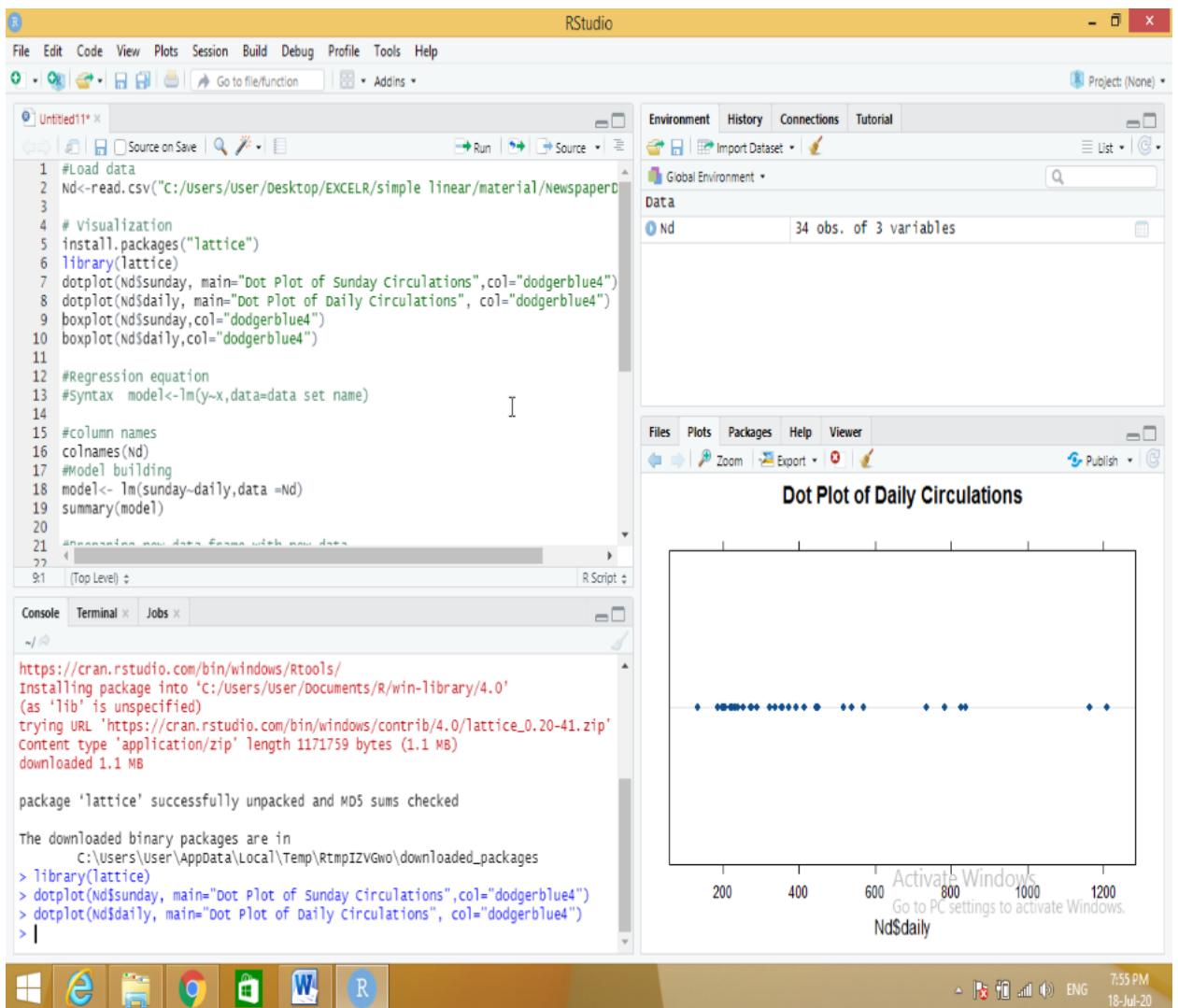
```
finaldata<-data.frame(Nd,pred,"Error"= Nd$sunday-pred)
```

output:

```
dotplot(Nd$sunday,main="Dot Plot of Sunday Circulations",col="dodgerblue4")
```



```
dotplot(Nd$daily, main="Dot Plot of Daily Circulations", col="dodgerblue4")
```



b) Multiple Regression

Source code:

```
mileage<-read.csv("C:/Users/User/Desktop/EXCELR/multi linear/Cars.csv")
```

```
#mileage<-Cars
#Scatter Plot Matrix:
pairs(mileage)

#Correlation Matrix:#####
cor(mileage)
#Regression Model and Summary
model.car<-lm(MPG~HP+VOL+SP+WT,data = mileage)
summary(model.car)
```

```
#####Experiment#####
reg_vol<-lm(MPG~VOL,data = mileage)
summary(reg_vol)
reg_wt<-lm(MPG~WT,data = mileage)
summary(reg_wt)
reg_wt_vol<-lm(MPG~WT+VOL,data = mileage)
summary(reg_wt_vol)
#####
```

```
####stepAIC to find out variable that removes multicollinearity
#Regression Model and Summary
model.car<-lm(MPG~.,data = mileage)
summary(model.car)
#Multi-colinearity
install.packages("car")
library(car)
#Variance Inflation Factor - Multi collinearity values
car::vif(model.car)
##Subset selection
install.packages("MASS")
library(MASS)
stepAIC(model.car)
#####
```

```
#####
#Full Model Building process

#mileage<-Cars
#Scatter Plot Matrix:
pairs(mileage)

#Correlation Matrix:#####
cor(mileage)

#Regression Model and Summary
model.car<-lm(MPG~.,data = mileage)

#Multi-colinearity
install.packages("car")
library(car)
#Variance Inflation Factor
car::vif(model.car)
#Diagnostic Plots:
  #Residual Plots, QQ-Plos, Std. Residuals vs Fitted
plot(model.car)
#Residuals vs Regressors

residualPlots(model.car)
#Added Variable Plots
avPlots(model.car)
#QQ plots of studentized residuals
qqPlot(model.car)
#Deletion Diagnostics
influenceIndexPlot(model.car) # Index Plots of the influence measures

####Iteration 1
#Remove 77th observation
mileage["HP2"] <- mileage$HP*mileage$HP
```

```
mileage["SP2"] <- mileage$SP*mileage$SP
mileage1<-mileage[-77,]
```

```
model1<-lm(MPG~.,data = mileage1)
summary(model1)
```

```
plot(model1)
residualPlots(model1)
qqPlot(model1)
influenceIndexPlot(model1)
```

```
#iteration2
mileage2<-mileage[-c(77,66,81,79),]
```

```
model2<-lm(MPG~.,data = mileage2[,c(5,4)])
summary(model2)
influenceIndexPlot(model2)
```

output:

```
cor(mileage)
      HP     MPG     VOL     SP     WT
HP 1.00000000 -0.7250383 0.07745947 0.9738481 0.07651307
MPG -0.72503835 1.0000000 -0.52905658 -0.6871246 -0.52675909
VOL 0.07745947 -0.5290566 1.00000000 0.1021700 0.99920308
SP 0.97384807 -0.6871246 0.10217001 1.0000000 0.10243919
WT 0.07651307 -0.5267591 0.99920308 0.1024392 1.00000000
>#Regression Model and Summary
>model.car<-lm(MPG~HP+VOL+SP+WT,data=mileage)
>summary(model.car)
```

Call:
`lm(formula = MPG ~ HP + VOL + SP + WT, data = mileage)`

Residuals:
`Min 1Q Median 3Q Max`

-8.6320 -2.9944 -0.3705 2.2149 15.6179

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	30.67734	14.90030	2.059	0.0429*
HP	-0.20544	0.03922	-5.239	1.4e-06***
VOL	-0.33605	0.56864	-0.591	0.5563
SP	0.39563	0.15826	2.500	0.0146*
WT	0.40057	1.69346	0.237	0.8136

Signif. codes: 0 ' ***' 0.001 ' **' 0.01 ' *' 0.05 ' .' 0.1 ' ' 1

Residual standard error: 4.488 on 76 degrees of freedom

Multiple R-squared: 0.7705, Adjusted R-squared: 0.7585

F-statistic: 63.8 on 4 and 76 DF, p-value:< 2.2e-16

> #####Experiment#####

> reg_vol<-lm(MPG~VOL,data = mileage)

> summary(reg_vol)

Call:

lm(formula = MPG ~ VOL, data = mileage)

Residuals:

Min	1Q	Median	3Q	Max
-25.3074	-5.2026	0.1902	5.4536	17.1632

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	55.81709	3.95696	14.106	<2e-16***
VOL	-0.21662	0.03909	-5.5413.82e-07***	

Signif. codes: 0 ' ***' 0.001 ' **' 0.01 ' *' 0.05 ' .' 0.1 ' ' 1

Residual standard error: 7.798 on 79 degrees of freedom

Multiple R-squared: 0.2799, Adjusted R-squared: 0.2708

F-statistic: 30.71 on 1 and 79 DF, p-value: 3.823e-07

> reg_wt<-lm(MPG~WT,data = mileage)

> summary(reg_wt)

Call:

lm(formula = MPG ~ WT, data = mileage)

Residuals:

Min	1Q	Median	3Q	Max
-25.3933	-5.4377	0.2738	5.2951	16.9351

Coefficients:

Estimate Std. Error t value Pr(>|t|)
 (Intercept) 55.2296 3.8761 14.249 < 2e-16 ***
 WT -0.6420 0.1165 -5.5084 3.8e-07 ***

 Signif. codes: 0 ‘ ***’ 0.001 ‘ **’ 0.01 ‘ *’ 0.05 ‘ .’ 0.1 ‘ ’ 1

Residual standard error: 7.811 on 79 degrees of freedom
 Multiple R-squared: 0.2775, Adjusted R-squared: 0.2683
 F-statistic: 30.34 on 1 and 79 DF, p-value: 4.383e-07

```
> reg_wt_vol<-lm(MPG~WT+VOL,data = mileage)
> #####stepAIC to find out variable that removes multicollinearity
> #Regression Model and Summary
> model.car<-lm(MPG~,data = mileage)
> summary(model.car)
```

Call:
 lm(formula = MPG ~ ., data = mileage)

Residuals:

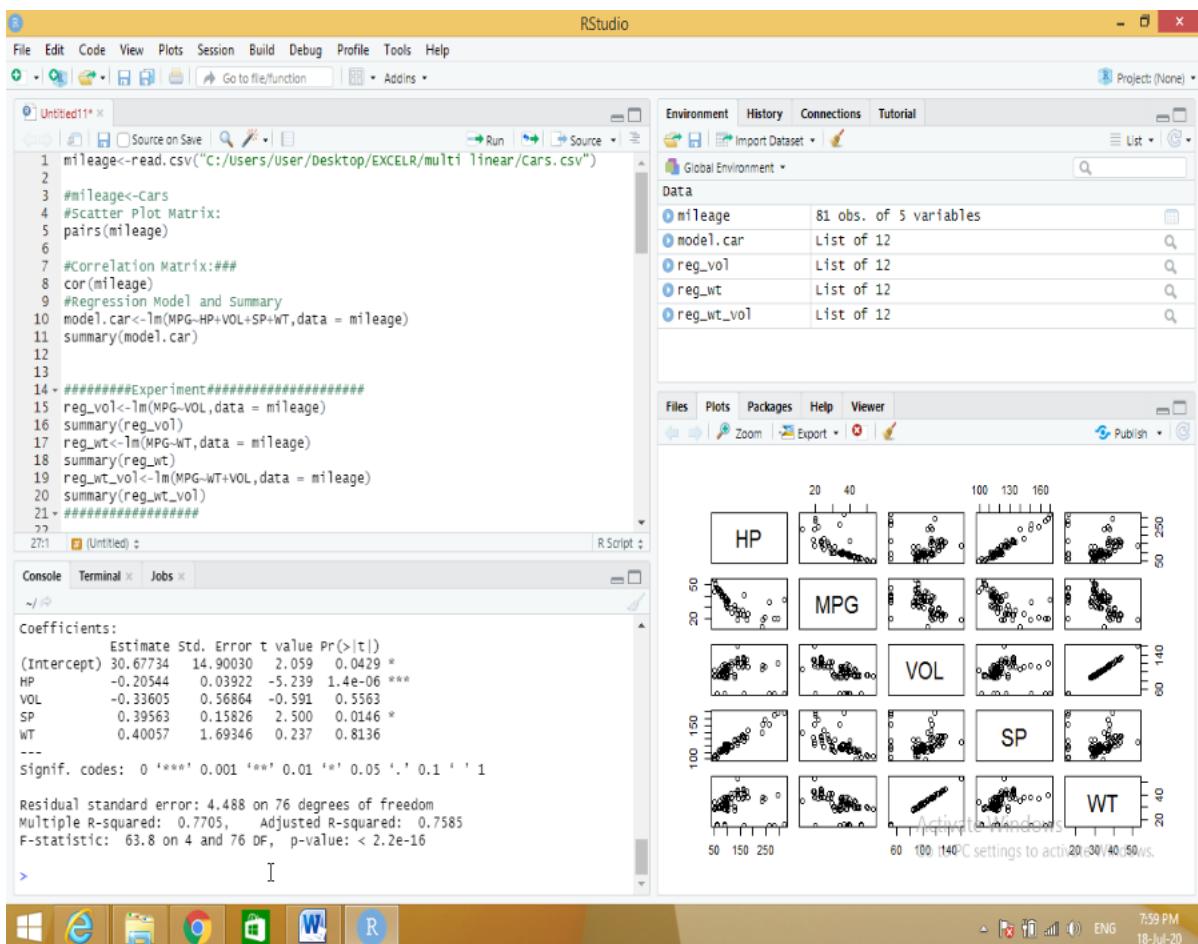
Min	1Q	Median	3Q	Max
-8.6320	-2.9944	-0.3705	2.2149	15.6179

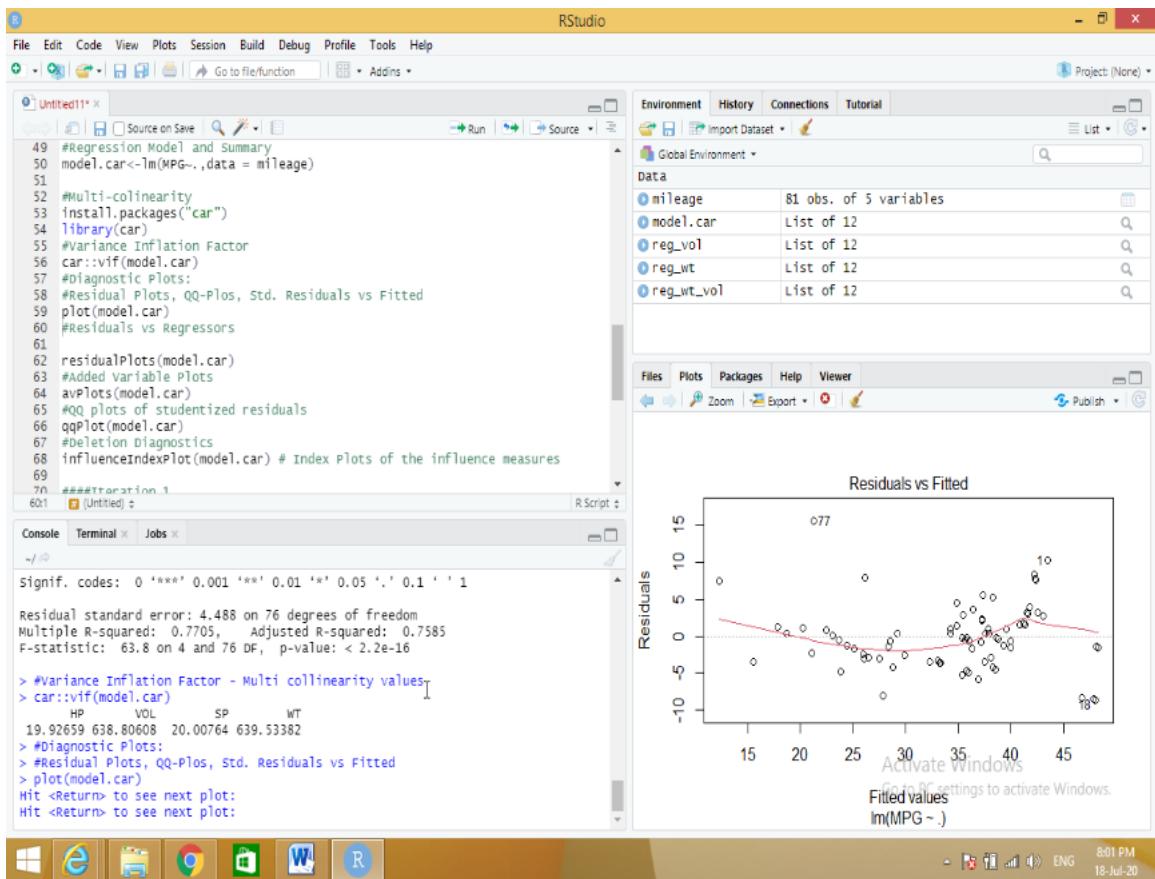
Coefficients:

Estimate Std. Error t value Pr(>|t|)
 (Intercept) 30.67734 14.90030 2.059 0.0429 *
 HP -0.20544 0.03922 -5.239 1.4e-06 ***
 VOL -0.33605 0.56864 -0.591 0.5563
 SP 0.39563 0.15826 2.500 0.0146 *
 WT 0.40057 1.69346 0.237 0.8136

 Signif. codes: 0 ‘ ***’ 0.001 ‘ **’ 0.01 ‘ *’ 0.05 ‘ .’ 0.1 ‘ ’ 1

Residual standard error: 4.488 on 76 degrees of freedom
 Multiple R-squared: 0.7705, Adjusted R-squared: 0.7585
 F-statistic: 63.8 on 4 and 76 DF, p-value: < 2.2e-16





c) Logistic Regression

Source code:

```
claimants<-read.csv("C:/Users/User/Desktop/EXCELR/logistic regression/claimants.csv")
#Finding null values
sum(is.na(claimants))
#Removing null values- na.omit(dataset)
claimants<- na.omit(claimants)
# Logistic Regression
#glm(y~x,family="bin....")
logit<-glm(ATTORNEY ~ factor(CLMSEX) + factor(CLMINSUR) + factor(SEATBELT)
          + CLMAGE + LOSS,family= "binomial",data=claimants)
summary(logit)

# Confusion Matrix Table
#predict(modelobject,testdataset)
```

```
prob=predict(logit,type=c("response"),claimants)
prob

#table(dataframe1,dataframe2) ..to create 2X2 matrix
confusion<-table(prob>0.5,claimants$ATTORNEY)
confusion

# Model Accuracy
#adding diagonal elements in the confusion matrix
Accuracy<-sum(diag(confusion))/sum(confusion)
Accuracy

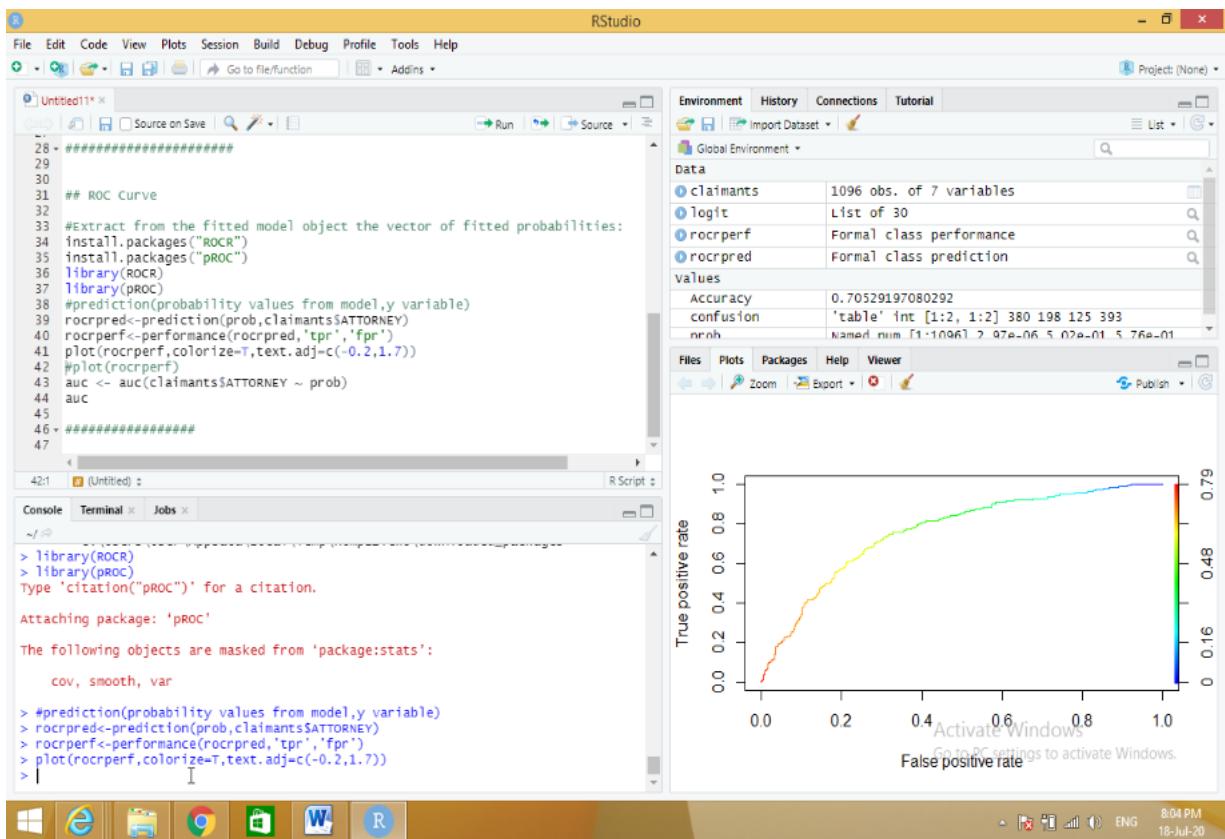
#####
#####

## ROC Curve

#Extract from the fitted model object the vector of fitted probabilities:
install.packages("ROCR")
install.packages("pROC")
library(ROCR)
library(pROC)
#prediction(probability values from model,y variable)
rocrpred<-prediction(prob,claimants$ATTORNEY)
rocrperf<-performance(rocrpred,'tpr','fpr')
plot(rocrperf,colorize=T,text.adj=c(-0.2,1.7))
#plot(rocrperf)
auc <- auc(claimants$ATTORNEY ~ prob)
auc

#####
```

output:



d) Poisson Regression.

Source code:

```
input <- warpbreaks
print(head(input))
```

```
output <-glm(formula = breaks ~ wool+tension, data = warpbreaks,  
family = poisson)  
print(summary(output))
```

Output:

```
input<-warpbreaks  
> print(head(input))  
breakswooltension  
1 26 A L  
2 30 A L  
3 54 A L  
4 25 A L  
5 70 A L  
6 52 A L  
>output <-glm(formula = breaks ~ wool+tension, data = warpbreaks,  
+ family = poisson)  
>print(summary(output))
```

Call:

```
glm(formula = breaks ~ wool + tension, family = poisson, data = warpbreaks)
```

Deviance Residuals:

	Min	1Q	Median	3Q	Max
(Intercept)	-3.6871	-1.6503	-0.4269	1.1902	4.2616

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	3.69196	0.04541	81.302	<2e-16 ***
woolB	-0.20599	0.05157	-3.994	6.49e-05 ***
tensionM	-0.32132	0.06027	-5.332	9.73e-08 ***
tensionH	-0.51849	0.06396	-8.107	5.21e-16 ***

Signif. codes: 0 ' ***' 0.001 ' **' 0.01 ' *' 0.05 ' .' 0.1 ' ' 1

(Dispersion parameter for poisson family taken to be 1)

Null deviance: 297.37 on 53 degrees of freedom
Residual deviance: 210.39 on 50 degrees of freedom
AIC: 493.06

Number of Fisher Scoring iterations: 4

3. a) Write R program to find Time Series Analysis with the sample data and visualize the regression graphically.

Source code:

```
# Getting the data points in form of a R vector.  
snowfall <- c(790,1170.8,860.1,1330.6,630.4,911.5,683.5,996.6,783.2,982,881.8,1021)  
# Convertting it into a time series object.  
snowfall_timeseries<- ts(snowfall,start = c(2013,1),frequency = 12)  
# Printing the timeseries data.  
print(snowfall_timeseries)  
# Giving a name to the chart file.  
png(file = "snowfall.png")  
# Plotting a graph of the time series.  
plot(snowfall_timeseries)  
# Saving the file.  
dev.off()
```

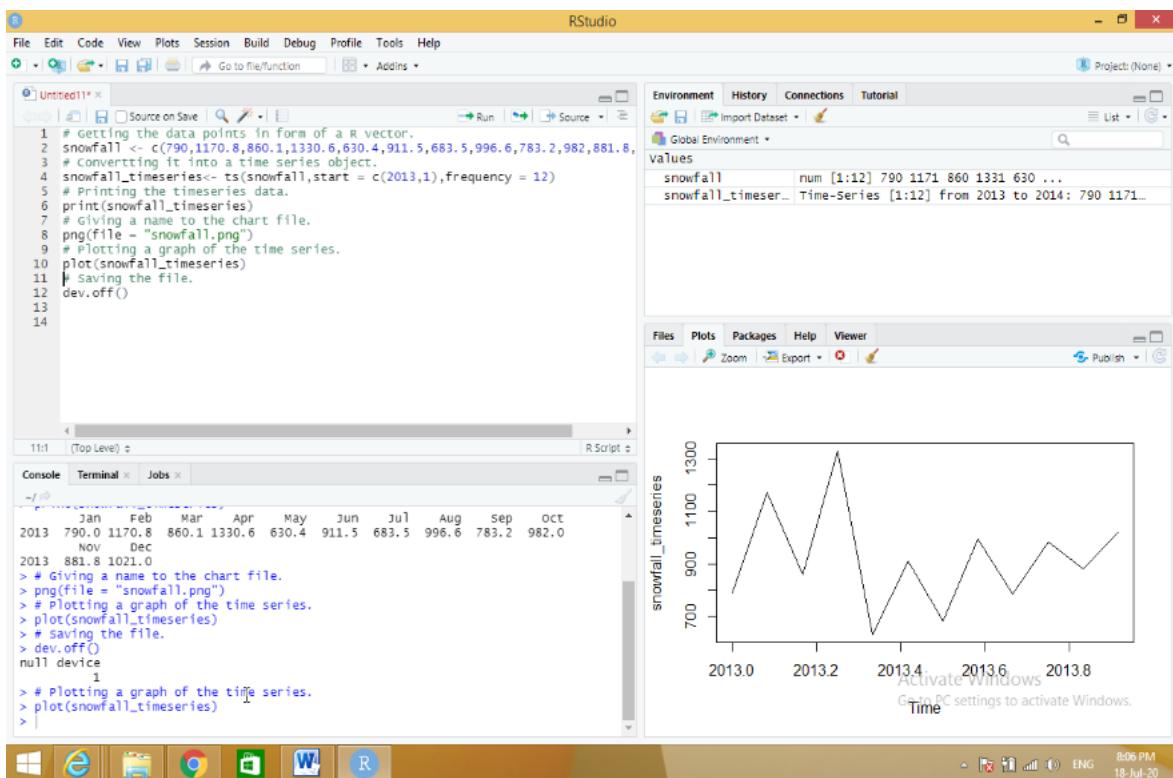
output:

```

snowfall<- c(790,1170.8,860.1,1330.6,630.4,911.5,683.5,996.6,783.2,982,881.8,1021)
> # Converting it into a time series object.
> snowfall_timeseries<- ts(snowfall,start = c(2013,1),frequency = 12)
> # Printing the timeseries data.
> print(snowfall_timeseries)

Jan Feb Mar Apr May Jun Jul Aug Sep Oct
2013 790.0 1170.8 860.1 1330.6 630.4 911.5 683.5 996.6 783.2 982.0
Nov Dec
2013 881.8 1021.0
> # Giving a name to the chart file.
> png(file = "snowfall.png")
> # Plotting a graph of the time series.
> plot(snowfall_timeseries)
> # Saving the file.
> dev.off()
null device
1
> # Plotting a graph of the time series.
> plot(snowfall_timeseries)

```



b) Write R program to find Non Linear Least Square with the sample data and visualize the regression graphically.

Source code:

```
xvalues <- c(1.6,2.1,2,2.23,3.71,3.25,3.4,3.86,1.19,2.21)
yvalues <- c(5.19,7.43,6.94,8.11,18.75,14.88,16.06,19.12,3.21,7.58)

# Give the chart file a name.
png(file = "nls.png")

# Plot these values.
plot(xvalues,yvalues)

# Take the assumed values and fit into the model.
model <- nls(yvalues ~ b1*xvalues^2+b2,start = list(b1 = 1,b2 = 3))

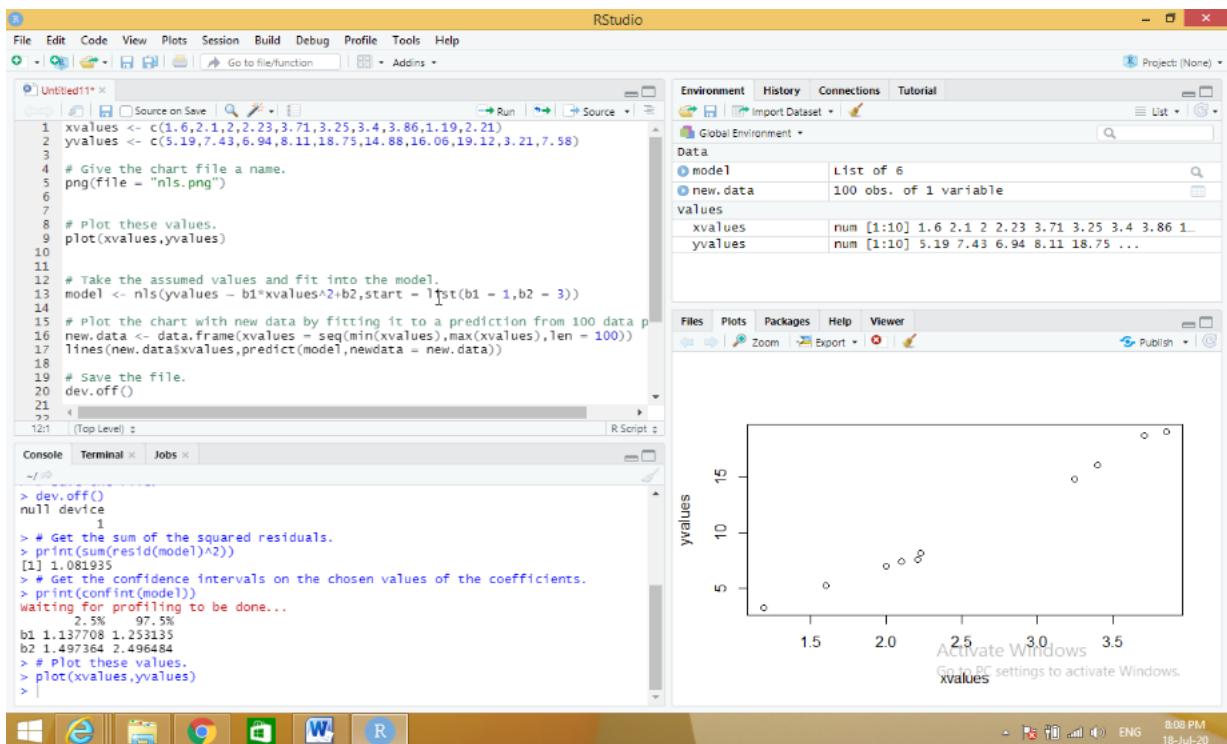
# Plot the chart with new data by fitting it to a prediction from 100 data points.
new.data <- data.frame(xvalues = seq(min(xvalues),max(xvalues),len = 100))
lines(new.data$xvalues,predict(model,newdata = new.data))

# Save the file.
dev.off()

# Get the sum of the squared residuals.
print(sum(resid(model)^2))

# Get the confidence intervals on the chosen values of the coefficients.
print(confint(model))
```

Output:



c) Write R program to find Decision Tree with the sample data and visualize the regression graphically.

Source code:

```

#Data Load
data("iris")

#Install the required packages
install.packages("caret")
install.packages("C50")

#Library invoke
library(caret)
library(C50)

#To make the results consistent across the runs

set.seed(7)

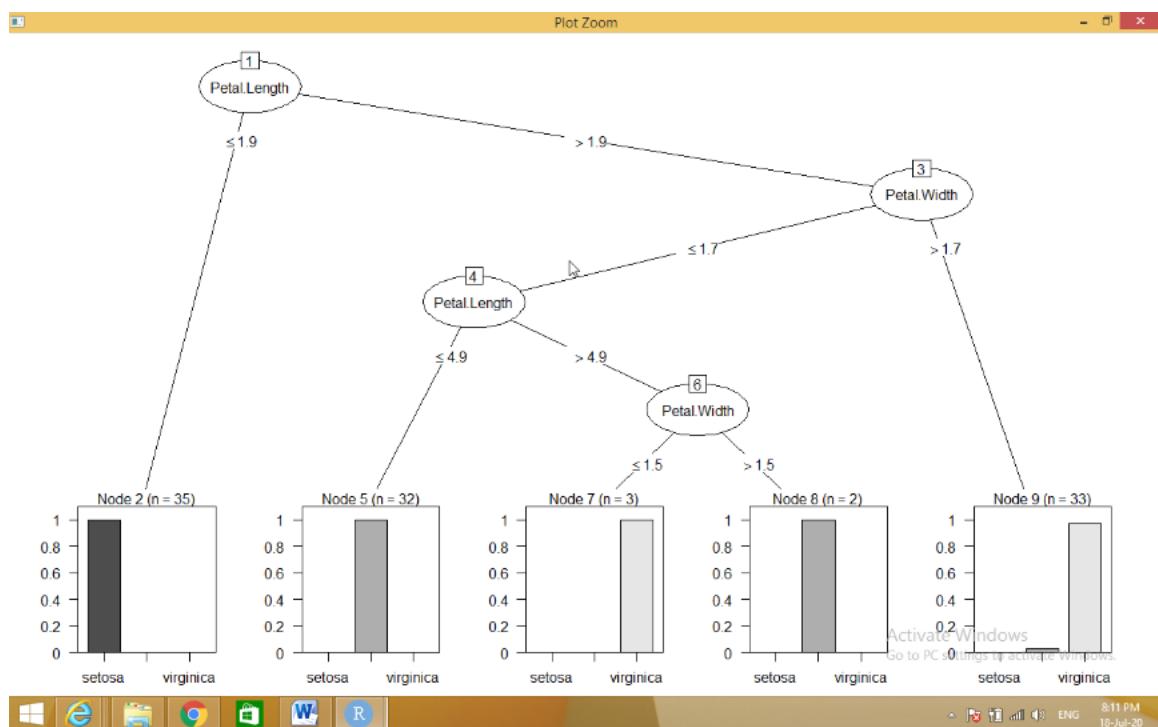
#Data Partition
inTraininglocal<-createDataPartition(iris$Species,p=.70,list = F)
training<-iris[inTraininglocal,]
testing<-iris[-inTraininglocal,]

#Model Building
model<-C5.0(Species~.,data = training)

```

```
#Generate the model summary
summary(model)
#Predict for test data set
pred<-predict.C5.0(model,testing[,-5]) #type ="prob"
#Accuracy of the algorithm
a<-table(testing$Species,pred)
sum(diag(a))/sum(a)
#Visualize the decision tree
plot(model)

output:
```



Source code:

a) Normal Distribution

dnorm

```
# Create a sequence of numbers between -10 and 10 incrementing by 0.1.
x <- seq(-10, 10, by = .1)
```

```
# Choose the mean as 2.5 and standard deviation as 0.5.
y <- dnorm(x, mean = 2.5, sd = 0.5)
```

```
# Give the chart file a name.
```

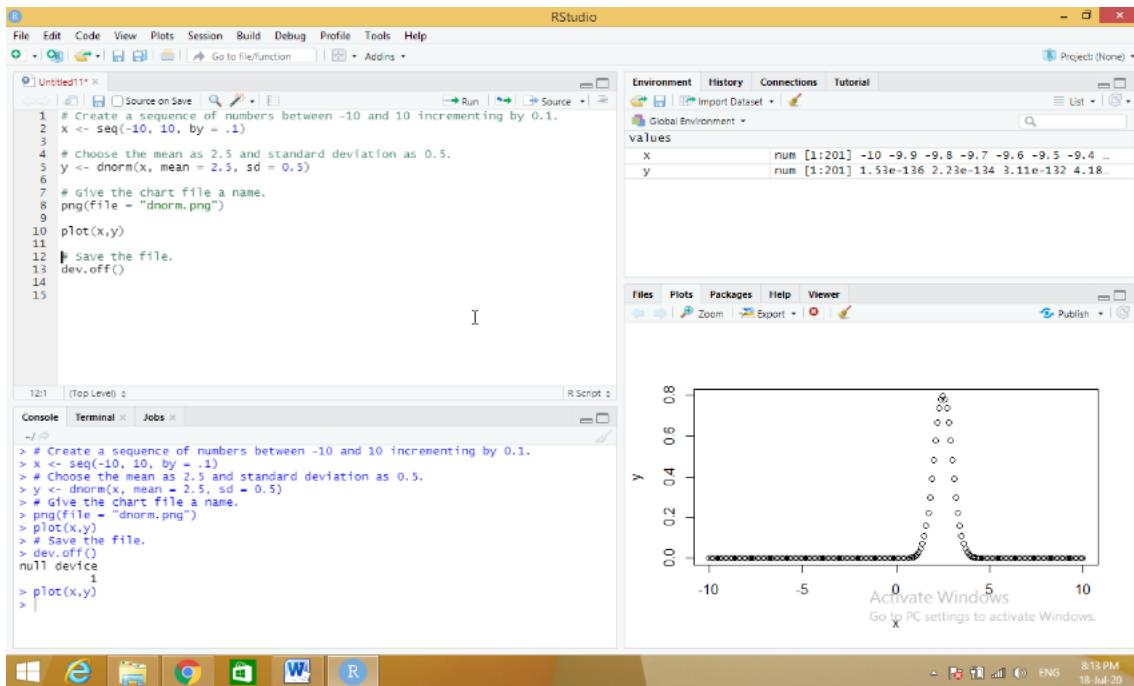
```
png(file = "dnorm.png")
```

```
plot(x,y)
```

```
# Save the file.
```

```
dev.off()
```

output:



Pnorm

```
# Create a sequence of numbers between -10 and 10 incrementing by 0.2.
x <- seq(-10,10,by = .2)

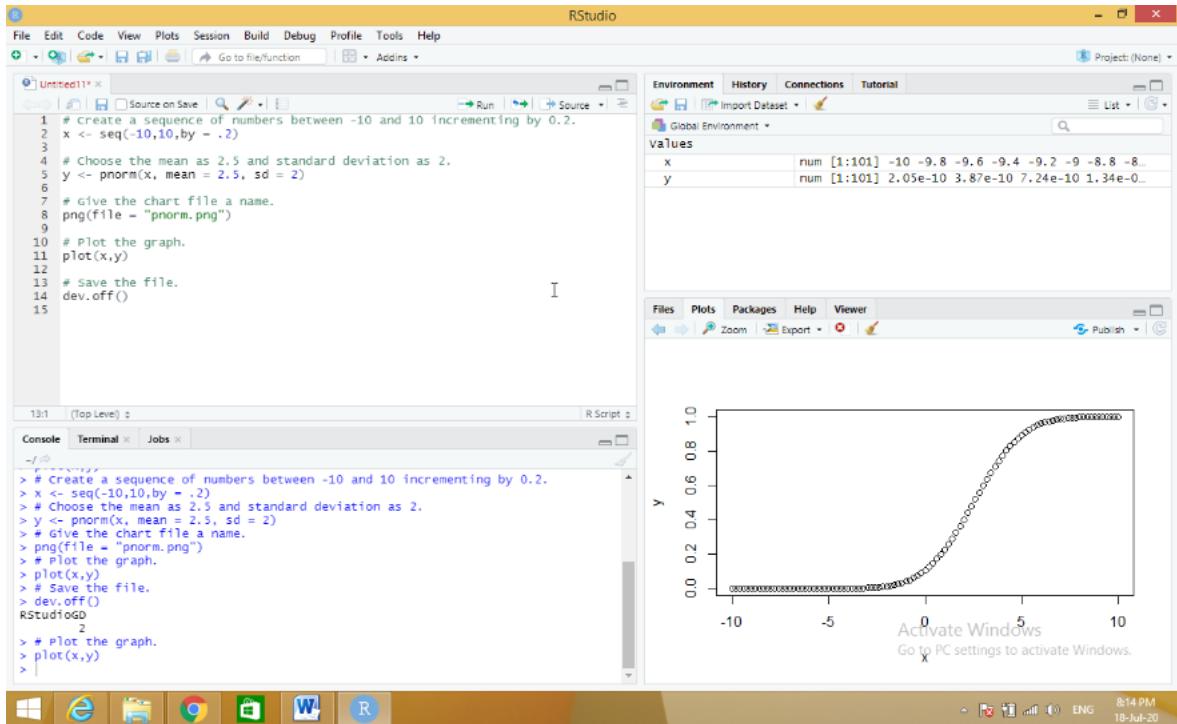
# Choose the mean as 2.5 and standard deviation as 2.
y <- pnorm(x, mean = 2.5, sd = 2)

# Give the chart file a name.
png(file = "pnorm.png")

# Plot the graph.
plot(x,y)

# Save the file.
dev.off()
```

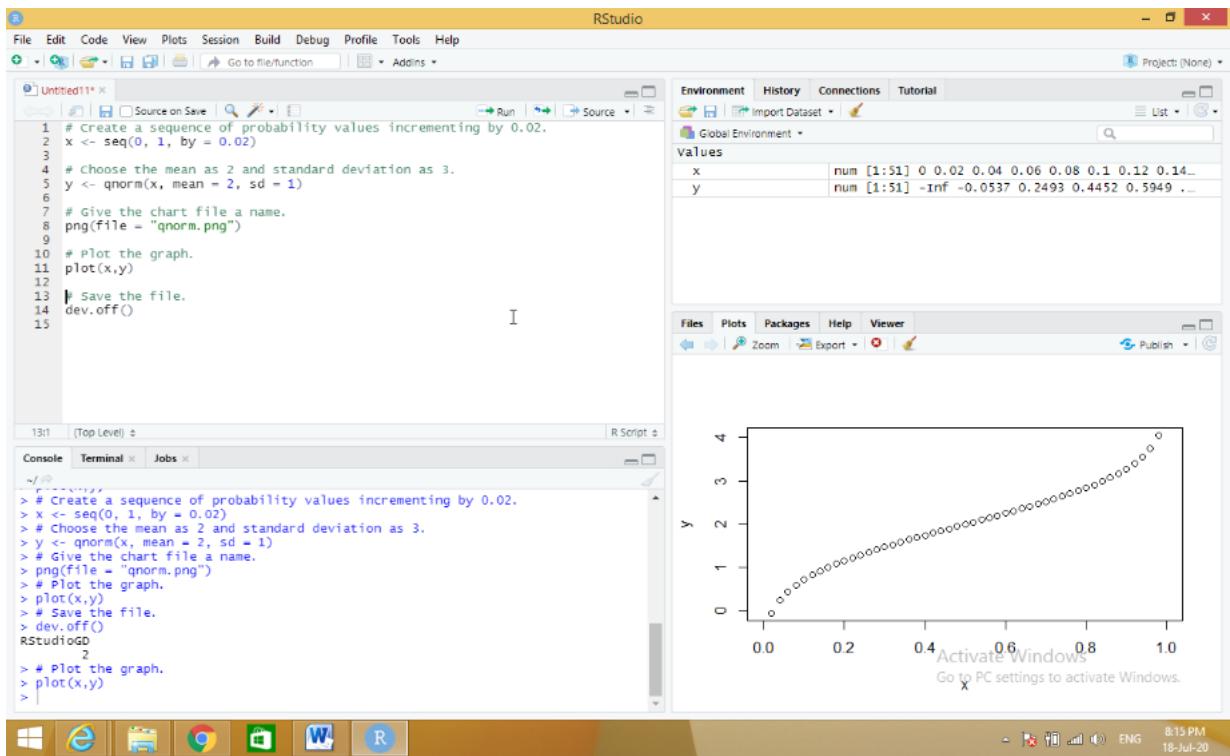
output:



qnorm

```
# Create a sequence of probability values incrementing by 0.02.  
x <- seq(0, 1, by = 0.02)  
  
# Choose the mean as 2 and standard deviation as 3.  
y <- qnorm(x, mean = 2, sd = 1)  
  
# Give the chart file a name.  
png(file = "qnorm.png")  
  
# Plot the graph.  
plot(x,y)  
  
# Save the file.  
dev.off()
```

output:



rnorm

Create a sample of 50 numbers which are normally distributed.

```
y <- rnorm(50)
```

Give the chart file a name.

```
png(file = "rnorm.png")
```

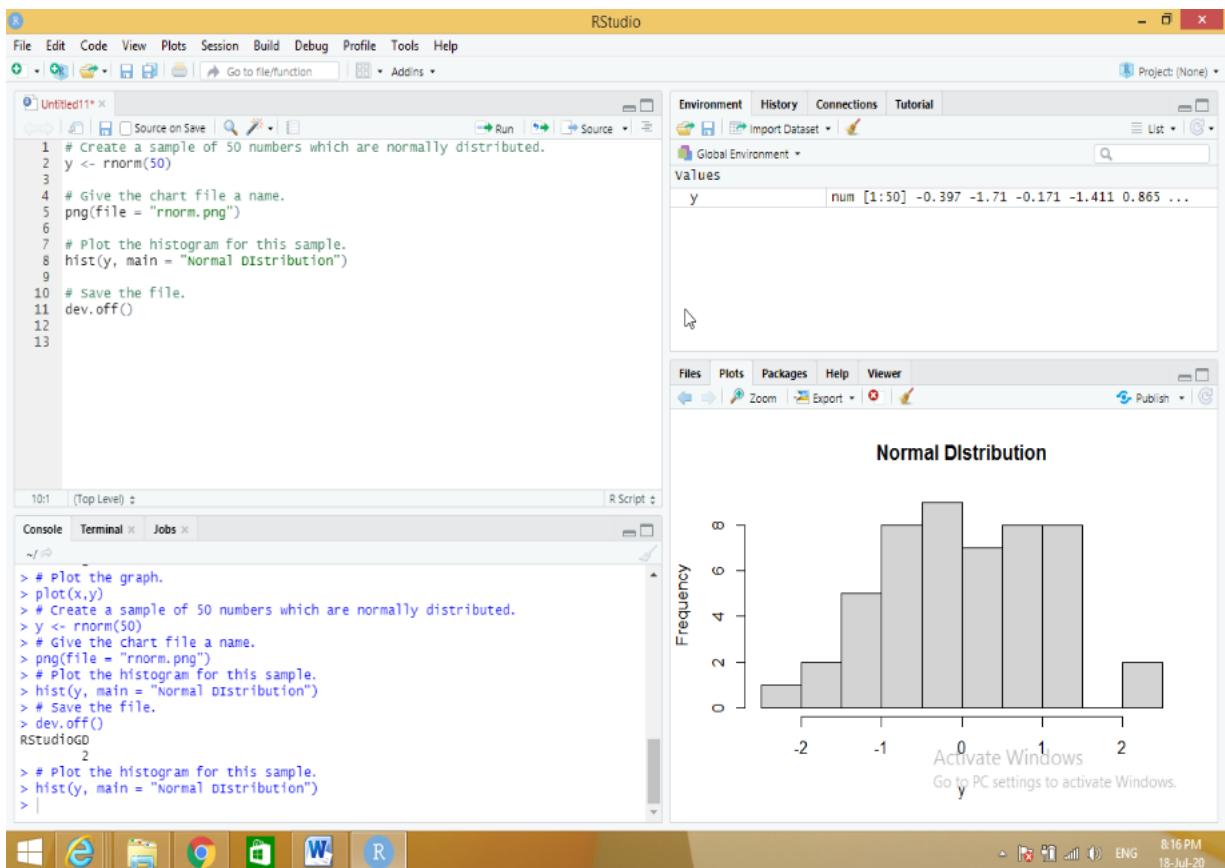
Plot the histogram for this sample.

```
hist(y, main = "Normal DIstribution")
```

Save the file.

```
dev.off()
```

output:



b) Binomial Distribution

dbinom

Source code:

```
# Create a sample of 50 numbers which are incremented by 1.
```

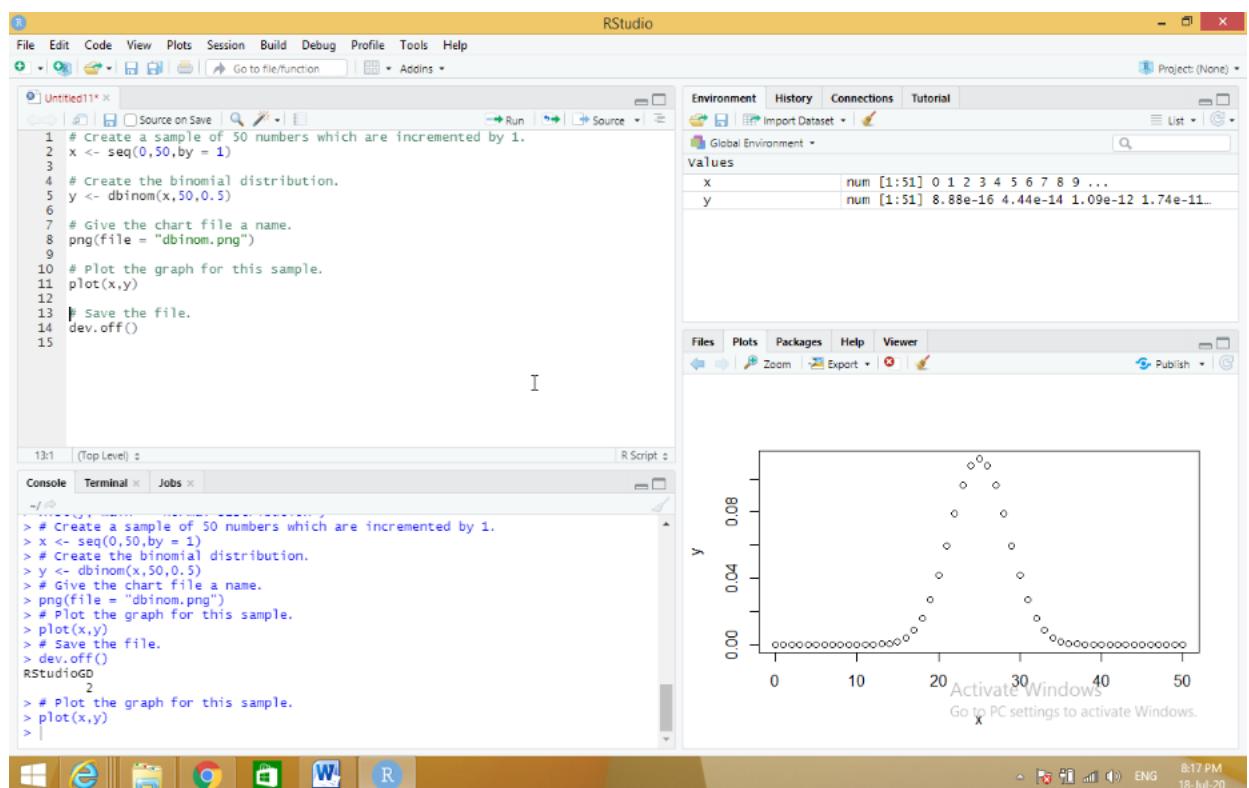
```
x <- seq(0,50,by = 1)
```

```
# Create the binomial distribution.
```

```
y <- dbinom(x,50,0.5)
```

```
# Give the chart file a name.  
png(file = "dbinom.png")  
  
# Plot the graph for this sample.  
plot(x,y)  
  
# Save the file.  
dev.off()
```

Output:



Pbinom

```
# Probability of getting 26 or less heads from a 51 tosses of a coin.  
x <- pbinom(26,51,0.5)  
  
print(x)
```

output:

```
print(x)
[1] 0.610116
```

qbinom

```
# How many heads will have a probability of 0.25 will come out when a coin
# is tossed 51 times.
x <- qbinom(0.25,51,1/2)
```

```
print(x)
```

output:

```
print(x)
[1] 23
```

rbinom

```
# Find 8 random values from a sample of 150 with probability of 0.4.
x <- rbinom(8,150,.4)
```

```
print(x)
```

output:

```
print(x)
[1] 68 59 55 49 51 59 53 55
```

5. Write R program to do the following tests with the sample data and visualize the results graphically.
Source code:

a) χ^2 -test

```
library("MASS")
print(str(Cars93))
# Loading the Mass library.

# Creating a data frame from the main data set.
car_data<- data.frame(Cars93$AirBags, Cars93$Type)
# Creating a table with the needed variables.
car_data = table(Cars93$AirBags, Cars93$Type)
print(car_data)
# Performing the Chi-Square test.
print(chisq.test(car_data))
```

output:

```
print(str(Cars93))
'data.frame':      93 obs. of  27 variables:
 $ Manufacturer :Factor w/ 32 levels "Acura", "Audi", ...:112 234 444 5...
 $ Model        :Factor w/ 93 levels "100", "190E", "240", ...:495 691 624 547 473 ...
 $ Type         :Factor w/ 6 levels "Compact", "Large", ...:4 313 332 232 ...
 $ Min.Price    :num 12.9 29.2 25.9 30.8 23.7 14.2 19.9 22.6 26.3 33 ...
 $ Price        :num 15.9 33.9 29.1 37.7 30.1 5.7 20.8 23.7 26.3 34.7 ...
 $ Max.Price   :num 18.8 38.7 32.3 44.6 36.2 17.3 21.7 24.9 26.3 36.3 ...
 $ MPG.city    :int 25 18 20 19 22 22 19 16 19 16 ...
 $ MPG.highway :int 31 25 26 26 30 31 28 25 27 25 ...
 $ AirBags     :Factor w/ 3 levels "Driver & Passenger", ...:31 21 22 22 22 ...
 $ DriveTrain  :Factor w/ 3 levels "4WD", "Front", ...:2 2 2 2 3 2 2 3 2 ...
 $ Cylinders   :Factor w/ 6 levels "3", "4", "5", "6", ...:2 4 4 4 2 2 4 4 5 ...
 $ EngineSize  :num 1.8 3.2 2.8 2.8 3.5 2.3 3.8 5.7 3.8 4.9 ...
 $ Horsepower  :int 140 200 172 172 208 110 170 180 170 200 ...
 $ RPM         :int 6300 5500 5500 5500 5700 5200 4800 4000 4800 4100 ...
 $ Rev.per.mile:int 2890 2335 2280 2535 2545 2565 1570 1320 1690 1510 ...
 $ Man.trans.avail :Factor w/ 2 levels "No", "Yes":2 2 2 2 1 1 1 1 ...
 $ Fuel.tank.capacity: num 13.2 18.9 21.121.116.4 18.2 318.8 18 ...
 $ Passengers  :int 5 5 5 6 4 6 6 6 5 6 ...
 $ Length       :int 177 195 180 193 186 189 200 216 198 206 ...
 $ Wheelbase    :int 102 115 102 106 109 105 111 116 108 114 ...
 $ Width        :int 687 16770 69 6974 78 73 73 ...
```

```
$ Turn.circle    :int 37383737394142454143...
$ Rear.seat.room :num 26.5302831272830.530.526.535...
$ Luggage.room   :int 11151417131617211418...
$ Weight         :int 2705356033753405364028803470410534953620...
$ Origin         :Factor w/ 2 levels "USA","non-USA":2 2 2 2 2 1 1 1 1...
$ Make           :Factor w/ 93 levels "Acura Integra",..:124 35679810...
NULL
> # Creating a data frame from the main data set.
> car_data<- data.frame(Cars93$AirBags, Cars93>Type)
> # Creating a table with the needed variables.
> car_data=table(Cars93$AirBags, Cars93>Type)
> print(car_data)
```

	Compact	Large	Midsize	Small	Sporty	Van
Driver & Passenger	2	4	7	0	3	0
Driver only	9	7	11	5	8	3
None	5	0	4	16	3	6

> # Performing the Chi-Square test.

> print(chisq.test(car_data))

Pearson's Chi-squared test

```
data: car_data
X-squared = 33.001, df = 10, p-value = 0.0002723
```

b) t-test

```
x <- c(0.593, 0.142, 0.329, 0.691, 0.231, 0.793, 0.519, 0.392, 0.418)
t.test(x, alternative="greater", mu=0.3)
```

output:

t.test(x, alternative="greater", mu=0.3)

One Sample t-test

```
data: x
t = 2.2051, df = 8, p-value = 0.02927
alternative hypothesis: true mean is greater than 0.3
95 percent confidence interval:
0.3245133 Inf
sample estimates:
mean of x
0.4564444
```

c) F-test

```
install.packages("randomForest")
# Load the party package. It will automatically load other
# required packages.
library(party)

# Print some records from data set readingSkills.
print(head(readingSkills))

# Load the party package. It will automatically load other
# required packages.
library(party)
library(randomForest)

# Create the forest.
output.forest <- randomForest(nativeSpeaker ~ age + shoeSize + score,
                                data = readingSkills)

# View the forest results.
print(output.forest)
```

Output:

print(output.forest)

Call:

randomForest(formula = nativeSpeaker ~ age + shoeSize + score, data = readingSkills)

Type of random forest: classification

Number of trees: 500

No. of variables tried at each split: 1

OOB estimate of error rate: 1.5%

Confusion matrix:

	no	yes	class.error
no	99	1	0.01
yes	2	98	0.02

Viva Questions:

1) What is the Main Goal of Data Mining?

Data mining is an interdisciplinary subfield of computer science and statistics with an overall goal to extract information (with intelligent methods) from a data set and transform the information into a comprehensible structure for further use.

2) What are the Two Types of Data Mining Tasks?

The data mining tasks can be classified generally into two types based on what a specific task tries to achieve. Those two categories are descriptive tasks and predictive tasks.

3) What are the Data Mining Functionalities?

- Class/Concept Description
- Mining of Frequent Patterns
- Mining of Associations
- Mining of Correlations
- Mining of Clusters
- Classification
- Prediction
- Outlier Analysis

Evolution Analysis

4) What are the Features of WEKA?

Weka features include machine learning, data mining, preprocessing, classification, regression, clustering, association rules, attribute selection, experiments, workflow and visualization.

5) Navigate the options available in the WEKA.

- Preprocess
- Classify
- Cluster
- Associate
- Select Attributes
- Visualize

6) What is ARFF file format?

An ARFF (Attribute-Relation File Format) file is an ASCII text file that describes a list of instances sharing a set of attributes.

7) Define Support and Confidence.

Support represents the popularity of that product of all the product transactions. Confidence can be interpreted as the likelihood of purchasing both the products A and B.

8) What are the frequent patterns?

Frequent patterns are itemsets, subsequences, or substructures that appear in a data set with frequency no less than a user-specified threshold. For example, a set of items, such as milk and bread, that appear frequently together in a transaction data set, is a frequent itemset.

9) Where we are using Apriori Algorithm in Real time scenario?

Usually, you operate this algorithm on a database containing a large number of transactions. One such example is the items customers buy at a supermarket. It helps the customers buy their items with ease, and enhances the sales performance of the departmental store.

10) Explain Association rule with a suitable example.

A classic example of association rule mining refers to a relationship between diapers and beers. The example, which seems to be fictional, claims that men who go to a store to buy diapers are also likely to buy beer.

11) What is Apriori Property?

Apriori is an algorithm for frequent item set mining and association rule learning over relational databases. It proceeds by identifying the frequent individual items in the database and extending them to larger and larger item sets as long as those item sets appear sufficiently often in the database.

12) How can we further improve the efficiency of Apriori-based mining?

Based on the inherent defects of Apriori algorithm, some related improvements are carried out: 1) using new database mapping way to avoid scanning the database repeatedly; 2) further pruning frequent itemsets and candidate itemsets in order to improve joining efficiency;

13) Define Classification.

Classification is a data mining function that assigns items in a collection to target categories or classes. The goal of classification is to accurately predict the target class for each case in the data.

14) Define Prediction.

The prediction, as its name implied, is one of a data mining techniques that discovers relationship between independent variables and relationship between dependent variables.

15) What are the classification techniques in data mining?

- Logistic Regression
- Naïve Bayes
- Stochastic Gradient Descent
- K-Nearest Neighbours
- Decision Tree
- Random Forest
- Support Vector Machine

16) What are the advantages of different classification algorithms

Mining Based Methods are cost effective and efficient. Helps in identifying criminal suspects. Helps in predicting risk of diseases.

17) What is the Kappa Statistic

“ The Kappa statistic (or value) is a metric that compares an Observed Accuracy with an Expected Accuracy (random chance). The kappa statistic is used not only to evaluate a single classifier, but also to evaluate classifiers amongst themselves.

18) Which classification algorithm is best for prediction and analysis?

Time Series Model. The time series model comprises a sequence of data points captured, using time as the input parameter. ...

Random Forest. Random Forest is perhaps the most popular classification algorithm, capable of both classification and regression.

19) What are the classification algorithms in data mining?

Six classification algorithms—Naive Bayes, Bayesian networks, J48, random forest, multilayer perceptron, and logistic regression

20) Which data mining algorithm provides best accuracy for classification?

The performances of the algorithms were compared according to accuracy, root mean squared error, ROC area, F-measure, precision, and recall criteria, and the logistic regression classification algorithm was found to be the best algorithm.

21) What are the best classification algorithms?

Top 5 Data Mining Algorithms for Classification

- Decision Trees.
- Logistic Regression.
- Naive Bayes Classification.
- k-nearest neighbors.
- Support Vector Machine.

22) Which are the best classifier algorithms for disease predictions?

Comparative analysis of classification techniques has shown that decision tree classifiers are simple and accurate [9]. Naïve Bayes was found to be the best algorithm, followed by neural networks and decision trees

23) Explain Decision Tree.

A decision tree is a diagram or chart that people use to determine a course of action or show a statistical probability. Each branch of the decision tree represents a possible decision, outcome, or reaction.

24) What is the Cross-Validation?

Cross validation is a technique for assessing how the statistical analysis generalizes to an independent data set. It is a technique for evaluating machine learning models by training several models on subsets of the available input data and evaluating them on the complementary subset of the data.

25) What is the Naïve-Bayes Classification?

It is a **classification** technique based on **Bayes' Theorem** with an assumption of independence among predictors. In simple terms, a **Naive Bayes classifier** assumes that the presence of a particular feature in a class is unrelated to the presence of any other feature.

26) What are the methods in classification methods in data mining?

There are many **techniques** for solving **classification** problems: **classification** trees, logistic regression, discriminant analysis, neural networks, boosted trees, random forests, deep learning **methods**, nearest neighbors, support vector machines.

27) Briefly explain the K-nearest neighbor algorithm

KNN works by finding the distances between a query and all the examples in the data, selecting the specified number examples (**K**) closest to the query, then votes for the most frequent label (in the case of classification) or averages the labels

28) How do you choose an algorithm for a classification problem?

Size of the training data. It is usually recommended to gather a good amount of data to get reliable predictions.

Accuracy and/or Interpretability of the output.

Speed or Training time.

Linearity.

Number of features.

29) What are the most useful algorithms used for data mining

K Nearest Neighbors Algorithm, Naïve Bayes Algorithm, SVM Algorithm, ANN Algorithm, 48 Decision Trees, Support Vector Machines, and SenseClusters.

30) What is the difference between classification and prediction?

Classification is the process of identifying the category or class label of the new observation to which it belongs. **Prediction** is the process of identifying the missing or unavailable numerical data for a new observation.

31) What is clustering with example?

Clustering is the task of dividing the population or data points into a number of groups such that data points in the same groups are more similar to other data points in the same group than those in other groups. In simple words, the aim is to segregate groups with similar traits and assign them into clusters.

32) What are the examples of clustering?

Identifying Fake News. Fake news is not a new phenomenon, but it is one that is becoming prolific.

Spam filter

Marketing and Sales

Classifying network traffic

Identifying fraudulent or criminal activity

Document analysis

Fantasy Football and Sports

33) Why Clustering is used in data mining

Clustering in Data Mining helps in the classification of animals and plants are done using similar functions or genes in the field of biology. It helps in gaining insight into the structure of the species. Areas are identified using the clustering in data mining.

34) Why we use K means clustering

The K-means clustering algorithm is used to find groups which have not been explicitly labeled in the data. This can be used to confirm business assumptions about what types of groups exist or to identify unknown groups in complex data sets.

35) What are the advantages and disadvantages of K means clustering

Advantages

Relatively simple to implement.

Scales to large data sets.

Guarantees convergence.

Disadvantages

Clustering data of varying sizes and density.

Clustering outliers.

36) Which is the best clustering algorithm?

K-means Clustering Algorithm. ...

Mean-Shift Clustering Algorithm. ...

DBSCAN – Density-Based Spatial Clustering of Applications with Noise.

37) Explain the Hierarchical Clustering

HCA is an unsupervised clustering algorithm which involves creating clusters that have predominant ordering from top to bottom.

38) State the other Clustering Techniques.

Connectivity-based Clustering (Hierarchical clustering)

Centroids-based Clustering (Partitioning methods)

- Distribution-based Clustering,
- Density-based Clustering (Model-based methods)
- Fuzzy Clustering.
- Constraint-based (Supervised Clustering)

39) What are the two types of hierarchical clustering?

Agglomerative clustering: It's also known as AGNES (Agglomerative Nesting). It works in a bottom-up manner. ...

Divisive hierarchical clustering: It's also known as DIANA (Divise Analysis) and it works in a top-down manner.

40) What are the disadvantages of agglomerative hierarchical clustering?

One drawback is that groups with close pairs can merge sooner than is optimal, even if those groups have overall dissimilarity. Complete Linkage: calculates similarity of the farthest away pair.

41) What is the use of hierarchical clustering?

Hierarchical clustering is the most popular and widely used method to analyze social network data. In this method, nodes are compared with one another based on their similarity. Larger groups are built by joining groups of nodes based on their similarity.

42) What is the difference between K means and hierarchical clustering?

Hierarchical clustering can't handle big data well but K Means clustering can. This is because the time complexity of K Means is linear i.e. $O(n)$ while that of hierarchical clustering is quadratic i.e. $O(n^2)$.

43) What is the mean, median and mode

The mean (average) of a data set is found by adding all numbers in the data set and then dividing by the number of values in the set. The median is the middle value when a data set is ordered from least to greatest. The mode is the number that occurs most often in a data set.

44) How do you find the mode in r?

To find the mode, or modal value, it is best to put the numbers in order. Then count how many of each number. A number that appears most often is the mode.

45) Define analysis of covariance

Analysis of Covariance (ANCOVA) is the inclusion of a continuous variable in addition to the variables of interest (i.e., the dependent and independent variable) as means for control.

46) What is analysis of covariance used for?

Analysis of covariance is used to test the main and interaction effects of categorical variables on a continuous dependent variable, controlling for the effects of selected other continuous variables, which co-vary with the dependent. The control variables are called the "covariates."

47) How do you analyze covariance?

The Analysis of covariance (ANCOVA) is done by using linear regression. This means that Analysis of covariance (ANCOVA) assumes that the relationship between the independent variable and the dependent variable must be linear in nature.

48) Why should we use R?

R is a programming language for statistical computing and graphics that you can use to clean, analyze, and graph your data. It is widely used by researchers from diverse disciplines to estimate and display results and by teachers of statistics and research methods.

49) Define regression and its types

Regression is a technique used to model and analyze the relationships between variables and often times how they contribute and are related to producing a particular outcome together. The two basic types of regression are simple linear regression and multiple linear regression.

50) What is regression in statistics with example?

Linear regression quantifies the relationship between one or more predictor variable(s) and one outcome variable. ... For example, it can be used to quantify the relative impacts of age, gender, and diet (the predictor variables) on height (the outcome variable).

51) What is meant by linear regression?

Linear regression attempts to model the relationship between two variables by fitting a linear equation to observed data. One variable is considered to be an explanatory variable, and the other is considered to be a dependent variable.

52) What is multiple linear regression explain with example

Multiple linear regression (MLR), also known simply as multiple regression, is a statistical technique that uses several explanatory variables to predict the outcome of a response variable. Multiple regression is an extension of linear (OLS) regression that uses just one explanatory variable.

53) What is the difference between linear regression and multiple regression?

Linear regression is one of the most common techniques of regression analysis. Multiple regression is a broader class of regressions that encompasses linear and nonlinear regressions with multiple explanatory variables.

54) What is difference between linear and logistic regression

The essential difference between these two is that

Logistic regression is used when the dependent variable is binary in nature.

In contrast, linear regression is used when the dependent variable is continuous and the nature of the regression line is linear.

55) What is time series analysis with example?

A time series is a sequence taken at successive equally spaced points in time. Thus it is a sequence of discrete-time data. Examples of time series are heights of ocean tides, counts of sunspots, and the daily closing value of the Dow Jones Industrial Average.

56) How do you analyze time series data in R?

The first thing that you will want to do to analyse your time series data will be to read it into R, and to plot the time series. You can read data into R using the scan() function, which assumes that your data for successive time points is in a simple text file with one column.

57) What is the purpose of time series analysis?

Time series analysis can be useful to see how a given asset, security, or economic variable changes over time. It can also be used to examine how the changes associated with the chosen data point compare to shifts in other variables over the same time period.

58) What is difference between linear and nonlinear?

While a linear equation has one basic form, nonlinear equations can take many different forms.

Thetas represent the parameters and X represents the predictor in the nonlinear functions. Unlike linear regression, these functions can have more than one parameter per predictor variable.

59) What is decision tree and example?

Decision Trees are a type of Supervised Machine Learning (that is you explain what the input is and what the corresponding output is in the training data) where the data is continuously split according to a certain parameter. ... An example of a decision tree can be explained using above binary tree.

60) How do you create a decision tree in R?

Step 1: Import the data.

Step 2: Clean the dataset.

Step 3: Create train/test set.

Step 4: Build the model.

Step 5: Make prediction.

Step 6: Measure performance.

Step 7: Tune the hyper-parameters.

61) What is the Rnorm function in R?

rnorm is the R function that simulates random variates having a specified normal distribution. As with pnorm , qnorm , and dnorm , optional arguments specify the mean and standard deviation of the distribution.

62) What is the Pnorm and Qnorm in R

pnorm: cumulative density function of the normal distribution. qnorm: quantile function of the normal distribution.

63) What is the normal distribution with example?

For example, heights, blood pressure, measurement error, and IQ scores follow the normal distribution. It is also known as the Gaussian distribution and the bell curve.

64) Where is normal distribution used?

Normal distribution, also called Gaussian distribution, the most common distribution function for independent, randomly generated variables. Its familiar bell-shaped curve is ubiquitous in statistical reports, from survey analysis and quality control to resource allocation.

65) What is Dbinom

The function dbinom returns this probability. There are three required arguments: the value(s) for which to compute the probability (j), the number of trials (n), and the success probability for each trial (p).

66) What is the difference between Pbinom and Dbinom

dbinom is a probability mass function of binomial distribution, while pbinom is a cumulative distribution function of this distribution.

67) Define chi square test and its application

The Chi Square test is a statistical hypothesis test in which the sampling distribution of the test statistic is a chi-square distribution when the null hypothesis is true. The Chi square test is used to compare a group with a value, or to compare two or more groups, always using categorical data.

68) When should we use chi square test

The Chi Square statistic is commonly used for testing relationships between categorical variables. The null hypothesis of the Chi-Square test is that no relationship exists on the categorical variables in the population; they are independent.

69) What are the limitations of chi square?

Chi-square, like any analysis has its limitations. One of the limitations is that all participants measured must be independent, meaning that an individual cannot fit in more than one category. If a participant can fit into two categories a chi-square analysis is not appropriate.

70) What is the difference between t test and F TEST?

T-test is used to test if two samples have the same mean. The assumptions are that they are samples from normal distribution. F-test is used to test if two samples have the same variance.

71) What is the difference between chi square and t test?

A t-test tests a null hypothesis about two means; most often, it tests the hypothesis that two means are equal, or that the difference between them is zero. A chi-square test tests a null hypothesis about the relationship between two variables.

72) Where do we use chi square test

The Chi Square statistic is commonly used for testing relationships between categorical variables. The null hypothesis of the Chi-Square test is that no relationship exists on the categorical variables in the population; they are independent.