# INDEX

| Prac.No. | Practical | Date | sing |
|---|---|---|---|
| 1 | Install, configure and run Hadoop and HDFS ad explore HDFS | | |
| 2 | Implement Decision tree classification techniquesb. | | |
| 3 | . Implement SVM classification techniques. | | |
| 4 | Implement of REGRESSION MODLE. | | |
| 5 | Implement of Simple Linear Regaression. | | |
| 6 | Implement of  Multiple Linear Regression. | | |
| 7 | Implement of Logistic regression. | | |
| 8 | Read a datafile grades_km_input.csv and apply k-means clustering. | | |
| 9 | Perform Apriori algorithm using Groceries dataset from the R arules package. | | |

**Shanti Chourasiya**

**Roll no.2023ITI1**

**Class-MSc.IT**

# Practical 1

**Aim:** -Install, configure and run Hadoop and HDFS and explore HDFS on Windows

**Code:**

**Steps to Install Hadoop**
1. Install Java JDK 1.8
2. Download Hadoop and extract and place under C drive
3. Set Path in Environment Variables
4. Config files under Hadoop directory
5. Create folder datanode and namenode under data directory
6. Edit HDFS and YARN files
7. Set Java Home environment in Hadoop environment
8. Setup Complete. Test by executing start-all.cmd

**There are two ways to install Hadoop, i.e.**
9. Single node
10. Multi node
Here, we use multi node cluster.
1. **Install Java**
    11.     – Java JDK Link to download
            https://www.oracle.com/java/technologies/javase-jdk8-downloads.html
    12.     – extract and install Java in C:\Java
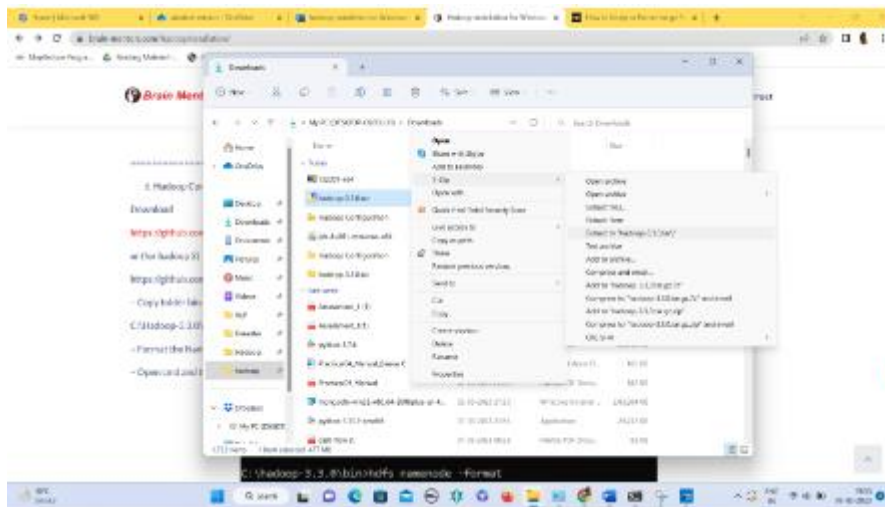    13.     – open cmd and type -> javac -version

```
C:\Users>cd Beena

C:\Users\Beena>java -version
java version "1.8.0_361"
Java(TM) SE Runtime Environment (build 1.8.0_361-b09)
Java HotSpot(TM) 64-Bit Server VM (build 25.361-b09, mixed mode)
```
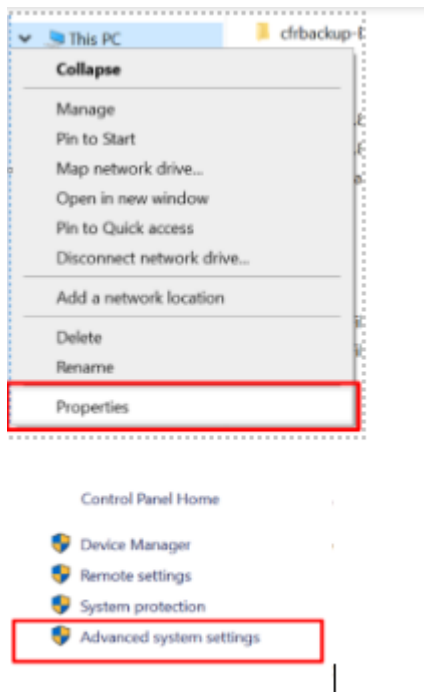
2. **Download Hadoop**
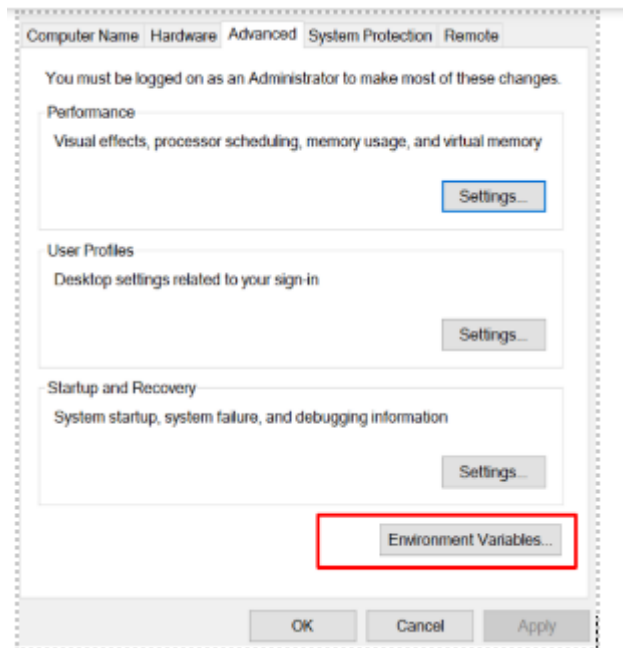    https://www.apache.org/dyn/closer.cgi/hadoop/common/hadoop-3.3.0/hadoop-3.3.0.tar.gz

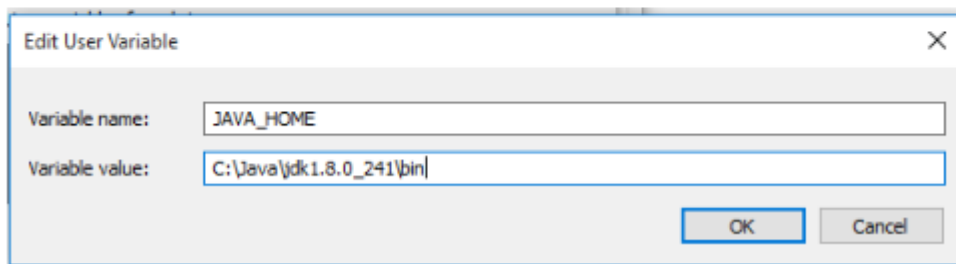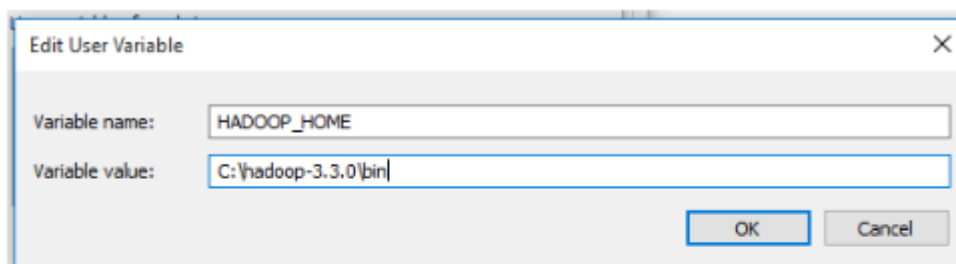    • right click .rar.gz file -> show more options -> 7-zip->and extract to C:\Hadoop-3.3.0\

3 **Set the path JAVA_HOME Environment variable**
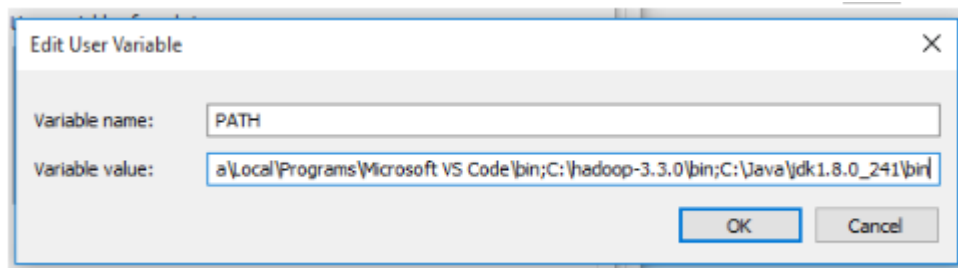
4 **Set the path HADOOP_HOME Environment variable**

Click on **New to both user variables and system variables.**





**Click on user variable -> path -> edit-> add path for Hadoop and java upto 'bin'**

**Edit User Variable**

| Variable name: | PATH |
| Variable value: | a\Local\Programs\Microsoft VS Code\bin;C:\hadoop-3.3.0\bin;C:\Java\jdk1.8.0_241\bin |

OK    Cancel

Click Ok, Ok, Ok.
5.          **Configurations**
**Edit file C:/Hadoop-3.3.0/etc/hadoop/core-site.xml,**
paste the xml code in folder and save

============================================================

```xml
<configuration>
<property>
    <name>fs.defaultFS</name>
    <value>hdfs://localhost:9000</value>
 </property>
</configuration>
```
============================================================

**Rename "mapred-site.xml.template" to "mapred-site.xml" and edit this file C:/Hadoop-3.3.0/etc/hadoop/mapred-site.xml, paste xml code and save this file.**
============================================================

```xml
<configuration>
  <property>
    <name>mapreduce.framework.name</name>
    <value>yarn</value>
  </property>
</configuration>
```
============================================================

**Create folder "data" under "C:\Hadoop-3.3.0"**

**Create folder "datanode" under "C:\Hadoop-3.3.0\data"**

**Create folder "namenode" under "C:\Hadoop-3.3.0\data"**

============================================================

**Edit file C:\Hadoop-3.3.0/etc/hadoop/hdfs-site.xml,**

**paste xml code and save this file.**

```xml
<configuration>
<property>
    <name>dfs.replication</name>
```

```
      <value>1</value>
  </property>

  <property>
     <name>dfs.namenode.name.dir</name>
     <value>/hadoop-3.3.0/data/namenode</value>
  </property>
  <property>
     <name>dfs.datanode.data.dir</name>
     <value>/hadoop-3.3.0/data/datanode</value>
  </property>
 </configuration>
```
=======================================================

**Edit file C:/Hadoop-3.3.0/etc/hadoop/yarn-site.xml,**

**paste xml code and save this file.**

```
<configuration>

</configuration>
```
=======================================================

6.      **Edit file C:/Hadoop-3.3.0/etc/hadoop/hadoop-env.cmd**
Find "JAVA_HOME=%JAVA_HOME%" and replace it as
set JAVA_HOME="C:\Java\jdk1.8.0_361"
=======================================================

7.      **Download "redistributable" package**


        **Download and run VC_redist.x64.exe**

This is a "redistributable" package of the Visual C runtime code for 64-bit applications, from
Microsoft. It contains certain shared code that every application written with Visual C expects to
have available on the Windows computer it runs on.

8.      **Hadoop Configurations**
**Download bin folder from**

**https://github.com/s911415/apache-hadoop-3.1.0-winutils**

**– Copy the bin folder to c:\hadoop-3.3.0. Replace the existing bin folder.**

9.      **copy "hadoop-yarn-server-timelineservice-3.0.3.jar" from ~\hadoop-3.0.3\share\hadoop\yarn\timelineservice to ~\hadoop-3.0.3\share\hadoop\yarn folder.**


10.      **Format the NameNode**
**– Open cmd 'Run as Administrator' and type command "hdfs namenode –format"**

**11. Testing**

**– Open cmd 'Run as Administrator' and change directory to** **C:\Hadoop-3.3.0\sbin**

**– type** **start-all.cmd**

 **OR**

 **- type** **start-dfs.cmd**

**– type** **start-yarn.cmd**



**– You will get 4 more running threads for Datanode, namenode, resouce manager and node manager**

**Output:**

12. Type JPS command to start-all.cmd command prompt, you will get following output.



13. Run http://localhost:9870/ from any browser
Or http://localhost:50070/

# Practical 2

**Aim: - Implement Decision tree classification techniquesb**.

install.packages('datasets')

install.packages('caTools')

install.packages('party')

install.packages('dplyr')

install.packages('magrittr')

library(datasets)

library(caTools)

library(party)

library(dplyr)

library(magrittr)

data("readingSkills")

head(readingSkills)

sample_data = sample.split(readingSkills, SplitRatio = 0.8)

train_data <- subset(readingSkills, sample_data == TRUE)

test_data <- subset(readingSkills, sample_data == FALSE)

model<- ctree(nativeSpeaker ~ ., train_data)

plot(model)

```
File   Edit   Code   View   Plots   Session   Build   Debug   Profile   Tools   Help
```

```
decision tree.R ×
```

⚠ Packages caTools , dplyr , and magrittr others required but are not installed. Install Don't Show Again

```
 2   install.packages('caTools')
 3   install.packages('party')
 4   install.packages('dplyr')
 5   install.packages('magrittr')
 6   library(datasets)
 7   library(caTools)
 8   library(party)
 9   library(dplyr)
10   library(magrittr)
11
12   data("readingSkills")
13   head(readingSkills)
14   sample_data = sample.split(readingSkills, SplitRatio = 0.8)
15   train_data <- subset(readingSkills, sample_data == TRUE)
16   test_data <- subset(readingSkills, sample_data == FALSE)
17   model<- ctree(nativeSpeaker ~ ., train_data)
18   plot(model)
19
```

```
18:12    (Top Level) ↕                                                          R Script
```

```
Console   Terminal ×   Background Jobs ×
```

```
R  R 4.4.1 · ~/
```

```
  nativeSpeaker age shoeSize    score
1           yes   5 24.83189 32.29385
2           yes   6 25.95238 36.63105
3            no  11 30.42170 49.60593
4           yes   7 28.66450 40.28456
5           yes  11 31.88207 55.46085
6           yes  10 30.07843 52.83124
> sample_data = sample.split(readingSkills, SplitRatio = 0.8)
> train_data <- subset(readingSkills, sample_data == TRUE)
> test_data <- subset(readingSkills, sample_data == FALSE)
> model<- ctree(nativeSpeaker ~ ., train_data)
> plot(model)
>
```
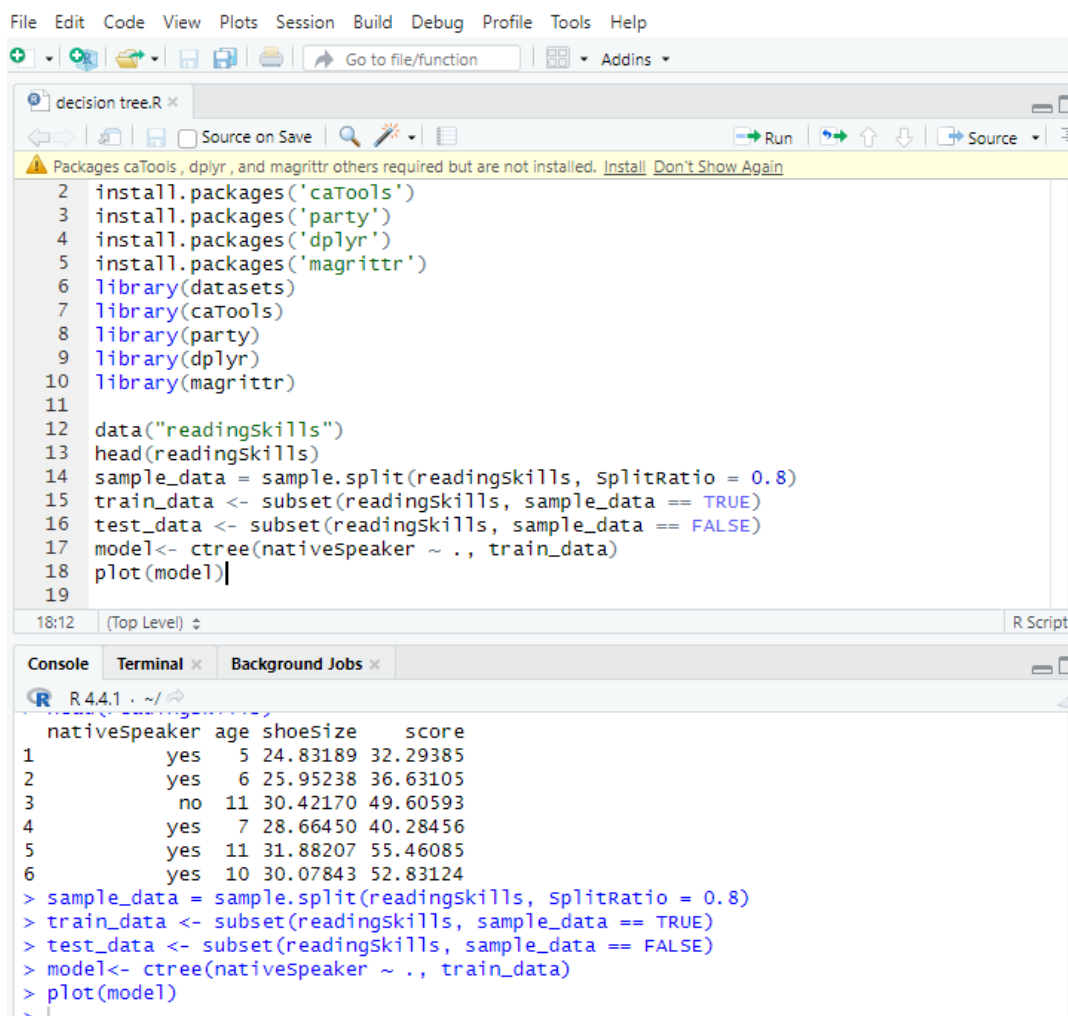
# Practical 3

**Aim:- Implement SVM classification techniques**

#Code for installation of all necessary packages

install.packages("caret")

install.packages("ggplot2")

install.packages("GGally")

install.packages("psych")

install.packages("ggpubr")

install.packages("reshape")


**# Code for importation of all necessary packages**

library(caret)

library(ggplot2)

library(GGally)

library(psych)

library(ggpubr)

library(reshape)

# Code

```
df <- read.csv("D:\\diabetes.csv")

head(df)
```

```
> # Code
> df <- read.csv("D:\\diabetes.csv")
> head(df)
  Pregnancies Glucose BloodPressure SkinThickness Insulin  BMI DiabetesPedigreeFunction Age Outcome
1           6     148            72            35       0 33.6                    0.627  50       1
2           1      85            66            29       0 26.6                    0.351  31       0
3           8     183            64             0       0 23.3                    0.672  32       1
4           1      89            66            23      94 28.1                    0.167  21       0
5           0     137            40            35     168 43.1                    2.288  33       1
6           5     116            74             0       0 25.6                    0.201  30       0
>
```

```
# Code

sum(is.na(df))

# Code

dim(df)
```

```
> # Code
> sum(is.na(df))
[1] 0
> # Code     .
> dim(df)
[1] 768    9
>
```

```
# Code

sapply(df, class)

# Code

summary(df) # to calculate the summary of our dataset
```

```
> # Code
> sapply(df, class)
              Pregnancies                  Glucose            BloodPressure             SkinThickness
                "integer"                "integer"                "integer"                "integer"
                  Insulin                      BMI DiabetesPedigreeFunction                      Age
                "integer"                "numeric"                "numeric"                "integer"
                  Outcome
                "integer"
> # Code
> summary(df) # to calculate the summary of our dataset
  Pregnancies        Glucose      BloodPressure    SkinThickness      Insulin           BMI
 Min.   : 0.000   Min.   :  0.0   Min.   :  0.00   Min.   : 0.00   Min.   :  0.0   Min.   : 0.00
 1st Qu.: 1.000   1st Qu.: 99.0   1st Qu.: 62.00   1st Qu.: 0.00   1st Qu.:  0.0   1st Qu.:27.30
 Median : 3.000   Median :117.0   Median : 72.00   Median :23.00   Median : 30.5   Median :32.00
 Mean   : 3.845   Mean   :120.9   Mean   : 69.11   Mean   :20.54   Mean   : 79.8   Mean   :31.99
 3rd Qu.: 6.000   3rd Qu.:140.2   3rd Qu.: 80.00   3rd Qu.:32.00   3rd Qu.:127.2   3rd Qu.:36.60
 Max.   :17.000   Max.   :199.0   Max.   :122.00   Max.   :99.00   Max.   :846.0   Max.   :67.10
 DiabetesPedigreeFunction      Age            Outcome
 Min.   :0.0780           Min.   :21.00   Min.   :0.000
 1st Qu.:0.2437           1st Qu.:24.00   1st Qu.:0.000
 Median :0.3725           Median :29.00   Median :0.000
 Mean   :0.4719           Mean   :33.24   Mean   :0.349
 3rd Qu.:0.6262           3rd Qu.:41.00   3rd Qu.:1.000
 Max.   :2.4200           Max.   :81.00   Max.   :1.000
>
```

# Code

```
a <- ggplot(data = df, aes(x = Pregnancies)) +

  geom_histogram( color = "red", fill = "blue", alpha = 0.1) +

  geom_density()

b <- ggplot(data = df, aes(x = Glucose)) +

  geom_histogram( color = "red", fill = "blue", alpha = 0.1) +

  geom_density()

c <- ggplot(data = df, aes(x = BloodPressure)) +

  geom_histogram( color = "red", fill = "blue", alpha = 0.1) +

  geom_density()


d <- ggplot(data = df, aes(x = SkinThickness)) +

  geom_histogram( color = "red", fill = "blue", alpha = 0.1) +

  geom_density()

e <- ggplot(data = df, aes(x = Insulin)) +

  geom_histogram( color = "red", fill = "blue", alpha = 0.1) +

  geom_density()

f <- ggplot(data = df, aes(x = BMI)) +

  geom_histogram( color = "red", fill = "blue", alpha = 0.1) +

  geom_density()

g <- ggplot(data = df, aes(x = DiabetesPedigreeFunction)) +

  geom_histogram( color = "red", fill = "blue", alpha = 0.1) +

  geom_density()

h <- ggplot(data = df, aes(x = Age)) +

  geom_histogram( color = "red", fill = "blue", alpha = 0.1) +geom_density()

ggarrange(a, b, c, d,e,f,g, h + rremove("x.text"),

       labels = c("a", "b", "c", "d","e", "f", "g", "h"),

       ncol = 3, nrow = 3)
```
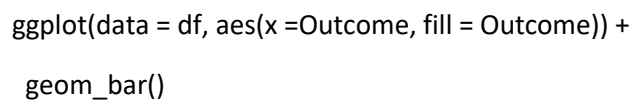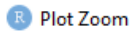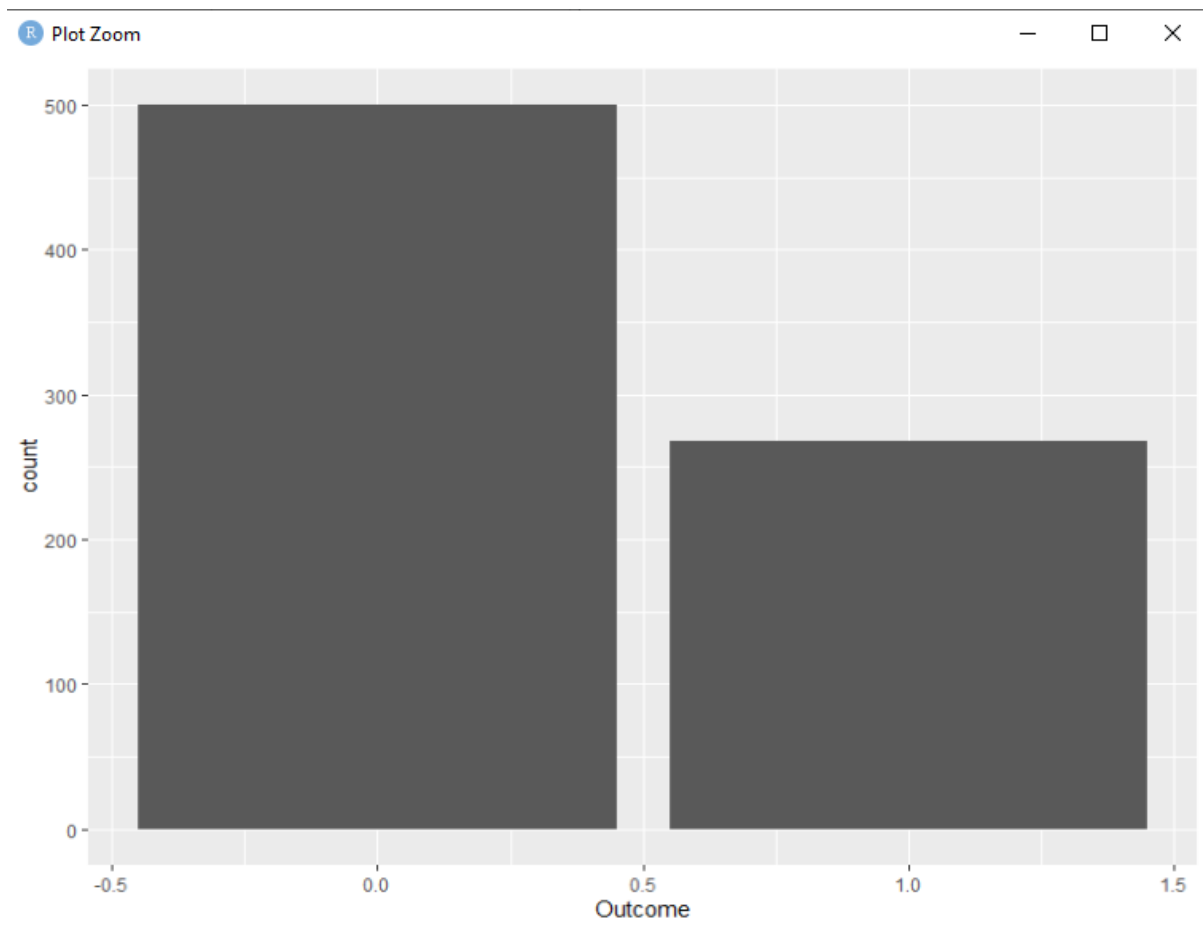
# Code

```
ggplot(data = df, aes(x =Outcome, fill = Outcome)) +
  geom_bar()
```

**# Code to label our categorical variable as a factor**

```
df$Outcome<- factor(df$Outcome,

        levels = c(0, 1),

        labels = c("Negative", "Positive"))

out <- subset(df,

      select = c(Pregnancies,Glucose,

            BloodPressure,SkinThickness,

            Insulin,BMI,

            DiabetesPedigreeFunction,Age))


# Code for boxplot

ggplot(data = melt(out),

    aes(x=variable, y=value)) +

  geom_boxplot(aes(fill=variable))
```

corPlot(df[, 1:8])

cutoff <- createDataPartition(df$Outcome, p=0.85, list=FALSE)

# select 15% of the data for validation

testdf <- df[-cutoff,]

# use the remaining 85% of data to training and testing the models

traindf <- df[cutoff,]

# Code to train the SVM

set.seed(1234)

# set the 10 fold crossvalidation with AU

# to pick for us what we call the best model

control <- trainControl(method="cv",number=10, classProbs = TRUE)

metric <- "Accuracy"

model <- train(Outcome ~., data = traindf, method = "svmRadial",

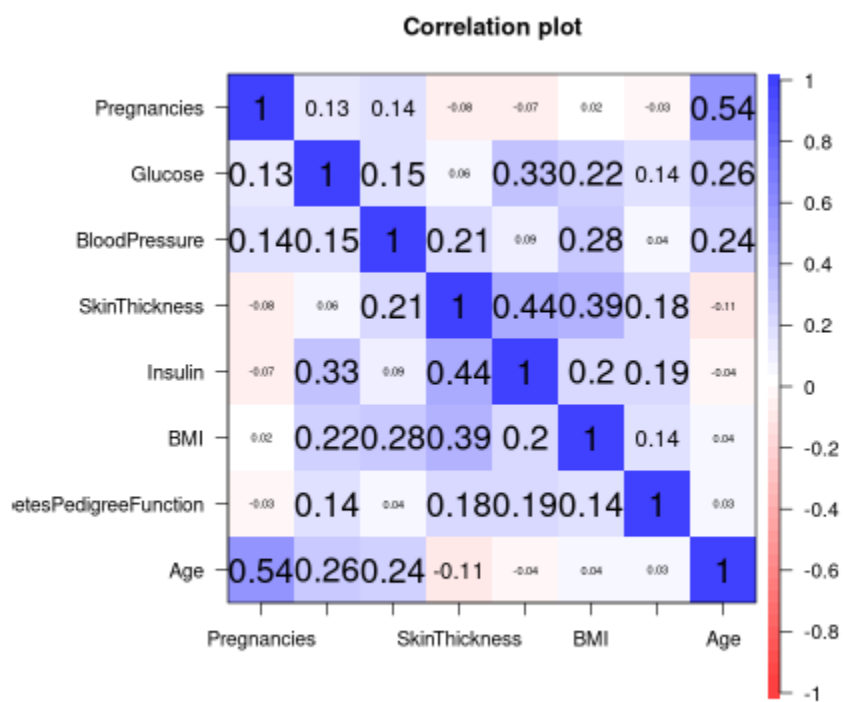        tuneLength = 8,preProc = c("center","scale"),

        metric=metric, trControl=control)

# Code for model summary

Model

```
Source

Console    Terminal ×    Background Jobs ×

R   R 4.4.1 · ~/
> set.seed(1234)
> # set the 10 fold crossvalidation with AU
> # to pick for us what we call the best model
> control <- trainControl(method="cv",number=10, classProbs = TRUE)
> metric <- "Accuracy"
> model <- train(Outcome ~., data = traindf, method = "svmRadial",
+                tuneLength = 8,preProc = c("center","scale"),
+                metric=metric, trControl=control)
> # Code for model summary
> model
Support Vector Machines with Radial Basis Function Kernel

653 samples
  8 predictor
  2 classes: 'Negative', 'Positive'

Pre-processing: centered (8), scaled (8)
Resampling: Cross-Validated (10 fold)
Summary of sample sizes: 587, 588, 587, 588, 587, 589, ...
Resampling results across tuning parameters:

  C       Accuracy   Kappa
   0.25   0.7578919  0.4410700
   0.50   0.7532044  0.4289572
   1.00   0.7594063  0.4383532
   2.00   0.7564001  0.4380336
   4.00   0.7456294  0.4150171
   8.00   0.7273317  0.3654897
  16.00   0.7227156  0.3481216
  32.00   0.7149752  0.3225678

Tuning parameter 'sigma' was held constant at a value of 0.1347287
Accuracy was used to select the optimal model using the largest value.
The final values used for the model were sigma = 0.1347287 and C = 1.
> |
```

# Code

plot(model)

# Practical 4

**Aim: - Implement of REGRESSION MODLE.**

```
# Generate random IQ values with mean = 30 and sd =2

IQ <- rnorm(40, 30, 2)

# Sorting IQ level in ascending order

IQ <- sort(IQ)

# Generate vector with pass and fail values of 40 students

result <- c(0, 0, 0, 1, 0, 0, 0, 0, 0, 1,
        1, 0, 0, 0, 1, 1, 0, 0, 1, 0,
        0, 0, 1, 0, 0, 1, 1, 0, 1, 1,
        1, 1, 1, 0, 1, 1, 1, 1, 0, 1)

# Data Frame

df <- as.data.frame(cbind(IQ, result))

# Print data frame

print(df)
```

```
> df <- as.data.frame(cbind(IQ, result))
>
> # Print data frame
> print(df)
          IQ result
1   26.06677      0
2   26.62661      0
3   27.46988      0
4   27.72373      1
5   27.86435      0
6   27.94799      0
7   28.54222      0
8   28.62629      0
9   28.74992      0
10  28.87905      1
11  28.88832      1
12  29.05442      0
13  29.10868      0
14  29.23906      0
15  29.38807      1
16  29.40986      1
17  29.53965      0
18  29.56405      0
19  29.87618      1
20  30.14102      0
21  30.22137      0
22  30.25858      0
23  30.30675      1
24  30.71963      0
25  30.80154      0
26  30.85293      1
27  30.92183      1
28  30.99570      0
29  31.10784      1
30  31.37728      1
```

# Plotting IQ on x-axis and result on y-axis

plot(IQ, result, xlab = "IQ Level",

    ylab = "Probability of Passing")

# Create a logistic model
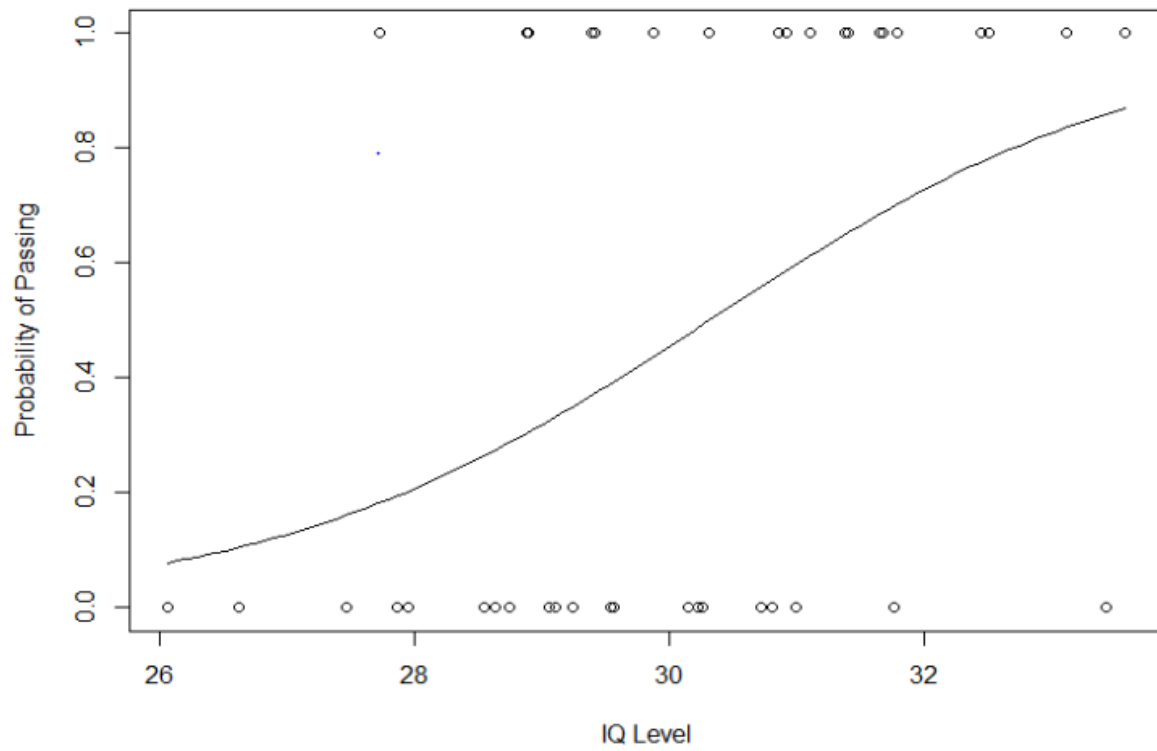
g = glm(result~IQ, family=binomial, df)

# Create a curve based on prediction using the regression model

curve(predict(g, data.frame(IQ=x), type="resp"), add=TRUE)

# Based on fit to the regression model

points(IQ, fitted(g), pch=30)

# Summary of the regression model

summary(g)

```
> # Summary of the regression model
> summary(g)

Call:
glm(formula = result ~ IQ, family = binomial, data = df)

Coefficients:
            Estimate Std. Error z value Pr(>|z|)
(Intercept) -17.6032     7.0531  -2.496   0.0126 *
IQ            0.5808     0.2335   2.488   0.0129 *
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

    Null deviance: 55.352  on 39  degrees of freedom
Residual deviance: 47.331  on 38  degrees of freedom
AIC: 51.331

Number of Fisher Scoring iterations: 3
```

# Practical 5

**Aim :- Implement of Simple Linear Regaression.**

years_of_exp = c(7,5,1,3)

salary_in_lakhs = c(21,13,6,8)

#employee.data = data.frame(satisfaction_score, years_of_exp, salary_in_lakhs)

employee.data = data.frame(years_of_exp, salary_in_lakhs)

employee.data

# Estimation of the salary of an employee, based on his year of experience and satisfaction score in his company.

model <- lm(salary_in_lakhs ~ years_of_exp, data = employee.data)

summary(model)

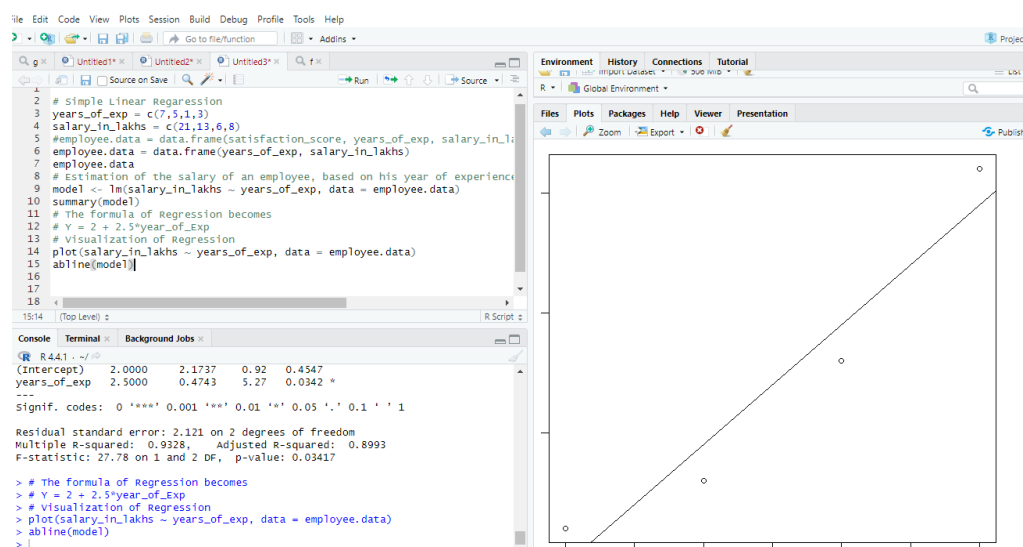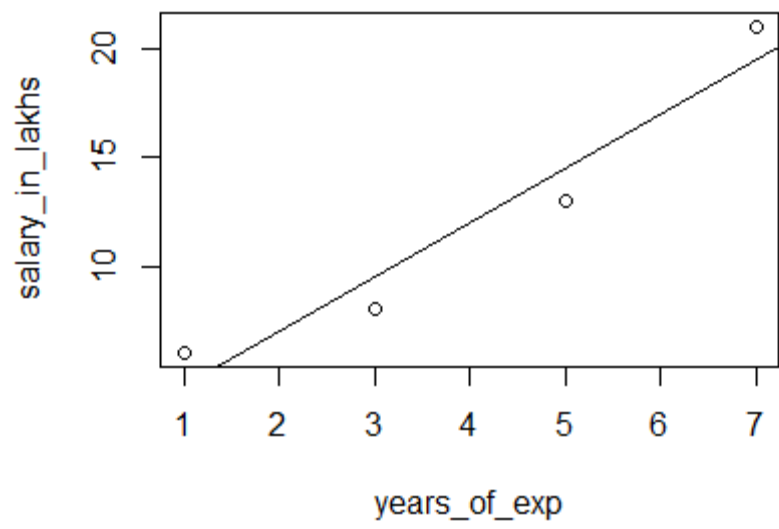# The formula of Regression becomes

# Y = 2 + 2.5*year_of_Exp

# Visualization of Regression

plot(salary_in_lakhs ~ years_of_exp, data = employee.data)

abline(model)

# Practical 6

**Aim :- Implement of Multiple Linear Regression.**

```
# Importing the dataset

dataset = read.csv('D:\\data2.csv')

# Encoding categorical data

dataset$State = factor(dataset$State,

            levels = c('New York', 'California', 'Florida'),

            labels = c(1, 2, 3))

dataset$State

# Splitting the dataset into the Training set and Test set

install.packages('caTools')

library(caTools)

set.seed(123)

split = sample.split(dataset$Profit, SplitRatio = 0.8)

training_set = subset(dataset, split == TRUE)

test_set = subset(dataset, split == FALSE)


# Feature Scaling

# training_set = scale(training_set)

# test_set = scale(test_set)


# Fitting Multiple Linear Regression to the Training set

regressor = lm(formula = Profit ~ .,

        data = training_set)


# Predicting the Test set results

y_pred = predict(regressor, newdata = test_set)
```

```
> regressor

Call:
lm(formula = Profit ~ ., data = training_set)

Coefficients:
    (Intercept)         R.D.Spend     Administration   Marketing.Spend
      2.816e+04         8.884e-01          5.670e-02         2.859e-02
         State2            State3
     -2.861e+03         9.172e+03
```

# Practical 7

**Aim :- Implement of Logistic regression.**

Source code:

install.packages("ISLR")

library(ISLR)

#load dataset

data <- ISLR::Default

print (head(ISLR::Default))

#view summary of dataset

summary(data)

#find total observations in dataset

nrow(data)

#Create Training and Test Samples

#split the dataset into a training set to train the model on and a testing set to test the model

set.seed(1)

#Use 70% of dataset as training set and remaining 30% as testing set

sample <- sample(c(TRUE, FALSE), nrow(data), replace=TRUE, prob=c(0.7,0.3))

print (sample)

train <- data[sample, ]

test <- data[!sample, ]

nrow(train)

nrow(test)

# Fit the Logistic Regression Model

# use the glm (general linear model) function and specify family="binomial"

#so that R fits a logistic regression model to the dataset

model <- glm(default~student+balance+income, family="binomial", data=train)
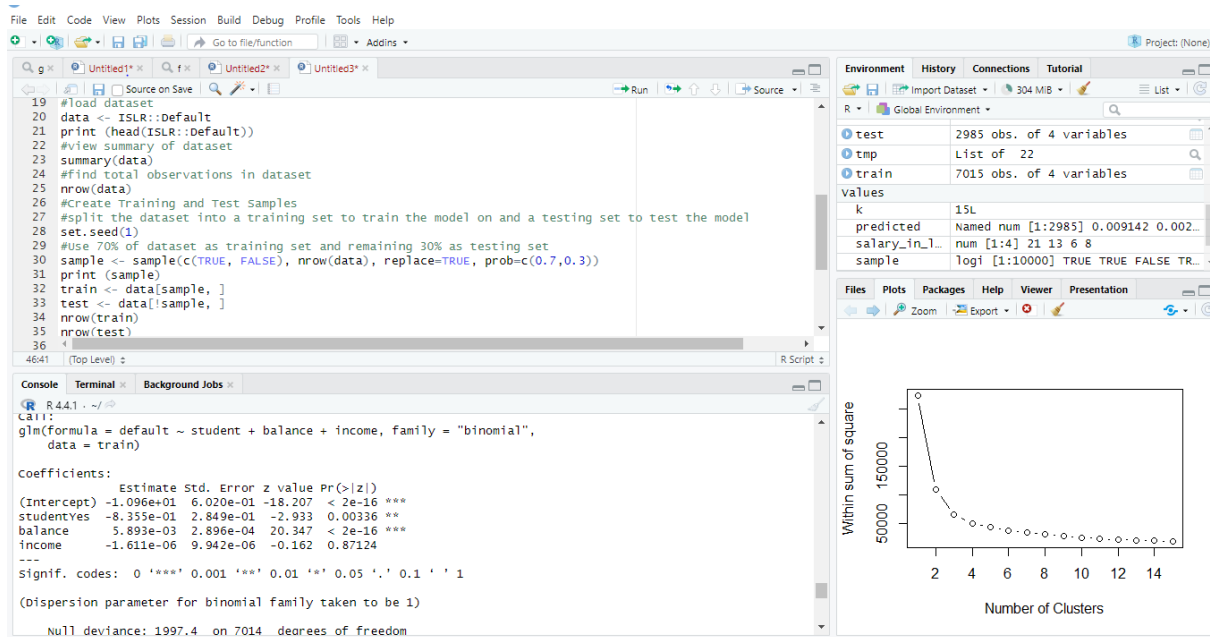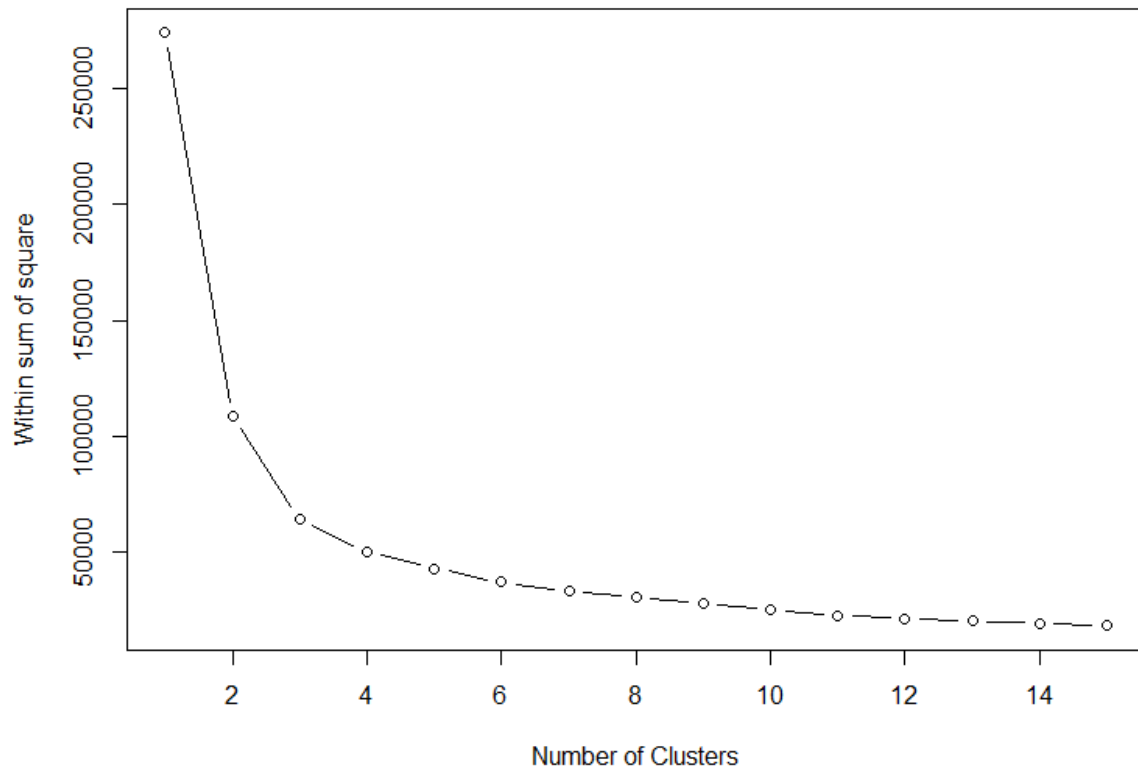
#view model summary

summary(model)

#Model Diagnostics

install.packages("InformationValue")

library(InformationValue)

predicted <- predict(model, test, type="response")

confusionMatrix(test$default, predicted)

# Practical 8

**Aim: Read a datafile grades_km_input.csv and apply k-means clustering.**

Datafile:

# install required packages

install.packages("plyr")

install.packages("ggplot2")

install.packages("cluster")

install.packages("lattice")

install.packages("grid")

install.packages("gridExtra")

# Load the package

library(plyr)

library(ggplot2)

library(cluster)

library(lattice)

library(grid)

library(gridExtra)

# A data frame is a two-dimensional array-like structure in which each column contains values of one variable and each row contains one set of values from each column.

grade_input=as.data.frame(read.csv("D:\\grades_km_input.csv"))

kmdata_orig=as.matrix(grade_input[, c ("Student","English","Math","Science")])

kmdata=kmdata_orig[,2:4]

kmdata[1:10,]

# the k-means algorithm is used to identify clusters for k = 1, 2, .. , 15. For each value of k, the WSS is calculated.

wss=numeric(15)

# the option n start=25 specifies that the k-means algorithm will be repeated 25 times, each starting with k random initial centroids

for(k in 1:15)wss[k]=sum(kmeans(kmdata,centers=k,nstart=25)$withinss)

plot(1:15,wss,type="b",xlab="Number of Clusters",ylab="Within sum of square")

#As can be seen, the WSS is greatly reduced when k increases from one to two. Another substantial reduction in WSS occurs at k = 3. However, the improvement in WSS is fairly linear fork > 3.

```
km = kmeans(kmdata,3,nstart=25)

km

c( wss[3] , sum(km$withinss))

df=as.data.frame(kmdata_orig[,2:4])

df$cluster=factor(km$cluster)

centers=as.data.frame(km$centers)

g1=ggplot(data=df, aes(x=English, y=Math, color=cluster )) +

  geom_point() + theme(legend.position="right") +

  geom_point(data=centers,aes(x=English,y=Math, color=as.factor(c(1,2,3))),size=10, alpha=.3, show.legend =FALSE)

g2=ggplot(data=df, aes(x=English, y=Science, color=cluster )) +

  geom_point () +geom_point(data=centers,aes(x=English,y=Science, color=as.factor(c(1,2,3))),size=10, alpha=.3, show.legend=FALSE)

g3 = ggplot(data=df, aes(x=Math, y=Science, color=cluster )) +

  geom_point () + geom_point(data=centers,aes(x=Math,y=Science, color=as.factor(c(1,2,3))),size=10, alpha=.3, show.legend=FALSE)

tmp=ggplot_gtable(ggplot_build(g1))

grid.arrange(arrangeGrob(g1 + theme(legend.position="none"),g2 + theme(legend.position="none"),g3 + theme(legend.position="none"),top ="High School Student Cluster Analysis" ,ncol=1))
```
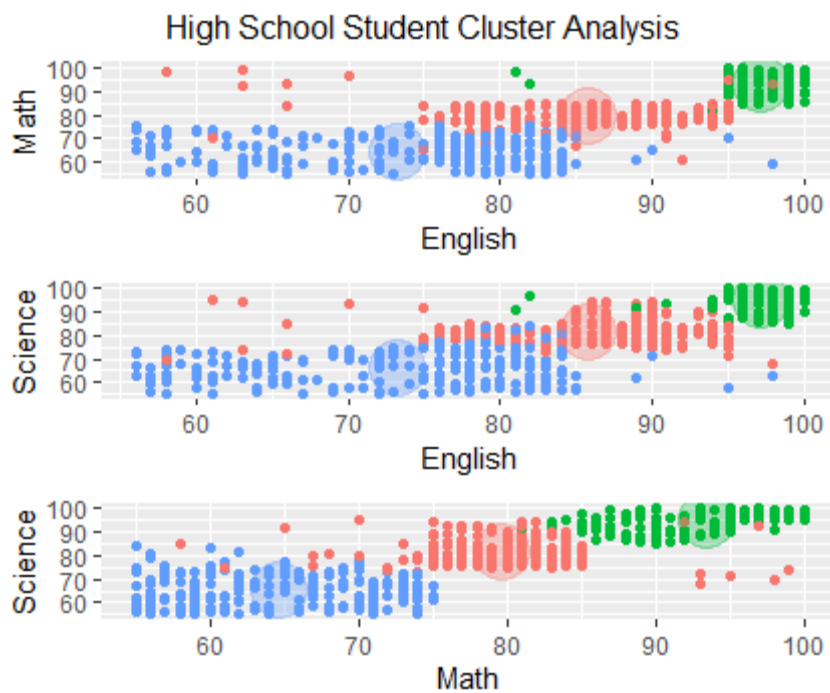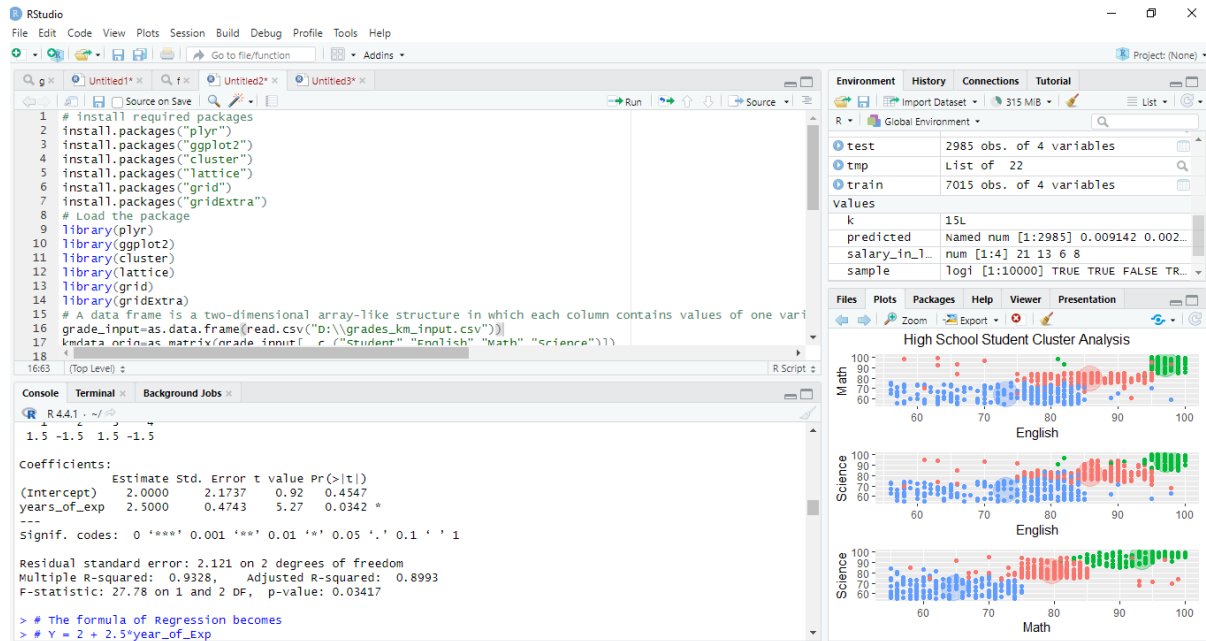
High School Student Cluster Analysis

# Practical 9

**Aim: Perform Apriori algorithm using Groceries dataset from the R arules package.**

install.packages("arules")

install.packages("arulesViz")

install.packages("RColorBrewer")

# Loading Libraries

library(arules)

library(arulesViz)

library(RColorBrewer)

# import dataset

data(Groceries)

Groceries

summary(Groceries)

class(Groceries)

# using apriori() function

rules = apriori(Groceries, parameter = list(supp = 0.02, conf = 0.2))

summary (rules)

# using inspect() function

inspect(rules[1:10])

# using itemFrequencyPlot() function

arules::itemFrequencyPlot(Groceries, topN = 20,

col = brewer.pal(8, 'Pastel2'),

main = 'Relative Item Frequency Plot',

type = "relative",

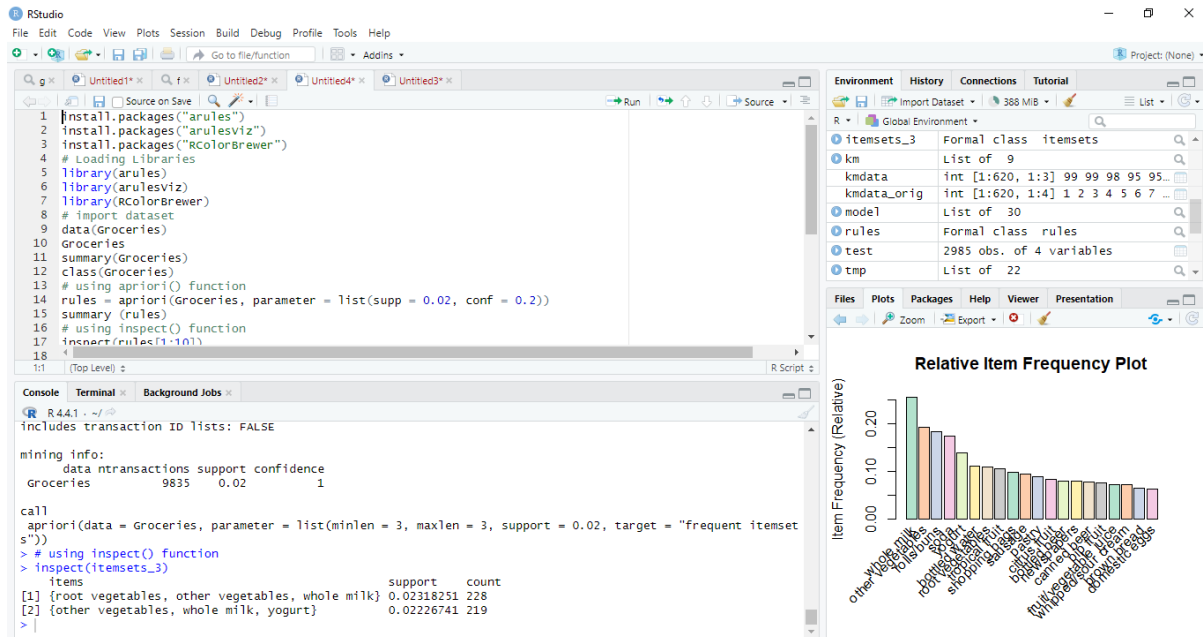ylab = "Item Frequency (Relative)")

itemsets = apriori(Groceries, parameter = list(minlen=2, maxlen=2,support=0.02, target="frequent itemsets"))

summary(itemsets)

# using inspect() function

inspect(itemsets[1:10])

itemsets_3 = apriori(Groceries, parameter = list(minlen=3, maxlen=3,support=0.02, target="frequent itemsets"))

summary(itemsets_3)

# using inspect() function

inspect(itemsets_3)

Relative Item Frequency Plot