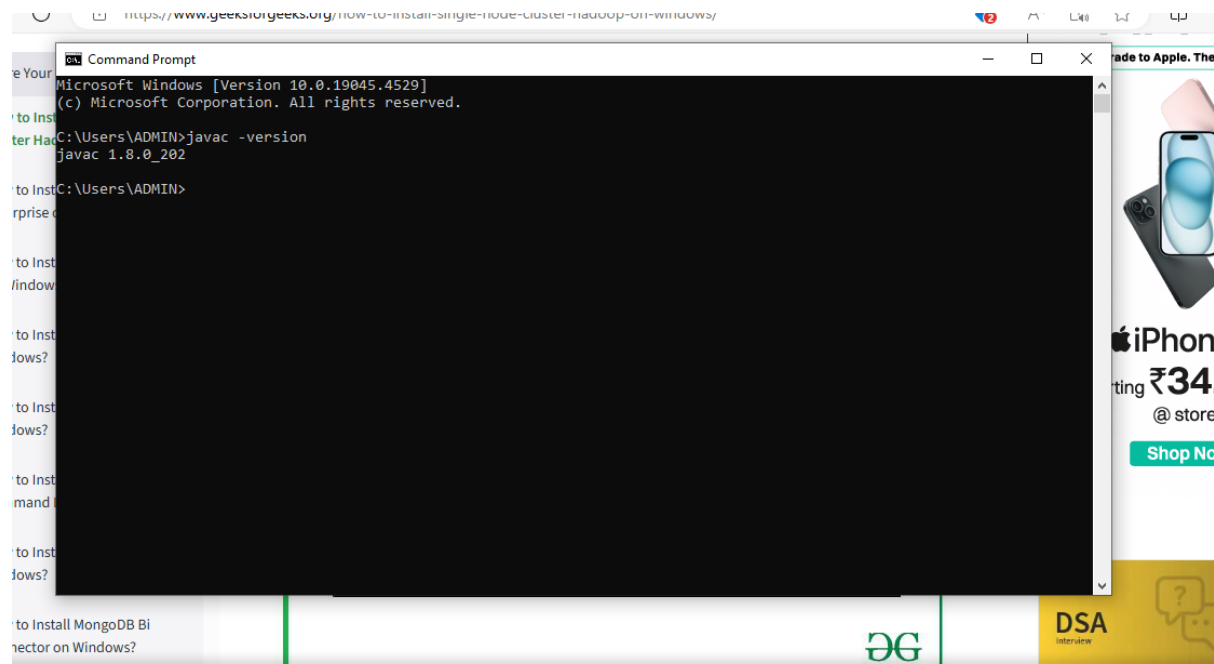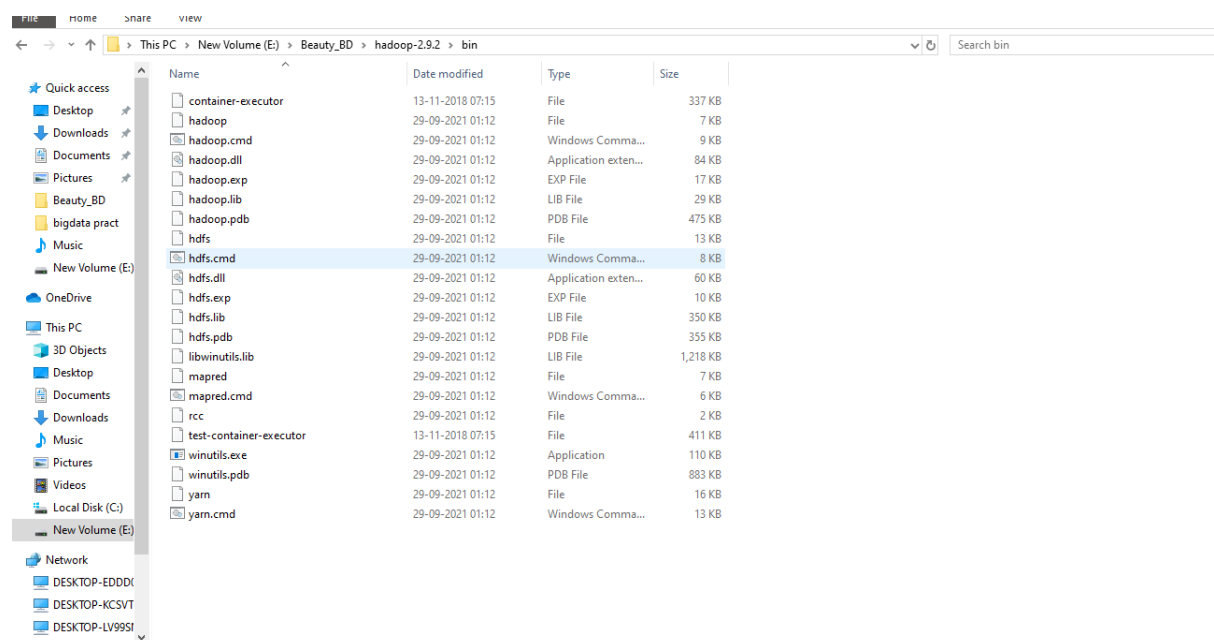## Practical.1

Step 1 :- **Verify the Java installed**
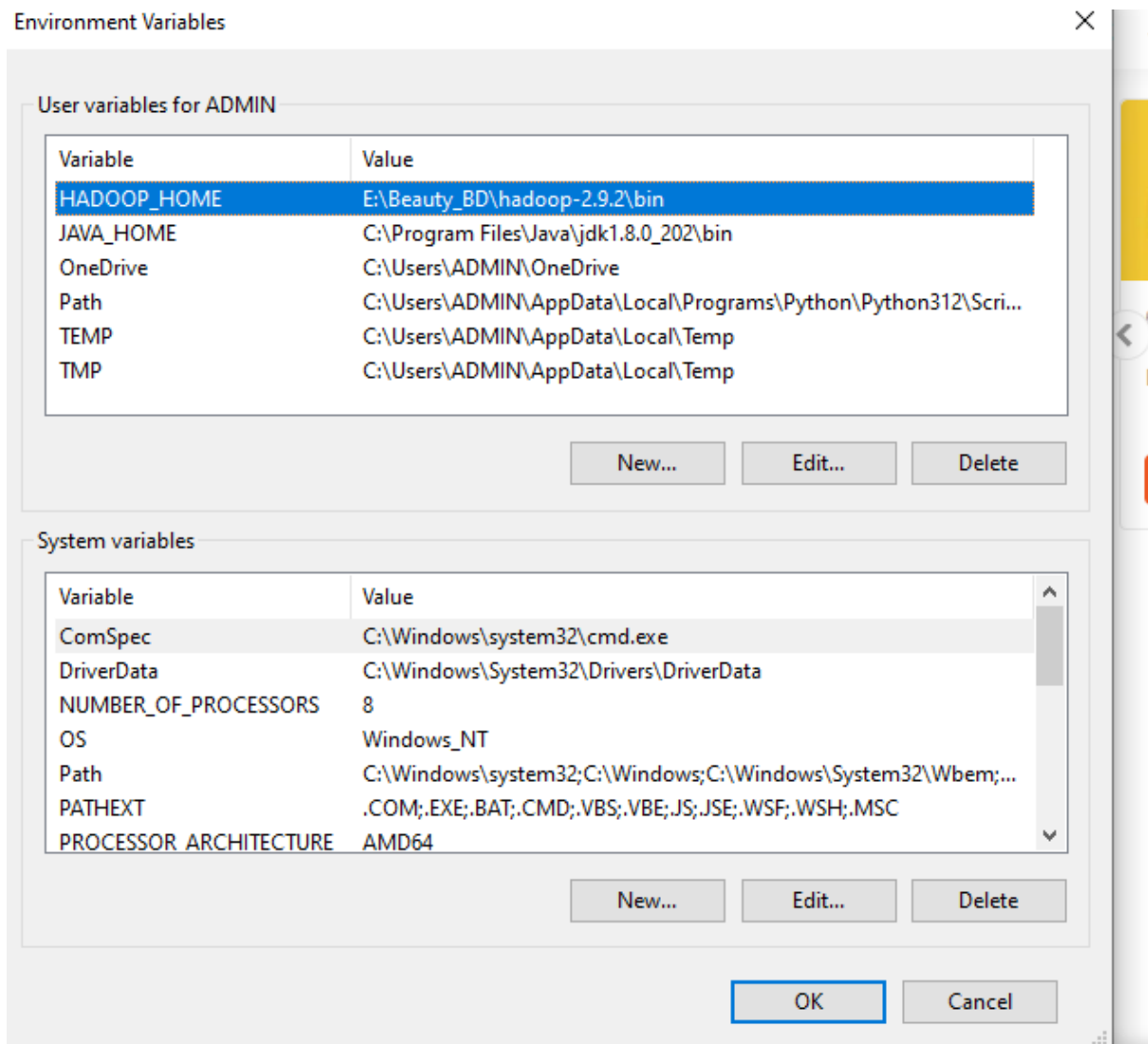


## Step 2: Extract Hadoop at C:\Hadoop



## Step 3: Setting up the HADOOP_HOME variable And

## Step 4: Set JAVA_HOME variable

## Step 5: Set Hadoop and Java bin directory path

## Step 6: Hadoop Configuration :

1. Core-site.xml

2. Mapred-site.xml

3. Hdfs-site.xml

4. Yarn-site.xml

5. Hadoop-env.cmd

6. Create two folders datanode and namenode

### Step 6.1: Core-site.xml configuration

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

### Step 6.2: Mapred-site.xml configuration

```xml
<configuration>

    <property>

        <name>mapreduce.framework.name</name>

        <value>yarn</value>

    </property>

</configuration>
```

### Step 6.3: Hdfs-site.xml configuration

```xml
<configuration>

    <property>

        <name>dfs.replication</name>

        <value>1</value>

    </property>

    <property>

        <name>dfs.namenode.name.dir</name>

        <value>C:\hadoop-2.8.0\data\namenode</value>

    </property>

    <property>

        <name>dfs.datanode.data.dir</name>

        <value>C:\hadoop-2.8.0\data\datanode</value>

    </property>

</configuration>
```

### Step 6.4: Yarn-site.xml configuration

```xml
<configuration>

    <property>

        <name>yarn.nodemanager.aux-services</name>

        <value>mapreduce_shuffle</value>

    </property>

    <property>


<name>yarn.nodemanager.auxservices.mapreduce.shuffle.class</name>

     <value>org.apache.hadoop.mapred.ShuffleHandler</value>

    </property>
```
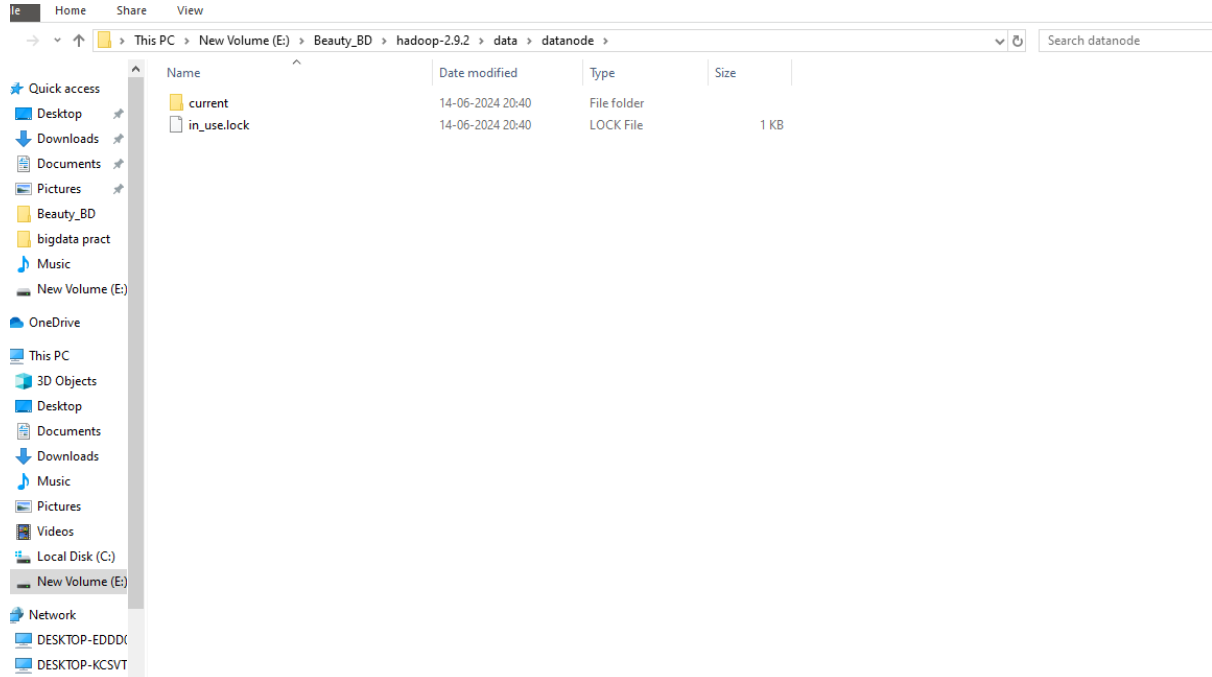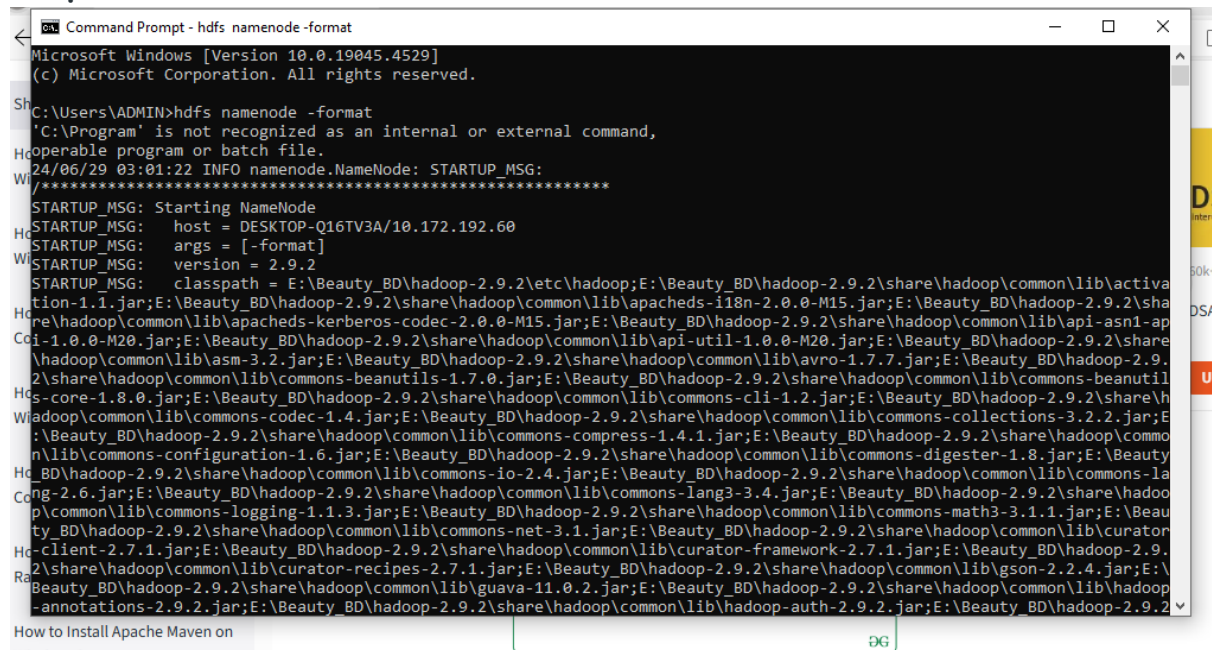
```
</configuration>
```

## Step 6.5: Hadoop-env.cmd configuration

```
Set "JAVA_HOME=C:\Java" (On C:\java this is path to file jdk.18.0)
```
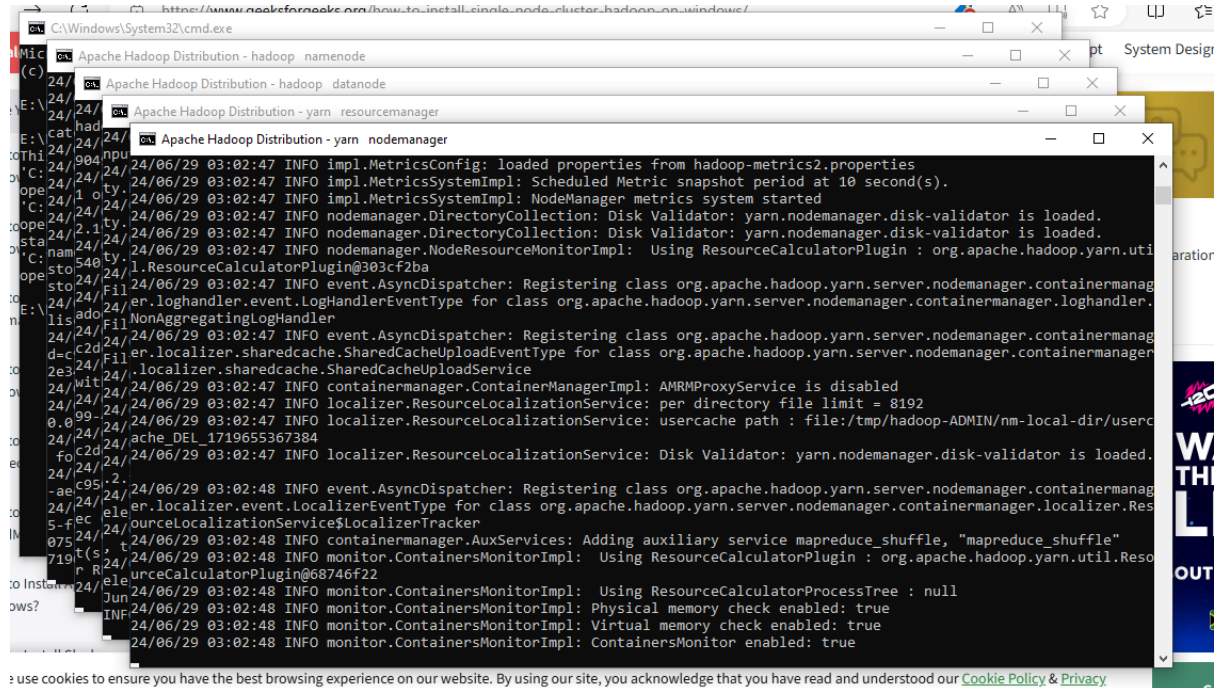
## Step 6.6: Create datanode and namenode folders
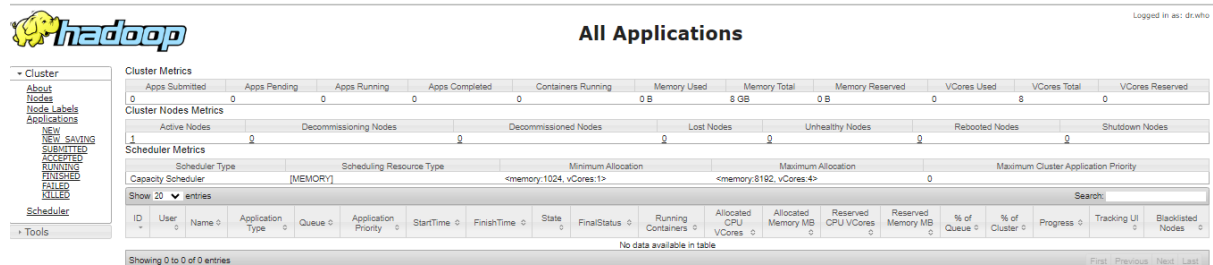


## Step 7: Format the namenode folder

# Step 8: Testing the setup



# Step 8.1: Testing the setup:

# Step 9: Open: http://localhost:8088

## Step 10: http://localhost:50070

| Hadoop | Overview | Datanodes | Datanode Volume Failures | Snapshot | Startup Progress | Utilities |

## Overview 'localhost:9000' (active)

| | |
|---|---|
| **Started:** | Sat Jun 29 03:02:46 -0700 2024 |
| **Version:** | 2.9.2, r826afbeae31ca687bc2f8471dc841b66ed2c6704 |
| **Compiled:** | Tue Nov 13 04:42:00 -0800 2018 by ajisaka from branch-2.9.2 |
| **Cluster ID:** | CID-f0ac2e34-c5ea-4b71-815c-46ab2d3b2027 |
| **Block Pool ID:** | BP-366826651-10.172.192.60-1719049540241 |

## Summary

Security is off.

Safemode is off.

1 files and directories, 0 blocks = 1 total filesystem object(s).

Heap Memory used 96.34 MB of 182.5 MB Heap Memory. Max Heap Memory is 889 MB.

Non Heap Memory used 41.84 MB of 43.38 MB Commited Non Heap Memory. Max Non Heap Memory is <unbounded>.

## Practical 2:

**Aim**: Classification using SVM

**Requirement:**
R tool

**Code:**
```
getwd()

read.csv()

ds=read.csv("E:/Rajdeep/bigdata pract/dataset/social.csv",TRUE,",")

ds

ds=ds[3:5]

ds

install("catools")

library(caTools)

set.seed(123)

split=sample.split(ds$Purchased, SplitRatio=0.75)

training_set=(subset(ds, split == TRUE))

test_set =(subset(ds, split == FALSE))

ds

test_set[-3]=scale(test_set[-3])

training_set[-3]=scale(training_set[-3])

test_set[-3]

training_set[-3]

install.packages('e1071')

library('e1071')

classifier=svm(formula=Purchased ~ ., data= training_set , type='C-classification',kernal='linear')

classifier

y_pred=predict(classifier, newdata=test_set[-3])

y_pred

cm=table(test_set[, 3],y_pred)

cm
```
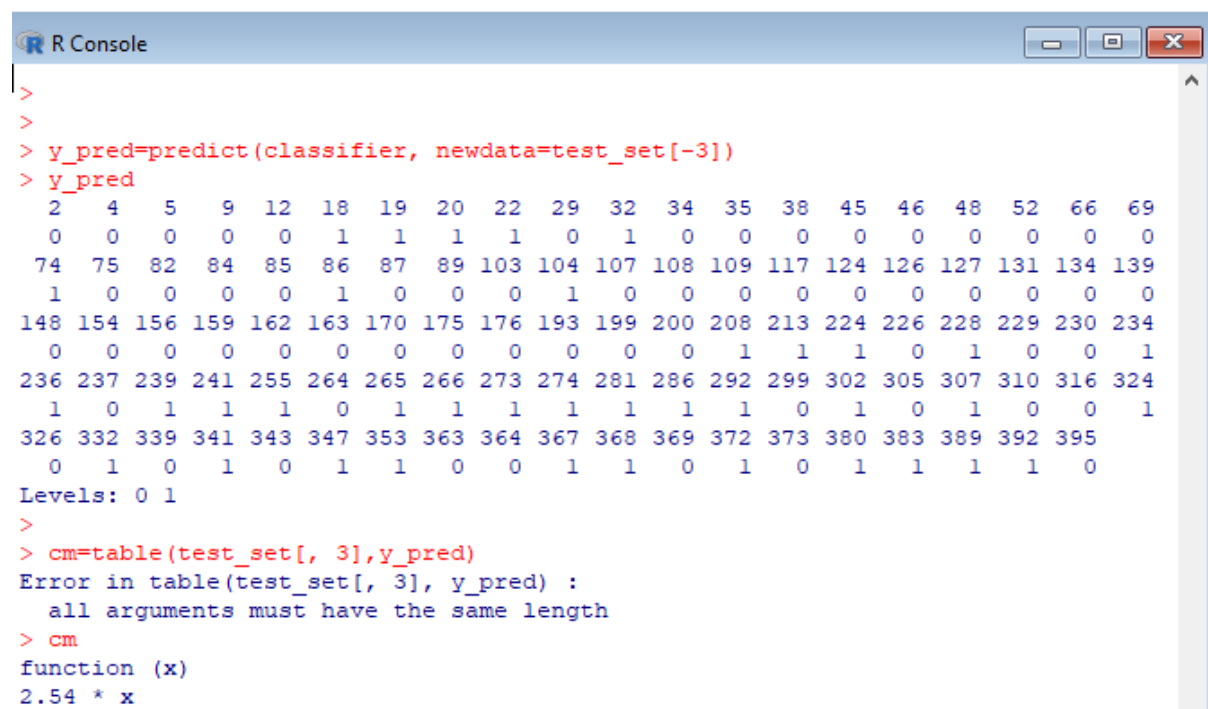
```
> set = training_set
> X1 = seq(min(set[, 1]) - 1, max(set[, 1]) + 1, by = 0.01)
> X2 = seq(min(set[, 2]) - 1, max(set[, 2]) + 1, by = 0.01)

> grid_set = expand.grid(X1, X2)
> colnames(grid_set) = c('Age', 'EstimatedSalary')
> y_grid = predict(classifier, newdata = grid_set)
> plot(set[, -3],
+       main = 'SVM (Training set)',
+       xlab = 'Age', ylab = 'Estimated Salary',
+       xlim = range(X1), ylim = range(X2))

> contour(X1, X2, matrix(as.numeric(y_grid), length(X1), length(X2)), add = TRUE)
> points(grid_set, pch = '.', col = ifelse(y_grid == 1, 'coral1', 'aquamarine'))
> points(set, pch = 21, bg = ifelse(set[, 3] == 1, 'green4', 'red3'))
```

**Output:**

```
R R Console                                                    [_][□][X]

>
>
> y_pred=predict(classifier, newdata=test_set[-3])
> y_pred
   2    4    5    9   12   18   19   20   22   29   32   34   35   38   45   46   48   52   66   69
   0    0    0    0    0    1    1    1    1    0    1    0    0    0    0    0    0    0    0    0
  74   75   82   84   85   86   87   89  103  104  107  108  109  117  124  126  127  131  134  139
   1    0    0    0    0    1    0    0    0    1    0    0    0    0    0    0    0    0    0    0
 148  154  156  159  162  163  170  175  176  193  199  200  208  213  224  226  228  229  230  234
   0    0    0    0    0    0    0    0    0    0    0    0    1    1    1    0    1    0    0    1
 236  237  239  241  255  264  265  266  273  274  281  286  292  299  302  305  307  310  316  324
   1    0    1    1    1    0    1    1    1    1    1    1    1    1    0    1    0    1    0    0    1
 326  332  339  341  343  347  353  363  364  367  368  369  372  373  380  383  389  392  395
   0    1    0    1    0    1    1    0    0    1    1    0    1    0    1    0    1    1    1    1    0
Levels: 0 1
>
> cm=table(test_set[, 3],y_pred)
Error in table(test_set[, 3], y_pred) :
  all arguments must have the same length
> cm
function (x)
2.54 * x
```

## Practical 3:

**Aim**: write program in R of Naive baye's theorem

**Requirement:**
R tool

**Code:**

```
data(iris)

str(iris)


install packages("e1071")

install packages("caTools")

install packages("caret")


library(e1071)

library(caTools)

library(caret)


split <- sample.split(iris,SplitRatio=0.7)

train_c1 <-subset(iris,split=="TRUE")

test_c1 <- subset(iris,split == "FALSE")

train_scale <- scale(train_c1[, 1:4])

test_scale <- scale(test_c1[,1:4])


set.seed(120)

classifier_c1 <- naiveBayes(Species ~ ., data = train_c1)

classifier_c1


y_pred <- predict(classifier_c1, newdata= test_c1)

cm <- table(test_c1$Species, y_pred)

cm


confusionMatrix(cm)
```

**Output:**

```
 virginica        0           1          19

Overall Statistics

              Accuracy : 0.9333
                95% CI : (0.838, 0.9815)
   No Information Rate : 0.3667
   P-Value [Acc > NIR] : < 2.2e-16

                 Kappa : 0.9

 Mcnemar's Test P-Value : NA

Statistics by Class:

                    Class: setosa Class: versicolor Class: virginica
Sensitivity                1.0000            0.9444           0.8636
Specificity                1.0000            0.9286           0.9737
Pos Pred Value             1.0000            0.8500           0.9500
Neg Pred Value             1.0000            0.9750           0.9250
Prevalence                 0.3333            0.3000           0.3667
Detection Rate             0.3333            0.2833           0.3167
Detection Prevalence       0.3333            0.3333           0.3333
Balanced Accuracy          1.0000            0.9365           0.9187
>
```

## Install python package:

1. You will need to make the hidden folder visible: go to "C:" drive on top click on tab "view"
2. Select "hidden Items" option:



3. Go to the below path:
   C:\Users\*Your Name*\AppData\Local\Programs\Python\Python36-32\Scripts
4. Set the below path in command prompt and then use the below command:
   python -m pip install pymongo

## Practical :4

**Aim:** Implement an application that stores big data in Hbase / MongoDB and manipulate it using R / Python

**Requirement:**
a. Python Package: PyMongo
b. Mongo Database

**Step A: Install Mongo database**
**Step 1)** Go to (https://www.mongodb.com/download-center/community)  and Download MongoDB Community Server. We will install the 64-bit version for Windows.



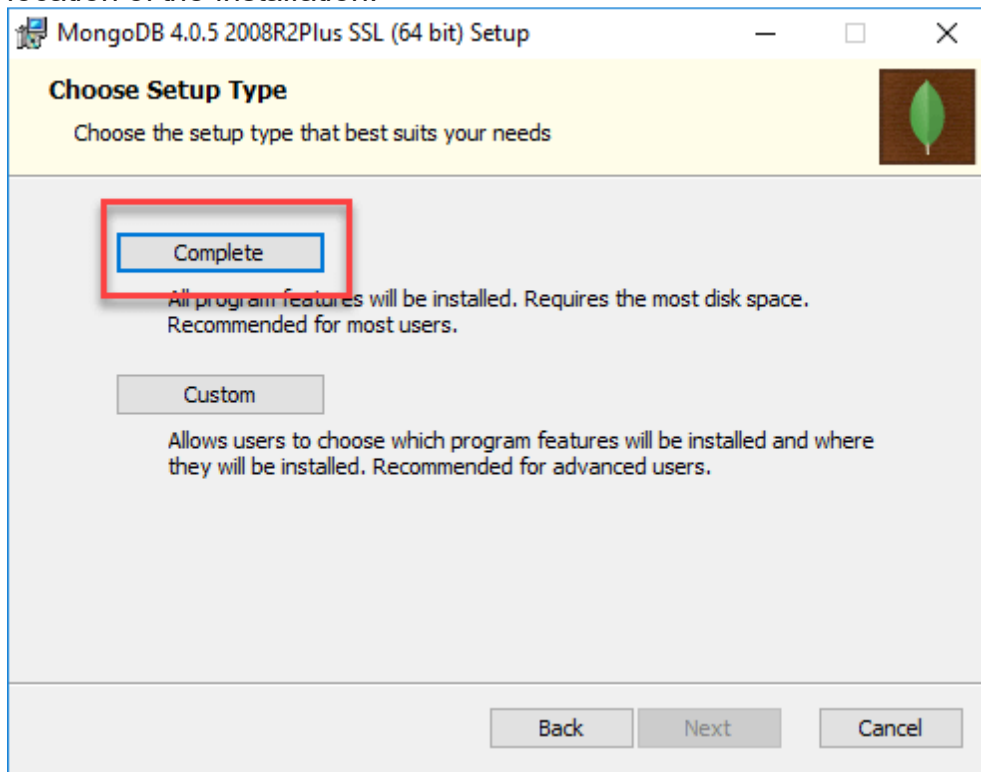**Step 2)** Once download is complete open the msi file. Click Next in the start up screen

**Step 3)**
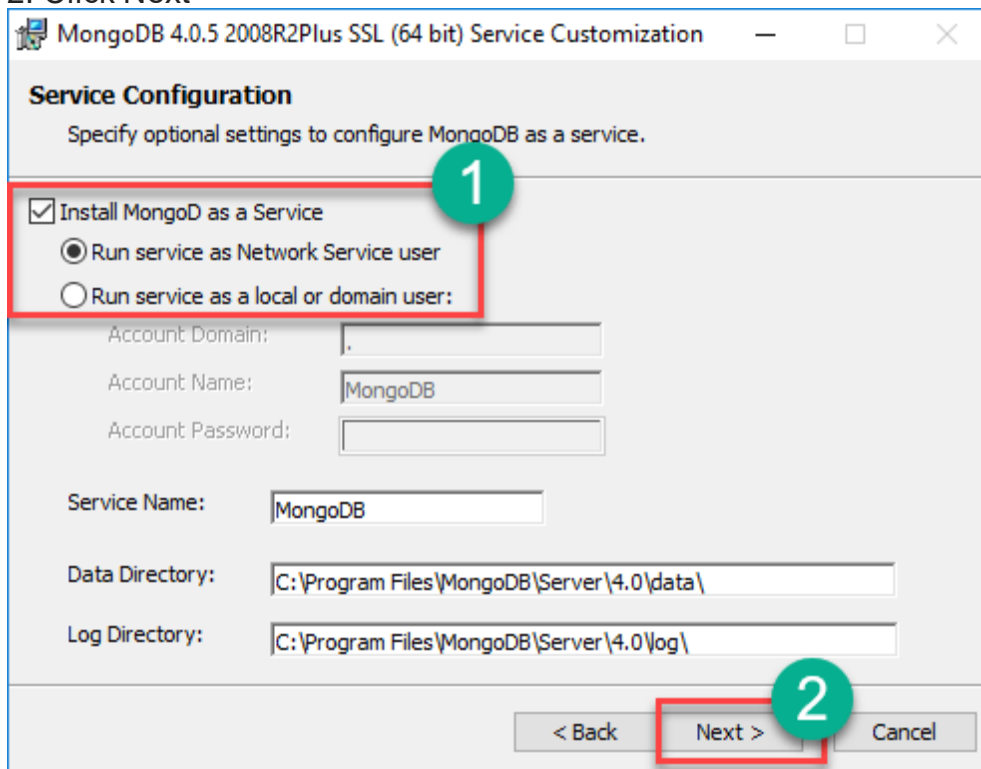1. Accept the End-User License Agreement
2. Click Next

**Step 4)** Click on the "complete" button to install all of the components. The custom option can be used to install selective components or if you want to change the location of the installation.
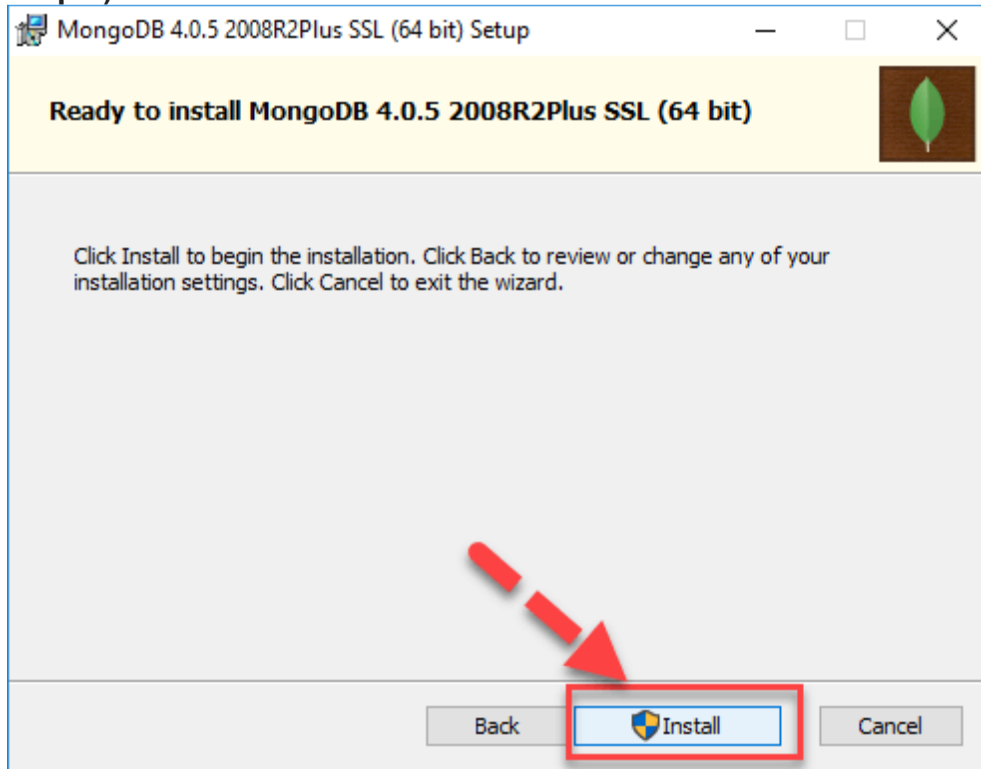


**Step 5)**
1. Select "Run service as Network Service user". make a note of the data directory, we'll need this later.
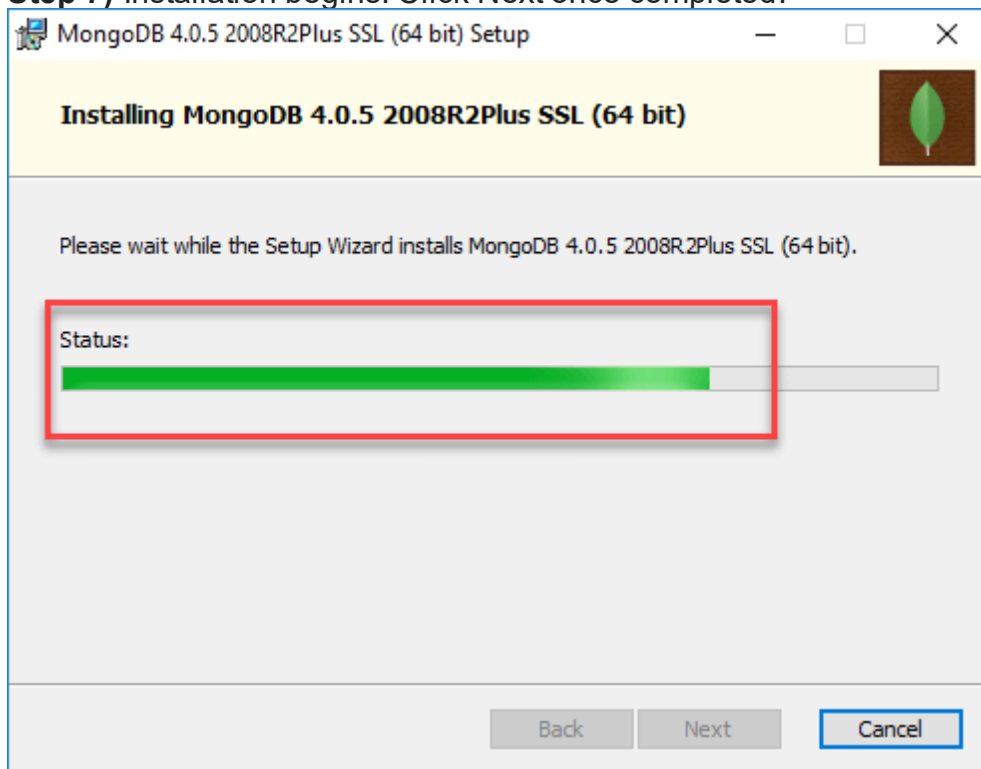2. Click Next

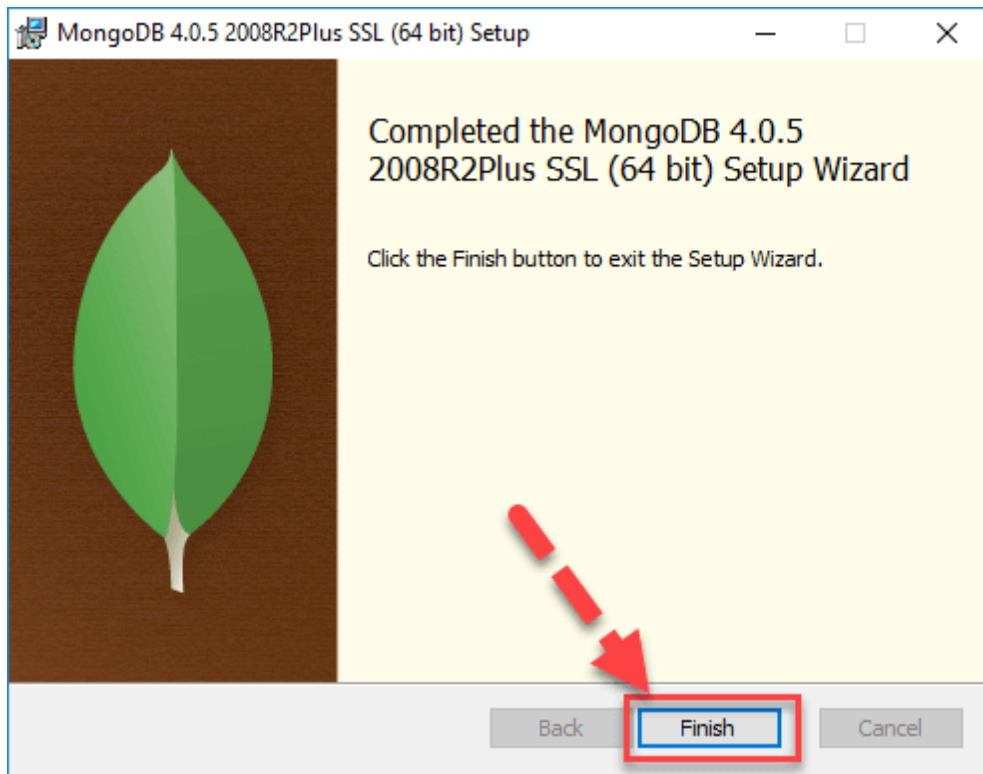**Step 6)** Click on the Install button to start the installation.



**Step 7)** Installation begins. Click Next once completed.



**Step 8) Click** on the Finish button to complete the installation

**Program 1:** Displaying the database name:

```
import pymongo
myclient = pymongo.MongoClient("mongodb://localhost:27017/")
mydb = myclient["mybigdata"]
print(myclient.list_database_names())
```
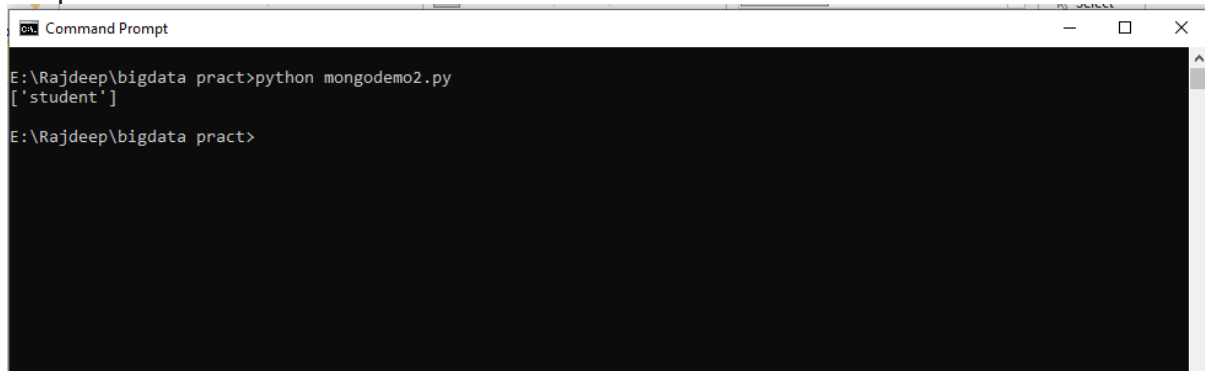
Output:



**Program 2:** Creating collection:

```
import pymongo
myclient = pymongo.MongoClient("mongodb://localhost:27017/")
mydb = myclient["mybigdata"]
mycol=mydb["student"]
print(mydb.list_collection_names())
```
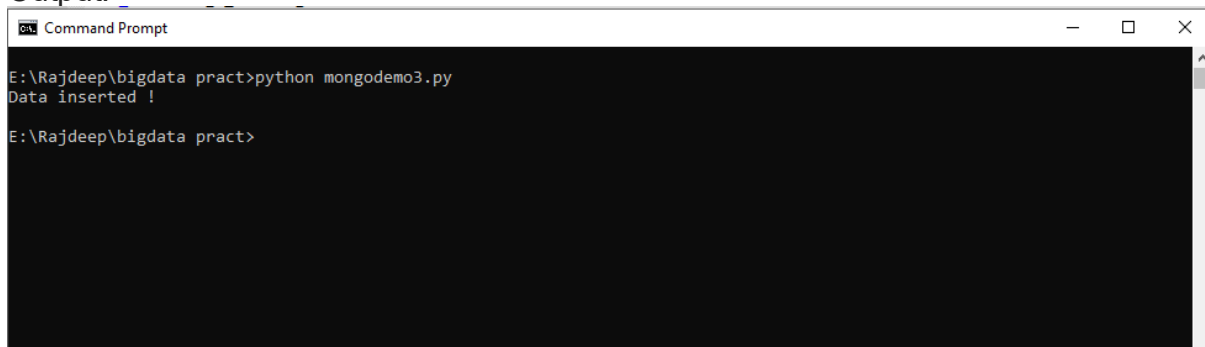
Output:



**Program 3**: Inserting Data

```
import pymongo
myclient = pymongo.MongoClient("mongodb://localhost:27017/")
mydb = myclient["mybigdata"]
mycol=mydb["student"]
mydict={"name":"vai", "address":"bhy"}
x=mycol.insert_one(mydict)
print("Data inserted !")
```
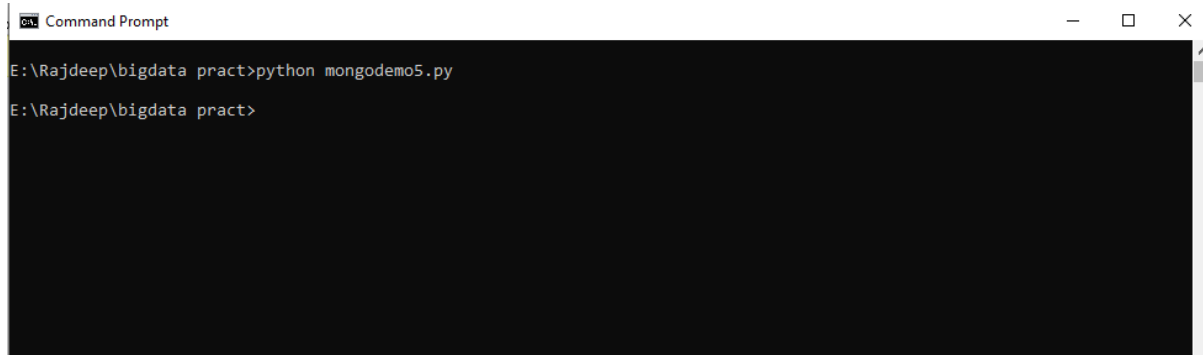
Output:



**Program 4**: Insert Multiple data into Collection

```
import pymongo
myclient = pymongo.MongoClient("mongodb://localhost:27017/")
mydb = myclient["mybigdata"]
mycol=mydb["student"]
mylist=[{"name":"Ganesh", "address":"Mumbai"}, {"name":"Varun",
"address":"Mumbai"},
{"name":"Prasoon", "address":"Pune"}, {"name":"Satish", "address":"Pune"},]
x=mycol.insert_many(mylist)
print("Data inserted !")
```

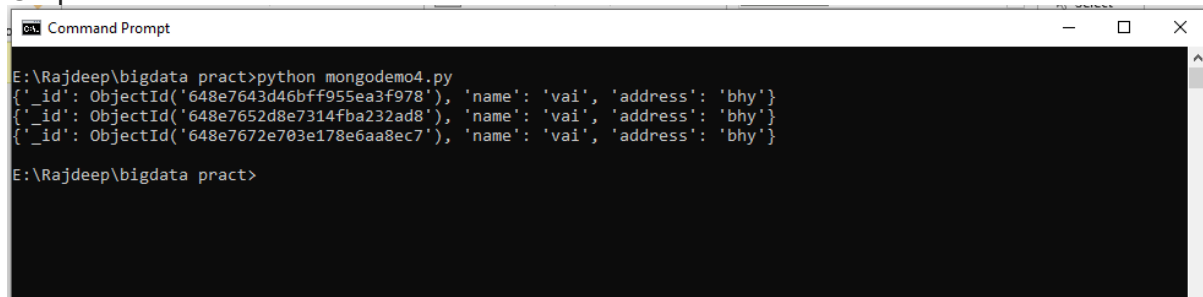Output:

**Program 5:** Displaying the collection data:

```
import pymongo
myclient = pymongo.MongoClient("mongodb://localhost:27017/")
mydb = myclient["mybigdata"]
mycol = mydb["student"]

myquery = { "name": "Vai" }

mydoc = mycol.find(myquery)

for x in mydoc:
  print(x)
```

Output:

**Practical 5:**

K means clustering.

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

**Requirement:**

R tool

**Code:**

```
install.packages("plyr")
install.packages("ggplot2")
install.packages("cluster")
install.packages("lattice")
install.packages("grid")
install.packages("gridExtra")

library(plyr)
library(ggplot2)
library(cluster)
library(lattice)
library(grid)
library(gridExtra)

grade_input=as.data.frame(read.csv("E:/Rajdeep/bigdata
pract/dataset/grades_km_input.csv"))

kmdata_orig=as.matrix(grade_input[, c ("Student","English","Math","Science")])
kmdata=kmdata_orig[,2:4]
kmdata[1:10,]
wss=numeric(15)

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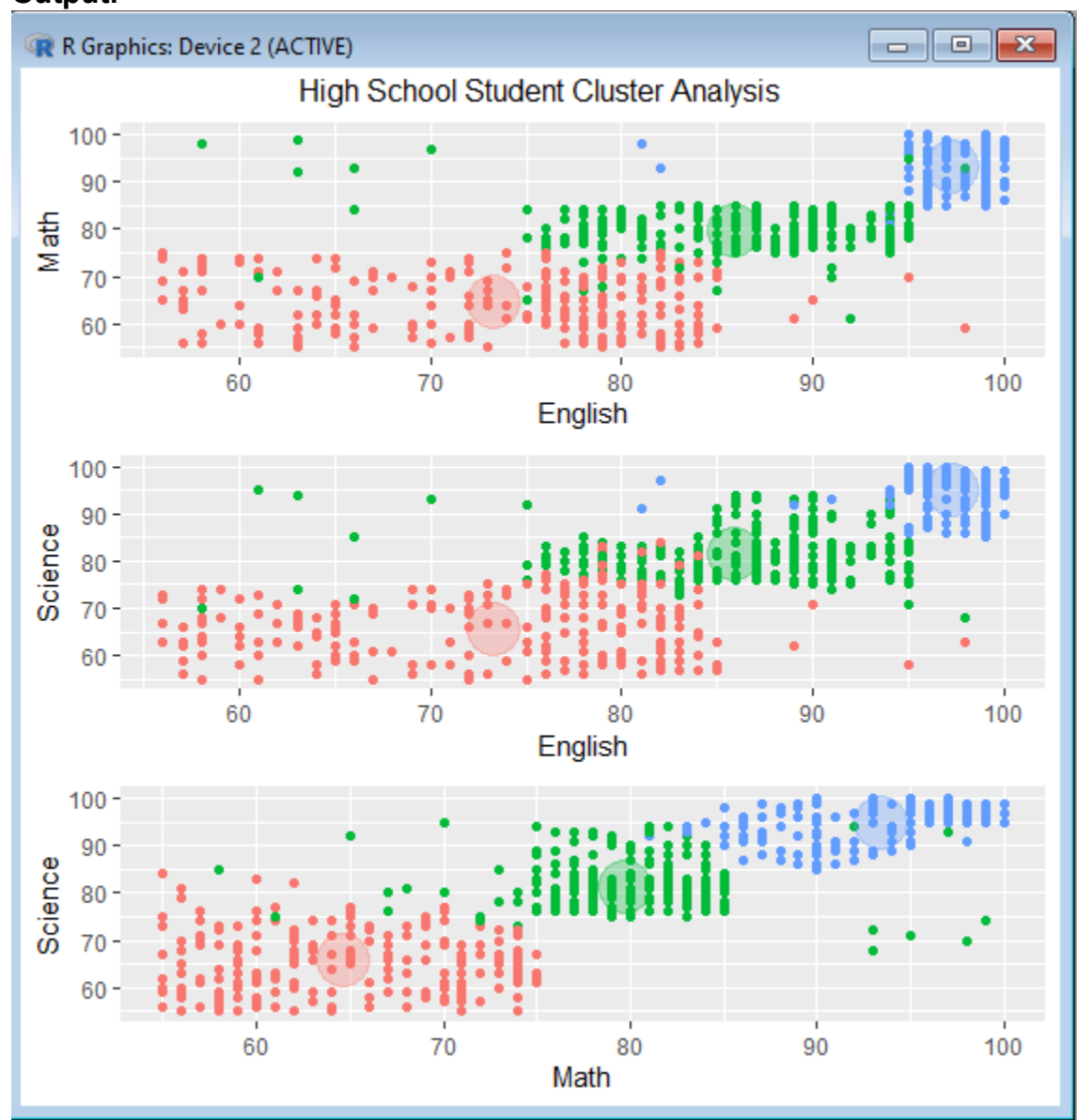
```
geom_point () + geom_point(data=centers,aes(x=Math,y=Science,
color=as.factor(c(1,2,3))),size=10,  alpha=.3, show.legend=FALSE)

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))
```

**Output:**

## Practical 6:
   a.  Simple Linear regression
**Aim:** Create your own data for years of experience and salary in lakhs and apply linear regression model to predict the salary
**Requirement:**
R tool

Code:
```
years_of_exp = c(7,5,1,3)
salary_in_lakhs = c(21,13,6,8)
employee.data = data.frame(years_of_exp, salary_in_lakhs)
employee.data

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

plot(salary_in_lakhs ~ years_of_exp, data = employee.data)
abline(model)
```
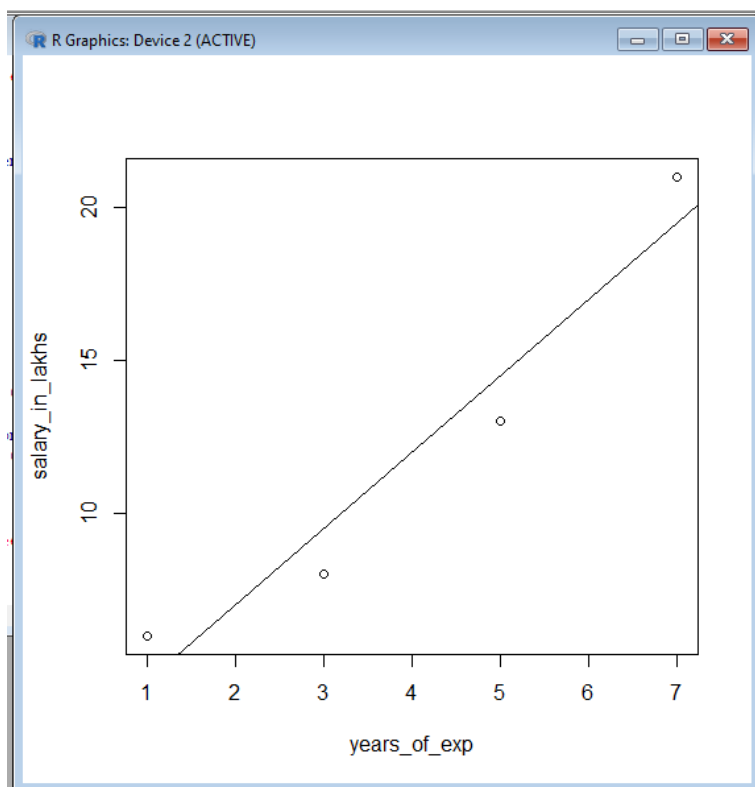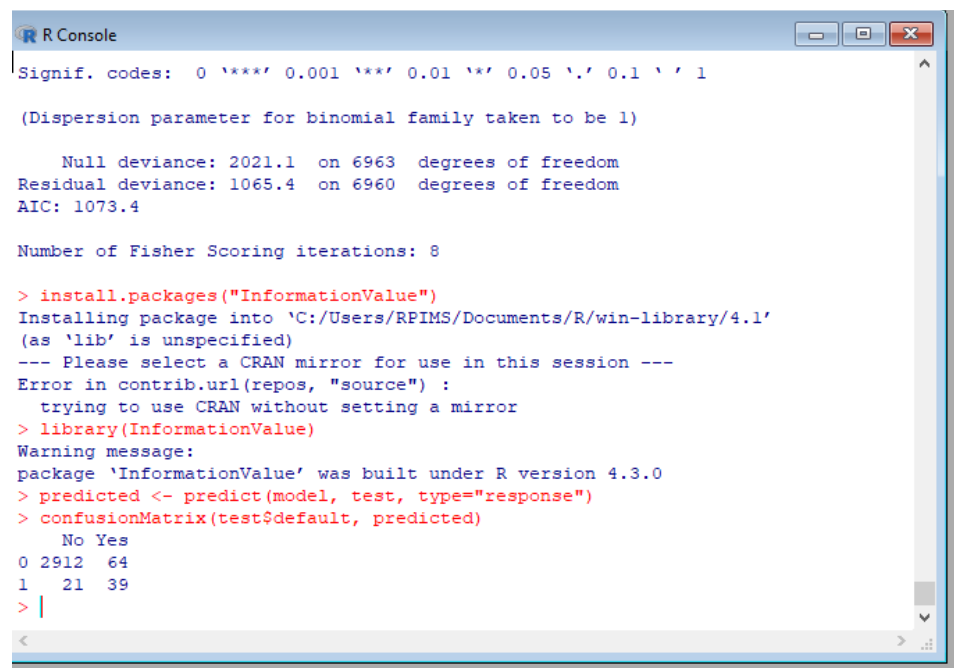
## Output:

b.: Logistic regression:
**Aim:** Take the in-built data from ISLR package and apply generalized logistic regression to find whether a person would be defaulter or not; considering input as student, income and balance.

Code:

```
install.packages("ISLR")
library(ISLR)
data <- ISLR::Default
print (head(ISLR::Default))
summary(data)
nrow(data)
set.seed(1)
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)
model <- glm(default~student+balance+income, family="binomial", data=train)
summary(model)
install.packages("InformationValue")
library(InformationValue)
predicted <- predict(model, test, type="response")
confusionMatrix(test$default, predicted)
```

Output:

**Practical 7:**
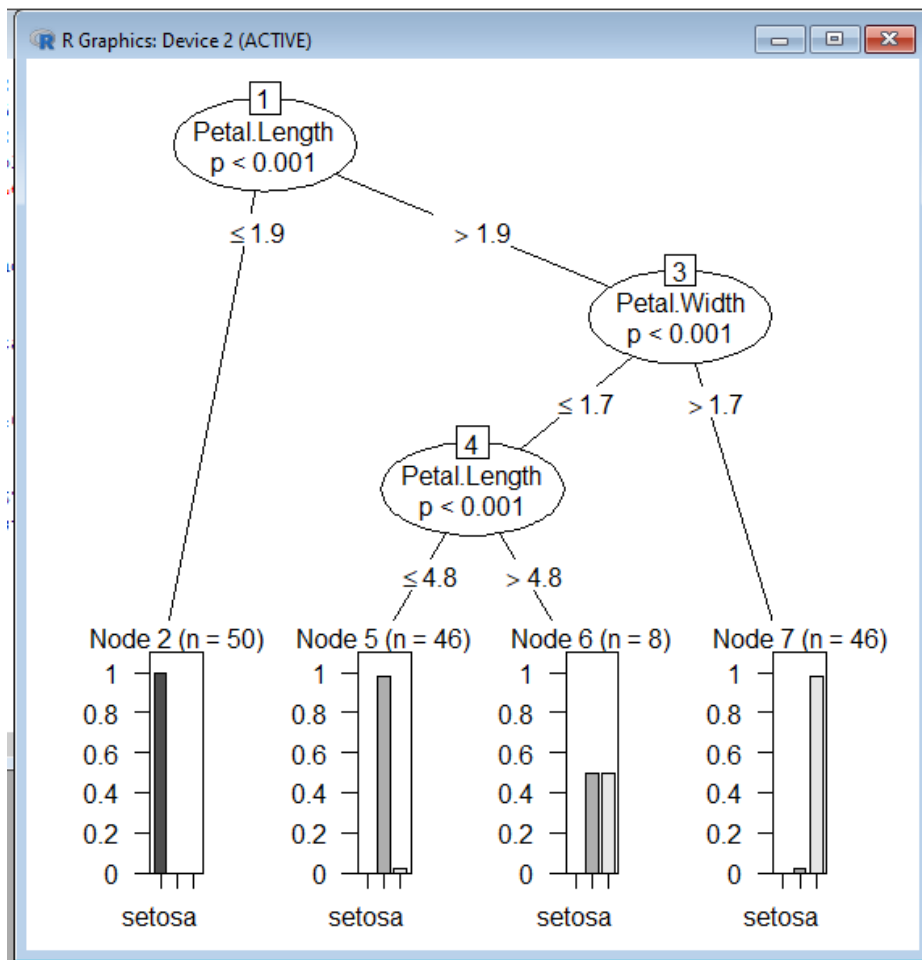**Aim:** Implement Decision tree classification techniques

**Requirement:**
R tool

**Code:**

```
library("party")
print(head(readingSkills))

str(iris)
iris_ctree <- ctree(Species ~ Sepal.Width + Sepal.Length + Petal.Length +
Petal.Width, data=iris

print (iris_ctree)
plot(iris_ctree)
```

**Output:**

### Practical 8:

Apriori algorithm

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

**Requirement:**
R tool

**Code:**

```
library(arules)
library(arulesViz)
library(RColorBrewer)

data(Groceries)
Groceries

summary(Groceries)
class(Groceries)

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

inspect(rules[1:10])

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


inspect(itemsets_3)
```
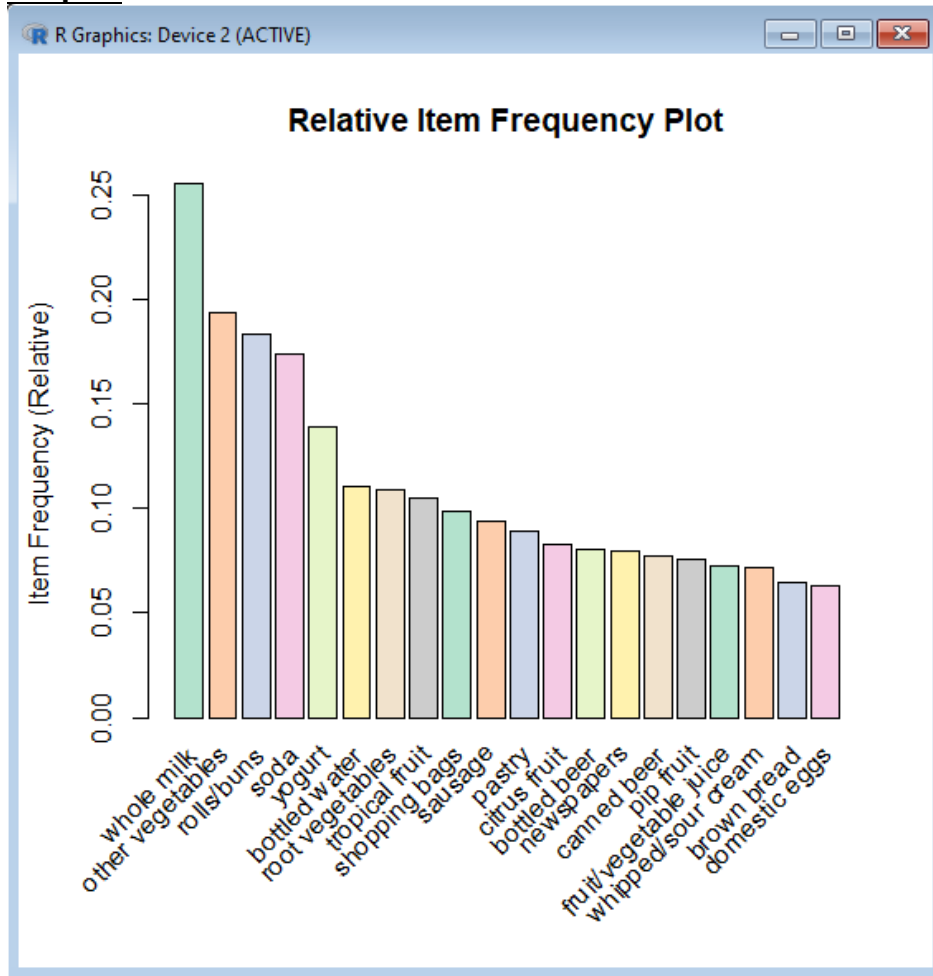
**Output:**



Relative Item Frequency Plot



```
        3       3       3       3       3       3

summary of quality measures:
    support            count
 Min.   :0.02227   Min.    :219.0
 1st Qu.:0.02250   1st Qu.:221.2
 Median :0.02272   Median :223.5
 Mean   :0.02272   Mean    :223.5
 3rd Qu.:0.02295   3rd Qu.:225.8
 Max.   :0.02318   Max.    :228.0

includes transaction ID lists: FALSE

mining info:
      data ntransactions support confidence
 Groceries          9835    0.02          1
                                                                  $
 apriori(data = Groceries, parameter = list(minlen = 3, maxlen = 3, support = 0$
>
>
> inspect(itemsets_3)
     items                                            support    count
[1] {root vegetables, other vegetables, whole milk} 0.02318251 228
[2] {other vegetables, whole milk, yogurt}          0.02226741 219
>
```