

hadoop@DESKTOP-24S8RFC:\$ sudo adduser --ingroup hadoop vishakha

Adding user 'vishakha' ---

Adding new 'vishakha' (1001) with group 'hadoop' ---

Creating home directory '/home/vishakha' ---

Enter the new value or press ENTER for the default

Full Name []: Viva

Room Number []:

Work Phone []:

Home Phone []:

Other []:

Is the information correct? [Y/n] Y

hadoop@DESKTOP-24S8RFC:\$

hadoop@DESKTOP-8HS8RFC:\$ sudo usermod -aG sudo

hadoop@DESKTOP-24S8RFC:\$ sudo visudo

Password

To run a command as administrator (use "root"),
use "sudo <command>"
See "man sudo_root" for details.

Generating public/private rsa key pair
The key fingerprint is

The keys randomart image is :

+---- [RSA 3072] ----+

----- 0 -----

0 - + +

B. + + =

+ x - - = S

= . + B + = + =

0 - - x + + + 0

.. = . E 0 " + 0 0

Shot on moto g⁴⁰ fusion

[PRANAV]TM

Practical 1

Tim: Install, configure and run Hadoop and HDFS.

Start TERMINAL and enter below commands to check version of java and Hadoop.

Command 1: Sudo addgroup hadoop

Command 2: Sudo adduser --ingroup hadoop Vishakha

Command 3: Sudo usermod -aG sudo Vishakha

Command 4: su Vishakha

Command 5: ssh-keygen -t rsa -P ""

Vishakha@DESKTOP-2458RFC:/home/hadoop\$ cat \$HOME/.ssh/id_rsa.
 pub > \$HOME/.ssh
 Vishakha@DESKTOP-2458RFC:/home/hadoop\$ sudo nano /etc/sysctl.conf
 [Sudo] password for vishakha:

for what other value do you want to change
 # kernel.sysrq = 438
 net.ipv6.conf.all.disable_ipv6 = 1 obj : 1 business
 net.ipv6.conf.default.disable_ipv6 = 1 obj : 3 business
 net.ipv6.conf.lo.disable_ipv6 = 1

^{^G1 Help} ^O Write out ^W WhereIS ^K Cut
^{^X Exit} ^R Read File ^V Replace ^U Paste

Resolving mirrors.estointernet.in ([mirrors.estointernet.net])
 Connecting to mirrors.estointernet.in
 HTTP request sent, awaiting response 200 OK
 Length: 605187279 (577m) [application/octet-stream]
 Saving to: 'hadoop-3.3.1.tar.gz'
 hadoop - 3.3.1.tar.gz

total 591044
 drwxr-xr-x 2 root root 4096 NOV 23 03:08
 drwxr-xr-x 2 root root 4096 NOV 23 03:08
 drwxr-xr-x 2 root root 4096 NOV 23 03:08
 -rw-r--r-- 1 Vishakha hadoop 625187279 JUN 15 2021

Command 6 : Cat \$HOME/.ssh/id_rsa.pub > \$HOME/.ssh/
 authorized_keys
 Command 7 : Sudo nano /etc/sysctl.conf

net.ipv6.conf.all.disable_ipv6 = 1
 net.ipv6.conf.default.disable_ipv6 = 1
 net.ipv6.conf.lo.disable_ipv6 = 1

Command 8 : cd /usr/local
 Command 9 : Sudo wget https://

Command 10 : Sudo tar xf hadoop-3.3.1.tar.gz
 Command 11 : Sudo mv hadoop-3.3.1 hadoop
 Command 12 : Sudo chown -R Vishakha:hadoop
 Command : ls -l



```
Vishalkha@DESKTOP-UNRAFS:~/usr/local$ source ~ / bashrc  
Vishalkha@DESKTOP-UNRAFS:~/usr/local$ cd /usr/local/hadoop/  
etc/hadoop  
Vishalkha@DESKTOP-UNRAFS:~/usr/local/hadoop/etc/hadoop
```

Hadoop 3.3.1
Source code repository <http://github.com/apache/hadoop>
git -> a3bac37a397ad188041dd80621
compilled with protoc 3.7.1
From Source with checksum 884uddb229gac054416d6d
This Source command was run using /usr/local/
hadoop/share/hadoop/common/hadoop-common-3.1.1/jar

```
[sudo] password for vishalkha  
Reading Package Lists... Done  
Building dependency tree... Done  
Reading state information... Done  
The following additional package will be installed  
libwrap ncurses-term openssh-server openssh  
sftp-server ssh-import-fd - suggested package  
molly-guard monkey sphere ssh-askpass  
The following NEW package will be installed  
libgssapi-ncurses-term openssh-server  
Need to get 803 kB of archives.  
After this operation, 6290 kB of additional disk space  
will be used.  
Do you want to continue? [y/N]y
```

Sudo Service ssh restart
• Restarting open BD secure Shell Server ssbd
vishalkha@DESKTOP user/local/hadoop/etc/hadoop.

Shot on moto g⁴⁰ fusion
[PRANAV]TM

DATE:

Command 13: nano \$HOME/.bashrc

Hadoop Related Options

export HADOOP_HOME=/usr/local/hadoop

export JAVA_HOME=INSTALL=\$HADOOP_HOME

export HADOOP_MAPRED_HOME=\$HADOOP_HOME

export HADOOP_COMMON_HOME=\$HADOOP_HOME

Command 14: Source ~/.bashrc

Command 15: cd /usr/local/hadoop/etc/hadoop

Command 16: nano hadoop-env.sh

Command 17: sudo mkdir -p /app/hadoop/tmp

Command 18: sudo chown vishalkha:hadoop /app/hadoop/
tmp

Command 19: Nano core-site.xml

<Property>

<name> hadoop.tmp.dir </name>

<Value> /app/hadoop/tmp </value>

</property>

<Property>

<name> fs.default.name </name>

<Value> hdfs://localhost:54310 </value>

</property>

Command 20: nano mapred-site.sh

<configuration>

<Property>

<name> mapred.job.tracker </name>

<Value> localhost:5434 </value>

</Property> </configuration>

Command 21: nano hdfs-site.xml

FOR EDUCATIONAL USE



WARNING! Permanently added 'localhost' (ED25519) to the
list of known hosts.
 Welcome to Ubuntu 22.04.2 LTS
 • Documentation: <https://help.ubuntu.com>
 • Management: <https://landscape.canonical.com>
 • Support: <https://ubuntu.com/advantage>

Welcome to Ubuntu 22.04.2 LTS (GNU/Linux 5.15.0-102-generic)
 Microsoft Standard WSL2 x86_64
 • Documentation: <https://help.ubuntu.com>
 • Management: <https://landscape.canonical.com>
 • Support: <https://ubuntu.com/advantage>

Vishakha@DESKTOP UNARAFS:

WARNING! Attempting to start all apache hadoop demons as
Vishakha in seconds. WARNING! This is not a recommended
production deployment configuration!

WARNING! USE CTRL-C to abort [CTRL-C]
 Starting namenodes on [localhost]
 Starting datanodes.
 Starting secondary namenode [DESKTOP-2458RFC]

DESKTOP-2458RFC WARNING! PARAMETER

2644 JPS
 1638 NameNode
 2311 NodeManager
 1976 Secondary NameNode. *New*
 1787 DataNode
 2187 ResourceManager

WARNING! Stopping all Apache Hadoop daemons [are
Vishakha in seconds]

WARNING! USE CTRL-C to abort

<property>
<name>dfs.replication </name>
<property>

Command 22: cd

Command 23: hadoop version

HDFS

Command 24: cd /usr/local/hadoop/etc/hadoop

Command 25: hadoop namenode -format

Command 26: ssh

Command 27: sudo apt-get install ssh

Command 28: sudo service ssh restart

Command 29: ssh localhost

Command 30: cd /usr/local/hadoop/sbin

Command 31: start-all.sh

Command 32: jps

Command 33: stop-all.sh

Thus we have successfully installed Hadoop.

o/p :-

```
import sys
for line in sys.stdin:
    line = line.strip()
    words = line.split()
    for word in words:
        print ('%s\t%s' % (word, 1))
```

o/p :-

bus 1
car 1
train 1
ship 1
Ship 1
train1
bus 1
Ship 1
car 1

*Next*Practical 2.

Aim: Implement word count program using Map Reduce.

Command : cd

Command 1: Mkdir MapReduce

Command 2: cd MapReduce

Command 3: touch word_count_data.txt

Command 4: nano word_count_data.txt

bus car train

ship Ship train

bus Ship Car

Command 5: cat word_count_data.txt

Command 6: touch mapper.py

Command 7: nano mapper.py

import sys

for line in sys.stdin:

line = line.strip()

words = line.split()

for word in words:

print ('%s\t%s' % (word, 1))

Command 8: cat word_count_data.txt | python: mapper.py

O/P =

```

Shipt 1 : bus
bus 2 : bananas
Car 2 : line = line.rstrip()
Shipt 2 : line = line.split(',')
train 2 : words

```

Next

O/P =

```

bus 1 : bus for : 2 bananas
car 1 : car for : 2 bananas
train 1 : train for : 2 bananas
Shipt 1 : shipt for : 2 bananas
bus 2 : bus for : 2 bananas
bus 3 : bus for : 2 bananas
Shipt 2 : shipt for : 2 bananas
Car 2 : line = line.rstrip()

```

Command 9: touch reducer.py
 Command 10: nano reducer.py

```

from operator import itemgetter
import sys
current_word = None
current_count = 0
word = None
for line in sys.stdin:
    line = line.rstrip()
    word, count = line.split(' ', 1)
    try:
        count = int(count)
    except ValueError:
        continue
    if current_word == word:
        current_word += count
    else:
        if current_word:
            print('%s\t%s' % (current_word, current_count))
        current_count = count
        current_word = word
    if current_word == word:
        print('%s\t%s' % (current_word, current_count))

```

Command 11: cat word-count-data.txt | python3
 mapper.py | sort -k1,1 | python3 reducer.py



vishakha@DESKTOP-UNARAF5:/usr/local/hadoop/etc/hadoop\$

Sudo service ssh restart

* Restarting OpenBSD Secure Shell server sshd

vishakha@DESKTOP-UNARAF5:/usr/local/hadoop/etc/hadoop\$

This key is not known by any other names
Are you sure you want to continue connecting (yes/no)?

Warning: Permanently added 'localhost' (RSA) to the list of known hosts.

Welcome to Ubuntu 22.04.2 LTS

* Documentation: <https://help.ubuntu.com>
* Management: <https://landscape.canonical.com>
* Support: <https://ubuntu.com/advantage>
Ubuntu comes with ABSOLUTE NO WARRANTY to the extent
permitted by application law.
* Documentation: <http://help.ubuntu.com>
* Management: <https://landscape.canonical.com>
* Support: <https://ubuntu.com/advantage>.

vishakha@DESKTOP-UNARAF5 \$

WARNING: Attempting to start all Apache Hadoop daemons
as vishakha in 10 seconds.

WARNING: This is not a recommended production deployment
configuration.

Starting NameNode on [Localhost]

Starting datanode

Starting Secondary namenode [DESKTOP - 245.8 RFC]

2644 JPS

1638 NameNode

1231 NodeManager

1787 DataNode

Shot on moto g⁴⁰ fusion

DATE:

Practical 3.

Exploring Hadoop Distributed file system (HDFS)

Command: su vishakha

Command: java -version

Command: hadoop version

Command 1: cd /usr/local/hadoop/etc/hadoop

Command 2: hadoop namenode -format

Command 3: ssh

Command 4: sudo apt-get install ssh

Command 5: sudo service ssh restart

Command 6: ssh localhost

Command 7: cd /usr/local/hadoop/sbin

Command 8: start-all.sh

Command 9: jps

Command 10: stop-all.sh

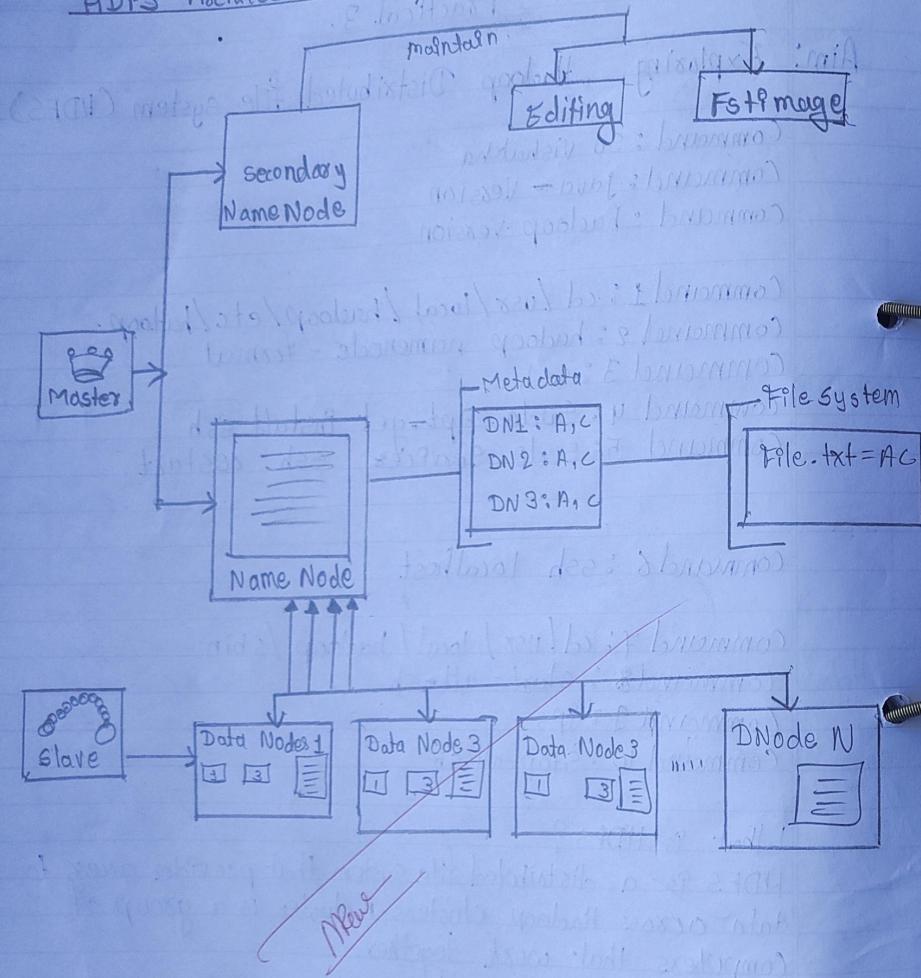
What is HDFS?

HDFS is a distributed file system that provides access to data across Hadoop clusters. A cluster is a group of computers that work together.

Why HDFS?

Before 2011, storing and retrieving petabytes or zettabytes of data had the following three more challenges:
Cost, Speed, Reliability.

HDFS Architecture



1. Cost: HDFS Hadoop clusters can read or write more than a terabyte of data per second.
 2. Speed: Large Hadoop clusters can read or write more than a terabyte of data per second.
 3. Reliability: HDFS copies the data multiple times & distributes the copies to individual nodes.
- Characteristics of HDFS.

1. HDFS has high fault-tolerance.
2. HDFS may consist of thousands of server machines.
3. HDFS has high throughput.
4. HDFS is designed to store & scan millions of rows of data.
5. HDFS is designed in such a way that it can be built on commodity hardware.

HDFS Architecture and Components

Broadly, HDFS architecture is known as the master & slave architecture which is shown below:

HDFS Components:

1. Namenode: The namenode server is the core component of an HDFS cluster. There can be only one NameNode server in an entire cluster.
2. Secondary NameNode: The Secondary NameNode server maintains the edit log and namespace image information in sync with the NameNode server.
3. File system: HDFS exposes a file system namespace & allows user data to be stored in files.
4. Metadata: When new DataNodes join a cluster, metadata loads the blocks.
5. DataNode: Data Node manages names & locations of file blocks. By default, each file block is 128 Megabytes.

Practical 4

Aim:

Implement a Program in Pig.

Command 1: Su Vishalkha

Command 2: java -version

Command: hadoop version

Command 3: Sudo Service ssh restart

Command 4: ssh localhost

Command 5: start-all.sh

Command 6: cd /home/vishalkha

Command 7: sudo nano sample.txt

1. Sanvi, 2000, lecturer

2. Manoj, 25000, Accountant

3. Akhil, 30000, IT

Command 8: cd /usr/local

Command 9: sudo wget https://mirrors.estointernet.in/

apache/pig/pig-0.16.0/pig-0.16.0.tar.gz

Command 10: sudo tar -xvzf pig-0.16.0.tar.gz

Command 11: sudo mv pig-0.16.0 pig

Command 12: ls -l

Command 13: cd /home/vishalkha

Command 14: nano hashxc

export PATH=\$PATH:\$HADOOP_HOME/bin

export PATH=\$PATH:\$HADOOP_HOME/sbin

export PIG_HOME = /usr/local/pig

export PIG_CLASSPATH = /usr/local/hadoop/etc/hadoop

Command 15: Source .bashrc

Command 16: Pig -version

Command 17: pig

Command 18: Employees = LOAD 'Sample.txt' USING PigStorage
(' ') AS (id : int, firstname : chararray, salary : int,
dept : chararray);

Command 19: Dump Employees;

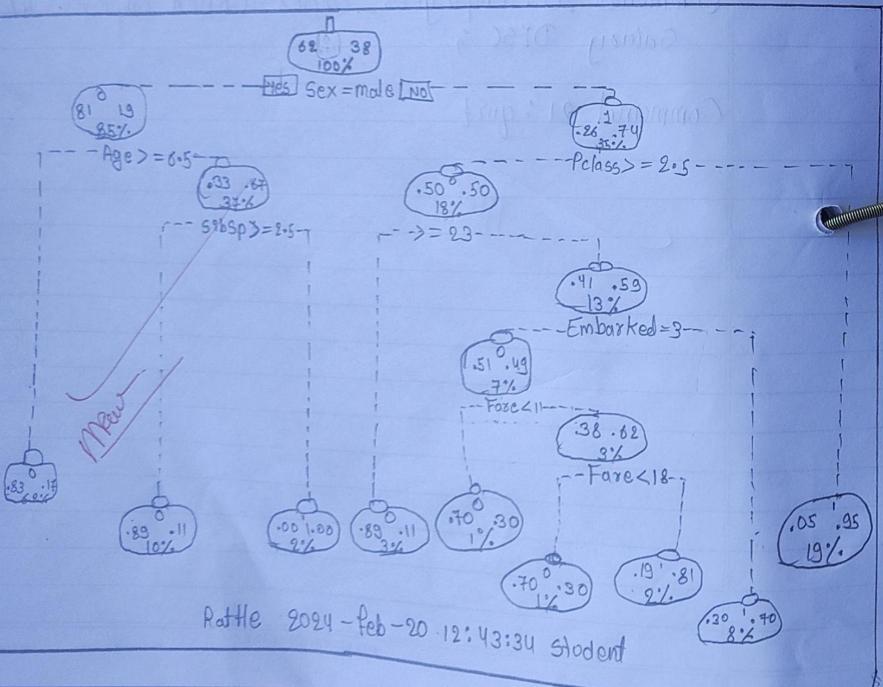
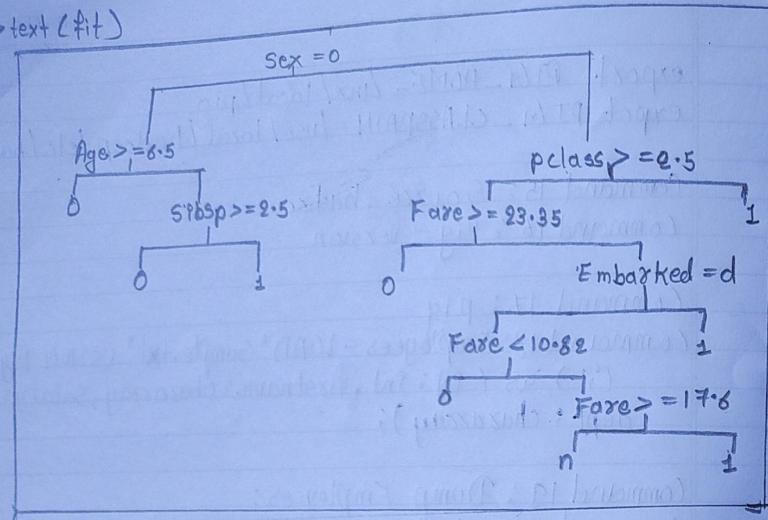
Command 20: Employees_order = ORDER Employees BY
Salary DESC ;

Command 21: quit

Never



>text(fit)

Practical 5.

- Aim :
- Implement Decision tree classification technique.
 - Implement SVM classification technique.

```
titanic <- read.csv(file.choose(), header=T, sep=",")
```

```
Summary(titanic)
```

```
name(titanic)
```

```
library(xpart) titanic
```

```
library(xpart)
```

```
fit <- xpart(C ~ survived ~ Pclass + Sex + Age + Sibsp + Parch +  
Fare + Embarked, data = titanic, method = "class")
```

```
plot(fit)
```

```
text(fit)
```

```
install.packages("xattle")
```

```
install.packages("xpart.plot")
```

```
install.packages("RColorBrewer")
```

```
library(xattle)
```

```
library(xpart.plot)
```

```
library(RColorBrewer)
```

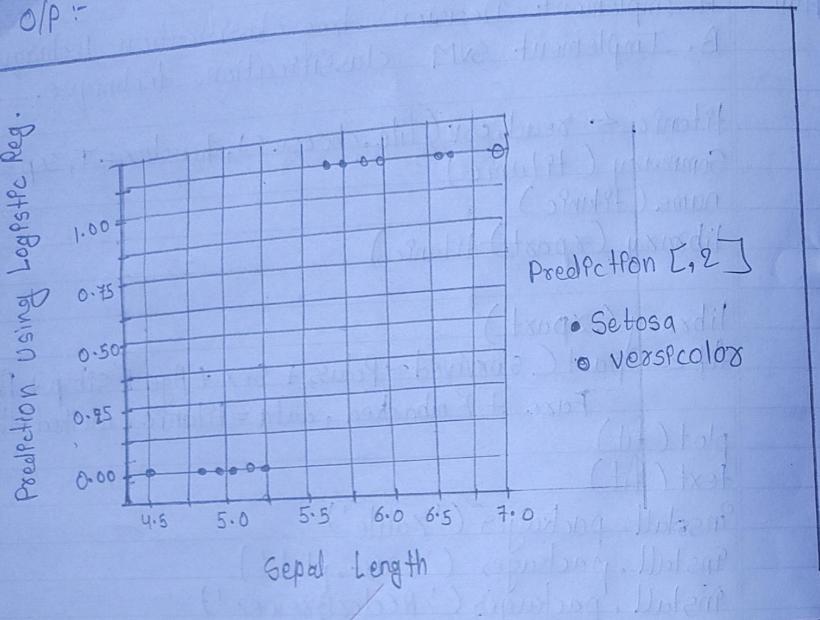
```
fancyRpartPlot(fit)
```

```
Prediction <- predict(fit, titanic, type = "class")
```

```
Prediction
```



O/P :-



More

Aim:-

Regression Model : Import a data from web storage. Name the dataset & do Logistic Regression to find out relation between variable that are affecting the admission of a student.

library (datasets)

iris_data <- iris

head (iris_data)

str (iris_data)

levels (iris_data\$Species)

sum (is.na(iris_data))

iris_data <- iris_data[1:100,]

set.seed (100)

samp <- sample (1:100, 30)

iris_test <- iris_data[samp,]

iris_ctrl <- iris_data[-samp,]

install.package ("ggplot2")

library (ggplot2)

install.package ("gridExtra")

library (gridExtra)

ggpairs (iris_test)

y <- iris_test\$Species; x <- iris_test\$Sepal.Length

glifit <- glm (y ~ x, family = "binomial") summary (glifit)

newdata <- data.frame (x = iris_ctrl\$Sepal.Length)

predicted_val <- predict (glifit, newdata, type = "response")

Prediction <- data.frame (iris_ctrl\$Sepal.Length, iris_ctrl\$Species,

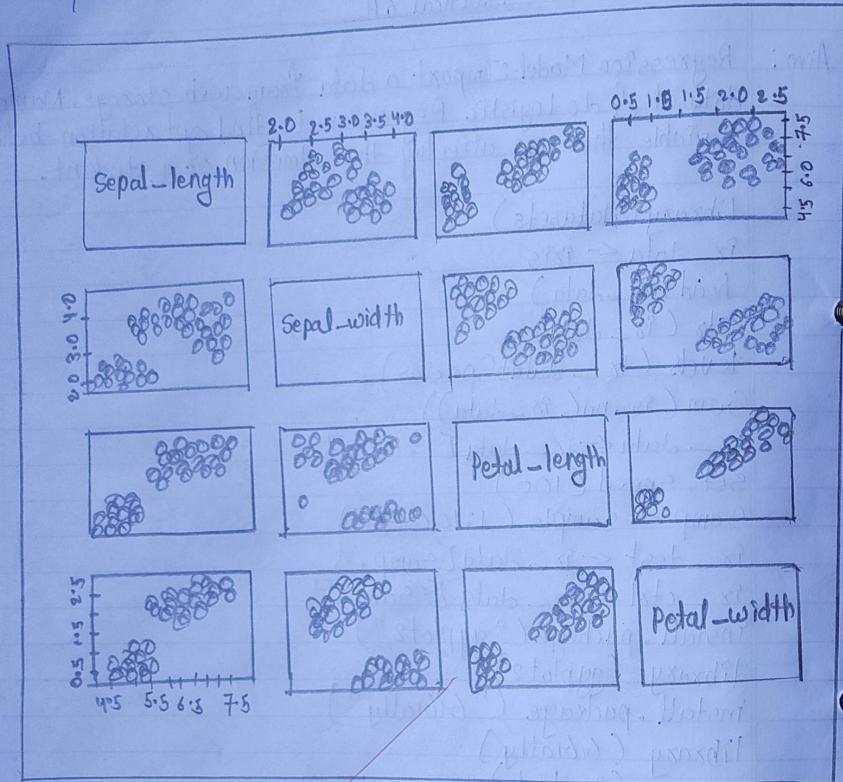
predicted_val)

Prediction

qplot (predicted_val[, 1], round (predicted_val[, 3]), col = prediction [, 2], xlab = 'Sepal.Length', ylab = 'prediction using Logistic Reg.')

Q/P :-

Ans / Answer

Practical 6B

Aim:

Multiple Regression model:- Apply multiple Regression, if data have a continuous independent variable.
Apply on above dataset.

"Linear Regression"

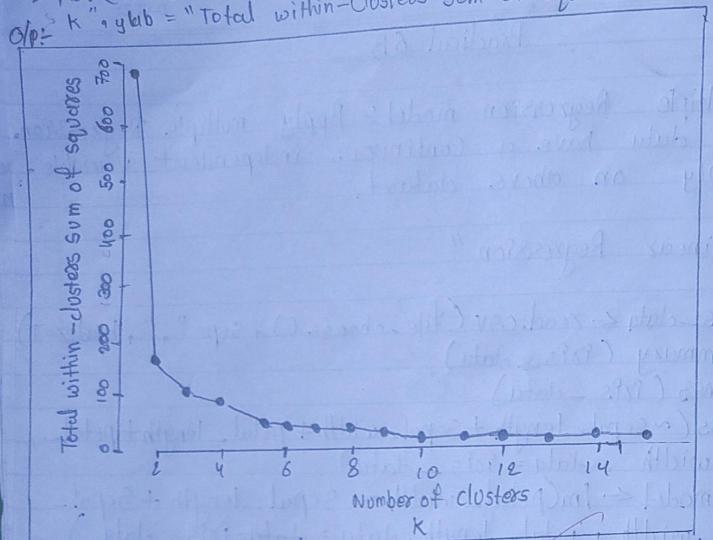
```
iris_data <- read.csv("file.choose()", sep = ", ", header = T)
summary(iris_data)
names(iris_data)

pairs(~ Sepal.length + Sepal.width + Petal.length + Petal.width, data = iris_data)
irismodel <- lm(Petal.width ~ Sepal.length + Sepal.width + Petal.length, data = iris_data)
```

Summary(irismodel)



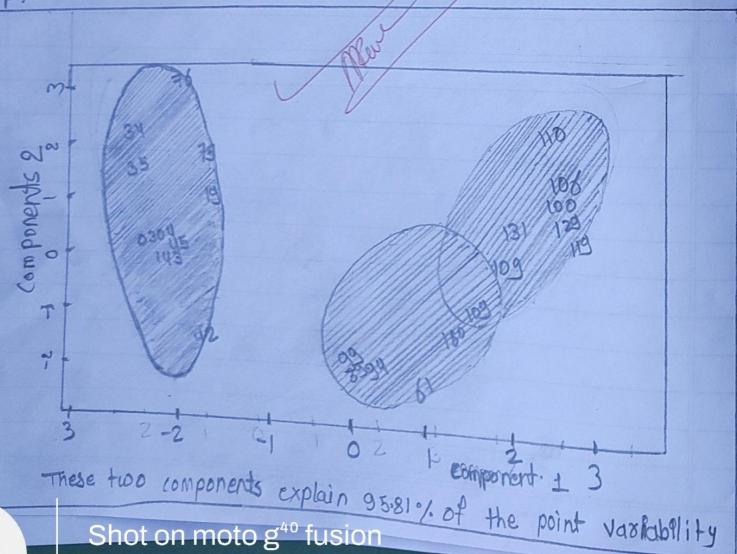
```
> plot(1:15, wss, type = "b", pch = 19, frame = FALSE, xlab = "Number of cluster  
k", ylab = "Total within-Clusters sum of square")
```



```
>clusplot(new_data, cl$cluster, color=TRUE, shade=TRUE, labels=2, line=0)
```

```
>cl$cluster
```

O/P:



Practical 7.

Aim:- a. classification Model:

- a] install genre-relevant package for classification
- b] choose classifier for classification problem.
- c] evaluate the performance of classifier

b. clustering Model:

- a] clustering algorithm for unsupervised classification
- b] plot the cluster data using R visualizations.

"k-mean clustering"

```
data("iris")
```

```
name(iris)
```

```
new_data <- subset(iris, select = c(~Species))
```

```
new_data
```

```
cl <- kmeans(new_data, 3)
```

```
cl
```

```
data <- new_data
```

"Simplify"

```
wss <- sapply(1:15, function(k) kmeans(data, k)$tot.wss)
```

```
wss
```

```
plot(1:15, wss, type = "b", pch = 19, frame = FALSE, xlab = "Number of clusters k", ylab = "Total within-clusters sum of squares")
```

Install.package("cluster")

library(cluster)

```
clusplot(new_data, cl$cluster, color = TRUE, shade = TRUE, labels = 9, line = 0)
```

cl\$cluster



Practical 8.

Aim:

Implement an application that store big data in MongoDB.

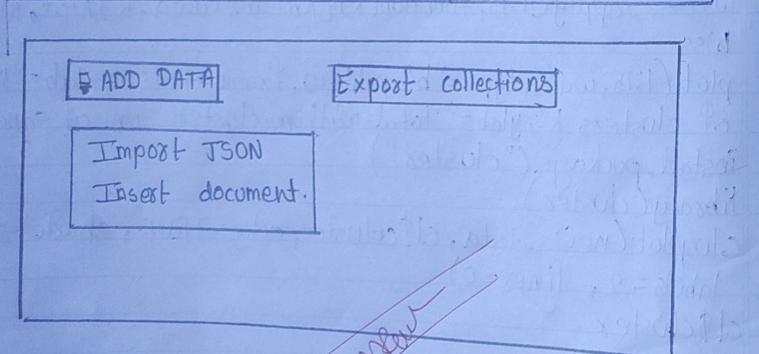
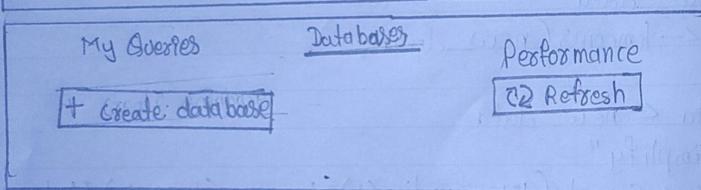
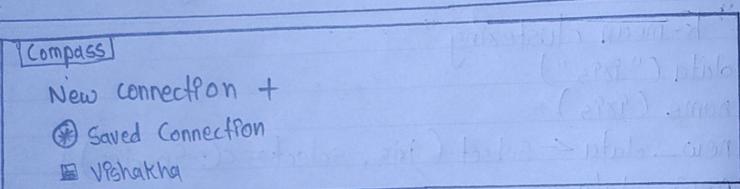
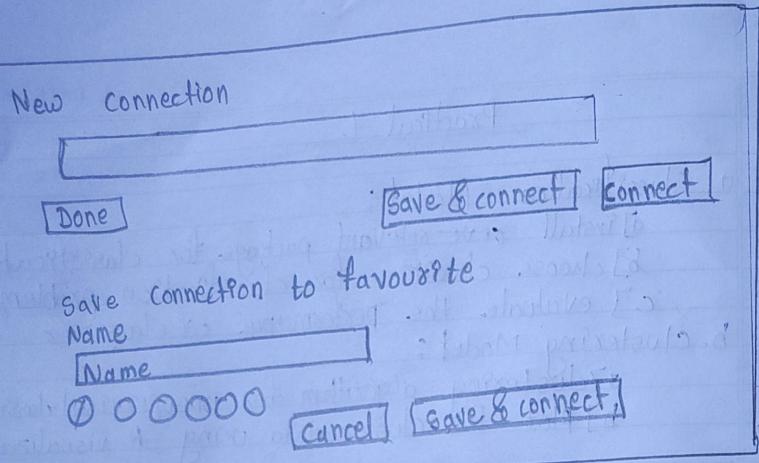
Step 1: Install MongoDB compass

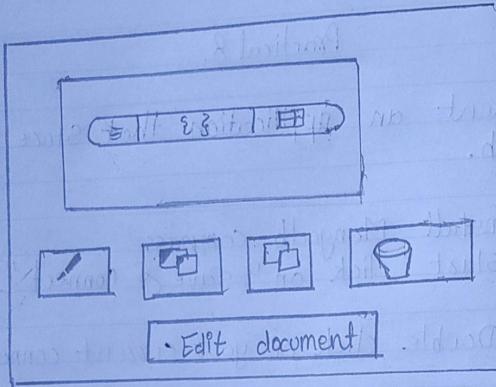
Step 2: Start click on "save & connect".

Step 3: Double click on your current connection.

Step 4: Move to database and click on → + create database.

Step 5: Now that the database and collection have been created. you can insert a list of unsorted document into the cities collection. click on ADD DATA button. Select the Insert Document option.





ADD DATA ▾ EXPORT COLLECTION 1 - 2 of 2 ↗

id : objectId ("65effc19350117b84ab")
 name: "Mexico City"
 Country: "Mexico"
 Continent: "North America"
 population: 21.581

Note: This document has been modified.

Step 6: Enter the following set of document onto the files.

```
{
  "name": "Karachi", "country": "Pakistan", "continent": "Asia",
  "population": 15.400
},
{
  "name": "Dhaka", "country": "Bangladesh", "continent": "Asia",
  "population": 19.578
},
{
  "name": "New York", "country": "U.S", "continent": "North America",
  "population": 18.819
},
{
  "name": "Manila", "country": "Philippines",
  "continent": "Asia",
  "population": 14.838
}
```

Step 7: Viewing, Editing and Filtering the data by choosing from.

Step 8: You can also delete data & connections.