

1.HypothesisTest

```
x<-rnorm(100) t.test(x,mu=5)
x1<-rnorm(100) y1<-rnorm(100) t.test(x1,y1)
x2<-rnorm(100) t.test(x2,mu=2,alternative='greater')
x3<-rnorm(100) wilcox.test(x3,exact=FALSE)
cor.test(mtcars$mpg,mtcars$hp)
```

.....

2.KMeans Clustering

```
library(cluster) library(ggplot2)
set.seed(20)
irisCluster<-kmeans(iris[,3:4],3,nstart=20) irisCluster
irisCluster$cluster<-as.factor(irisCluster$cluster)
ggplot(iris,aes(Petal.Length,Petal.Width,color=irisCluster$cluster))+geom_p
oint()
d<-dist(as.matrix(mtcars)) hc<-hclust(d)
plot(hc)
x<-rbind(cbind(rnorm(10,0,0.5),rnorm(10,0,0.5)),cbind(rnorm(15,5,0.5),rnorm(15,5,0.5)))
clusplot(pam(x,2))
x4<-cbind(x,rnorm(25),rnorm(25)) clusplot(pam(x4,2))
```

.....

3.Implement Linearand Logistic Regression

```
dataset=read.csv("/content/data-marketing-budget-12mo.csv",header=T,
colClasses=c("numeric","numeric","numeric"))
head(dataset)
simple.fit = lm(Sales~Spend,data=dataset) summary(simple.fit)
multi.fit=lm(Sales~Spend+Month,data=dataset) summary(multi.fit)
input<-mtcars[,c("am","cyl","hp","wt")] print(head(input))
input<-mtcars[,c("am","cyl","hp","wt")]
am.data=glm(formula=am~cyl+hp+wt,data=input,family=binomial)
print(summary(am.data))
```

.....

4.Time Series Analysis.

```
kings <- scan("http://robjhyndman.com/tsdldata/misc/kings.dat",skip=3)
kings
kingstimeseries<-ts(kings)
kingstimeseries
births<-scan("http://robjhyndman.com/tsdldata/data/nybirths.dat")birth
stimeseries<-ts(births,frequency=12,start=c(1946,1))
birthstimeseries
souvenir<-
scan("http://robjhyndman.com/tsdldata/data/fancy.dat")
souvenirtimeseries<-
ts(souvenir,frequency=12,start=c(1987,1))
Souvenirtimeseries
```

5.Data Visualizationtoexplore Various

CodeForHistogram:

```
library(RColorBrewer) data(VADeaths) par(mfrow=c(2,3))
hist(VADeaths,breaks=10, col=brewer.pal(3,"Set3"),main="Set3 3 colors")
hist(VADeaths,breaks=3 ,col=brewer.pal(3,"Set2"),main="Set2 3 colors")
hist(VADeaths,breaks=7, col=brewer.pal(3,"Set1"),main="Set1 3 colors")
hist(VADeaths,,breaks= 2, col=brewer.pal(8,"Set3"),main="Set3 8 colors")
hist(VADeaths,col=brewer.pal(8,"Greys"),main="Greys8colors")
hist(VADeaths,col=brewer.pal(8,"Greens"),main="Greens8colors")
```

CodeForLineChart:

```
data(AirPassengers)
plot(AirPassengers,type="l")
```

CodeforBarChart:

```
data("iris")
barplot(iris$Petal.Length) #Creating simple Bar Graph barplot(iris$Sepal.Length,col =
brewer.pal(3,"Set1")) barplot(table(iris$Species,iris$Sepal.Length),col
=brewer.pal(3,"Set1"))
```

CodeforBoxplot:

```
data(iris) par(mfrow=c(2,2))
boxplot(iris$Sepal.Length,col="red") boxplot(iris$Sepal.Length~iris$Species,col="red")
boxplot(iris$Sepal.Length~iris$Species,col=heat.colors(3))boxplot(iris$Sepal.Length~iris$Sp
ecies,col=topo.colors(3)) boxplot(iris$Petal.Length~iris$Species)
Code for Scatter Plot: plot(x=iris$Petal.Length) plot(x=iris$Petal.Length,y=iris$Species)
```

6.Install and ConfigureHadoop

SETUP:

SettingupHadooponWindows

PREREQUISITES:

Windows64bitOS

JavaJDK

AdministratorAccess

INSTALLATION

1.DownloadHadoopforwindowsfromtheofficialsite

2.ExtracttheZIPtoyourChosendirectory(Hadoopinstallatio)

3.SetEnvironmentVariables: 'HADOOP_HOME' 'JAVA_HOME'

Edit the "Path" variable and add: '%HADOOP_HOME%\bin' '%HADOOP_HOME%\sbin'
'JAVA_HOME%\bin'

4.ConfigureHadoop

a)Open'hadoop-env.cmd',set'JAVA_HOME'.

Create/edit'core-site.xml'with:

xml

<property>

<name>fs.defaultFS</name>

<value>hdfs:localhost:9000</value>

</property>

b) Create/edit 'hdfs-site.xml' with: xml

```
<property>
<name>dfs.replication</name>
<value>1</value>
</property>
```

d) create/edit 'yarn-site.xml' with: xml

```
<property>
<name>yarn.nodemanager.aux-services</name>
<value>mapreduce_shuffle</value>
</property>
<property>
<name>yarn.nodemanager.env-whitelist</name>
<value>HADOOP_HOME,YARN_HOME</value>
</property>
```

Formatting HDFS:

5. Open a command prompt and run: `hdfs namenode format`

Starting Hadoop Services:

6. In Command prompt run

`cd %HADOOP_HOME%`

`start-dfs.cmd start-yarn.cmd`

Testing Hadoop:

7. Open a browser and visit 'http://localhost:9870' to see Hadoop namenode Web interface.

Running a MapReduce Example:

8. To test, run a MapReduce example (replace jar file)

Hadoop jar

`share/hadoop/mapreduce/hadoop*-mapreduce-examples.jar` `PI 16 1000`

.....

7. MapReduce

PREPARE:

Download `MapReduceClient.jar`

Download `Input_file.txt` Place both files in "C:/"

HADOOP OPERATIONS:

1. Open cmd in administrative mode and move to "C:/Hadoop-2.8.0/sbin" and start the cluster.

2. `Start-all.cmd`

3. Create an input directory in HDFS

`hadoopfs-mkdir/input_dir`

4. Copy the input text file named `input_file.txt` in the input directory of HDFS

`hadoopfs-put C:/input_file.txt/input_dir`

Verify the file `input_file.txt` is available in HDFS input directory.

`Hadoopfs-ls/input_dir/`

5. Run `MapReduceClient.jar` and also provide input and output directories.

`hadoopjar C:/MapReduceClient.jar wordcount/input_dir/output_dir`

6. Verify the content for generated output file

`hadoopdfs-cat/output_dir/`

Some other useful commands:

7. To leave Safe mode:

`hadoopdfsadmin-safemodeleave`

8. To delete file from HDFS directory

`hadoop fs -rm -r /input_dir/input_file.txt`

9. To delete directory from HDFS directory

`hadoop fs -rm -r /input_dir`

.....

8. implement an application that stores big data

HBase/MongoDB/Pig using Hadoop/R/Cassandra

Setup the Environment:

You should have Hadoop, HBase, MongoDB, Pig, R, and Cassandra installed and Configured.

Sample dataset:

Let us assume you have a simple dataset named "sample_data.csv" like this: Name, Age, City

John, 30, New York Alice, 25, Los Angeles Bob, 35, Chicago

Hadoop MapReduce (Python)

Python mapreduce script to process the dataset and store it in HBase: #mapper

```
import sys
```

```
for line in sys.stdin:
```

```
    line = line.strip()
    fields = line.split(",")
```

```
    Name, age, city = fields
    print(f'{Name}\t{age}\t{city}')
```

```
#reducer
```

```
import happybase
```

```
connection = happybase.Connection(host='localhost', port=9090)
table =
```

```
connection.table('my_table')
```

```
for line in sys.stdin:
```

```
    name, age, city = line.strip().split("\t")
    table.put(name, {'info:age': age, 'info:city': city})
```

Pig Script:

APig script to perform some transformations and store the data in MongoDB:

```
data = LOAD 'sample_data.csv' USING PigStorage(',') AS (name:chararray, age:int,
```

```
city:chararray);
```

```
filtered_data = FILTER data BY age >= 30;
```

```
STORE filtered_data INTO 'mongodb://localhost:27017/mydb.mycollection' USING
```

```
org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.PigMongoStorage();
```

R Script:

```
An R script to analyze the data: library(rhbase)
```

```
hbase <- HBase$new(host="localhost", port=9090)
data <- hbase$getTable("my_table")
```

```
print(data)
```

Cassandra (CQL):

You would create a Cassandra keyspace, define a table schema, and insert data using CQL

commands. Below is a simplified example.

```
cql
```

```
CREATE KEYSPACE mykeyspace WITH replication = ('class': 'SimpleStrategy',
```

```
'replication_factor': 1);
```

```
USE mykeyspace;
```

```
CREATE TABLE mytable (name TEXT PRIMARY KEY, age INT, city TEXT);
```

```
INSERT INTO mytable (name, age, city) VALUES (John, 30, 'New York');
```

Running the Code:

Run the Python MapReduce script using Hadoop, execute the Pig script, run the R script, and execute the CQL commands in Cassandra.

9. Use apach spark for data Analytics

Creating SparkSession:

```
from pyspark import SparkContext
from pyspark.sql import SparkSession
from pyspark.sql.types import *
from pyspark.sql.functions import *
from pyspark.sql.types import Row
from datetime import datetime
```

Initializing SparkSession:

```
sc = SparkContext()
spark = SparkSession.builder.appName("PythonSparkSQLbasic example").config("spark.some.config.option", "some-value").getOrCreate()
```

Creation of Spark RDD:

```
srecord = sc.parallelize([Row(roll_no=1, name="john", passed=True, marks={'Math':89, 'Physics':87, 'Chemistry':96}, sports=['chess', 'football'], DoB=datetime(2012,5,1,12,1,5)),
row(roll_no=2, name="Vignesh", passed=False, marks={'Math':95, 'Physics':66, 'Chemistry':77}, sports=['carrom', 'tennis'], DoB=datetime(2012,5,12,14,2,5)),
Row(roll_no=3, name="Sidharth", passed=True, marks={'Math':95, 'Physics':100, 'Chemistry':95}, sports=['football', 'kabadi'], DoB=datetime(2012,5,14,12,2,5))])
```

Creating a DataFrame:

```
srdf = srecord.toDF()
srdf.show()
```

Create Temporary View:

```
srdf.createOrReplaceTempView('records')
spark.sql("SELECT * FROM records").show()
re = spark.sql("SELECT * FROM records")
type(re)
```

Accessing Elements of a List or Dictionary within DataFrame:

```
spark.sql('SELECT roll_no, marks["Physics"], sports[1] FROM records').show()
```

Usage of Where Clause:

```
spark.sql("SELECT * FROM records where passed= True").show()
spark.sql('SELECT * FROM records where marks["Chemistry"] < 40').show()
```

Creating Global View:

```
srdf.createGlobalTempView('globalrecord')
spark.sql("SELECT * FROM global_temp.globalrecord").show()
```

Dropping Columns From DataFrame:

```
srdf.columns
srdf = srdf.drop('passed')
```

Few more queries:

```
spark.sql("SELECT round((marks.Physics+marks.Chemistry+marks.Math)/3) avg_marks FROM records").show()
srdf = spark.sql("SELECT *, round((marks.Physics+marks.Chemistry+marks.Math)/3) avg_marks FROM records")
srdf.show()
```

10. Filter and sort the sales data using Pig latin in Pig

1. Prepare the sales dataset in the format sales.csv

Transaction_id, customer, amount

1,	Alice,	500
2,	Bob,	1500
3,	Charlie,	700

4, Diana, 2000
5, Eve, 1200

2. Upload this file to HDFS:

Upload this csv file to the hdfs by using the Commands:

`hdfsdfs-mkdir/input`

`hdfsdfs-put sales.csv/input/`

3. Write the script named `filter_and_sort` and save this file in the extension .pig

After uploading data set to hdfs, write pig script in save it in the extension .pig

PigScript:

Load the sales data from the HDFS:

```
sales_data=LOAD '/input/sales.csv' USING PigStorage(',') AS (transaction_id:int,  
customer:chararray, amount:float);
```

Filter transactions with amounts greater than 1000:

```
filtered_data=FILTER sales_data BY amount>1000;
```

Sort the filtered data by amount in descending order:

```
sorted_data=ORDER filtered_data BY amount DESC;
```

Store the result in HDFS:

```
STORE sorted_data INTO '/output/sorted_sales' USING PigStorage(',');
```

4. Run the PigScript:

Run the above pig script by using the command in the CMD: Verify Output directory:

`hdfsdfs-ls/output/sorted_sales`

View the output: