

VISVESVARAYA TECHNOLOGICAL UNIVERSITY

“JnanaSangama”, Belgaum -590014, Karnataka.



LAB REPORT
on

BIG DATA ANALYTICS **(20CS6PEBDA)**

Submitted by

CHEETHANA D (1BM20CS405)

in partial fulfillment for the award of the degree of
BACHELOR OF ENGINEERING
in
COMPUTER SCIENCE AND ENGINEERING



B.M.S. COLLEGE OF ENGINEERING

(Autonomous Institution under VTU)

BENGALURU-560019

May-2022 to July-2022

B. M. S. College of Engineering,
Bull Temple Road, Bangalore 560019
(Affiliated To Visvesvaraya Technological University, Belgaum)
Department of Computer Science and Engineering



CERTIFICATE

This is to certify that the Lab work entitled “**BIG DATA ANALYTICS**” carried out by **CHETHANA D (1BM20CS405)**, who is bonafide student of **B. M. S. College of Engineering**. It is in partial fulfillment for the award of **Bachelor of Engineering in Computer Science and Engineering** of the Visvesvaraya Technological University, Belgaum during the year 2022. The Lab report has been approved as it satisfies the academic requirements in respect of a **BIG DATA ANALYTICS - (20CS6PEBDA)** work prescribed for the said degree.

Antara Roy Choudhary
Assistant Professor
Department of CSE
BMSCE, Bengaluru

Dr. Jyothi S Nayak
Professor and Head
Department of CSE
BMSCE, Bengaluru

Index Sheet

Sl. No.	Experiment Title	Page No.
1	MongoDB- CRUD Demonstration	4-11
2	Perform the following DB operations using Cassandra (Employee)	12-15
3	Perform the following DB operations using Cassandra (Library)	16-19
4	Execution of HDFS Commands for interaction with Hadoop Environment.	20-22
5	Screenshot of Hadoop Installed	23
6	Create a Map Reduce program to a) find average temperature for each year from NCDC data set. b) find the mean max temperature for every month	24-30
7	For a given Text file, Create a Map Reduce program to sort the content in an alphabetic order listing only top 10 maximum occurrences of words.	31-36
8	Create a Map Reduce program to demonstrating join operation	37-47
9	Program to print word count on scala shell and print “Hello world” on scala IDE	48
10	Using RDD and FlatMap count how many times each word appears in a file and write out a list of words whose count is strictly greater than 4 using Spark	49-51

Course Outcome

CO1	Apply the concept of NoSQL, Hadoop or Spark for a given task
CO2	Analyze the Big Data and obtain insight using data analytics mechanisms.
CO3	Design and implement Big data applications by applying NoSQL, Hadoop or Spark

```

root@kali:~# ssh root@10.10.10.10
Warning: Permanently added host 10.10.10.10 (SSH-2.0-PuTTY_Release_0.76)
root@mohamed-macbook-pro:~# mongo --quiet
MongoDB shell version v2.6.8
connecting to mongodb://107.8.4.1:27017
Implicit session: session ["id" = 905Dc875c7810-dbc7-4d73-b251-a8e88d18332"]
MongoDB server version: 2.6.8
Server has startup warnings:
2012-04-11T14:31:56.326+0000 I STORAGE [initandlisten]
2012-04-11T14:31:56.326+0000 I STORAGE [initandlisten] ** WARNING: Using the XFS filesystem is strongly recommended with the WiredTiger storage engine
2012-04-11T14:31:56.326+0000 I STORAGE [initandlisten] ** See http://dochub.mongodb.org/core/production-filesystem
2012-04-11T14:32:02.009+0000 I CONTROL [initandlisten]
2012-04-11T14:32:02.009+0000 I CONTROL [initandlisten] ** WARNING: Access control is not enabled for the database.
2012-04-11T14:32:02.009+0000 I CONTROL [initandlisten] ** Read and write access to data and configuration is unrestricted.
> use vmsh01db
switched to db vmsh01db
> db;
vmsh01db
> db.createCollection("Students");
{ "ok" : 1 }
> db.Student.insert({$date:"2013-04-28T00:00Z"}, {"Student": "Vanshi", "Grade": "X", "Hobbies":["Watching Movies", "Playing Games"]}];
writeResult({ "nInserted" : 1 })
> db.Student.find()
[ { "_id" : ObjectId("625f2354ab9c48287fec12f"), "Student" : "Vanshi", "Grade" : "X", "Hobbies" : [ "watching Movies", "Playing Games" ] } ]
> db.Student.update({'_id':"$date":"2013-04-28T00:00Z"}, {"Student": "Vanshi", "Grade": "X", "Hobbies":["Watching Movies", "Playing Games"]}];
2012-04-11T14:34:04.049+0000 E QUERY [thread1] error: need an object :
DBCollection.prototype._parseupdatearg/mongo/shell/collection.js:44:1
DBCollection.prototype.updatearg/mongo/shell/collection.js:48:18
gshell():1:1
> db.Student.find()
[ { "_id" : ObjectId("625f2354ab9c48287fec12f"), "Student" : "Vanshi", "Grade" : "X", "Hobbies" : [ "Watching Movies", "Playing Games" ] } ]
> db.Student.$eq()
true
> db.Student.find()
> db.Student.find()
> db.createCollection("Students");
{ "ok" : 1 }
> db.Student.update({'_id':"$date":"2013-04-28T00:00Z"}, {"Student": "Vanshi", "Grade": "X", "Hobbies":["Watching Movies", "Playing Games"]}];
2012-04-11T14:36:15.223+0000 E QUERY [thread1] error: need an object :
DBCollection.prototype._parseupdatearg/mongo/shell/collection.js:44:1
DBCollection.prototype.updatearg/mongo/shell/collection.js:48:18
gshell():1:1
> db.Student.find()
> db.Student.insert({'_id':"$date":"2013-04-28T00:00Z"}, {"Student": "Vanshi", "Grade": "X", "Hobbies":["Watching Movies", "Playing Games"]}];
writeResult({ "nInserted" : 1 })
> db.Student.find()
[ { "_id" : 1, "Student" : "Vanshi", "Grade" : "X", "Hobbies" : [ "Watching Movies", "Playing Games" ] } ]
> db.Student.update({'_id':"$date":"2013-04-28T00:00Z"}, {"Student": "Vanshi", "Grade": "X", "Hobbies":["Football"]}];
2012-04-11T14:38:11.959+0000 E QUERY [thread1] syntaxerror: Missing } after property list gshell():1:18
> db.Student.update({'_id':"$date":"2013-04-28T00:00Z"}, {"Student": "Vanshi", "Grade": "X", "Hobbies":["Football"]}];
writeResult({ "updated" : 1, "upserted" : ObjectId("625f2354ab9c48287fec12f") })
> db.Student.find()
[ { "_id" : 1, "Student" : "Vanshi", "Grade" : "X", "Hobbies" : "Football" } ]
> db.Student.find().pretty()
[ { "_id" : 1, "Student" : "Vanshi", "Grade" : "X", "Hobbies" : "Football" } ]
> db.Student.find()
[ { "_id" : 1, "Student" : "Vanshi", "Grade" : "X", "Hobbies" : "Football" } ]

```


2. Perform the following DB operations using Cassandra.

1. Create a keyspace by name Employee
2. Create a column family by name Employee-Info
with attributes
Emp_Id Primary Key, Emp_Name,
Designation, Date_of_Joining, Salary,
Dept_Name
3. Insert the values into the table in batch
4. Update Employee name and Department of Emp-Id 121
5. Sort the details of Employee records based on salary
6. Alter the schema of the table Employee_Info to add a column Projects which stores a set of Projects done by the corresponding Employee.
7. Update the altered table to add project names.
8. Create a TTL of 15 seconds to display the values of Employee

1. Create a key space by name Employee

```
cqlsh> describe keyspaces;

system_schema  system          system_distributed
system_auth     test_keyspace  system_traces

cqlsh> create keyspace Employee with replication={'class':'SimpleStrategy','replication_factor':2};
cqlsh> describe keyspace Employee;

CREATE KEYSPACE employee WITH replication = {'class': 'SimpleStrategy', 'replication_factor': '2'} AND durable_writes = true;

cqlsh> use employee;
```

2. Create a column family by name Employee-Info with attributes Emp_Id Primary Key, Emp_Name, Designation, Date_of_Joining, Salary, Dept_Name

```
cqlsh:employee> CREATE TABLE Employee_Info(cluster_col text,Emp_Id int,Emp_Name text,Designation text,Date_of_Joining timestamp,Salary int,Dept_Name text,primary key(cluster_col,Salary)) WITH CLUSTERING ORDER BY(Salary DESC);
cqlsh:employee> select*from employee;
InvalidRequest: Error from server: code=1200 [Invalid query] message="unconfigured table employee"
cqlsh:employee> select*from employee_info;
```

```
cluster_col | salary | date_of_joining | dept_name | designation | emp_id | emp_name
-----
```

3. Insert the values into the table in batch

```
(4 rows)
cqlsh:employee> BEGIN BATCH
... INSERT INTO Employee_Info(cluster_col,emp_id,emp_name,designation,Date_of_Joining,Salary,Dept_Name) VALUES ('xyz',1,'Ravi','MTS','2020-05-24',12000,'TESTING');
... INSERT INTO Employee_Info(cluster_col,emp_id,emp_name,designation,Date_of_Joining,Salary,Dept_Name) VALUES ('xyz',2,'Vanshi','MANAGER','2021-05-28',20000,'DEVELOPMENT');
... INSERT INTO Employee_Info(cluster_col,emp_id,emp_name,designation,Date_of_Joining,Salary,Dept_Name) VALUES ('xyz',123,'Kiran','SDE','2019-04-21',10000,'PRODUCTION');
... INSERT INTO Employee_Info(cluster_col,emp_id,emp_name,designation,Date_of_Joining,Salary,Dept_Name) VALUES ('xyz',3,'Ramesh','ANALYST','2020-05-07',18000,'QUALITY');
... APPLY BATCH;
cqlsh:employee> SELECT*FROM Employee_Info;

cluster_col | salary | date_of_joining | dept_name | designation | emp_id | emp_name
-----
xyz | 20000 | 2021-05-28 18:38:00.000000-0800 | DEVELOPMENT | MANAGER | 2 | Vanshi
xyz | 18000 | 2020-05-07 18:38:00.000000-0800 | QUALITY | ANALYST | 3 | Ramesh
xyz | 12000 | 2019-04-21 18:38:00.000000-0800 | PRODUCTION | SDE | 123 | Kiran
xyz | 10000 | 2020-05-24 18:38:00.000000-0800 | TESTING | MTS | 1 | Ravi
(4 rows)
```

4. Update Employee name and Department of Emp-id 121

```
sqlsh:employee> update Employee_Info SET emp_name='karthik',dept_name='Compliance' where cluster_col='xyz' and salary=10000 IF emp_id=121;
```

```
[applied]
```

```
True
```

```
sqlsh:employee> SELECT*FROM Employee_Info;
```

cluster_col	salary	date_of_joining	dept_name	designation	emp_id	emp_name	projects
xyz	50000	2021-03-19 18:30:00.000000+0000	DEVELOPEMENT	MANAGER	2	Vanshi	['AI', 'DS']
xyz	20000	2020-05-06 18:30:00.000000+0000	QUALITY	ANALYST	3	Ramesh	['DEVOPS']
xyz	12000	2020-08-23 18:30:00.000000+0000	TESTING	MTS	1	Ravi	['%']
xyz	10000	2019-04-20 18:30:00.000000+0000	Compliance	SDE	121	karthik	['QUANTUM COMPUTING']

5. Sort the details of Employee records based on salary

```
sqlsh:employee> SELECT*FROM Employee_Info;
```

cluster_col	salary	date_of_joining	dept_name	designation	emp_id	emp_name	projects
xyz	50000	2021-03-19 18:30:00.000000+0000	DEVELOPEMENT	MANAGER	2	Vanshi	null
xyz	20000	2020-05-06 18:30:00.000000+0000	QUALITY	ANALYST	3	Ramesh	null
xyz	12000	2020-08-23 18:30:00.000000+0000	TESTING	MTS	1	Ravi	null
xyz	10000	2019-04-20 18:30:00.000000+0000	Finance	SDE	121	Signesh	null

6. Alter the schema of the table Employee_Info to add a column Projects which stores a set of Projects done by the corresponding Employee.

```
sqlsh:employee> alter table Employee_Info add Projects set(text);
```

```
sqlsh:employee> SELECT*FROM Employee_Info;
```

cluster_col	salary	date_of_joining	dept_name	designation	emp_id	emp_name	projects
xyz	50000	2021-03-19 18:30:00.000000+0000	DEVELOPEMENT	MANAGER	2	Vanshi	null
xyz	20000	2020-05-06 18:30:00.000000+0000	QUALITY	ANALYST	3	Ramesh	null
xyz	12000	2020-08-23 18:30:00.000000+0000	TESTING	MTS	1	Ravi	null
xyz	10000	2019-04-20 18:30:00.000000+0000	Finance	SDE	121	Signesh	null

```
(4 rows)
```

- Update the altered table to add project names.

```

sqlsh>employee> update Employee_Info SET projects=projects('M') where cluster_col='xyz' and salary<10000 IF emp_id=1;
[applied]
-----
True

sqlsh>employee> update Employee_Info SET projects=projects('AI','DS') where cluster_col='xyz' and salary<50000 IF emp_id=2;
[applied]
-----
True

sqlsh>employee> update Employee_Info SET projects=projects('DEVOPS') where cluster_col='xyz' and salary<20000 IF emp_id=3;
[applied]
-----
True

sqlsh>employee> update Employee_Info SET projects=projects('QUANTUM COMPUTING') where cluster_col='xyz' and salary<10000 IF emp_id=121;
[applied]
-----
True

sqlsh>employee> SELECT*FROM Employee_Info;
cluster_col | salary | date_of_joining | dept_name | designation | emp_id | emp_name | projects
-----
xyz | 50000 | 2021-01-19 18:30:00.000000+0000 | DEVELOPMENT | MANAGER | 2 | Vanshi | ('AI', 'DS')
xyz | 20000 | 2020-05-05 18:30:00.000000+0000 | QUALITY | ANALYST | 3 | Ramesh | ('DEVOPS')
xyz | 12000 | 2020-04-23 18:30:00.000000+0000 | TESTING | HTS | 1 | Ravi | ('M')
xyz | 10000 | 2019-04-20 18:30:00.000000+0000 | Finance | SDE | 121 | Jignesh | ('QUANTUM COMPUTING')

```

- Create a TTL of 15 seconds to display the values of Employees.

```

(4 row)
sqlsh>employee> INSERT INTO Employee_Info(cluster_col,emp_id,emp_name,designation,date_of_joining,salary,dept_name) VALUES ('xyz',121,'Rodi','SDE','2022-04-
TIDM') using TTL 15;
sqlsh>employee> SELECT*FROM Employee_Info;

cluster_col | salary | date_of_joining | dept_name | designation | emp_id | emp_name | projects
-----
xyz | 50000 | 2021-01-19 18:30:00.000000+0000 | DEVELOPMENT | MANAGER | 2 | Vanshi | ('AI', 'DS')
xyz | 20000 | 2020-05-05 18:30:00.000000+0000 | QUALITY | ANALYST | 3 | Ramesh | ('DEVOPS')
xyz | 12000 | 2020-04-23 18:30:00.000000+0000 | TESTING | HTS | 1 | Ravi | ('M')
xyz | 10000 | 2019-04-20 18:30:00.000000+0000 | Finance | SDE | 121 | Jignesh | ('QUANTUM COMPUTING')
xyz | 1000 | 2022-04-20 18:30:00.000000+0000 | PRODUCTION | SDE | 121 | Rodi | null

(5 row)
sqlsh>employee> SELECT*FROM Employee_Info;

cluster_col | salary | date_of_joining | dept_name | designation | emp_id | emp_name | projects
-----
xyz | 50000 | 2021-01-19 18:30:00.000000+0000 | DEVELOPMENT | MANAGER | 2 | Vanshi | ('AI', 'DS')
xyz | 20000 | 2020-05-05 18:30:00.000000+0000 | QUALITY | ANALYST | 3 | Ramesh | ('DEVOPS')
xyz | 12000 | 2020-04-23 18:30:00.000000+0000 | TESTING | HTS | 1 | Ravi | ('M')
xyz | 10000 | 2019-04-20 18:30:00.000000+0000 | Finance | SDE | 121 | Jignesh | ('QUANTUM COMPUTING')

(4 row)

```


3. Perform the following DB operations using Cassandra.

1. Create a keyspace by name Library
2. Create a column family by name Library-Info with attributes
Stud_Id Primary Key,
Counter_value of type Counter,
Stud_Name, Book-Name, Book-Id,
Date_of_issue
3. Insert the values into the table in batch
4. Display the details of the table created and increase the value of the counter
5. Write a query to show that a student with id 112 has taken a book “BDA” 2 times.
6. Export the created column to a csv file
7. Import a given csv dataset from local file system into Cassandra column family

1. Create a key space by name Library

```
cqlsh> describe keyspaces;

system_schema  system          system_distributed  system_traces
system_auth    test_keyspace  employee

cqlsh> CREATE KEYSPACE library WITH replication={'class':'SimpleStrategy','replication_factor':'1'};
cqlsh> describe keyspace library;

CREATE KEYSPACE library WITH replication = {'class': 'SimpleStrategy', 'replication_factor': '1'} AND durable_writes = true;

cqlsh> use library;
```

2. Create a column family by name Library-Info with attributes Stud_Id Primary Key, Counter_value of type Counter, Stud_Name, Book-Name, Book-Id, Date_of_issue

```
cqlsh> use library;
cqlsh:library> create table library_info (Stud_Id int,counter_value counter,Stud_Name text,Book_Name text,Book_Id int,Date_of_issue timestamp,PRIMARY KEY(Stud_Id,Stud_Name,Book_Name,Book_Id,Date_of_issue));
cqlsh:library> describe table library_info;

CREATE TABLE library.library_info (
  stud_id int,
  stud_name text,
  book_name text,
  book_id int,
  date_of_issue timestamp,
  counter_value counter,
  PRIMARY KEY (stud_id, stud_name, book_name, book_id, date_of_issue)
) WITH CLUSTERING ORDER BY (stud_name ASC, book_name ASC, book_id ASC, date_of_issue ASC)
AND bloom_filter_fp_chance = 0.01
AND caching = {'keys': 'ALL', 'rows_per_partition': 'NONE'}
AND comment = ''
AND compaction = {'class': 'org.apache.cassandra.db.compaction.SizeTieredCompactionStrategy', 'max_threshold': '32', 'min_threshold': '4'}
AND compression = {'chunk_length_in_kb': '64', 'class': 'org.apache.cassandra.io.compress.LZ4Compressor'}
AND crc_check_chance = 1.0
AND dclocal_read_repair_chance = 0.1
AND default_time_to_live = 0
AND gc_grace_seconds = 8640000
AND max_index_interval = 1000
AND memtable_flush_period_in_ms = 0
AND min_index_interval = 128
AND read_repair_chance = 0.0
AND speculative_retry = 'WHENEXHAUSTED';
```

3. Insert the values into the table in batch

4. Display the details of the table created and increase the value of the counter

```
cdsh:library> update library_info set counter_value=counter_value+1 where stud_id=1 and stud_name='Vamsi' and book_name='OVD' and book_id=100 and date_of_issue='2022-03-13';
cdsh:library> update library_info set counter_value=counter_value+1 where stud_id=2 and stud_name='Ravi' and book_name='OIS' and book_id=101 and date_of_issue='2022-03-13';
cdsh:library> update library_info set counter_value=counter_value+1 where stud_id=112 and stud_name='Ramesh' and book_name='BDA' and book_id=102 and date_of_issue='2022-03-13';
cdsh:library> select*from library_info;
```

stud_id	stud_name	book_name	book_id	date_of_issue	counter_value
1	Vamsi	OVD	100	2022-04-17 18:38:00.000000+0000	1
2	Ravi	OIS	101	2022-03-14 18:38:00.000000+0000	1
112	Ramesh	BDA	102	2022-03-13 18:38:00.000000+0000	1

(3 rows)

5. Write a query to show that a student with id 112 has taken a book "BDA" 2 times.

```
cdsh:library> update library_info set counter_value=counter_value+1 where stud_id=112 and stud_name='Ramesh' and book_name='BDA' and book_id=102 and date_of_issue='2022-03-13';
cdsh:library> select*from library_info;
```

stud_id	stud_name	book_name	book_id	date_of_issue	counter_value
1	Vamsi	OVD	100	2022-04-17 18:38:00.000000+0000	1
2	Ravi	OIS	101	2022-03-14 18:38:00.000000+0000	1
112	Ramesh	BDA	102	2022-03-13 18:38:00.000000+0000	2

(3 rows)

6. Export the created column to a csv file

```
cdsh:library> copy library_info (stud_id,stud_name,book_name,book_id,date_of_issue,counter_value) to 'C:/Users/shari/OneDrive/Documents/library_info.csv' with header=true using 7 child processes
```

Starting copy of library.library_info with columns [stud_id, stud_name, book_name, book_id, date_of_issue, counter_value].
Processed: 3 rows; Rate: 2 row/s; Avg. rate: 1 row/s
3 rows exported to 1 files in 3.164 seconds.

	A	B	C	D	E	F	G	H
1	stud_id	stud_name	book_name	book_id	date_of_issue	counter_value		
2	112	Ramesh	BDA	102	2022-03-1	2		
3	1	Vamshi	OOMD	100	2022-04-1	1		
4	2	Ravi	CNS	101	2022-03-1	1		
5								

7. Import a given csv dataset from local file system into Cassandra column family

```
cqlsh:library> create table library_test(stud_id int,counter_value counter,stud_name text,book_name text,book_id int,date_of_issue timestamp,primary key(stud_id,stud_name,book_name,book_id,date_of_issue));
```

```
cqlsh:library> COPY library_test(stud_id,stud_name,book_name,book_id,date_of_issue,counter_value) FROM 'C:\Users\shana\OneDrive\Documents\Library_Info.csv' with header=true;
```

Using 7 child processes

Starting copy of library.library_test with columns (stud_id, stud_name, book_name, book_id, date_of_issue, counter_value).

Process ImportProcess-0: 1 rows/s; Avg. rate: 1 rows/s

```
AttributeError: 'NoneType' object has no attribute 'add_timer'
Processed: 3 rows; Rate: 1 rows/s; Avg. rate: 0 rows/s
3 rows imported from 1 files in 6.268 seconds (0 skipped).
```

```
cqlsh:library> select*from library_test;
```

stud_id	stud_name	book_name	book_id	date_of_issue	counter_value
1	Vamshi	OOMD	100	2022-04-17 18:30:00.000000+0000	1
2	Ravi	CNS	101	2022-03-14 18:30:00.000000+0000	1
112	Ramesh	BDA	102	2022-03-18 18:30:00.000000+0000	2

(3 rows)

```
cqlsh:library>
```

4. Execution of HDFS Commands for interaction with Hadoop Environment.

1. Successful installation proof

```
hadoop@sharat-VirtualBox:~/Desktop$ hadoop
Usage: hadoop [OPTIONS] SUBCOMMAND [SUBCOMMAND OPTIONS]
or:   hadoop [OPTIONS] CLASSNAME [CLASSNAME OPTIONS]
      where CLASSNAME is a user-provided Java class

      OPTIONS is none or any of:
      buildpaths      attempt to add class files from build tree
      --config-dir    Hadoop config directory
      --debug         turn on shell script debug mode
      --help          usage information
      hostnames list[,of,host,names] hosts to use in slave mode
      hosts filename  list of hosts to use in slave mode
      loglevel level  set the log4j level for this command
      workers         turn on worker mode

      SUBCOMMAND is one of:

      Admin Commands:
      daemonlog       get/set the log level for each daemon

      Client Commands:
```

```
hadoop@sharat-VirtualBox:~/hadoop-3.2.3$ jps
5809 NodeManager
5170 DataNode
5650 ResourceManager
5015 NameNode
5415 SecondaryNameNode
6442 Jps
```

2. Mkdir

```
hadoop@sharat-VirtualBox:~/hadoop-3.2.3$ hdfs dfs -mkdir /lab5
2022-06-08 09:57:10,719 WARN util.NativeCodeLoader: Unable to load native-hadoop
p library for your platform... using builtin-java classes where applicable
```

3. Is

```
hadoop@sharat-VirtualBox:~/hadoop-3.2.3$ hadoop fs -ls /
2022-06-08 09:57:33,445 WARN util.NativeCodeLoader: Unable to load native-hadoop
p library for your platform... using builtin-java classes where applicable
Found 2 items
-rw-r--r-- 1 hadoop supergroup 3745 2022-06-06 23:56 /Hadoop_Installatio
n_Commands.txt
drwxr-xr-x - hadoop supergroup 0 2022-06-08 09:57 /lab5
```

4. put

```
hadoop@sharat-VirtualBox:~/hadoop-3.2.3$ hdfs dfs -put /home/hadoop/Desktop/a.txt /lab5
2022-06-08 10:02:57,394 WARN util.NativeCodeLoader: Unable to load native-hadoop
p library for your platform... using builtin-java classes where applicable
hadoop@sharat-VirtualBox:~/hadoop-3.2.3$ hdfs fs -ls /
2022-06-08 10:03:08,676 WARN util.NativeCodeLoader: Unable to load native-hadoop
p library for your platform... using builtin-java classes where applicable
Found 2 items
-rw-r--r-- 1 hadoop supergroup 3745 2022-06-06 23:56 /Hadoop_Installatio
n_Commands.txt
drwxr-xr-x - hadoop supergroup 0 2022-06-08 10:02 /lab5
```

5. copyFromLocal

```
hadoop@sharat-VirtualBox:~/hadoop-3.2.3$ hdfs dfs -copyFromLocal /home/hadoop/Des
ktop/b.txt /home/hadoop/Desktop/c.txt /lab5
2022-06-08 10:07:13,375 WARN util.NativeCodeLoader: Unable to load native-hadoop
p library for your platform... using builtin-java classes where applicable
hadoop@sharat-VirtualBox:~/hadoop-3.2.3$ hdfs dfs -cat /lab52022-06-08 10:07:36,
122 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your p
latform... using builtin-java classes where applicable
cat: '/lab5': Is a directory
hadoop@sharat-VirtualBox:~/hadoop-3.2.3$ hdfs dfs -ls /lab5
2022-06-08 10:08:26,214 WARN util.NativeCodeLoader: Unable to load native-hadoop
p library for your platform... using builtin-java classes where applicable
Found 3 items
-rw-r--r-- 1 hadoop supergroup 15 2022-06-08 10:02 /lab5/a.txt
-rw-r--r-- 1 hadoop supergroup 0 2022-06-08 10:07 /lab5/b.txt
-rw-r--r-- 1 hadoop supergroup 0 2022-06-08 10:07 /lab5/c.txt
```

6. Get

i. get

```
hadoop@sharat-VirtualBox:~/hadoop-3.2.3$ hdfs dfs -get /lab5/a.txt /home/hadoop/D
esktop/test.txt
2022-06-08 10:10:29,649 WARN util.NativeCodeLoader: Unable to load native-hadoop
p library for your platform... using builtin-java classes where applicable
hadoop@sharat-VirtualBox:~/hadoop-3.2.3$ hdfs dfs -ls /home/hadoop/Desktop
2022-06-08 10:11:08,053 WARN util.NativeCodeLoader: Unable to load native-hadoop
p library for your platform... using builtin-java classes where applicable
ls: '/home/hadoop/Desktop': No such file or directory
hadoop@sharat-VirtualBox:~/hadoop-3.2.3$ ls /home/hadoop/Desktop
a.txt b.txt c.txt test.txt
```

ii. getmerge

```
hadoop@sharat-VirtualBox:~/hadoop-3.2.3$ hdfs dfs -getmerge /lab5/b.txt /lab5/c.
txt /home/hadoop/Desktop/merge.txt
2022-06-08 10:15:24,732 WARN util.NativeCodeLoader: Unable to load native-hadoop
p library for your platform... using builtin-java classes where applicable
hadoop@sharat-VirtualBox:~/hadoop-3.2.3$ ls /home/hadoop/Desktop
a.txt b.txt c.txt merge.txt test.txt
hadoop@sharat-VirtualBox:~/hadoop-3.2.3$ cat /home/hadoop/Desktop/merge.txt
```

iii. getfacl

```
hadoop@sharat-VirtualBox:~/hadoop-3.2.3$ hdfs dfs -getfacl /lab5
2022-06-08 10:17:20,801 WARN util.NativeCodeLoader: Unable to load native-hadoop
p library for your platform... using builtin-java classes where applicable
# file: /lab5
# owner: hadoop
# group: supergroup
user::rwx
group::r-x
other::r-x
```

7. copyToLocal

```
hadoop@sharat-VirtualBox:~/hadoop-3.2.3$ hdfs dfs -copyToLocal /lab5/a.txt /home/hadoop/Documents
2022-06-08 10:19:08,660 WARN util.NativeCodeLoader: Unable to load native-hadoop
p library for your platform... using builtin-java classes where applicable
hadoop@sharat-VirtualBox:~/hadoop-3.2.3$ ls /home/hadoop/Documents
a.txt
```

8.cat

```
hadoop@sharat-VirtualBox:~/hadoop-3.2.3$ hdfs dfs -cat /lab5/a.txt
2022-06-08 10:19:48,077 WARN util.NativeCodeLoader: Unable to load native-hadoop
p library for your platform... using builtin-java classes where applicable
this is a test
```

9.mv

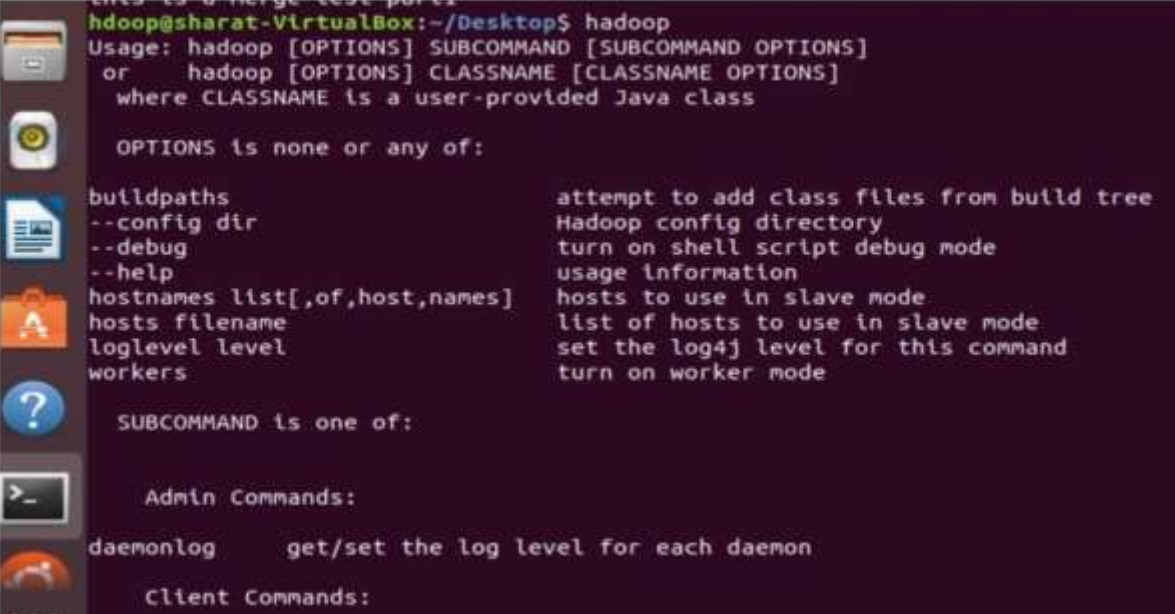
```
hadoop@sharat-VirtualBox:~/hadoop-3.2.3$ hadoop fs -mv /lab5/a.txt /lab5_part2
2022-06-08 10:22:18,644 WARN util.NativeCodeLoader: Unable to load native-hadoop
p library for your platform... using builtin-java classes where applicable
```

10.cp

```
hadoop@sharat-VirtualBox:~/hadoop-3.2.3$ hadoop fs -cp /lab5/b.txt /lab5_part2
2022-06-08 10:23:16,644 WARN util.NativeCodeLoader: Unable to load native-hadoop
p library for your platform... using builtin-java classes where applicable
hadoop@sharat-VirtualBox:~/hadoop-3.2.3$ hadoop fs -ls /lab5_part2
2022-06-08 10:23:21,944 WARN util.NativeCodeLoader: Unable to load native-hadoop
p library for your platform... using builtin-java classes where applicable
Found 2 items
-rw-r--r-- 1 hadoop supergroup      15 2022-06-08 10:02 /lab5_part2/a.txt
-rw-r--r-- 1 hadoop supergroup       0 2022-06-08 10:23 /lab5_part2/b.txt
```


5. Screenshot of Hadoop installed

1. Successful installation proof



```
hadoop@sharat-VirtualBox:~/Desktop$ hadoop
Usage: hadoop [OPTIONS] SUBCOMMAND [SUBCOMMAND OPTIONS]
or      hadoop [OPTIONS] CLASSNAME [CLASSNAME OPTIONS]
       where CLASSNAME is a user-provided Java class

      OPTIONS is none or any of:

buildpaths          attempt to add class files from build tree
--config dir        Hadoop config directory
--debug            turn on shell script debug mode
--help            usage information
hostnames list[,of,host,names] hosts to use in slave mode
hosts filename      list of hosts to use in slave mode
loglevel level      set the log4j level for this command
workers            turn on worker mode

      SUBCOMMAND is one of:

      Admin Commands:

daemonlog          get/set the log level for each daemon

      Client Commands:
```

6. Create a Map Reduce program to

a) find average temperature for each year from NCDC data set.

b) find the mean max temperature for every month

a)

CODE:

AverageDriver

```
package temp;

import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;

public class AverageDriver {
    public static void main(String[] args) throws Exception {
        if (args.length != 2) {
            System.err.println("Please Enter the input and output parameters");
            System.exit(-1);
        }
        Job job = new Job();
        job.setJarByClass(AverageDriver.class);
        job.setJobName("Max temperature");
        FileInputFormat.addInputPath(job, new Path(args[0]));
        FileOutputFormat.setOutputPath(job, new Path(args[1]));
        job.setMapperClass(AverageMapper.class);
        job.setReducerClass(AverageReducer.class);
        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(IntWritable.class);
        System.exit(job.waitForCompletion(true) ? 0 : 1);
    }
}
```

AverageMapper

```
package temp;
```

```
import java.io.IOException;
```

```
import org.apache.hadoop.io.IntWritable;
```

```
import org.apache.hadoop.io.LongWritable;
```

```
import org.apache.hadoop.io.Text;
```

```
import org.apache.hadoop.mapreduce.Mapper;
```

```
public class AverageMapper extends Mapper<LongWritable,  
Text, Text, IntWritable> {
```

```
    public static final int MISSING = 9999;
```

```
    public void map(LongWritable key, Text value,  
Mapper<LongWritable, Text, Text, IntWritable>.Context  
context) throws IOException, InterruptedException {
```

```
        int temperature;
```

```
        String line = value.toString();
```

```
        String year = line.substring(15, 19);
```

```
        if (line.charAt(87) == '+') {
```

```
            temperature = Integer.parseInt(line.substring(88,  
92));
```

```
        } else {
```

```
            temperature = Integer.parseInt(line.substring(87,  
92));
```

```
        }
```

```
        String quality = line.substring(92, 93);
```

```
        if (temperature != 9999 && quality.matches("[01459]"))
```

```
            context.write(new Text(year), new  
IntWritable(temperature));
```

```
    }
```

```
}
```

AverageReducer

```
package temp;
```

```
import java.io.IOException;
```

```
import org.apache.hadoop.io.IntWritable;
```

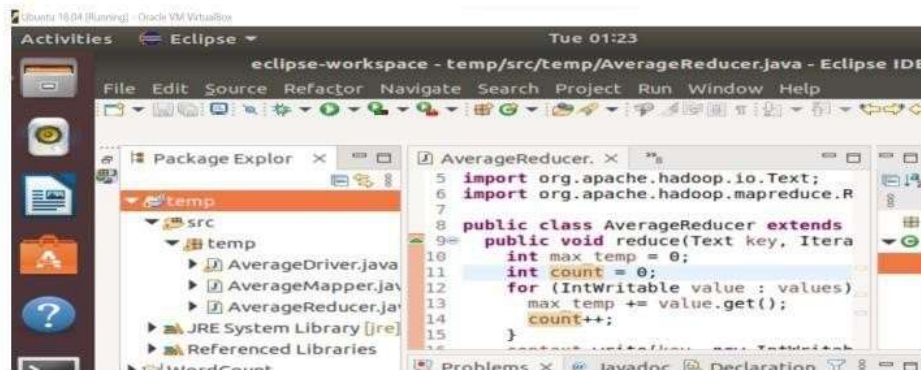
```

import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Reducer;

public class AverageReducer extends Reducer<Text,
IntWritable, Text, IntWritable> {
    public void reduce(Text key, Iterable<IntWritable>
values, Reducer<Text, IntWritable, Text,
IntWritable>.Context context) throws IOException,
InterruptedException {
        int max_temp = 0;
        int count = 0;
        for (IntWritable value : values) {
            max_temp += value.get();
            count++;
        }
        context.write(key, new IntWritable(max_temp / count));
    }
}

```

OUTPUT:



```

put: /home/hadoop/Desktop/1901.txt: No such file or directory
hadoop@sharat-VirtualBox:~$ hdfs dfs -put /home/hadoop/Desktop/1901 /inputt
2022-06-28 01:12:47,278 WARN util.NativeCodeLoader: Unable to load native-hadoop
p library for your platform... using builtin-java classes where applicable
hadoop@sharat-VirtualBox:~$ hdfs dfs -ls /inputt
2022-06-28 01:13:05,646 WARN util.NativeCodeLoader: Unable to load native-hadoop
p library for your platform... using builtin-java classes where applicable
Found 4 items
-rw-r--r-- 1 hadoop supergroup      888190 2022-06-28 01:12 /inputt/1901
-rw-r--r-- 1 hadoop supergroup       15 2022-06-20 16:51 /inputt/a.txt
-rw-r--r-- 1 hadoop supergroup       38 2022-06-27 22:01 /inputt/b.txt
drwxr-xr-x - hadoop supergroup       0 2022-06-20 16:52 /inputt/output

```

```

hadoop@sharat-VirtualBox:~$ hadoop jar weathertwo.jar temp.AverageDriver /inputt
/1901 /inputt/outputweather
2022-06-28 01:21:32,366 WARN util.NativeCodeLoader: Unable to load native-hadoo
p library for your platform... using builtin-java classes where applicable
2022-06-28 01:21:33,696 INFO client.RMProxy: Connecting to ResourceManager at /
127.0.0.1:8032
2022-06-28 01:21:34,100 WARN mapreduce.JobResourceUploader: Hadoop command-line
option parsing not performed. Implement the Tool interface and execute your ap
plication with ToolRunner to remedy this.
2022-06-28 01:21:34,131 INFO mapreduce.JobResourceUploader: Disabling Erasure C
oding for path: /tmp/hadoop-yarn/staging/hadoop/.staging/job_1656358828291_0001
2022-06-28 01:21:35,309 INFO input.FileInputFormat: Total input files to proces
s : 1
2022-06-28 01:21:35,410 INFO mapreduce.JobSubmitter: number of splits:1
2022-06-28 01:21:35,589 INFO mapreduce.JobSubmitter: Submitting tokens for job:
job_1656358828291_0001
2022-06-28 01:21:35,590 INFO mapreduce.JobSubmitter: Executing with tokens: []
2022-06-28 01:21:36,346 INFO conf.Configuration: resource-types.xml not found
2022-06-28 01:21:36,346 INFO resource.ResourceUtils: Unable to find 'resource-t
ypes.xml'.
2022-06-28 01:21:37,378 INFO impl.YarnClientImpl: Submitted application applica
tion_1656358828291_0001
2022-06-28 01:21:38,336 INFO mapreduce.Job: The url to track the job: http://sh
arat-VirtualBox:8088/proxy/application_1656358828291_0001/
2022-06-28 01:21:38,338 INFO mapreduce.Job: Running job: job_1656358828291_0001
2022-06-28 01:21:48,759 INFO mapreduce.Job: Job job_1656358828291_0001 running
in uber mode : false
2022-06-28 01:21:48,760 INFO mapreduce.Job: map 0% reduce 0%

```

```

Reduce input groups=1
Reduce shuffle bytes=72210
Reduce input records=6564
Reduce output records=1
Spilled Records=13128
Shuffled Maps =1
Failed Shuffles=0
Merged Map outputs=1
GC time elapsed (ms)=754
CPU time spent (ms)=1840
Physical memory (bytes) snapshot=645009408
Virtual memory (bytes) snapshot=5166370816
Total committed heap usage (bytes)=658505728
Peak Map Physical memory (bytes)=450666496
Peak Map Virtual memory (bytes)=2579943424
Peak Reduce Physical memory (bytes)=194342912

```

```

Bytes Written=8
hadoop@sharat-VirtualBox:~$ hdfs dfs -ls /inputt/outputweather
2022-06-28 01:22:16,506 WARN util.NativeCodeLoader: Unable to load native-hadoo
p library for your platform... using builtin-java classes where applicable
Found 2 items
-rw-r--r-- 1 hadoop supergroup 0 2022-06-28 01:21 /inputt/outputweath
er/_SUCCESS
-rw-r--r-- 1 hadoop supergroup 8 2022-06-28 01:21 /inputt/outputweath
er/part-r-000000

```

```

hadoop@sharat-VirtualBox:~$ hdfs dfs -cat /inputt/outputweather/part-r-000000
2022-06-28 01:23:07,585 WARN util.NativeCodeLoader: Unable to load native-hadoo
p library for your platform... using builtin-java classes where applicable
1901 46

```

b)

CODE:

MeanMaxDriver.class

```
package meanmax;
```

```
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import
org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import
org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;

public class MeanMaxDriver {
    public static void main(String[] args) throws Exception {
        if (args.length != 2) {
            System.err.println("Please Enter the input and output
parameters");
            System.exit(-1);
        }
        Job job = new Job();
        job.setJarByClass(MeanMaxDriver.class);
        job.setJobName("Max temperature");
        FileInputFormat.addInputPath(job, new Path(args[0]));
        FileOutputFormat.setOutputPath(job, new Path(args[1]));
        job.setMapperClass(MeanMaxMapper.class);
        job.setReducerClass(MeanMaxReducer.class);
        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(IntWritable.class);
        System.exit(job.waitForCompletion(true) ? 0 : 1);
    }
}
```

MeanMaxMapper.class

```
package meanmax;
```

```
import java.io.IOException;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
```



```

import org.apache.hadoop.mapreduce.Mapper;

public class MeanMaxMapper extends Mapper<LongWritable,
Text, Text, IntWritable> {
    public static final int MISSING = 9999;

    public void map(LongWritable key, Text value,
Mapper<LongWritable, Text, Text, IntWritable>.Context
context) throws IOException, InterruptedException {
        int temperature;
        String line = value.toString();
        String month = line.substring(19, 21);
        if (line.charAt(87) == '+') {
            temperature = Integer.parseInt(line.substring(88,
92));
        } else {
            temperature = Integer.parseInt(line.substring(87,
92));
        }
        String quality = line.substring(92, 93);
        if (temperature != 9999 && quality.matches("[01459]"))
            context.write(new Text(month), new
IntWritable(temperature));
    }
}

```

MeanMaxReducer.class

```

package meanmax;

import java.io.IOException;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Reducer;

public class MeanMaxReducer extends Reducer<Text,
IntWritable, Text, IntWritable> {
    public void reduce(Text key, Iterable<IntWritable>
values, Reducer<Text, IntWritable, Text,
IntWritable>.Context context) throws IOException,
InterruptedException {
        int max_temp = 0;
        int total_temp = 0;
    }
}

```

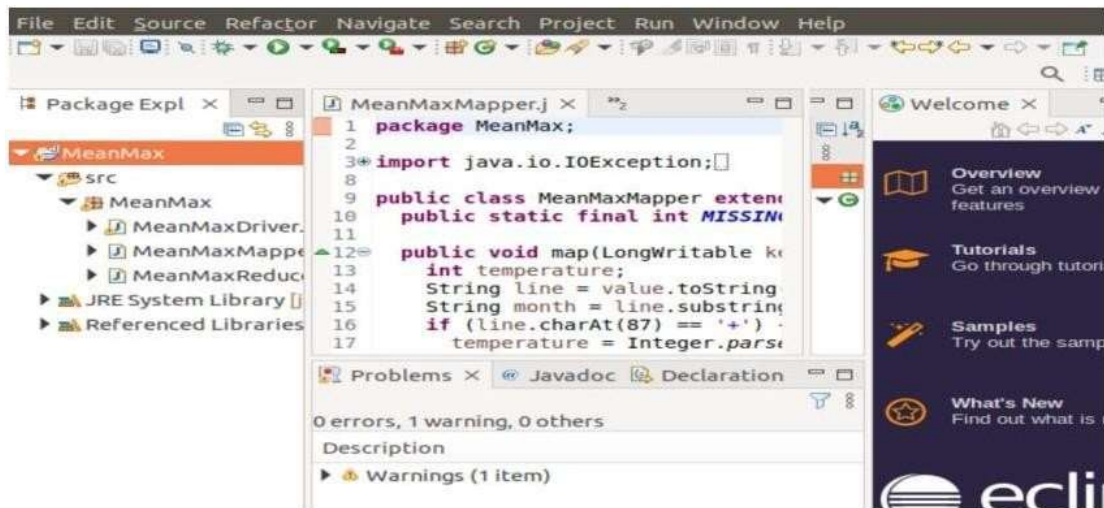


```

int count = 0;
int days = 0;
for (IntWritable value : values) {
    int temp = value.get();
    if (temp > max_temp)
        max_temp = temp;
    count++;
    if (count == 3) {
        total_temp += max_temp;
        max_temp = 0;
        count = 0;
        days++;
    }
}
context.write(key, new IntWritable(total_temp / days));
}
}

```

OUTPUT:



```

2022-06-28 02:35:15,863 WARN util.NativeCodeLoader: Unable to load native-hadoop
hadoop@sharat-VirtualBox:~$ hadoop jar MeanMaxweather2.jar MeanMax.MeanMaxDriver
/inputt/1901 /inputt/outputmeanmax
p library for your platform... using builtin-java classes where applicable
2022-06-28 02:35:16,403 INFO client.RMProxy: Connecting to ResourceManager at /
127.0.0.1:8032
2022-06-28 02:35:16,741 WARN mapreduce.JobResourceUploader: Hadoop command-line
option parsing not performed. Implement the Tool interface and execute your ap
plication with ToolRunner to remedy this.
2022-06-28 02:35:16,774 INFO mapreduce.JobResourceUploader: Disabling Erasure C
oding for path: /tmp/hadoop-yarn/staging/hadoop/.staging/job_1656363425892_0001
2022-06-28 02:35:17,464 INFO input.FileInputFormat: Total input files to proces
s : 1
2022-06-28 02:35:17,959 INFO mapreduce.JobSubmitter: number of splits:1
2022-06-28 02:35:18,176 INFO mapreduce.JobSubmitter: Submitting tokens for job:
job_1656363425892_0001
2022-06-28 02:35:18,177 INFO mapreduce.JobSubmitter: Executing with tokens: []
2022-06-28 02:35:18,417 INFO conf.Configuration: resource-types.xml not found
2022-06-28 02:35:18,418 INFO resource.ResourceUtils: Unable to find 'resource-t
ypes.xml'.
2022-06-28 02:35:18,932 INFO impl.YarnClientImpl: Submitted application applica
tion 1656363425892_0001

```

```

hadoop@sharat-VirtualBox:~$ hdfs dfs -ls /inputt/outputmeanmax
2022-06-28 02:36:40,638 WARN util.NativeCodeLoader: Unable to loa
p library for your platform... using builtin-java classes where a
Found 2 items
-rw-r--r-- 1 hadoop supergroup 0 2022-06-28 02:35 /inputt/outputmeanmax/_SUCCESS
-rw-r--r-- 1 hadoop supergroup 74 2022-06-28 02:35 /inputt/outputmeanmax/part-r-000000
hadoop@sharat-VirtualBox:~$ hdfs dfs -cat /inputt/outputmeanmax/part-r-000000
2022-06-28 02:36:57,109 WARN util.NativeCodeLoader: Unable to loa
p library for your platform... using builtin-java classes where a
01 4
02 0
03 7
04 44
05 100
06 168
07 219
08 198
09 141
10 100
11 19
12 3

```

7. For a given Text file, Create a Map Reduce program to sort the content in an alphabetic order listing only top 10 maximum occurrences of words.

CODE:

Driver-TopN.class

```
package samples.topn;

import java.io.IOException;
import java.util.StringTokenizer;
import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.Mapper;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
import org.apache.hadoop.util.GenericOptionsParser;

public class TopN {
    public static void main(String[] args) throws Exception {
        Configuration conf = new Configuration();
        String[] otherArgs = (new GenericOptionsParser(conf,
args)).getRemainingArgs();
        if (otherArgs.length != 2) {
            System.err.println("Usage: TopN <in> <out>");
            System.exit(2);
        }
        Job job = Job.getInstance(conf);
        job.setJobName("Top N");
        job.setJarByClass(TopN.class);
        job.setMapperClass(TopNMapper.class);
        job.setReducerClass(TopNReducer.class);
        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(IntWritable.class);
        FileInputFormat.addInputPath(job, new
Path(otherArgs[0]));
```

```

        FileOutputFormat.setOutputPath(job, new
Path(otherArgs[1]));
        System.exit(job.waitForCompletion(true) ? 0 : 1);
    }

    public static class TopNMapper extends Mapper<Object,
Text, Text, IntWritable> {
        private static final IntWritable one = new
IntWritable(1);

        private Text word = new Text();

        private String tokens =
"[_]|$#<>\\^=\\[\\]\\|\\*/\\\\\\\\,;,.\\|-:()?!\\\"'"]";

        public void map(Object key, Text value, Mapper<Object,
Text, Text, IntWritable>.Context context) throws
IOException, InterruptedException {
            String cleanLine =
value.toString().toLowerCase().replaceAll(this.tokens, "
");
            StringTokenizer itr = new StringTokenizer(cleanLine);
            while (itr.hasMoreTokens()) {
                this.word.set(itr.nextToken().trim());
                context.write(this.word, one);
            }
        }
    }
}

```

TopNCombiner.class

```

package samples.topn;

import java.io.IOException;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Reducer;

public class TopNCombiner extends Reducer<Text,
IntWritable, Text, IntWritable> {

```

```

    public void reduce(Text key, Iterable<IntWritable>
values, Reducer<Text, IntWritable, Text,
IntWritable>.Context context) throws IOException,
InterruptedException {
    int sum = 0;
    for (IntWritable val : values)
        sum += val.get();
    context.write(key, new IntWritable(sum));
}
}

```

TopNMapper.class

```

package samples.topn;

```

```

import java.io.IOException;
import java.util.StringTokenizer;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;

```

```

public class TopNMapper extends Mapper<Object, Text, Text,
IntWritable> {
    private static final IntWritable one = new
IntWritable(1);

```

```

    private Text word = new Text();

```

```

    private String tokens = "[_!$#<>\\^=\\[\\]\\|\\*\\/\\\\\\\\\\,;\\.\\-
:()?!\\\"'"]";

```

```

    public void map(Object key, Text value,
Mapper<Object, Text, Text, IntWritable>.Context context)
throws IOException, InterruptedException {
        String cleanLine =
value.toString().toLowerCase().replaceAll(this.tokens, "
");
        StringTokenizer itr = new StringTokenizer(cleanLine);
        while (itr.hasMoreTokens()) {
            this.word.set(itr.nextToken().trim());
            context.write(this.word, one);
        }
    }
}

```

```
}
```

TopNReducer.class

```
package samples.topn;
```

```
import java.io.IOException;
```

```
import java.util.HashMap;
```

```
import java.util.Map;
```

```
import org.apache.hadoop.io.IntWritable;
```

```
import org.apache.hadoop.io.Text;
```

```
import org.apache.hadoop.mapreduce.Reducer;
```

```
import utils.MiscUtils;
```

```
public class TopNReducer extends Reducer<Text, IntWritable,  
Text, IntWritable> {
```

```
    private Map<Text, IntWritable> countMap = new  
    HashMap<>();
```

```
    public void reduce(Text key, Iterable<IntWritable>  
values, Reducer<Text, IntWritable, Text,  
IntWritable>.Context context) throws IOException,  
InterruptedException {  
        int sum = 0;  
        for (IntWritable val : values)  
            sum += val.get();  
        this.countMap.put(new Text(key), new IntWritable(sum));  
    }
```

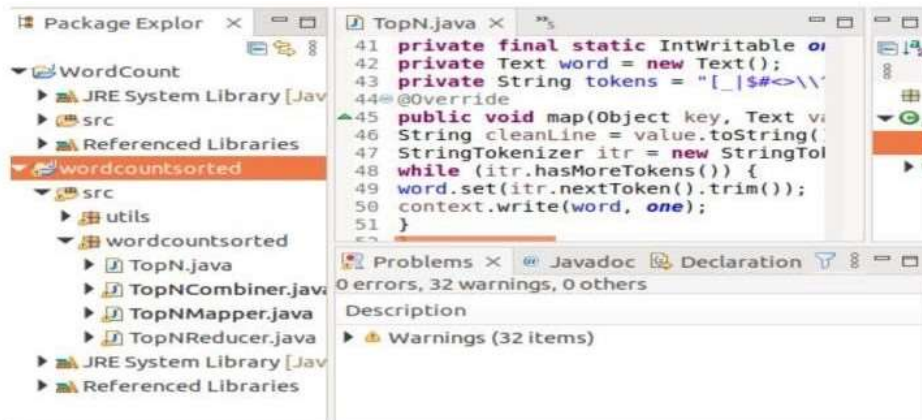
```
    protected void cleanup(Reducer<Text, IntWritable, Text,  
IntWritable>.Context context) throws IOException,  
InterruptedException {
```

```
        Map<Text, IntWritable> sortedMap =  
MiscUtils.sortByValues(this.countMap);  
        int counter = 0;  
        for (Text key : sortedMap.keySet()) {  
            if (counter++ == 20)  
                break;  
            context.write(key, sortedMap.get(key));  
        }
```



```
}  
}
```

OUTPUT:



```
hadoop@sharat-VirtualBox:~/hadoop-3.2.3/sbin$ hdfs dfs -mkdir /input
2022-06-27 21:59:42,586 WARN util.NativeCodeLoader: Unable to load native-hadoop
p library for your platform... using builtin-java classes where applicable
mkdir: '/input': File exists
hadoop@sharat-VirtualBox:~/hadoop-3.2.3/sbin$ hdfs dfs -put /home/hadoop/Docum
s/b.txt /input
2022-06-27 22:00:59,014 WARN util.NativeCodeLoader: Unable to load native-hadoop
p library for your platform... using builtin-java classes where applicable
put: '/input/b.txt': File exists
hadoop@sharat-VirtualBox:~/hadoop-3.2.3/sbin$ hdfs dfs -put /home/hadoop/Docum
s/b.txt /inputt
2022-06-27 22:01:16,095 WARN util.NativeCodeLoader: Unable to load native-hadoop
p library for your platform... using builtin-java classes where applicable
hadoop@sharat-VirtualBox:~/hadoop-3.2.3/sbin$ hdfs dfs -ls /inputt
2022-06-27 22:01:33,726 WARN util.NativeCodeLoader: Unable to load native-hadoop
p library for your platform... using builtin-java classes where applicable
Found 3 items
-rw-r--r-- 1 hadoop supergroup      15 2022-06-20 16:51 /inputt/a.txt
-rw-r--r-- 1 hadoop supergroup     38 2022-06-27 22:01 /inputt/b.txt
drwxr-xr-x - hadoop supergroup      0 2022-06-20 16:52 /inputt/output
```



```
hadoop@sharat-VirtualBox:~/hadoop-3.2.3/sbin$ hdfs dfs -ls inputt/outputword
2022-06-27 22:08:26,995 WARN util.NativeCodeLoader: Unable to load native-hadoop
p library for your platform... using builtin-java classes where applicable
Found 2 items
-rw-r--r-- 1 hadoop supergroup 0 2022-06-27 22:05 inputt/outputword/_
SUCCESS
-rw-r--r-- 1 hadoop supergroup 35 2022-06-27 22:05 inputt/outputword/p
art-r-00000
hadoop@sharat-VirtualBox:~/hadoop-3.2.3/sbin$ hdfs dfs -cat inputt/outputword/pa
rt-r-00000
2022-06-27 22:09:12,199 WARN util.NativeCodeLoader: Unable to load native-hadoop
p library for your platform... using builtin-java classes where applicable
test 2
is 2
this 2
a 1
important 1
```

8. Create a Map Reduce program to demonstrating join operation

CODE:

```
// JoinDriver.java
import org.apache.hadoop.conf.Configured;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapred.*;
import org.apache.hadoop.mapred.libMultipleInputs;
import org.apache.hadoop.util.*;

public class JoinDriver extends Configured implements Tool {

    public static class KeyPartitioner implements Partitioner<TextPair,
    Text> {
        @Override
        public void configure(JobConf job) {}

        @Override
        public int getPartition(TextPair key, Text value, int numPartitions) {
            return (key.getFirst().hashCode() & Integer.MAX_VALUE) %
            numPartitions;
        }
    }

    @Override
    public int run(String[] args) throws Exception {
        if (args.length != 3) {
            System.out.println("Usage: <Department Emp Strength input>
            <Department Name input> <output>");
            return -1;
        }

        JobConf conf = new JobConf(getConf(), getClass());
```

```
conf.setJobName("Join 'Department Emp Strength input' with  
'Department Name  
input'");  
  
Path AInputPath = new Path(args[0]);  
Path BInputPath = new Path(args[1]);  
Path outputPath = new Path(args[2]);  
  
MultipleInputs.addInputPath(conf, AInputPath, TextInputFormat.class,  
Posts.class);  
  
MultipleInputs.addInputPath(conf, BInputPath, TextInputFormat.class,  
User.class);  
  
FileOutputFormat.setOutputPath(conf, outputPath);  
  
conf.setPartitionerClass(KeyPartitioner.class);  
  
conf.setOutputValueGroupingComparator(TextPair.FirstComparator.cl  
ass);  
  
conf.setMapOutputKeyClass(TextPair.class);  
  
conf.setReducerClass(JoinReducer.class);  
  
conf.setOutputKeyClass(Text.class);  
  
JobClient.runJob(conf);  
  
return 0;  
}  
  
public static void main(String[] args) throws Exception {  
  
int exitCode = ToolRunner.run(new JoinDriver(), args);  
System.exit(exitCode);
```

```
}}
```

```
// JoinReducer.java
```

```
import java.io.IOException;
```

```
import java.util.Iterator;
```

```
import org.apache.hadoop.io.Text;
```

```
import org.apache.hadoop.mapred.*;
```

```
public class JoinReducer extends MapReduceBase implements  
Reducer<TextPair, Text, Text,  
Text> {
```

```
@Override
```

```
public void reduce (TextPair key, Iterator<Text> values,  
OutputCollector<Text, Text>  
output, Reporter reporter)
```

```
throws IOException
```

```
{
```

```
Text nodeId = new Text(values.next());
```

```
while (values.hasNext()) {
```

```
Text node = values.next();
```

```
Text outValue = new Text(nodeId.toString() + "\t\t" + node.toString());
```

```
output.collect(key.getFirst(), outValue);
```

```
}
```

```
}
```

```
}
```

```
// User.java
```

```
import java.io.IOException;
```

```
import java.util.Iterator;
```

```
import org.apache.hadoop.conf.Configuration;
```

```
import org.apache.hadoop.fs.FSDataInputStream;
```

```
import org.apache.hadoop.fs.FSDataOutputStream;
```

```

import org.apache.hadoop.fs.FileSystem;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapred.*;

import org.apache.hadoop.io.IntWritable;

public class User extends MapReduceBase implements
Mapper<LongWritable, Text, TextPair,
Text> {

    @Override
    public void map(LongWritable key, Text value,
OutputCollector<TextPair, Text> output,
Reporter reporter)

        throws IOException

    {

        String valueString = value.toString();

        String[] SingleNodeData = valueString.split("\t");
        output.collect(new TextPair(SingleNodeData[0], "1"), new
Text(SingleNodeData[1]));
    }
}

//Posts.java
import java.io.IOException;

import org.apache.hadoop.io.*;
import org.apache.hadoop.mapred.*;

```

```

public class Posts extends MapReduceBase implements
Mapper<LongWritable, Text, TextPair,
Text> {

@Override
public void map(LongWritable key, Text value,
OutputCollector<TextPair, Text> output,
Reporter reporter)
throws IOException
{
String valueString = value.toString();
String[] SingleNodeData = valueString.split("\t");
output.collect(new TextPair(SingleNodeData[3], "0"), new
Text(SingleNodeData[9]));
}
}

```

```

// TextPair.java
import java.io.*;

import org.apache.hadoop.io.*;

public class TextPair implements WritableComparable<TextPair> {

private Text first;
private Text second;

public TextPair() {
set(new Text(), new Text());
}

public TextPair(String first, String second) {
set(new Text(first), new Text(second));
}

public TextPair(Text first, Text second) {

```

```
set(first, second);  
}
```

```
public void set(Text first, Text second) {  
    this.first = first;  
    this.second = second;  
}
```

```
public Text getFirst() {  
    return first;  
}
```

```
public Text getSecond() {  
    return second;  
}
```

```
@Override  
public void write(DataOutput out) throws IOException {  
    first.write(out);  
    second.write(out);  
}
```

```
@Override  
public void readFields(DataInput in) throws IOException {  
    first.readFields(in);  
    second.readFields(in);  
}
```

```
@Override  
public int hashCode() {  
    return first.hashCode() * 163 + second.hashCode();  
}
```

```
@Override  
public boolean equals(Object o) {  
    if (o instanceof TextPair) {
```



```
TextPair tp = (TextPair) o;  
return first.equals(tp.first) && second.equals(tp.second);  
}  
return false;  
}
```

```
@Override  
public String toString() {  
return first + "\t" + second;  
}
```

```
@Override  
public int compareTo(TextPair tp) {  
int cmp = first.compareTo(tp.first);  
if (cmp != 0) {  
return cmp;  
}  
return second.compareTo(tp.second);  
}  
// ^^ TextPair
```

```
// vv TextPairComparator  
public static class Comparator extends WritableComparator {  
  
private static final Text.Comparator TEXT_COMPARATOR = new  
Text.Comparator();  
  
public Comparator() {  
super(TextPair.class);  
}
```

```
@Override  
public int compare(byte[] b1, int s1, int l1,  
byte[] b2, int s2, int l2) {  
  
try {
```

```

int firstL1 = WritableUtils.decodeVIntSize(b1[s1]) + readVInt(b1, s1);
int firstL2 = WritableUtils.decodeVIntSize(b2[s2]) + readVInt(b2, s2);
int cmp = TEXT_COMPARATOR.compare(b1, s1, firstL1, b2, s2,
firstL2);
if (cmp != 0) {
return cmp;
}
return TEXT_COMPARATOR.compare(b1, s1 + firstL1, l1 - firstL1,

b2, s2 + firstL2, l2 - firstL2);
} catch (IOException e) {
throw new IllegalArgumentException(e);
}
}
}

static {
WritableComparator.define(TextPair.class, new Comparator());
}
public static class FirstComparator extends WritableComparator {

private static final Text.Comparator TEXT_COMPARATOR = new
Text.Comparator();

public FirstComparator() {
super(TextPair.class);
}

@Override
public int compare(byte[] b1, int s1, int l1,
byte[] b2, int s2, int l2) {

try {
int firstL1 = WritableUtils.decodeVIntSize(b1[s1]) + readVInt(b1, s1);
int firstL2 = WritableUtils.decodeVIntSize(b2[s2]) + readVInt(b2, s2);
return TEXT_COMPARATOR.compare(b1, s1, firstL1, b2, s2, firstL2);

```

```

} catch (IOException e) {
throw new IllegalArgumentException(e);
}
}

```

@Override

```

public int compare(WritableComparable a, WritableComparable b) {
if (a instanceof TextPair && b instanceof TextPair) {
return ((TextPair) a).first.compareTo(((TextPair) b).first);
}
return super.compare(a, b);
}
} }

```

OUTPUT:

```

hadoop@sharat-VirtualBox:~$ hdfs dfs -copyFromLocal DeptName.txt DeptStrength.txt /
2022-06-28 01:49:34,172 WARN util.NativeCodeLoader: Unable to load native-hadoop
p library for your platform... using builtin-java classes where applicable
copyFromLocal: 'DeptStrength.txt': No such file or directory
hadoop@sharat-VirtualBox:~$ hdfs dfs -copyFromLocal DeptName.txt DeptEmpStrength
.txt /
2022-06-28 01:50:03,670 WARN util.NativeCodeLoader: Unable to load native-hadoop
p library for your platform... using builtin-java classes where applicable
copyFromLocal: '/DeptName.txt': File exists
hadoop@sharat-VirtualBox:~$ hdfs dfs -copyFromLocal DeptEmpStrength.txt /
2022-06-28 01:50:14,698 WARN util.NativeCodeLoader: Unable to load native-hadoop
p library for your platform... using builtin-java classes where applicable
copyFromLocal: '/DeptEmpStrength.txt': File exists

```

```

hadoop@sharat-VirtualBox:~$ hadoop jar MapReduceJoin.jar /DeptEmpStrength.txt /D
eptName.txt /output_mapreducejoin
2022-06-28 01:54:22,260 WARN util.NativeCodeLoader: Unable to load native-hadoop
p library for your platform... using builtin-java classes where applicable
2022-06-28 01:54:22,634 INFO client.RMProxy: Connecting to ResourceManager at /
127.0.0.1:8032
2022-06-28 01:54:22,756 INFO client.RMProxy: Connecting to ResourceManager at /
127.0.0.1:8032
2022-06-28 01:54:22,936 INFO mapreduce.JobResourceUploader: Disabling Erasure C
oding for path: /tmp/hadoop-yarn/staging/hadoop/.staging/job_1656358828291_0002
2022-06-28 01:54:23,108 INFO mapred.FileInputFormat: Total input files to proce
ss : 1
2022-06-28 01:54:23,121 INFO mapred.FileInputFormat: Total input files to proce
ss : 1
2022-06-28 01:54:23,607 INFO mapreduce.JobSubmitter: number of splits:4
2022-06-28 01:54:23,771 INFO mapreduce.JobSubmitter: Submitting tokens for job:
job_1656358828291_0002
2022-06-28 01:54:23,772 INFO mapreduce.JobSubmitter: Executing with tokens: []
2022-06-28 01:54:23,909 INFO conf.Configuration: resource-types.xml not found
2022-06-28 01:54:23,909 INFO resource.ResourceUtils: Unable to find 'resource-t
ypes.xml'.
2022-06-28 01:54:23,967 INFO impl.YarnClientImpl: Submitted application applica
tion_1656358828291_0002

```

```

Bytes Written=85
hadoop@sharat-VirtualBox:~$ hdfs dfs -ls /output/output_mapreducejoin
2022-06-28 01:55:29,436 WARN util.NativeCodeLoader: Unable to load native-hadoop
p library for your platform... using builtin-java classes where applicable
ls: '/output/output_mapreducejoin': No such file or directory
hadoop@sharat-VirtualBox:~$ hdfs dfs -ls /output_mapreducejoin
2022-06-28 01:55:36,422 WARN util.NativeCodeLoader: Unable to load native-hadoop
p library for your platform... using builtin-java classes where applicable
Found 2 items
-rw-r--r-- 1 hadoop supergroup 0 2022-06-28 01:54 /output_mapreducejo
ln/_SUCCESS
-rw-r--r-- 1 hadoop supergroup 85 2022-06-28 01:54 /output_mapreducejo
ln/part-00000
hadoop@sharat-VirtualBox:~$ hdfs dfs -cat /output_mapreducejoin/part-00000
2022-06-28 01:56:01,186 WARN util.NativeCodeLoader: Unable to load native-hadoop
p library for your platform... using builtin-java classes where applicable
A11 50 Finance
B12 100 HR
C13 250 Manufacturing
Dept_ID Total_Employee Dept_Name

```

9. Program to print word count on scala shell and print “Hello world” on scala IDE

CODE:

```
package wordcount
```

```
import org.apache.spark.SparkConf import
org.apache.spark.SparkContext
import org.apache.spark.rdd.RDD._, PairRDDFunctions
```

```
object WordCount {
  def
  main(args: Array[String]) = {
    //Start the Spark context
    val conf = new SparkConf().setAppName("WordCount").setMaster("local") val
    sc = new SparkContext(conf)
    //Read some example file to a test RDD val
    test = sc.textFile("input.txt") test.flatMap {
      line => //for
      each line
      line.split(" ") //split
      the line in word by word.
      } .map {
      word => //for
      each word
      (word, 1) //Return a key/value tuple, with the word as key and 1 as value
      } .reduceByKey(_ + _) //Sum
      all of the value with same key
      .saveAsTextFile("output.txt") //Save to a
      text file
      //Stop the Spark context
      sc.stop
    }
  }
```

OUTPUT:

```
scala> val test=sc.textFile("/home/hadoop/spark_word_count.txt")
test: org.apache.spark.rdd.RDD[String] = /home/hadoop/spark_word_count.txt MapPartitionsRDD[1] at textFile at <console>:23

scala> test.collect;
[Stage 0:>                                     (0 + 0) /
[Stage 0:>                                     (0 + 2) /

res4: Array[String] = Array(This is a test, This is an evaluation, do you want
a test, why do you want a test)

scala> val count=test.flatMap(line=>line.split(" "))
count: org.apache.spark.rdd.RDD[String] = MapPartitionsRDD[2] at flatMap at <console>:23

scala> count.collect
res5: Array[String] = Array(This, is, a, test, This, is, an, evaluation, do,
u, want, a, test, why, do, you, want, a, test)

scala> val map_frequency=count.map(entry=>(entry,1))
map_frequency: org.apache.spark.rdd.RDD[(String, Int)] = MapPartitionsRDD[3]
map at <console>:23

scala> map_frequency.collect
res6: Array[(String, Int)] = Array((This,1), (is,1), (a,1), (test,1), (This,1),
(is,1), (an,1), (evaluation,1), (do,1), (you,1), (want,1), (a,1), (test,1),
hy,1), (do,1), (you,1), (want,1), (a,1), (test,1))
```

```
scala> map_frequency.reduceByKey(_+_ )
res7: org.apache.spark.rdd.RDD[(String, Int)] = ShuffledRDD[4] at reduceByKey
t <console>:24

scala> map_frequency.collect
res8: Array[(String, Int)] = Array((This,1), (is,1), (a,1), (test,1), (This,1),
(is,1), (an,1), (evaluation,1), (do,1), (you,1), (want,1), (a,1), (test,1),
hy,1), (do,1), (you,1), (want,1), (a,1), (test,1))

scala>

scala> val final_output=map_frequency.reduceByKey(_+_ )
final_output: org.apache.spark.rdd.RDD[(String, Int)] = ShuffledRDD[5] at reduceByKey at <console>:23

scala> final_output.collect
[Stage 4:>                                     (0 + 2) /

res9: Array[(String, Int)] = Array((is,2), (evaluation,1), (This,2), (why,1),
want,2), (test,3), (you,2), (a,3), (do,2), (an,1))
```


10. Using RDD and FlatMap count how many times each word appears in a file and write out a list of words whose count is strictly greater than 4 using Spark

CODE:

```
val textFile = sc.textFile("/home/Desktop/test.txt")
val counts = textFile.flatMap(line => line.split(" ")).map(word => (word, 1)).reduceByKey(_ + _)
import scala.collection.immutable.ListMap
val sorted = ListMap(counts.collect.sortWith(_._2 > _._2):_*) // sort in descending order based on values
println(sorted)
for((k,v) <- sorted)
{
  if(v > 4)
  {
    print(k+",")
    print(v)
    println()
  }
}
```

OUTPUT:

```
sorted: scala.collection.immutable.ListMap[String,Int] = ListMap(test -> 5, is -> 2, This -> 2, want -> 2, do -> 2, why -> 1, you -> 1)

scala> for((k,v) <- sorted)
| {
|   if(v > 4)
|   {
|     print(k+",")
|     print(v)
|     println()
|   }
| }
test,5
```

```
scala> val word_count = sc.textFile("/home/hadoop/spark_word_count.txt")
word_count: org.apache.spark.rdd.RDD[String] = /home/hadoop/spark_word_count.txt MapPartitionsRDD[1] at textFile at <console>:23

scala> val frequency = word_count.flatMap((line) => line.split(" ")).map(word => (word, 1)).reduceByKey(_ + _)
frequency: org.apache.spark.rdd.RDD[(String, Int)] = ShuffledRDD[4] at reduceByKey at <console>:23

scala> val sorted = ListMap(frequency.collect.sortWith(_._2 > _._2):_*)
<console>:23: error: not found: value ListMap
      val sorted = ListMap(frequency.collect.sortWith(_._2 > _._2):_*)
                        ^

scala> import scala.collection.immutable.ListMap
import scala.collection.immutable.ListMap

scala> val sorted = ListMap(frequency.collect.sortWith(_._2 > _._2):_*)
[Stage 0:>] (0 + 2)

sorted: scala.collection.immutable.ListMap[String,Int] = ListMap(test -> 5, is -> 2, This -> 2, want -> 2, do -> 2, why -> 1, you -> 1, an -> 1)
```



```
sorted: scala.collection.immutable.ListMap[String,Int] = ListMap(  
> 3, is -> 2, This -> 2, want -> 2, do -> 2, why -> 1, you -> 3  
  
scala> for((k,v)<-sorted)  
      | {  
      |   if(v>4)  
      |   {  
      |     print(k+",")  
      |     print(v)  
      |     println()  
      |   }  
      | }  
test,5
```