# B.M.S. COLLEGE OF ENGINEERING BENGALURU
Autonomous Institute, Affiliated to VTU

Lab Record

## Big Data Analytics

*Submitted in partial fulfillment for the 6th Semester Laboratory*

Bachelor of Technology
in
Computer Science and Engineering

*Submitted by:*

## Ankit Kesar

1BM18CS150

Department of Computer Science and Engineering
B.M.S. College of Engineering
Bull Temple Road, Basavanagudi, Bangalore 560 019
Mar-June 2021

# B.M.S. COLLEGE OF ENGINEERING

# DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING



## *CERTIFICATE*

This is to certify that the Big Data Analytics (20CS6PEBDA) laboratory has been carried out by Ankit Kesar  (1BM18CS150) during the 6<sup>th</sup> Semester Mar-June-2021.

Signature of the Faculty Incharge:

Bhoomika :

Department of Computer Science and Engineering
B.M.S. College of Engineering, Bangalore

# Table of Contents

**Program 1.** 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 Employees.

```
cqlsh> create keyspace employee with replication = {'class': 'SimpleStrategy', 'replication_factor': 1};
cqlsh>  use employee;
cqlsh:employee> create table employeeinfo(emp_id int primary key, emp_name text, designation text, doj timestamp, salary double, dept_name text);
cqlsh:employee> describe table employeeinfo;

CREATE TABLE employee.employeeinfo (
    emp_id int PRIMARY KEY,
    dept_name text,
    designation text,
    doj timestamp,
    emp_name text,
    salary double
) WITH 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 = 864000
    AND max_index_interval = 2048
    AND memtable_flush_period_in_ms = 0
    AND min_index_interval = 128
    AND read_repair_chance = 0.0
    AND speculative_retry = '99PERCENTILE';

cqlsh:employee> begin batch
        ... insert into employeeinfo(emp_id, emp_name, designation, doj, salary, dept_name) values (2, 'Akanksha', 'Data analyst', '2010-05-15', 23456.90, 'Corporate');
        ... insert into employeeinfo(emp_id, emp_name, designation, doj, salary, dept_name) values (3, 'Abhinay', 'Manager', '2012-09-05', 33333, 'Web development');
        ... insert into employeeinfo(emp_id, emp_name, designation, doj, salary, dept_name) values (8, 'Akshita', 'Software developer', '2003-05-05', 123123, 'Data
    analytics');
        ... insert into employeeinfo(emp_id, emp_name, designation, doj, salary, dept_name) values (4, 'Anmol', 'Corporate', '2003-06-05', 242, 'IT');
        ... apply batch;
cqlsh:employee> select * from employeeinfo;
```

Fig 1.1

```
 emp_id | dept_name        | designation       | doj                              | emp_name | salary
--------+------------------+-------------------+----------------------------------+----------+-----------
      8 |   Data analytics | Software developer | 2003-05-04 18:30:00.000000+0000 |  Akshita | 1.2312e+05
      2 |        Corporate |      Data analyst | 2010-05-14 18:30:00.000000+0000 | Akanksha |    23456.9
      4 |               IT |         Corporate | 2003-06-04 18:30:00.000000+0000 |    Anmol |       242
      3 |  Web development |           Manager | 2012-09-04 18:30:00.000000+0000 |  Abhinay |     33333

(4 rows)
cqlsh:employee> begin batch insert into employeeinfo(emp_id, emp_name, designation, doj, salary, dept_name) values (121, 'Akash', 'HR', '2012-09-05', 111111, 'Corporate');
apply batch;
cqlsh:employee> select * from employeeinfo;

 emp_id | dept_name        | designation       | doj                              | emp_name | salary
--------+------------------+-------------------+----------------------------------+----------+-----------
      8 |   Data analytics | Software developer | 2003-05-04 18:30:00.000000+0000 |  Akshita | 1.2312e+05
      2 |        Corporate |      Data analyst | 2010-05-14 18:30:00.000000+0000 | Akanksha |    23456.9
      4 |               IT |         Corporate | 2003-06-04 18:30:00.000000+0000 |    Anmol |       242
    121 |        Corporate |                HR | 2012-09-04 18:30:00.000000+0000 |    Akash | 1.1111e+05
      3 |  Web development |           Manager | 2012-09-04 18:30:00.000000+0000 |  Abhinay |     33333

(5 rows)
cqlsh:employee> update employeeinfo set emp_name = 'Jinny', dept_name = 'Management' where emp_id = 121;
cqlsh:employee> select * from employeeinfo;

 emp_id | dept_name        | designation       | doj                              | emp_name | salary
--------+------------------+-------------------+----------------------------------+----------+-----------
      8 |   Data analytics | Software developer | 2003-05-04 18:30:00.000000+0000 |  Akshita | 1.2312e+05
      2 |        Corporate |      Data analyst | 2010-05-14 18:30:00.000000+0000 | Akanksha |    23456.9
      4 |               IT |         Corporate | 2003-06-04 18:30:00.000000+0000 |    Anmol |       242
    121 |       Management |                HR | 2012-09-04 18:30:00.000000+0000 |    Jinny | 1.1111e+05
      3 |  Web development |           Manager | 2012-09-04 18:30:00.000000+0000 |  Abhinay |     33333

(5 rows)
cqlsh:employee> create index on employeeinfo(salary);
```

Fig 1.2

```
cqlsh:employee> create index on employeeinfo(salary);
cqlsh:employee> update employeeinfo set projects = {'Test', 'Start'} where emp_id in(8,2,4,121,3);
cqlsh:employee> alter table employeeinfo add projects set<text>;
cqlsh:employee>  update employeeinfo set projects = {'Test', 'Start'} where emp_id in(8,2,4,121,3);
cqlsh:employee> select * from employeeinfo;

 emp_id | dept_name        | designation       | doj                              | emp_name | projects          | salary
--------+------------------+-------------------+----------------------------------+----------+-------------------+-----------
      8 |   Data analytics | Software developer | 2003-05-04 18:30:00.000000+0000 |  Akshita | {'Start', 'Test'} | 1.2312e+05
      2 |        Corporate |      Data analyst | 2010-05-14 18:30:00.000000+0000 | Akanksha | {'Start', 'Test'} |    23456.9
      4 |               IT |         Corporate | 2003-06-04 18:30:00.000000+0000 |    Anmol | {'Start', 'Test'} |       242
    121 |       Management |                HR | 2012-09-04 18:30:00.000000+0000 |    Jinny | {'Start', 'Test'} | 1.1111e+05
      3 |  Web development |           Manager | 2012-09-04 18:30:00.000000+0000 |  Abhinay | {'Start', 'Test'} |     33333

(5 rows)
cqlsh:employee> begin batch insert into employeeinfo(emp_id, emp_name, designation, doj, salary, dept_name) values (121, 'Boris', 'MTO', '2001-08-05', 12212, 'Corporate') using
ttl 120; apply batch;
cqlsh:employee>  select ttl(designation) from employeeinfo where emp_id = 121;

 ttl(designation)
------------------
              103

(1 rows)
cqlsh:employee>
```

Fig 1.3

**Program 2.** 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

```
cqlsh> create keyspace library  with replication = { 'class' : 'SimpleStrategy','replication_factor':1};
cqlsh> use library;
cqlsh:library> create table library_info( id int, counter_val counter, stud_name text, book_name text, book_id int, issue_date timestamp,primary
key(id,stud_name,book_name,book_id,issue_date));
cqlsh:library> update library_info SET counter_val = counter_val +1 where id = 1 and stud_name = 'Akanksha' and book_name = 'DBMS' and book_id = 121 and issue_date='2017-10-
08';
cqlsh:library> update library_info SET counter_val = counter_val +1 where id = 3 and stud_name = 'Akshay' and book_name = 'BDA' and book_id = 112 and issue_date='2011-12-20';
cqlsh:library> update library_info SET counter_val = counter_val +1 where id = 5 and stud_name = 'Akshat' and book_name = 'Java' and book_id = 114 and issue_date='2009-08-27';
cqlsh:library> update library_info SET counter_val = counter_val +1 where id = 10 and stud_name = 'Akash' and book_name = 'Operating system' and book_id = 118 and
issue_date='2005-12-03';
cqlsh:library> select * from library_info;

 id | stud_name | book_name       | book_id | issue_date                     | counter_val
----+-----------+-----------------+---------+--------------------------------+-------------
  5 |    Akshat |            Java |     114 | 2009-08-26 18:30:00.000000+0000 |           1
 10 |     Akash | Operating system |     118 | 2005-12-02 18:30:00.000000+0000 |           1
  1 |  Akanksha |            DBMS |     121 | 2017-10-07 18:30:00.000000+0000 |           1
  3 |    Akshay |             BDA |     112 | 2011-12-19 18:30:00.000000+0000 |           1

(4 rows)
cqlsh:library> update library_info SET counter_val = counter_val +1 where id = 3 and stud_name = 'Akshay' and book_name = 'BDA' and book_id = 112 and issue_date='2011-12-20';
cqlsh:library> select * from library_info where counter_val = 2 allow filtering;

 id | stud_name | book_name | book_id | issue_date                     | counter_val
----+-----------+-----------+---------+--------------------------------+-------------
  3 |    Akshay |       BDA |     112 | 2011-12-19 18:30:00.000000+0000 |           2

(1 rows)
```

Fig 2.1

```
6 cqlsh:library> copy employee_info(id,counter_val,stud_name,book_name,book_id,issue_date) to '/home/bmsce/Desktop/week2_library_data.csv';
7 Column family 'employee_info' not found
8 cqlsh:library> copy library_info(id,counter_val,stud_name,book_name,book_id,issue_date) to '/home/bmsce/Desktop/week2_library_data.csv';
9 Using 11 child processes
0
1 Starting copy of library.library_info with columns [id, counter_val, stud_name, book_name, book_id, issue_date].
2 Processed: 4 rows; Rate:      24 rows/s; Avg. rate:      24 rows/s
3 4 rows exported to 1 files in 0.174 seconds.
4 cqlsh:library> copy library_info(id,counter_val,stud_name,book_name,book_id,issue_date) from '/home/bmsce/Desktop/week2_library_data.csv';
5 Using 11 child processes
6
7 Starting copy of library.library_info with columns [id, counter_val, stud_name, book_name, book_id, issue_date].
8 Processed: 4 rows; Rate:       7 rows/s; Avg. rate:      10 rows/s
9 4 rows imported from 1 files in 0.399 seconds (0 skipped).
0 cqlsh:library>
1
```

Fig 2.2

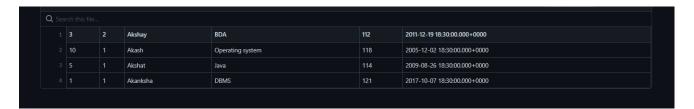| 1 | 3 | 2 | Akshay | BDA | 112 | 2011-12-19 18:30:00.000+0000 |
| 2 | 10 | 1 | Akash | Operating system | 118 | 2005-12-02 18:30:00.000+0000 |
| 3 | 5 | 1 | Akshat | Java | 114 | 2009-08-26 18:30:00.000+0000 |
| 4 | 1 | 1 | Akanksha | DBMS | 121 | 2017-10-07 18:30:00.000+0000 |

Fig 2.3

**Program 3**. Perform the following DB operations using MongoDB.1. Create a database "Student" with the following attributes Rollno, Age, ContactNo, Email-Id.2. Insert appropriate values3. Write query to update Email-Id of a student with rollno 10.4. Replace the student name from "ABC" to "FEM" of rollno 11.5. Export the created table into local file system6. Drop the table7. Import a given csv dataset from local file system into mongodb collection.



```
> show databases;
admin    0.000GB
config   0.000GB
local    0.000GB
student  0.000GB
> use Studentdb
switched to db Studentdb
> var information = [
... {
...      "Name" : "Akanksha",
...      "Age" : 21,
...      "Contact" : 1290345678,
...      "Email" : "abc@gmail.com",
...      "Rollno" : 7
... },
... {
...      "Name" : "Akash",
...      "Age" : 18,
...      "Conatct" : 2310458955,
...      "Email" : "xyz@yahoo.com",
...      "Rollno" : 10
... },
... {
...      "Name" : "Mohit",
...      "Age" : 25,
...      "Conatct" : 124567830,
...      "Email" : "qwer@hike.com",
...      "Rollno" : 11
... },
... {
...      "Name" : "Ayush",
...      "Age" : 12,
...      "Contact" : 0987654321,
...      "Email" : "qazx@gmail.com",
...      "Rollno" : 15
... },
... ];
> db.student_database.insert(information);
```

Fig 3.1

7

```
> db.student_database.insert(information);
BulkWriteResult({
        "writeErrors" : [ ],
        "writeConcernErrors" : [ ],
        "nInserted" : 4,
        "nUpserted" : 0,
        "nMatched" : 0,
        "nModified" : 0,
        "nRemoved" : 0,
        "upserted" : [ ]
})
> db.student_database.find().pretty()
{
        "_id" : ObjectId("606768e8719e10fb5c03819d"),
        "Name" : "Akanksha",
        "Age" : 21,
        "Contact" : 1290345678,
        "Email" : "abc@gmail.com",
        "Rollno" : 7
}
{
        "_id" : ObjectId("606768e8719e10fb5c03819e"),
        "Name" : "Akash",
        "Age" : 18,
        "Conatct" : 2310458955,
        "Email" : "xyz@yahoo.com",
        "Rollno" : 10
}
{
        "_id" : ObjectId("606768e8719e10fb5c03819f"),
        "Name" : "Mohit",
        "Age" : 25,
        "Conatct" : 124567830,
        "Email" : "qwer@hike.com",
        "Rollno" : 11
}
```

Fig 3.2

```
{
        "_id" : ObjectId("606768e8719e10fb5c0381a0"),
        "Name" : "Ayush",
        "Age" : 12,
        "Contact" : 987654321,
        "Email" : "qazx@gmail.com",
        "Rollno" : 15
}
> db.student_database.update({"Rollno":10},{$set:{"Email":"xyz@gmail.com"}})
WriteResult({ "nMatched" : 1, "nUpserted" : 0, "nModified" : 1 })
> db.student_database.find({"Rollno":10}).pretty()
{
        "_id" : ObjectId("606768e8719e10fb5c03819e"),
        "Name" : "Akash",
        "Age" : 18,
        "Conatct" : 2310458955,
        "Email" : "xyz@gmail.com",
        "Rollno" : 10
}
> db.student_database.update({"Rollno":11},{$set:{"Name":"Piyush"}})
WriteResult({ "nMatched" : 1, "nUpserted" : 0, "nModified" : 1 })
> db.student_database.find({"Rollno":11}).pretty()
{
        "_id" : ObjectId("606768e8719e10fb5c03819f"),
        "Name" : "Piyush",
        "Age" : 25,
        "Conatct" : 124567830,
        "Email" : "qwer@hike.com",
        "Rollno" : 11
}
```

Fig 3.3

```
> db.student_database.replaceOne({"Rollno":11},{"Name":"FEM","Age":25,"Contact":"124567830","Email" : "qwer@hike.com","Rollno" : 11});
{ "acknowledged" : true, "matchedCount" : 1, "modifiedCount" : 1 }
> db.student_database.find({"Rollno":11}).pretty()
{
        "_id" : ObjectId("606768e8719e10fb5c03819f"),
        "Name" : "FEM",
        "Age" : 25,
        "Contact" : "124567830",
        "Email" : "qwer@hike.com",
        "Rollno" : 11
}

> show tables
student_database
student_db
> db.student_db.find().pretty()
{
        "_id" : ObjectId("606817866de84a3417e07a9b"),
        "Name" : "Akash",
        "Age" : 18,
        "Contact" : NumberLong("2310458955"),
        "Email" : "xyz@gmail.com",
        "Rollno" : 10
}
{
        "_id" : ObjectId("606817866de84a3417e07a9c"),
        "Name" : "Akanksha",
        "Age" : 21,
        "Contact" : 1290345678,
        "Email" : "abc@gmail.com",
        "Rollno" : 7
}
{
        "_id" : ObjectId("606817866de84a3417e07a9d"),
        "Name" : "Name",
        "Age" : "Age",
        "Contact" : "Contact",
        "Email" : "Email",
        "Rollno" : "Rollno"
}
```

Fig 3.4

```
{
        "_id" : ObjectId("606817866de84a3417e07a9e"),
        "Name" : "FEM",
        "Age" : 25,
        "Contact" : 124567830,
        "Email" : "qwer@hike.com",
        "Rollno" : 11
}
{
        "_id" : ObjectId("606817866de84a3417e07a9f"),
        "Name" : "Ayush",
        "Age" : 12,
        "Contact" : 987654321,
        "Email" : "qazx@gmail.com",
        "Rollno" : 15
}
> db.student_db.find().pretty()
{
        "_id" : ObjectId("606817866de84a3417e07a9d"),
        "Name" : "Name",
        "Age" : "Age",
        "Contact" : "Contact",
        "Email" : "Email",
        "Rollno" : "Rollno"
}
{
        "_id" : ObjectId("606817866de84a3417e07a9b"),
        "Name" : "Akash",
        "Age" : 18,
        "Contact" : NumberLong("2310458955"),
        "Email" : "xyz@gmail.com",
        "Rollno" : 10
}
```

Fig 3.5

```
{
        "_id" : ObjectId("606817866de84a3417e07a9c"),
        "Name" : "Akanksha",
        "Age" : 21,
        "Contact" : 1290345678,
        "Email" : "abc@gmail.com",
        "Rollno" : 7
}
{
        "_id" : ObjectId("606817866de84a3417e07a9e"),
        "Name" : "FEM",
        "Age" : 25,
        "Contact" : 124567830,
        "Email" : "qwer@hike.com",
        "Rollno" : 11
}
{
        "_id" : ObjectId("606817866de84a3417e07a9f"),
        "Name" : "Ayush",
        "Age" : 12,
        "Contact" : 987654321,
        "Email" : "qazx@gmail.com",
        "Rollno" : 15
}
> db.student_database.drop()
true
> show tables
student_db
>
```

Fig 3.6



Fig 3.7

| | Name | Age | Contact | Email | Rollno |
|---|---|---|---|---|---|
| 2 | Akanksha | 21 | 1290345678 | abc@gmail.com | 7 |
| 3 | Akash | 18 | 2310458955 | xyz@gmail.com | 10 |
| 4 | FEM | 25 | 124567830 | qwer@hike.com | 11 |
| 5 | Ayush | 12 | 987654321 | qazx@gmail.com | 15 |

Fig 3.8

10

**Program 4.** Screenshot of Hadoop installed



```
8083 SecondaryNameNode
7908 DataNode
8485 NodeManager
10054 Jps
8361 ResourceManager
7759 NameNode
```

Fig 4.1

**Program 5.** Execution of HDFS Commands for interaction with Hadoop Environment. (Minimum 10 commands to be executed).



```
hduser@lab-VirtualBox:/usr/local/sbin$ hadoop fs -cat /mydir/file1.txt
21/04/19 23:38:07 WARN util.NativeCodeLoader: Unable to load native-hadoop libr
ary for your platform... using builtin-java classes where applicable
I am using Hadoop
line1
line2
```

```
hduser@lab-VirtualBox:/usr/local/sbin$ hadoop fs -ls /
21/04/19 22:58:36 WARN util.NativeCodeLoader: Unable to load native-hadoop libr
ary for your platform... using builtin-java classes where applicable
Found 2 items
drwxr-xr-x   - hduser supergroup          0 2021-04-19 22:58 /mydir
drwxr-xr-x   - hduser supergroup          0 2021-04-18 19:27 /mydr
```

Fig 5.1

```
hduser@lab-VirtualBox:/usr/local/sbin$ hadoop fs -copyFromLocal ~/file1.txt /my
dir
21/04/19 23:19:36 WARN util.NativeCodeLoader: Unable to load native-hadoop libr
ary for your platform... using builtin-java classes where applicable
hduser@lab-VirtualBox:/usr/local/sbin$ hadoop fs -ls /mydir
21/04/19 23:20:13 WARN util.NativeCodeLoader: Unable to load native-hadoop libr
ary for your platform... using builtin-java classes where applicable
Found 1 items
-rw-r--r--   1 hduser supergroup         30 2021-04-19 23:19 /mydir/file1.txt
hduser@lab-VirtualBox:/usr/local/sbin$

hduser@lab-VirtualBox:/usr/local/sbin$ hadoop fs -ls /
21/04/19 22:58:36 WARN util.NativeCodeLoader: Unable to load native-hadoop libr
ary for your platform... using builtin-java classes where applicable
Found 2 items
drwxr-xr-x   - hduser supergroup          0 2021-04-19 22:58 /mydir
drwxr-xr-x   - hduser supergroup          0 2021-04-18 19:27 /mydr
```

Fig 5.2

```
hduser@lab-VirtualBox:/usr/local/sbin$ hadoop fs -copyToLocal /mydir ~/hadoopco
py
21/04/19 23:29:39 WARN util.NativeCodeLoader: Unable to load native-hadoop libr
ary for your platform... using builtin-java classes where applicable
hduser@lab-VirtualBox:/usr/local/sbin$

hduser@lab-VirtualBox:/usr/local/sbin$ hadoop fs -ls /
21/04/19 22:58:36 WARN util.NativeCodeLoader: Unable to load native-hadoop libr
ary for your platform... using builtin-java classes where applicable
Found 2 items
drwxr-xr-x   - hduser supergroup          0 2021-04-19 22:58 /mydir
drwxr-xr-x   - hduser supergroup          0 2021-04-18 19:27 /mydr
```

Fig 5.3

```
hduser@lab-VirtualBox:/usr/local/sbin$ hadoop fs -ls /
21/04/19 23:48:41 WARN util.NativeCodeLoader: Unable to load native-hadoop libr
ary for your platform... using builtin-java classes where applicable
Found 2 items
drwxr-xr-x   - hduser supergroup          0 2021-04-19 23:45 /mydir
drwxr-xr-x   - hduser supergroup          0 2021-04-19 23:41 /newdir
hduser@lab-VirtualBox:/usr/local/sbin$ hadoop fs -cp /mydir/sample.txt /newdir
21/04/19 23:48:56 WARN util.NativeCodeLoader: Unable to load native-hadoop libr
ary for your platform... using builtin-java classes where applicable
hduser@lab-VirtualBox:/usr/local/sbin$ hadoop fs -ls /newdir
21/04/19 23:49:22 WARN util.NativeCodeLoader: Unable to load native-hadoop libr
ary for your platform... using builtin-java classes where applicable
Found 2 items
drwxr-xr-x   - hduser supergroup          0 2021-04-19 23:21 /newdir/mydr
-rw-r--r--   1 hduser supergroup         13 2021-04-19 23:48 /newdir/sample.txt
hduser@lab-VirtualBox:/usr/local/sbin$
```

Fig 5.4

```
hduser@lab-VirtualBox:/usr/local/sbin$ hadoop fs -get /mydr ~/copyfromhadoop
21/04/19 23:25:49 WARN util.NativeCodeLoader: Unable to load native-hadoop libr
ary for your platform... using builtin-java classes where applicable

hduser@lab-VirtualBox:/usr/local/sbin$ hadoop fs -ls /
21/04/19 22:58:36 WARN util.NativeCodeLoader: Unable to load native-hadoop libr
ary for your platform... using builtin-java classes where applicable
Found 2 items
drwxr-xr-x   - hduser supergroup          0 2021-04-19 22:58 /mydir
drwxr-xr-x   - hduser supergroup          0 2021-04-18 19:27 /mydr
```

Fig 5.5



```
hduser@lab-VirtualBox:/usr/local/sbin$ hadoop fs -ls /
21/04/19 22:58:36 WARN util.NativeCodeLoader: Unable to load native-hadoop libr
ary for your platform... using builtin-java classes where applicable
Found 2 items
drwxr-xr-x   - hduser supergroup          0 2021-04-19 22:58 /mydir
drwxr-xr-x   - hduser supergroup          0 2021-04-18 19:27 /mydr
```

Fig 5.6

**Program 6.** From the following link extract the weather data

https://github.com/tomwhite/hadoop-book/tree/master/input/ncdc/all .Create a Map Reduce program to

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

ii) find the mean max temperature for every month.

```
hduser@lab-VirtualBox:/home/lab$ hadoop dfs -cat /tempmax/part-r-00000
DEPRECATED: Use of this script to execute hdfs command is deprecated.
Instead use the hdfs command for it.

WARNING: An illegal reflective access operation has occurred
WARNING: Illegal reflective access by org.apache.hadoop.security.authentication
.util.KerberosUtil (file:/usr/local/hadoop/share/hadoop/common/lib/hadoop-auth-
2.6.0.jar) to method sun.security.krb5.Config.getInstance()
WARNING: Please consider reporting this to the maintainers of org.apache.hadoop
.security.authentication.util.KerberosUtil
WARNING: Use --illegal-access=warn to enable warnings of further illegal reflec
tive access operations
WARNING: All illegal access operations will be denied in a future release
21/05/10 16:08:48 WARN util.NativeCodeLoader: Unable to load native-hadoop libr
ary for your platform... using builtin-java classes where applicable
03      111
05      22
```

Fig 6.1.1

```
hduser@lab-VirtualBox:/home/lab$ hadoop jar /home/lab/temperaturemax.jar temper
atureMax.TempDriver /input/sample_temp.txt /tempmax
WARNING: An illegal reflective access operation has occurred
WARNING: Illegal reflective access by org.apache.hadoop.security.authentication
.util.KerberosUtil (file:/usr/local/hadoop/share/hadoop/common/lib/hadoop-auth-
2.6.0.jar) to method sun.security.krb5.Config.getInstance()
WARNING: Please consider reporting this to the maintainers of org.apache.hadoop
.security.authentication.util.KerberosUtil
WARNING: Use --illegal-access=warn to enable warnings of further illegal reflec
tive access operations
WARNING: All illegal access operations will be denied in a future release
21/05/10 16:07:56 WARN util.NativeCodeLoader: Unable to load native-hadoop libr
ary for your platform... using builtin-java classes where applicable
21/05/10 16:07:57 INFO Configuration.deprecation: session.id is deprecated. Ins
tead, use dfs.metrics.session-id
21/05/10 16:07:57 INFO jvm.JvmMetrics: Initializing JVM Metrics with processNam
e=JobTracker, sessionId=
21/05/10 16:07:58 WARN mapreduce.JobSubmitter: Hadoop command-line option parsi
ng not performed. Implement the Tool interface and execute your application wit
h ToolRunner to remedy this.
21/05/10 16:07:58 INFO input.FileInputFormat: Total input paths to process : 1
21/05/10 16:07:58 INFO mapreduce.JobSubmitter: number of splits:1
21/05/10 16:07:59 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_l
ocal701083544_0001
21/05/10 16:08:00 INFO mapreduce.Job: The url to track the job: http://localhos
t:8080/
21/05/10 16:08:00 INFO mapreduce.Job: Running job: job_local701083544_0001
21/05/10 16:08:00 INFO mapred.LocalJobRunner: OutputCommitter set in config nul
l
```

Fig 6.1.2

Fig 6.2

**Program 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.



Fig 7.1

Fig 7.2

**Program 8.** Create a Map Reduce program to combine information from the users file along with Information from the posts file by using the concept of join and display user_id, Reputation and Score**.**



Fig 8.1

Fig 8.2



Fig 8.3

**Program 9.**  Scala programs and Screenshot of Spark Installed



Fig 9.1



Fig 9.2

Fig 9.3



Fig 9.4

Fig 9.5



Fig 9.6

Fig 9.7



Fig 9.8

22

Fig 9.9



Fig 9.10

**Program 10.** Using RDD and Flat Map 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.



```
Using Scala version 2.12.10 (OpenJDK 64-Bit Server VM, Java 11.0.11)
Type in expressions to have them evaluated.
Type :help for more information.

scala> val textFile = sc.textFile("/home/akanksha/Desktop/test1.txt")
textFile: org.apache.spark.rdd.RDD[String] = /home/akanksha/Desktop/test1.txt MapPartitionsRDD[1] at textFile at <console>:24

scala> val counts = textFile.flatMap(line => line.split(" ")).map(word => (word, 1)).reduceByKey(_+_)
counts: org.apache.spark.rdd.RDD[(String, Int)] = ShuffledRDD[4] at reduceByKey at <console>:25

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

scala> val sorted = ListMap(counts.collect.sortWith(_._2>_._2):_*)
sorted: scala.collection.immutable.ListMap[String,Int] = ListMap(Dog -> 5, Cow -> 5, Akanksha -> 5, Cat -> 5, Hello -> 4, World -> 1
, world -> 1)

scala> println(sorted)
ListMap(Dog -> 5, Cow -> 5, Akanksha -> 5, Cat -> 5, Hello -> 4, World -> 1, world -> 1)

scala> for((k,v)<-sorted)
     | {
     | if(v>4)
     | {
     | print(k+",")
     | print(v)
     | println()
     | }
     | }
Dog,5
Cow,5
Akanksha,5
Cat,5

scala>
```

Fig 10.1