

VISVESVARAYA TECHNOLOGICAL UNIVERSITY

“JnanaSangama”, Belgaum -590014, Karnataka.



LAB REPORT on

Big Data Analytics (23CS6PCBDA)

Submitted by:

Pranav Anantha Rao (1BM22CS201)

**Under the Guidance of
Vikranth B.M.
Assistant Professor, BMSCE**

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

March 2024 - June 2024

**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 **Pranav Anantha Rao (1BM22CS201)**, 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 2024. The Lab report has been approved as it satisfies the academic requirements in respect of **Big Data Analytics –(23CS6PCBDA)** work prescribed for the said degree.

Vikranth B.M.
Associate Professor
Department of CSE
BMSCE, Bengaluru

Dr. Kavitha sooda
Professor and Head
Department of CSE
BMSCE, Bengaluru

Table Of Contents

Sl.no	Program details	Pg no
1	MongoDB- CRUD Operations Demonstration (Practice and Self Study)	5
2	Perform the following DB operations using Cassandra.	9
3	Perform the following DB operations using Cassandra	13
4	Execution of HDFS Commands for interaction with Hadoop Environment. (Minimum 10 commands to be executed)	17
5	Implement Wordcount program on Hadoop framework	28
6	Create a MapReduce program to find average temperature for each year from NCDC data set. b) find the mean max temperature for every month. 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.	33
7	Write a Scala program to print numbers from 1 to 100 using for loop.	44
8	Using RDD and FlatMap count how many times each word appears in a file and write out a list of words whose count is	46

	strictly greater than 4 using Spark.	
9	Write a simple streaming program in Spark to receive text data streams on a particular port, perform basic text cleaning (like white space removal, stop words removal, lemmatization, etc.), and print the cleaned text on the screen. (Open Ended Question).	47

Course Outcomes

CO1: Apply the concepts of NoSQL, Hadoop, Spark for a given task

CO2: Analyse data analytic techniques for a given problem .

CO3: Conduct experiments using data analytics mechanisms for a given problem.

1. Experiments

Experiment - 1

Question:

Perform the following DB operations using Cassandra.

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

Lab 04

Cassandra DB

→ Create Keyspace:

CREATE KEYSPACE students WITH REPLICATION = { 'class': 'SimpleStrategy', 'replication_factor': 1 };

→ Describe existing keyspaces:

DESCRIBE KEYSPACES;

→ More details:

SELECT * FROM system.schema.keyspaces

→ Use keyspace students:

USE students;

→ Create table Student_Info

CREATE TABLE students_Info (Roll_No int
PRIMARY KEY, StudName text, DateOfJoining
timestamp, last_exam_Percent double);

→ Lookup the names of tables:

DESCRIBE TABLES;

1.1.2 Code with Output:

```
bscsece@bscsece-HP-Elite-Tower-800-G9-Desktop-PC: ~ $ cqlsh
Connected to Test Cluster at 127.0.0.1:9042
[cqlsh 6.1.0 | Cassandra 4.1.4 | CQL spec 3.4.6 | Native protocol v5]
Use HELP for help.
cqlsh> create keyspace Employee with replication = {'class':'SimpleStrategy',;replicationfactor':1};
syntaxException: line 1:89 mismatched input ';' expecting ')'
...with replication = {'class':'SimpleStrategy',;replicationfactor':1...)
cqlsh> create keyspace Employee WITH replication={'class': 'SimpleStrategy', 'replicationfactor':1};
ConfigurationException: Unrecognized strategy option [replicationfactor] passed to SimpleStrategy for keyspace employee
cqlsh> create keyspace Employee WITH replication={'class': 'SimpleStrategy', 'replication_factor':1};
cqlsh> DESCRIBE KEYSPACES
employee      system_auth      system_schema      system_views
system      system_distributed      system_traces      system_virtual_schema

cqlsh> CREATE TABLE IF NOT EXISTS Employee_Info(
    ... Emp_Id INT PRIMARY KEY,
    ... Emp_name TEXT,
    ... designation TEXT,
    ... date_of_joining DATE,
    ... Salary FLOAT,
    ... Dep_name TEXT,
    ... Projects SET<TEXT>);
InvalidRequest: Error from server: code=2200 [Invalid query] message="No keyspace has been specified. USE a keyspace, or explicitly specify keyspace.tablename"
cqlsh> USE Employee
...
cqlsh> USE Employee

cqlsh> USE Employee;
cqlsh> Employee> CREATE TABLE IF NOT EXISTS Employee_Info( Emp_Id INT PRIMARY KEY, Emp_name TEXT, designation TEXT, date_of_joining DATE, Salary FLOAT, Dep_name TEXT, Projects SET<TEXT>);
cqlsh> Employee> describe keyspace Employee

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

CREATE TABLE employee.employee_info (
    emp_id int PRIMARY KEY,
    date_of_joining date,
    dep_name text,
    designation text,
    emp_name text,
    salary float,
    projects set<text>
) WITH additional.write_policy = '99p'
AND bloom_filter_fp_chance = 0.01
AND caching = {'keys': 'ALL', 'rows_per_partition': 'NONE'}
AND cdc = false
AND comment = ''
AND compaction = {'class': 'org.apache.cassandra.db.compaction.SizeTieredCompactionStrategy', 'max_threshold': '32', 'min_threshold': '4'}
AND compression = {'chunk_length_in_kb': '16', 'class': 'org.apache.cassandra.io.compress.LZ4Compressor'}
AND memtable = 'default'
AND crc_check_chance = 1.0
AND default_time_to_live = 0
AND extensions = {}
AND gc_grace_seconds = 864000
AND max_index_interval = 2048
AND memtable_flush_period_in_ms = 0
AND min_index_interval = 128
```

```
cqlsh:employee> update employee_info using ttl 15 set salary = 0 where emp_id = 121;
cqlsh:employee> select * from employee_info;
+-----+-----+-----+-----+-----+-----+-----+-----+
| emp_id | bonus | date_of_joining | dep_name | designation | emp_name | projects | salary |
+-----+-----+-----+-----+-----+-----+-----+-----+
| 120 | 12000 | 2024-05-06 | Engineering | Developer | Priyanka GH | {'Project B', 'ProjectA'} | 1e+06
| 123 | null | 2024-05-07 | Engineering | Engineer | Sadhana | {'Project M', 'Project P'} | 1.2e+06
| 122 | null | 2024-05-06 | Management | HR | Rachana | {'Project C', 'Project M'} | 9e+05
| 121 | 11000 | 2024-05-06 | Management | Developer | Shreya | {'Project C', 'ProjectA'} | 0
+-----+-----+-----+-----+-----+-----+-----+-----+
(4 rows)
cqlsh:employee> select * from employee_info;
+-----+-----+-----+-----+-----+-----+-----+-----+
| emp_id | bonus | date_of_joining | dep_name | designation | emp_name | projects | salary |
+-----+-----+-----+-----+-----+-----+-----+-----+
| 120 | 12000 | 2024-05-06 | Engineering | Developer | Priyanka GH | {'Project B', 'ProjectA'} | 1e+06
| 123 | null | 2024-05-07 | Engineering | Engineer | Sadhana | {'Project M', 'Project P'} | 1.2e+06
| 122 | null | 2024-05-06 | Management | HR | Rachana | {'Project C', 'Project M'} | 9e+05
| 121 | 11000 | 2024-05-06 | Management | Developer | Shreya | {'Project C', 'ProjectA'} | null
+-----+-----+-----+-----+-----+-----+-----+-----+
(4 rows)
cqlsh:employee>
```

```
AND speculative_retry = '99p';
cqlsh:employee> select * from employee_info;
+-----+-----+-----+-----+-----+-----+-----+-----+
| emp_id | date_of_joining | dep_name | designation | emp_name | projects | salary |
+-----+-----+-----+-----+-----+-----+-----+-----+
| 120 | 2024-05-06 | Engineering | Developer | Priyanka | {'Project B', 'ProjectA'} | 1e+06
| 123 | 2024-05-07 | Engineering | Engineer | Sadhana | {'Project M', 'Project P'} | 1.2e+06
| 122 | 2024-05-06 | Management | HR | Rachana | {'Project C', 'Project M'} | 9e+05
| 121 | 2024-05-06 | Management | Developer | Shreya | {'Project C', 'ProjectA'} | 9e+05
+-----+-----+-----+-----+-----+-----+-----+-----+
(4 rows)
cqlsh:employee> update employee_info set emp_name = 'Priyanka GH' Where emp_id = '120';
InvalidRequest: Error from server: code=2200 [Invalid query] message="Invalid STRING constant (120) for "emp_id" of type int"
cqlsh:employee> update employee_info set emp_name = 'Priyanka GH' Where emp_id=120;
cqlsh:employee> select * from employee_info;
+-----+-----+-----+-----+-----+-----+-----+-----+
| emp_id | date_of_joining | dep_name | designation | emp_name | projects | salary |
+-----+-----+-----+-----+-----+-----+-----+-----+
| 120 | 2024-05-06 | Engineering | Developer | Priyanka GH | {'Project B', 'ProjectA'} | 1e+06
| 123 | 2024-05-07 | Engineering | Engineer | Sadhana | {'Project M', 'Project P'} | 1.2e+06
| 122 | 2024-05-06 | Management | HR | Rachana | {'Project C', 'Project M'} | 9e+05
| 121 | 2024-05-06 | Management | Developer | Shreya | {'Project C', 'ProjectA'} | 9e+05
+-----+-----+-----+-----+-----+-----+-----+-----+
(4 rows)
cqlsh:employee> select * from employee_info order by salary;
InvalidRequest: Error from server: code=2200 [Invalid query] message="ORDER BY is only supported when the partition key is restricted by an EQ or an IN."
cqlsh:employee> alter table employee_info add bonus INT;
cqlsh:employee> select * from employee_info;
+-----+-----+-----+-----+-----+-----+-----+-----+
| emp_id | bonus | date_of_joining | dep_name | designation | emp_name | projects | salary |
+-----+-----+-----+-----+-----+-----+-----+-----+
| 120 | null | 2024-05-06 | Engineering | Developer | Priyanka GH | {'Project B', 'ProjectA'} | 1e+06
| 123 | null | 2024-05-07 | Engineering | Engineer | Sadhana | {'Project M', 'Project P'} | 1.2e+06
| 122 | null | 2024-05-06 | Management | HR | Rachana | {'Project C', 'Project M'} | 9e+05
| 121 | null | 2024-05-06 | Management | Developer | Shreya | {'Project C', 'ProjectA'} | 9e+05
+-----+-----+-----+-----+-----+-----+-----+-----+
(4 rows)
cqlsh:employee> update employee_info set bonus = 12000 where emp_id = 120;
cqlsh:employee> select * from employee_info;
+-----+-----+-----+-----+-----+-----+-----+-----+
| emp_id | bonus | date_of_joining | dep_name | designation | emp_name | projects | salary |
+-----+-----+-----+-----+-----+-----+-----+-----+
| 120 | 12000 | 2024-05-06 | Engineering | Developer | Priyanka GH | {'Project B', 'ProjectA'} | 1e+06
| 123 | null | 2024-05-07 | Engineering | Engineer | Sadhana | {'Project M', 'Project P'} | 1.2e+06
| 122 | null | 2024-05-06 | Management | HR | Rachana | {'Project C', 'Project M'} | 9e+05
| 121 | null | 2024-05-06 | Management | Developer | Shreya | {'Project C', 'ProjectA'} | 9e+05
+-----+-----+-----+-----+-----+-----+-----+-----+
(4 rows)
cqlsh:employee> update employee_info set bonus = 11000 where emp_id = 121;
cqlsh:employee> select * from employee_info using ttl 15 where emp_id = 123;
SyntaxException: line 1:28 mismatched input 'using' expecting EOF (select * from employee_info [using] ttl...)
cqlsh:employee> select * from employee_info where emp_id = 121 using ttl 15;
SyntaxException: line 1:47 no viable alternative at input 'using' (...employee_info where emp_id = 121 [using]...)
cqlsh:employee> update employee_info using ttl 15 set salary = 0 where emp_id = 121;
cqlsh:employee> select * from employee_info;
```

1.2 Experiment - 2

1.2.1 Question:

Perform the following DB operations using Cassandra:

- Create a keyspace by name Library
- 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
- Insert the values into the table in batch
- Display the details of the table created and increase the value of the counter
- Write a query to show that a student with id 112 has taken a book "BDA" 2 times.
- Export the created column to a csv file
- Import a given csv dataset from local file system into Cassandra column family.

02/04/25

Lab 05

Cassandra Exercise :

1. Create a keyspace by name library

```
> create keyspace library with replication = {'class': 'SimpleStrategy', 'replication-factor': 1};
```

2. Create a column family by name Library_Info with attributes Stud_Id, Counter_value, Stud_Name, Book_Name, Book_Id, Date_of_Issue

```
> create table Library_Info (Stud_Id int, Counter_Value counter, Stud_Name text, Book_Name text, Book_Id int, Date_of_Issue date, PRIMARY KEY ((Stud_Id, Book_Id), Stud_Name, Book_Name, Date_of_Issue));
```

3. Insert values into the table

```
> Update Library_Info SET Counter_value = Counter_value + 1 WHERE Stud_Id = 112 AND Book_Id = 101 AND Stud_Name = 'John Doe' AND Date_of_Issue =
```

'2025-04-08';

4. Display the details of the table created

> select * from Library_Info;

Stud_Id	Book_id	Stud_Name	Book_Name	Date_Of_Issue	Counter_Value
113	102	Jane Doe	DEV	2025-04-08	1
112	101	John Doe	BDA	2025-04-08	2
114	103	Joey Doe	ML	2025-03-08	1

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

> select counter_value from Library_Info where Stud_Id = 112 AND Book_Name = 'BDA';

counter_value
2

6. Export the created column to a csv file

> copy Library_Info to 'lib.csv';

3 rows exported to 1 files in 0.110 seconds

1.2.2 Code with Output:

```

bmscecse@bmscecse-HP-Elite-Tower-800-G9-Desktop-PC: ~ cqlsh
Connected to Test Cluster at 127.0.0.1:9042
[cqlsh 6.1.0 | Cassandra 4.1.4 | CQL spec 3.4.6 | Native protocol v5]
Use HELP for help.
cqlsh> CREATE KEYSPACE Students WITH REPLICATION={
    ... 'class':'SimpleStrategy','replication_factor':1};
cqlsh> DESCRIBE KEYSPACES

students      system_auth      system_schema      system_views
system       system_distributed   system_traces      system_virtual_schema

cqlsh> SELECT * FROM system.schema_keyspaces;
InvalidRequest: Error from server: code=2200 [Invalid query] message="table schema_keyspaces does not exist"
cqlsh> use Students;
cqlsh:Students> create table Students_info(Roll_No int Primary key,StudName text,DateOfJoining timestamp,last_exam_Percent double);
cqlsh:Students> describe tables;

students_info

cqlsh:Students> describe table students;
Table 'students' not found in keyspace 'Students'
cqlsh:Students> describe table Students_info;

CREATE TABLE students.Students_info (
    roll_no int PRIMARY KEY,
    dateofjoining timestamp,
    last_exam_percent double,
    studname text
) WITH additional write_policy = '99p'
    AND bloom_filter_fp_chance = 0.01
    AND caching = {'keys': 'ALL', 'rows_per_partition': 'NONE'}
    AND cdc = false
    AND comment = ''
    AND compaction = {'class': 'org.apache.cassandra.db.compaction.SizeTieredCompactionStrategy', 'max_threshold': '32', 'min_threshold': '4'}
    AND compression = {'chunk_length_in_kb': '16', 'class': 'org.apache.cassandra.io.compress.LZ4Compressor'}
    AND memtable = 'default'
    AND crc_check_chance = 1.0
    AND default_time_to_live = 0
    AND extensions = {}
    AND gc_grace_seconds = 8640000
    AND max_index_interval = 2048
    AND memtable_flush_period_in_ms = 0
    AND min_index_interval = 128
    AND read_repair = 'BLOCKING'
    AND speculative_retry = '99p';

```



```

cqlsh:Students> Begin batch insert into Students_info(Roll_no,StudName,DateOfJoining, last_exam_Percent) values(1,'Sadhana','2023-10-09', 98) insert into Students_info(Roll_no,StudName,DateOfJoining, last_exam_Percent) values(2,'Ritu','2023-10-10', 97) insert into Students_info(Roll_no,StudName,DateOfJoining, last_exam_Percent) values(3,'Rachana','2023-10-06', 96.5) apply batch;
cqlsh:Students> select * from Students_info;

roll_no | dateofjoining | last_exam_percent | studname
-----+-----+-----+-----
  1 | 2023-10-09 18:30:00.000000+0000 |      98 | Sadhana
  2 | 2023-10-09 18:30:00.000000+0000 |      97 | Ritu
  4 | 2023-10-05 18:30:00.000000+0000 |     96.5 | Charu
  3 | 2023-10-09 18:30:00.000000+0000 |     97.5 | Rachana

(4 rows)
cqlsh:Students> select * from Students_info where roll_no in (1,2,3);

roll_no | dateofjoining | last_exam_percent | studname
-----+-----+-----+-----
  1 | 2023-10-09 18:30:00.000000+0000 |      98 | Sadhana
  2 | 2023-10-09 18:30:00.000000+0000 |      97 | Ritu
  3 | 2023-10-09 18:30:00.000000+0000 |     97.5 | Rachana

(3 rows)
cqlsh:Students> select * from Students_info where Studname='Charu';
InvalidRequest: Error from server: code=2200 [Invalid query] message="Cannot execute this query as it might involve data filtering and thus may have unpredictable performance. If you want to execute this query despite the performance unpredictability, use ALLOW FILTERING"
cqlsh:Students> create index on Students_info(Studname);
cqlsh:Students> select * from Students_info where Studname='Charu';

roll_no | dateofjoining | last_exam_percent | studname
-----+-----+-----+-----
  4 | 2023-10-05 18:30:00.000000+0000 |     96.5 | Charu

(1 rows)
cqlsh:Students> select Roll_no,StudName from Students_info LIMIT 2;

```

```
(4 rows)
cqlsh:students> select * from students_info where roll_no in (1,2,3);
+-----+-----+-----+
| roll_no | dateofjoining | last_exam_percent | studname |
+-----+-----+-----+
| 1 | 2023-10-08 18:30:00.000000+0000 | 98 | Sadhana |
| 2 | 2023-10-09 18:30:00.000000+0000 | 97 | Rutu |
| 3 | 2023-10-09 18:30:00.000000+0000 | 97.5 | Rachana |
+-----+-----+-----+
(3 rows)
cqlsh:students> select * from students_info where Studname='Charu';
invalid query. Errno: 0 [code=2200] message="Cannot execute this query as it might involve data filtering and thus may have unpredictable performance. If you want to execute this query despite the performance unpredictability, use ALLOW FILTERING"
cqlsh:students> create index on Students_info(StudName);
cqlsh:students> select * from students_info where Studname='Charu';
+-----+-----+-----+
| roll_no | dateofjoining | last_exam_percent | studname |
+-----+-----+-----+
| 4 | 2023-10-05 18:30:00.000000+0000 | 96.5 | Charu |
+-----+-----+-----+
(1 rows)
cqlsh:students> select Roll_no,Studname from students_info LIMIT 2;
+-----+-----+
| roll_no | studname |
+-----+-----+
| 1 | Sadhana |
| 2 | Rutu |
+-----+-----+
(2 rows)
cqlsh:students> SELECT Roll_no as "USN" from Students_info;
+-----+
| USN |
+-----+
| 1 |
| 2 |
| 4 |
| 3 |
+-----+
```

1.3 Experiment - 3

1.3.1 Question:

MongoDB - CRUD Demonstration.

Lab - 1
Working with mongodb

I Create database in mongodb
> use myDB;

Confirm the existence of your database
> db;

To list all databases:
> show dbs;

II CRUD Operations

1 To create a collection by the name "Student":
> db.createCollection("Student");

2 To drop a collection by the name "Student":
> db.Student.drop();

3 Create a table and store some data
> db.Student.insert({id: 1, StudentName: "Michelle", Grade: "VII", Hobbies: ["Cooking"]});

4 Update a document with new values
> db.Student.update({StudentName: "Aryan", id: 1, Grade: "VII"}, { \$set: {Hobbies: ["Skating"]}, upsert: true});

5. FIND Method
A: Search based on certain criteria:

> db.Student.find({ StudName: "Aryan" })

[{

-id: 3,

StudName: "Aryan",

Grade: "VII",

Hobbies: "Skating"

}]

B. To display only student name & grade

> db.Student.find({ \$: { StudName: 1, Grade: 1, _id: 0 } }) ;

If StudName: "Michelle", Grade: "VII",

{ StudName: "Aryan", Grade: "VII" }]

C. Find documents where grade is set to 'VII'

> db.Student.find({ grade: { \$eq: "VII" } }).pretty();

D. Find where hobbies is set to chess or skating

> db.Student.find({ \$in: ["Chess", "Skating"] }).pretty();

E. To find documents from student collection where the Student name begins with 'M'.

> db.Student.find({ StudName: /[^]M/ }).pretty();

F. Student name has "e"

> db.Student.find({ StudName: /e/ }).pretty();

G. Find number of documents

> db.Student.count();

5

H To sort the documents

> db.Student.find().sort({StudentName: -1}).pretty();

III Import data from a csv file

mongoimport --db Student --collection airlines --type csv --headerline --file /home/haider/Desktop/output.csv

IV Export data to a CSV file

mongoexport --host localhost --db Student --collection airlines --csv --out /home/haider/Desktop/output.txt --fields "Year", "Quarter"

V Save method:

> db.Student.save({StudentName: "Virus", Grade: "V"})

VI Add a new field to existing document.

> db.Students.update({l_id: 4}, { \$set: {Education: "Network"} })

VII Remove the field in an existing document:

> db.Students.update({l_id: 4}, { \$unset: {Education: "Network"} })

IX Set a particular value to NULL

> db.Student.update({l_id: 3}, { \$set: {Education: null} })

X Aggregate function

```
> db. Student.aggregate(  
  { group: {  
    - id: '$Grade',  
    total: { $sum: '$Fees' }  
  }  
});  
[  
  { -id: 'vii', total: 4000 },  
  { -id: 'I', total: 1000 },  
  { -id: 'VI', total: 2000 },  
  { -id: 'IX', total: 2000 },  
  { -id: 'VIII', total: 3000 }  
]
```

Q8
a/3/2

Customer data Lab 2

1 Customer data base:

1. Create a collection by name Customers

>use lab2

>db.createCollection("Customers")

2. Insert atleast 5 values into the table.

>db.Customers.insertMany([{"Cust_id": "001", "Acc_Bal": 1500, "Acc_Type": "s"}, {"Cust_id": "002", "Acc_Bal": 1000, "Acc_Type": "s"}, {"Cust_id": "003", "Acc_Bal": 8500, "Acc_Type": "c"}, {"Cust_id": "004", "Acc_Bal": 1800, "Acc_Type": "s"}, {"Cust_id": "005", "Acc_Bal": 950, "Acc_Type": "s"}])

3. Write a query to display those records where total account balance is greater than 1000 of account type 's' for each customer.

> db.Customers.find({\$Acc_Bal: { \$gt: 1000 }, Acc_Type: "s"})

O/P:

[{"Cust_id": "001", "Acc_Bal": 1500, "Acc_Type": "s"}, {"Cust_id": "002", "Acc_Bal": 1000, "Acc_Type": "s"}, {"Cust_id": "005", "Acc_Bal": 950, "Acc_Type": "s"}]

[5]

6. Determine min and max acc balance for each customer

> db. Customers.aggregate([

{ \$group: {

-id: "Customer id"

min_balance: { \$min: "\$Acc_Bal" },

max_balance: { \$max: "\$Acc_Bal" }

}]])

O/P:

[{ id: '005', min: 950, max: 950 },

{ id: '001', min: 1500, max: 1500 },

{ id: '002', min: 800, max: 800 },

{ id: '003', min: 2000, max: 2000 },

{ id: '004', min: 100, max: 100 }]

7. You are developing an e-commerce platform where users can browse and purchase products. Each product has a unique identifier, a name, a category, a price, and available quantity.

1. Retrieve all products:

> db. Products.find()

2. Retrieve products in a specific category:

> db. Products.find({ category: "Electronics" })

3. Retrieve products with quantity greater than 0.
 - > db.Products.find({ quantity: { \$gt: 0 } })
4. Retrieve products with price less than or equal to \$100
 - > db.Products.find({ price: { \$lte: 100 } })
5. Retrieve products sorted by price in ascending order
 - > db.Products.find().sort({ price: 1 })
6. Retrieve products added to a user's cart
 - > db.Carts.find({ user_id: "123abc" })
7. Retrieve orders placed by a user
 - > db.Orders.find({ user_id: "123abc" })
8. Retrieve total price of orders placed by a user
 - > db.Orders.aggregate([
 { \$match: { user_id: "123abc" } },
 { \$group: { _id: "\$user_id", total_spent: { \$sum: "\$total_price" } } }
]);

O/P:

[{ _id: "123abc", total_spent: 1798 }]

11 additional aggregation queries

1. Calculate total number of products in each category.

> db.Products.aggregate([

{ \$group : { category_id: "category" },
total_products : { \$sum : 1 } }]);

2. Calculate total price of products in each category

> db.Products.aggregate([

{ \$group : { _id: "category" }, total_price :
{ \$sum : "\$price" } }]);

3. Find average price of products

> db.Products.aggregate([

{ \$group : { _id: null, avg_price : { \$avg : "\$price" } } }]);

4. Find products with quantity less than 10

> db.Products.find({ \$quantity : { \$lt : 10 } });

5. Sort products by price in descending order

> db.Products.find().sort({ price : -1 });

6 Calculate Total price of orders placed by each user

```
> db.orders.aggregate([  
  { $group:  
    { id: "$user_id",  
     total_spent: {  
       $sum: "$total_price"  
     }  
   }  
]);
```

O/P:

```
[ { id: '456def', total_spent: 403,  
  { id: '123abc', total_spent: 1298 },  
  { id: '789ghi', total_spent: 1848 }  
 ]
```

7 Find user with the highest total price of orders

```
> db.orders.aggregate([  
  { $group:  
    { id: "$user_id",  
     total_spent: {  
       $sum: "$total_price"  
     }  
   }  
]);
```

```
{ $sort: { total_spent: -1 }  
 { $limit: 1 }  
});
```

O/P:

```
[ { id: '789ghi', total_spent: 1848 } ]
```

8. Find average total price of orders

> db.Orders aggregate ([
 \$group : {
 - id: null, avg_total_price: \$total
 },
 \$avg : "\$total-price"
 }]);

OIP:
[{ id: null, avg_total_price: 737.5 }]

~~return db.Orders aggregate ([
 \$group : {
 - id: null, avg_total_price: \$total
 },
 \$avg : "\$total-price"
 }]);~~

1.3.2 Code with Output:

1. Create a database “Student” with the following attributes Rollno, Name , Age, ContactNo, Email-Id, grade, hobby:
use Students

2. Insert 5 appropriate values according to the below queries.

```
db.students.insertMany([
  { "Rollno": 10, "Name": "John", "Age": 20, "ContactNo": "1234567890", "Email-Id": "john@example.com", "grade": "A", "hobby": "Reading" },
  { "Rollno": 11, "Name": "Alice", "Age": 21, "ContactNo": "9876543210", "Email-Id": "alice@example.com", "grade": "B", "hobby": "Painting" },
  { "Rollno": 12, "Name": "Bob", "Age": 22, "ContactNo": "2345678901", "Email-Id": "bob@example.com", "grade": "C", "hobby": "Cooking" },
  { "Rollno": 13, "Name": "Eve", "Age": 23, "ContactNo": "3456789012", "Email-Id": "eve@example.com", "grade": "A" },
  { "Rollno": 14, "Name": "Charlie", "Age": 24, "ContactNo": "4567890123", "Email-Id": "charlie@example.com", "hobby": "Gardening" }
])
```

```
Atlas atlas-wanmtx-shard-0 [primary] Student> use Students
switched to db Students
Atlas atlas-wanmtx-shard-0 [primary] Students> show collections

Atlas atlas-wanmtx-shard-0 [primary] Students> db.students.insertMany([
...   { "Rollno": 10, "Name": "John", "Age": 20, "ContactNo": "1234567890", "Email-Id": "john@example.com", "grade": "A", "hobby": "Reading" },
...   { "Rollno": 11, "Name": "Alice", "Age": 21, "ContactNo": "9876543210", "Email-Id": "alice@example.com", "grade": "B", "hobby": "Painting" },
...   { "Rollno": 12, "Name": "Bob", "Age": 22, "ContactNo": "2345678901", "Email-Id": "bob@example.com", "grade": "C", "hobby": "Cooking" },
...   { "Rollno": 13, "Name": "Eve", "Age": 23, "ContactNo": "3456789012", "Email-Id": "eve@example.com", "grade": "A" },
...   { "Rollno": 14, "Name": "Charlie", "Age": 24, "ContactNo": "4567890123", "Email-Id": "charlie@example.com", "hobby": "Gardening" }
... ])
{
  acknowledged: true,
  insertedIds: {
    '0': ObjectId("661ce9dc76a00ff8cc51dae1"),
    '1': ObjectId("661ce9dc76a00ff8cc51dae2"),
    '2': ObjectId("661ce9dc76a00ff8cc51dae3"),
    '3': ObjectId("661ce9dc76a00ff8cc51dae4"),
    '4': ObjectId("661ce9dc76a00ff8cc51dae5")
  }
}
```

3. Write query to update Email-Id of a student with rollno 10.

```
db.students.updateOne(
  { "Rollno": 10 },
  { $set: { "Email-Id": "john.doe@example.com" } }
)
```

```

Atlas atlas-wanmtx-shard-0 [primary] Students> db.students.updateOne(
...   { "Rollno": 10 },
...   { $set: { "Email-Id": "john.doe@example.com" } }
...
{
  acknowledged: true,
  insertedId: null,
  matchedCount: 1,
  modifiedCount: 1,
  upsertedCount: 0
}

```

4. Replace the student name from “Alice” to “Alicee” of rollno 11

```

db.students.updateOne(
  { "Rollno": 11 },
  { $set: { "Name": "Alicee" } }
)
Atlas atlas-wanmtx-shard-0 [primary] Students> db.students.updateOne(
...   { "Rollno": 11 },
...   { $set: { "Name": "Alicee" } }
...
{
  acknowledged: true,
  insertedId: null,
  matchedCount: 1,
  modifiedCount: 1,
  upsertedCount: 0
}

```

5. Display Student Name and grade(Add if grade is not present)where the _id column is 1.

```

db.students.find({}, { "Name": 1, "grade": { $ifNull: ["$grade", "Not available"] }, "_id": 0 })
Atlas atlas-wanmtx-shard-0 [primary] Students> db.students.find({}, { "Name": 1, "grade": {
  { $ifNull: ["$grade", "Not available"] }, "_id": 0 })
[
  { Name: 'John', grade: 'A' },
  { Name: 'Alicee', grade: 'B' },
  { Name: 'Bob', grade: 'C' },
  { Name: 'Eve', grade: 'A' },
  { Name: 'Charlie', grade: 'Not available' }
]

```

6. Update to add hobbies

```

db.students.updateMany(
  { "Name": "Eve" },
  { $set: { "hobby": "Dancing" } }
)

```

```
Atlas atlas-wanmtx-shard-0 [primary] Students> db.students.updateMany(  
...     { "Name": "Eve" },  
...     { $set: { "hobby": "Dancing" } }  
... )  
{  
    acknowledged: true,  
    insertedId: null,  
    matchedCount: 1,  
    modifiedCount: 1,  
    upsertedCount: 0  
}
```

7. Find documents where hobbies is set neither to Chess nor to Skating

```
db.students.find({ "hobby": { $nin: ["Chess", "Skating"] } })
```

```
Atlas atlas-wanmtx-shard-0 [primary] Students> db.students.find({ "hobby": { $nin: ["Chess", "Skating"] } })  
[  
  {  
    _id: ObjectId("661ce9dc76a00ff8cc51dae1"),  
    Rollno: 10,  
    Name: 'John',  
    Age: 20,  
    ContactNo: '1234567890',  
    'Email-Id': 'john.doe@example.com',  
    grade: 'A',  
    hobby: 'Reading'  
  },  
  {  
    _id: ObjectId("661ce9dc76a00ff8cc51dae2"),  
    Rollno: 11,  
    Name: 'Alicee',  
    Age: 21,  
    ContactNo: '9876543210',  
    'Email-Id': 'alice@example.com',  
    grade: 'B',  
    hobby: 'Painting'  
  },  
  {  
    _id: ObjectId("661ce9dc76a00ff8cc51dae3"),  
    Rollno: 12,  
    Name: 'Bob',  
    Age: 22,  
    ContactNo: '2345678901',  
    'Email-Id': 'bob@example.com',  
    grade: 'C',  
    hobby: 'Cooking'  
  },  
]
```

8. Find documents whose name begins with A

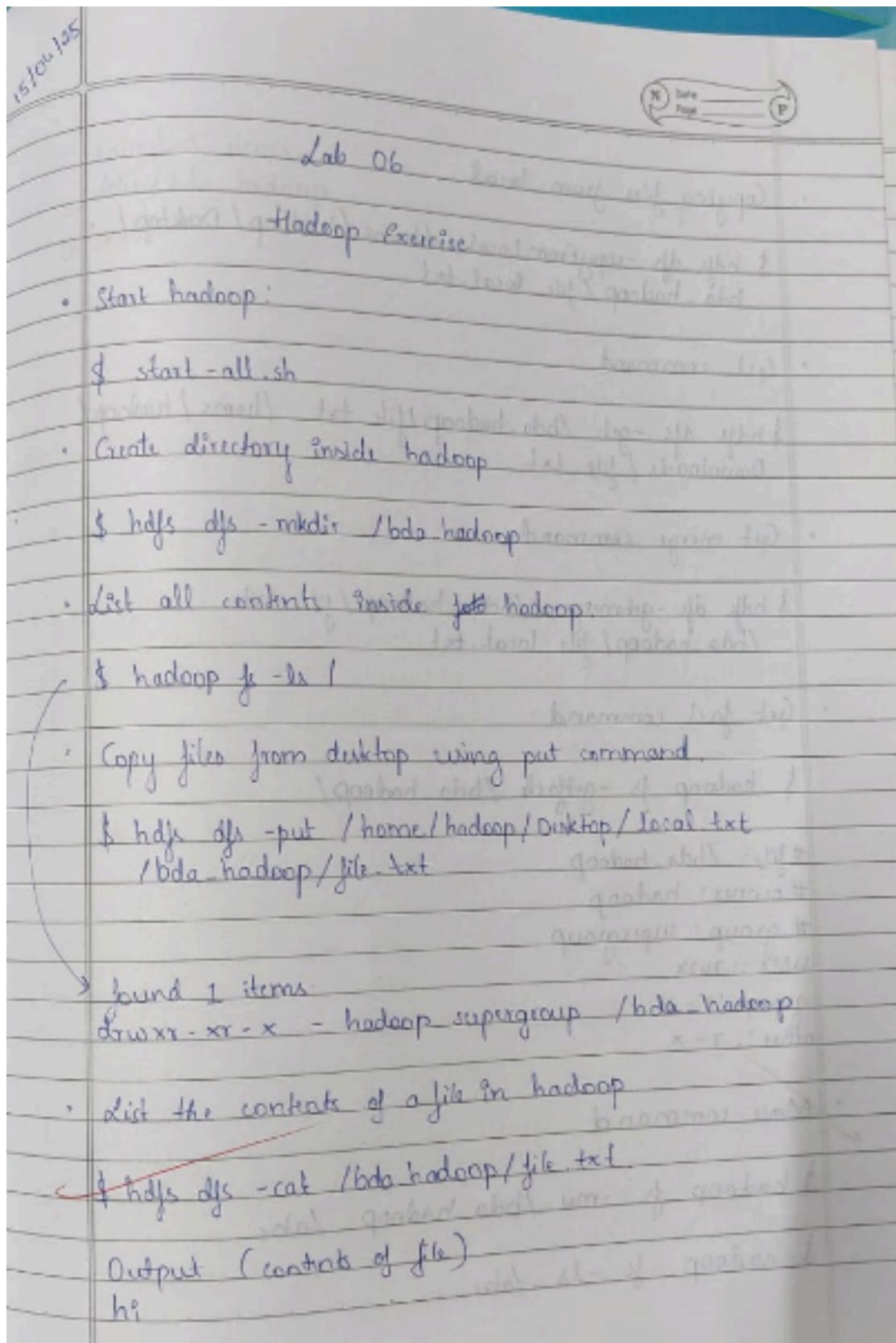
```
db.students.find({ "Name": /^A/ })
```

```
Atlas atlas-wanmtx-shard-0 [primary] Students> db.students.find({ "Name": /^A/ })
[
  {
    _id: ObjectId("661ce9dc76a00ff8cc51dae2"),
    Rollno: 11,
    Name: 'Alicee',
    Age: 21,
    ContactNo: '9876543210',
    'Email-Id': 'alice@example.com',
    grade: 'B',
    hobby: 'Painting'
  }
]
```

Experiment - 5

1.3.3 Question:

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



- Copying files from local

\$ hdfs dfs -copyFromLocal /home/hadoop/Desktop/
bda-hadoop/file local.txt

- Get command

\$ hdfs dfs -get /bda-hadoop/file.txt /home/hadoop/
Downloads/file.txt

- Get merge command

\$ hdfs dfs -getmerge /bda-hadoop/file.txt /home/
bda-hadoop/file.local.txt

- Get fsck command

\$ hadoop fs -getfacl /bda-hadoop/

#file: /bda-hadoop

#owner: hadoop

#group: supergroup

user::rwx

group::r-x

other::r-x

- Move command

\$ hadoop fs -mv /bda-hadoop/abc

\$ hadoop fs -ls /abc

~~forward 1 items~~

~~use 1 hda hadoop~~

~~only 1 item~~

1.3.4 Code with Output:

```
hadoop@bmscecse-HP-Elite-Tower-800-G9-Desktop-PC:~$ cd ./Desktop/
hadoop@bmscecse-HP-Elite-Tower-800-G9-Desktop-PC:~/Desktop$ start-all.sh
WARNING: Attempting to start all Apache Hadoop daemons as hadoop in 10 seconds.
WARNING: This is not a recommended production deployment configuration.
WARNING: Use CTRL-C to abort.
Starting namenodes on [localhost]
Starting datanodes
Starting secondary namenodes [bmscecse-HP-Elite-Tower-800-G9-Desktop-PC]
Starting resourcemanager
Starting nodemanagers
hadoop@bmscecse-HP-Elite-Tower-800-G9-Desktop-PC:~/Desktop$ hdfs dfs -mkdir /Lab05
```

```
hadoop@bmscecse-HP-Elite-Tower-800-G9-Desktop-PC:~/Desktop$ hadoop fs -ls /Hadoop
ls: '/Hadoop': No such file or directory
hadoop@bmscecse-HP-Elite-Tower-800-G9-Desktop-PC:~/Desktop$ hadoop fs -ls /Lab05
```

```
hadoop@bmscecse-HP-Elite-Tower-800-G9-Desktop-PC:~/Desktop$ touch test.txt
hadoop@bmscecse-HP-Elite-Tower-800-G9-Desktop-PC:~/Desktop$ nano text.txt
hadoop@bmscecse-HP-Elite-Tower-800-G9-Desktop-PC:~/Desktop$ hdfs dfs -put ./text.txt /Lab05/text.txt
hadoop@bmscecse-HP-Elite-Tower-800-G9-Desktop-PC:~/Desktop$ hadoop fs -ls /Lab05
Found 1 items
-rw-r--r-- 1 hadoop supergroup 19 2024-05-13 14:33 /Lab05/text.txt
hadoop@bmscecse-HP-Elite-Tower-800-G9-Desktop-PC:~/Desktop$ hdfs dfs -cat /Lab05/text.txt
Hello
How are you?
```

```
hadoop@bmscecse-HP-Elite-Tower-800-G9-Desktop-PC:~/Desktop$ hadoop fs -ls /Lab05
Found 2 items
-rw-r--r-- 1 hadoop supergroup 15 2024-05-13 14:40 /Lab05/test.txt
-rw-r--r-- 1 hadoop supergroup 19 2024-05-13 14:33 /Lab05/text.txt
hadoop@bmscecse-HP-Elite-Tower-800-G9-Desktop-PC:~/Desktop$ hdfs dfs -getmerge /Lab05 /text.txt /Lab05 /test.txt ..
Downloads/Merged.txt
getmerge: '/text.txt': No such file or directory
getmerge: '/test.txt': No such file or directory
hadoop@bmscecse-HP-Elite-Tower-800-G9-Desktop-PC:~/Desktop$ hdfs dfs -getmerge /Lab05/text.txt /Lab05/test.txt ..
Downloads/Merged.txt
hadoop@bmscecse-HP-Elite-Tower-800-G9-Desktop-PC:~/Desktop$ hadoop fs -getfacl /Lab05
# file: /Lab05
# owner: hadoop
# group: supergroup
user::rwx
group::r-x
other::r-x
```

```
hadoop@bmscecse-HP-Elite-Tower-800-G9-Desktop-PC:~/Desktop$ hdfs dfs -copyToLocal /Lab05/text.txt ..//Documents
hadoop@bmscecse-HP-Elite-Tower-800-G9-Desktop-PC:~/Desktop$ hdfs dfs -copyToLocal /Lab05/test.txt ..//Documents
```

```
hadoop@bmscecse-HP-Elite-Tower-800-G9-Desktop-PC:~/Desktop$ hdfs dfs -cat /Lab05/text.txt
Hello
How are you?
hadoop@bmscecse-HP-Elite-Tower-800-G9-Desktop-PC:~/Desktop$ hdfs dfs -mv /Lab05 /test_Lab05
```

```
hadoop@bmscecse-HP-Elite-Tower-800-G9-Desktop-PC:~/Desktop$ hdfs dfs -ls /test_Lab05
Found 2 items
-rw-r--r-- 1 hadoop supergroup      15 2024-05-13 14:40 /test_Lab05/test.txt
-rw-r--r-- 1 hadoop supergroup      19 2024-05-13 14:33 /test_Lab05/text.txt
hadoop@bmscecse-HP-Elite-Tower-800-G9-Desktop-PC:~/Desktop$ hdfs dfs -cp /test_Lab05/ /Lab05
hadoop@bmscecse-HP-Elite-Tower-800-G9-Desktop-PC:~/Desktop$ hdfs dfs -ls /Lab05
Found 2 items
-rw-r--r-- 1 hadoop supergroup      15 2024-05-13 14:51 /Lab05/test.txt
-rw-r--r-- 1 hadoop supergroup      19 2024-05-13 14:51 /Lab05/text.txt
hadoop@bmscecse-HP-Elite-Tower-800-G9-Desktop-PC:~/Desktop$ hdfs dfs -ls /test_Lab05
Found 2 items
-rw-r--r-- 1 hadoop supergroup      15 2024-05-13 14:40 /test_Lab05/test.txt
-rw-r--r-- 1 hadoop supergroup      19 2024-05-13 14:33 /test_Lab05/text.txt
```

Experiment - 6

1.3.5 Question:

Implement WordCount Program on Hadoop framework.

Lab 07

task 1: hadoop
granted what

Wordcount - MapReduce

→ WCMapper.java:

```
import java.io.IOException;  
import org.apache.hadoop.io;  
import org.apache.hadoop.mapred;  
  
public class WCMapper extends MapReduceBase implements  
Mapper<LongWritable, Text, Text, IntWritable> {  
    public void map(LongWritable key, Text value,  
    OutputCollector<Text, IntWritable> output,  
    Reporter reporter) throws IOException {  
        String line = value.toString();  
        for (String word : line.split(" ")) {  
            if (word.length() > 0) {  
                output.collect(new Text(word),  
                new IntWritable(1));  
            }  
        }  
    }  
}
```

→ WCReducer.java:

```
import java.io.IOException;  
import java.util.Iterator;  
import org.apache.hadoop.io;  
import org.apache.hadoop.mapred;
```

N Date _____ P

```

public class WCReducer extends MapReduceBase
    implements Reducer<Text, IntWritable, Text, IntWritable> {
    public void reduce(Text key, Iterator<IntWritable>
        value, OutputCollector<Text, IntWritable> output,
        Reporter rps) throws IOException {
        int count = 0;
        while (value.hasNext()) {
            IntWritable i = value.next();
            count += i.get();
        }
        output.collect(key, new IntWritable(count));
    }
}

```

→ WC Driver.java

```

import java.io.IOException;
import org.apache.hadoop.*;

public class WCDriver extends Configured implements Tool {
    public int run(String args[]) throws Exception {
        if (args.length < 2) {
            System.out.println("Please give valid inputs");
            return -1;
        }
        JobConf conf = new JobConf(WCDriver.class);
        FileInputFormat.setInputPath(conf,
            new Path(args[0]));
        FileOutputFormat.setOutputPath(conf,
            new Path(args[1]));
        conf.setMapperClass(WCMapper.class);

```

```
conf.setReducerClass(wcReducer.class);
conf.setMapOutputKeyClass(Text.class);
conf.setMapOutputValueClass(IntWritable.class);
conf.setOutputKeyClass(Text.class);
conf.setOutputValueClass(IntWritable.class);
JobClient.runJob(conf);
return 0;
```

```
public static void main(String args[])
throws Exception {
    int exitCode = ToolRunner.run(new
        wcDriver(), args);
    System.out.println(exitCode);
}
```

Aug 2011/5

```

Activities Terminal May 20 14:47
hadoop@bmscsece-HP-Elite-Tower-600-G9-Desktop-PC: ~
hadoop@bmscsece-HP-Elite-Tower-600-G9-Desktop-PC: ~
hadoop@bmscsece-HP-Elite-Tower-600-G9-Desktop-PC: $ start-all.sh
WARNING: Attempting to start all Apache Hadoop daemons as hadoop in 10 seconds.
WARNING: This is not a recommended production deployment configuration.
WARNING: Use CTRL-C to abort.
Starting namenodes on [localhost]
localhost: namenode is running as process 8499. Stop it first and ensure /tmp/hadoop-hadoop-namenode.pid file is empty before retry.
localhost: datanode is running as process 8673. Stop it first and ensure /tmp/hadoop-hadoop-datanode.pid file is empty before retry.
Starting secondary namenodes [bmscsece-HP-Elite-Tower-600-G9-Desktop-PC]
bmscsece-HP-Elite-Tower-600-G9-Desktop-PC: secondarynamenode is running as process 8959. Stop it first and ensure /tmp/hadoop-hadoop-secondarynamenode.pid file is empty before retry.
Starting resourcemanager
resourcemanager is running as process 9238. Stop it first and ensure /tmp/hadoop-hadoop-resourcemanager.pid file is empty before retry.
Starting nodemanagers
localhost: nodemanager is running as process 9399. Stop it first and ensure /tmp/hadoop-hadoop-nodemanager.pid file is empty before retry.
Starting datanodes
hadoop@bmscsece-HP-Elite-Tower-600-G9-Desktop-PC: ~
nano /home/hadoop/hadoop/etc/hadoop/mapred-site.xml
hadoop@bmscsece-HP-Elite-Tower-600-G9-Desktop-PC: ~ stop-all.sh
WARNING: Stopping all Apache Hadoop daemons as hadoop in 10 seconds.
WARNING: Use CTRL-C to abort.
Stopping namenodes on [localhost]
Stopping datanodes
Stopping secondary namenodes [bmscsece-HP-Elite-Tower-600-G9-Desktop-PC]
Stopping nodemanagers
Stopping resourcemanager
hadoop@bmscsece-HP-Elite-Tower-600-G9-Desktop-PC: ~ start-all.sh
WARNING: Attempting to start all Apache Hadoop daemons as hadoop in 10 seconds.
WARNING: This is not a recommended production deployment configuration.
Starting namenodes on [localhost]
Starting datanodes
Starting secondary namenodes [bmscsece-HP-Elite-Tower-600-G9-Desktop-PC]
Starting resourcemanager
Starting nodemanagers
hadoop@bmscsece-HP-Elite-Tower-600-G9-Desktop-PC: ~ Jps
14785 DataNode
15107 SecondaryNameNode
15989 Jps
15990 ResourceManager
15741 NodeManager
0270 org.eclipse.equinox.launcher_1.6.1000.v20250227-1734.jar
14591 NameNode
hadoop@bmscsece-HP-Elite-Tower-600-G9-Desktop-PC: ~ hadoop fs -ls /
Found 3 items
drwxr-xr-x - hadoop supergroup 0 2025-05-20 13:48 /folder1
drwxr-xr-x - hadoop supergroup 0 2025-05-20 13:48 /folder2
drwxr-xr-x - hadoop supergroup 0 2025-05-20 13:43 /tmp
hadoop@bmscsece-HP-Elite-Tower-600-G9-Desktop-PC: ~ hadoop fs -mkdirr /rgs
hadoop@bmscsece-HP-Elite-Tower-600-G9-Desktop-PC: ~ hadoop fs -copyFromLocal /home/hadoop/Desktop/sample.txt /rgs/test.txt
hadoop@bmscsece-HP-Elite-Tower-600-G9-Desktop-PC: ~ hadoop jar /home/hadoop/Desktop/wordcount.jar WordCount.WCDriver /rgs/test.txt /rgs/output
2025-05-20 14:45:00,274 INFO impl.MetricsConfig: Loaded properties from hadoop-metrics2.properties
2025-05-20 14:45:00,315 INFO impl.MetricSystemImpl: Scheduled Metric snapshot period at 10 second(s).
2025-05-20 14:45:00,315 INFO impl.MetricSystemImpl: JobTracker metrics system started
2025-05-20 14:45:00,321 WARN impl.MetricSystemImpl: JobTracker metrics system already initialized!
2025-05-20 14:45:00,436 INFO Mapred.FileInputFormat: Total input files to process : 1
2025-05-20 14:45:00,469 INFO mapreduce.JobSubmitter: number of splits:1
Activities Terminal May 20 14:48
hadoop@bmscsece-HP-Elite-Tower-600-G9-Desktop-PC: ~
hadoop@bmscsece-HP-Elite-Tower-600-G9-Desktop-PC: ~
hadoop@bmscsece-HP-Elite-Tower-600-G9-Desktop-PC: ~
HDFS: Number of bytes written=86
HDFS: Number of read operations=15
HDFS: Number of large read operations=0
HDFS: Number of write operations=4
HDFS: Number of bytes read erasure-coded=0
Map-Reduce Framework
  Map Input records=4
  Map output records=13
  Map output bytes=116
  Map output materialized bytes=148
  Input split bytes=80
  Combine input records=0
  Combine output records=0
  Reduce input groups=2
  Redundant output bytes=148
  Reduce input records=13
  Reduce output records=12
  Spilled Records=26
  Shuffled Maps =1
  Failed Shuffles=0
  Merged Map outputs=1
  GC time elapsed (ms)=0
  Total committed heap usage (bytes)=1375731712
Shuffle Errors
  BAD_ID=0
  CONNECTION=0
  IO_ERROR=0
  WRONG_LENGTH=0
  WRONG_MAP=0
  WRONG_PARTITION=0
  File Input Format Counters
    Bytes Read=66
    File Output Format Counters
      Bytes Written=86
hadoop@bmscsece-HP-Elite-Tower-600-G9-Desktop-PC: ~ hadoop fs -ls /output/
ls: '/output/': No such file or directory
hadoop@bmscsece-HP-Elite-Tower-600-G9-Desktop-PC: ~ hadoop fs -ls /rgs/output/
Found 2 items
-rw-r--r-- 1 hadoop supergroup 0 2025-05-20 14:45 /rgs/output/_SUCCESS
-rw-r--r-- 1 hadoop supergroup 86 2025-05-20 14:45 /rgs/output/part-00000
hadoop@bmscsece-HP-Elite-Tower-600-G9-Desktop-PC: ~ hadoop fs -cat /rgs/output/part-00000
am 1
are 1
beez 1
executed 1
feeling 1
good 1
hiiii 1
how 1
i 2
program 1
the 1
you 1
hadoop@bmscsece-HP-Elite-Tower-600-G9-Desktop-PC: ~ ss

```

1.3.6 Code with Output:

Mapper Code:

```
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.mapred.MapReduceBase;
import org.apache.hadoop.mapred.Mapper;
import org.apache.hadoop.mapred.OutputCollector;
import org.apache.hadoop.mapred.Reporter;

public class WCMapper extends MapReduceBase implements Mapper<LongWritable, Text, Text,
IntWritable> {
    public void map(LongWritable key, Text value, OutputCollector<Text, IntWritable> output, Reporter rep)
throws IOException
    {
        String line = value.toString();
        for (String word : line.split(" "))
        {
            if (word.length() > 0)
            {
                output.collect(new Text(word), new IntWritable(1));
            }
        }
    }
}
```

Reducer Code:

```
// Importing libraries import
java.io.IOException; import
java.util.Iterator;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapred.MapReduceBase;
import org.apache.hadoop.mapred.OutputCollector;
import org.apache.hadoop.mapred.Reducer;
import org.apache.hadoop.mapred.Reporter;
public class WCReducer extends MapReduceBase implements Reducer<Text, IntWritable, Text,
IntWritable> {
    // Reduce function
    public void reduce(Text key, Iterator<IntWritable> value,
OutputCollector<Text, IntWritable> output,
Reporter rep) throws IOException
    {
        int count = 0;
        // Counting the frequency of each words
        while (value.hasNext())
        {
            IntWritable i = value.next();
            count += i.get();
        }
        output.collect(key, new IntWritable(count));
    }
}
```

} }

Driver Code: WCDriver Java Class file.

```
import java.io.IOException;
import org.apache.hadoop.conf.Configured;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapred.FileInputFormat;
import org.apache.hadoop.mapred.FileOutputFormat;
import org.apache.hadoop.mapred.JobClient;
import org.apache.hadoop.mapred.JobConf;
import org.apache.hadoop.util.Tool;
import org.apache.hadoop.util.ToolRunner;
public class WCDriver extends Configured implements Tool {
    public int run(String args[]) throws IOException
    {
        if (args.length < 2)
        {
            System.out.println("Please give valid inputs");
            return -1;
        }
        JobConf conf = new JobConf(WCDriver.class);
        FileInputFormat.setInputPaths(conf, new Path(args[0]));
        FileOutputFormat.setOutputPath(conf, new Path(args[1]));
        conf.setMapperClass(WCMapper.class);
        conf.setReducerClass(WCReducer.class);
        conf.setMapOutputKeyClass(Text.class);
        conf.setMapOutputValueClass(IntWritable.class);
        conf.setOutputKeyClass(Text.class);
        conf.setOutputValueClass(IntWritable.class);
        JobClient.runJob(conf);
        return 0;
    }
    public static void main(String args[]) throws Exception
    {
        int exitCode = ToolRunner.run(new WCDriver(), args);
        System.out.println(exitCode);
    }
}
```

1.4 Experiment - 7

1.4.1 Question:

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:

- c) Find average temperature for each year from NCDC data set.
- d) Find the mean max temperature for every month.

1.4.2 Code with Output:

a) Find average temperature for each year from NCDC

data set. 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));
}}

```

```

.:~$ hadoop jar C:/avgtemp.jar temp.AverageDriver /input_dir/temp.txt /avgtemp_outputdir
2021-05-15 14:52:50,635 INFO Client.DefaultNamenodeFailoverProxyProvider: Connecting to ResourceManager at /0.0.0.0:8032
2021-05-15 14:52:51,005 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to remedy this.
2021-05-15 14:52:51,111 INFO mapreduce.JobResourceUploader: Disabling Erasure Coding for path: /tmp/hadoop-yarn/staging/Anusree/.staging/job_1621060230696_0005
2021-05-15 14:52:51,735 INFO input.FileInputFormat: Total input files to process : 1
2021-05-15 14:52:52,751 INFO mapreduce.JobSubmitter: number of splits:1
2021-05-15 14:52:53,073 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1621060230696_0005
2021-05-15 14:52:53,073 INFO mapreduce.JobSubmitter: Executing with tokens: []
2021-05-15 14:52:53,237 INFO conf.Configuration: resource-types.xml not found
2021-05-15 14:52:53,238 INFO resource.ResourceUtils: Unable to find 'resource-types.xml'.
2021-05-15 14:52:53,312 INFO impl.YarnClientImpl: Submitted application application_1621060230696_0005
2021-05-15 14:52:53,352 INFO mapreduce.Job: The url to track the job: http://LAPTOP-JG329ESD:8088/proxy/application_1621060230696_0005/
2021-05-15 14:52:53,353 INFO mapreduce.Job: Running job: job_1621060230696_0005
2021-05-15 14:53:06,640 INFO mapreduce.Job: Job job_1621060230696_0005 running in uber mode : false
2021-05-15 14:53:06,643 INFO mapreduce.Job: map 0% reduce 0%
2021-05-15 14:53:12,758 INFO mapreduce.Job: map 100% reduce 0%
2021-05-15 14:53:19,880 INFO mapreduce.Job: map 100% reduce 100%
2021-05-15 14:53:25,967 INFO mapreduce.Job: Job job_1621060230696_0005 completed successfully
2021-05-15 14:53:26,896 INFO mapreduce.Job: Counters: 54
File System Counters
    FILE: Number of bytes read=72210
    FILE: Number of bytes written=674341
    FILE: Number of read operations=0
    FILE: Number of large read operations=0
    FILE: Number of write operations=0
    HDFS: Number of bytes read=894860
    HDFS: Number of bytes written=8
    HDFS: Number of read operations=8
    HDFS: Number of large read operations=0
    HDFS: Number of write operations=2
    HDFS: Number of bytes read erasure-coded=0
Job Counters
    Launched map tasks=1
    Launched reduce tasks=1
    Data-local map tasks=1
    Total time spent by all maps in occupied slots (ms)=3782

```

```
C:\hadoop-3.3.0\sbin>hdfs dfs -ls /avgtemp_outputdir
Found 2 items
-rw-r--r-- 1 Anusree supergroup          0 2021-05-15 14:53 /avgtemp_outputdir/_SUCCESS
-rw-r--r-- 1 Anusree supergroup          8 2021-05-15 14:53 /avgtemp_outputdir/part-r-00000

C:\hadoop-3.3.0\sbin>hdfs dfs -cat /avgtemp_outputdir/part-r-00000
1901    46

C:\hadoop-3.3.0\sbin>
```

b) find the mean max temperature for every month 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));
}
}

```

```

C:\hadoop-3.3.0\sbin>hadoop jar C:\meanmax.jar meanmax.MeanMaxDriver /input_dir/temp.txt /meanmax_output
2021-05-21 20:28:05,258 INFO client.DefaultNoHARNFailoverProxyProvider: Connecting to ResourceManager at /0.0.0.0:8088
2021-05-21 20:28:06,662 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to remedy this.
2021-05-21 20:28:06,916 INFO mapreduce.JobResourceUploader: Disabling Erasure Coding for path: /tmp/hadoop-yarn/staging/Anusree/.staging/job_1621688943095_0001
2021-05-21 20:28:08,425 INFO input.FileInputFormat: Total input files to process :1
2021-05-21 20:28:09,167 INFO mapreduce.JobSubmitter: number of splits:1
2021-05-21 20:28:09,741 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1621688943095_0001
2021-05-21 20:28:09,741 INFO mapreduce.JobSubmitter: Executing with tokens: []
2021-05-21 20:28:10,029 INFO conf.Configuration: resource-types.xml not found
2021-05-21 20:28:10,030 INFO resource.ResourcesUtils: Unable to find 'resource-types.xml'
2021-05-21 20:28:10,676 INFO impl.YarnClientImpl: Submitted application application_1621688943095_0001
2021-05-21 20:28:11,005 INFO mapreduce.Job: The url to track the job: http://LAPTOP-JG329ESD:8088/proxy/application_1621688943095_0001/
2021-05-21 20:28:11,006 INFO mapreduce.Job: Running job: job_1621688943095_0001
2021-05-21 20:28:29,385 INFO mapreduce.Job: Job job_1621688943095_0001 running in uber mode : false
2021-05-21 20:28:29,389 INFO mapreduce.Job: map 0% reduce 0%
2021-05-21 20:28:40,664 INFO mapreduce.Job: map 100% reduce 0%
2021-05-21 20:28:50,832 INFO mapreduce.Job: map 100% reduce 100%
2021-05-21 20:28:58,965 INFO mapreduce.Job: Job job_1621688943095_0001 completed successfully
2021-05-21 20:28:59,178 INFO mapreduce.Job: Counters: 54
File System Counters
FILE: Number of bytes read=59882
FILE: Number of bytes written=648891
FILE: Number of read operations=0
FILE: Number of large read operations=0
FILE: Number of write operations=0
HDFS: Number of bytes read=834860
HDFS: Number of bytes written=74
HDFS: Number of read operations=8
HDFS: Number of large read operations=0
HDFS: Number of write operations=2
HDFS: Number of bytes read erasure-coded=0
Job Counters
Launched map tasks=1
Launched reduce tasks=1
Data-local map tasks=1
Total time spent by all maps in occupied slots (ms)=8877
Total time spent by all reduces in occupied slots (ms)=7511
Total time spent by all map tasks (ms)=8877
Total time spent by all reduce tasks (ms)=7511
Total vcore-milliseconds taken by all map tasks=8877
Total vcore-milliseconds taken by all reduce tasks=7511
Total megabyte-milliseconds taken by all map tasks=8270848
Total megabyte-milliseconds taken by all reduce tasks=7691264

```

```

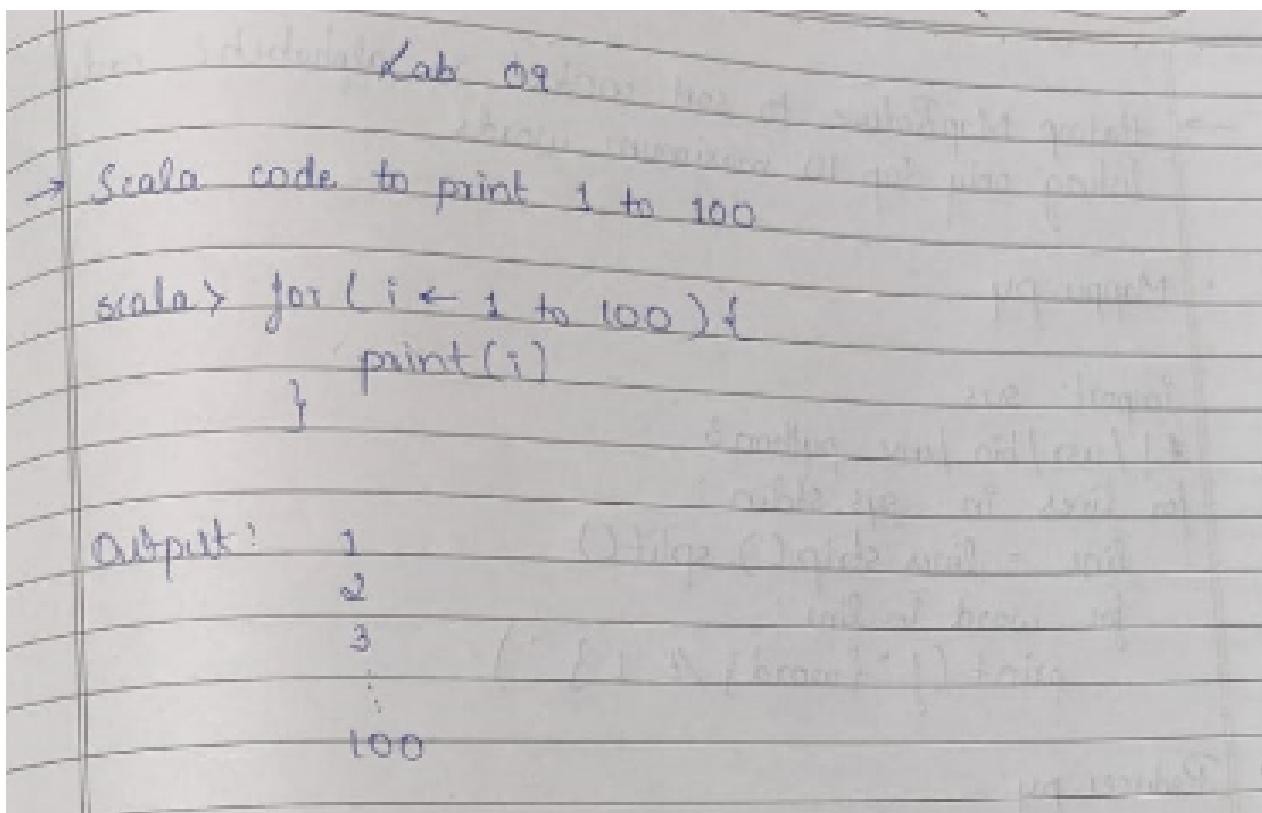
C:\hadoop-3.3.0\sbin>hdfs dfs -cat /meanmax_output/*
01      4
02      0
03      7
04     44
05    100
06    168
07    219
08    198
09    141
10    100
11     19
12      3

C:\hadoop-3.3.0\sbin>

```

Experiment – 8

Write a Scala program to print numbers from 1 to 100 using for loop.



Experiment – 9

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.

→ Using RDD and FlatMap counts how many words appears in a file and write out words which are more than 4.

scala > val textfile = sc.textfile("home/input.txt")

scala > val words = textfile.flatMap(line => line.split(" ")).map(word => (word, 1))

scala > val wordPairs = words.map(word => (word, 1)).reduceByKey((a, b) => a + b)

scala > val wordCounts = wordPairs.collect()

scala > val filtered = wordCounts.filter(case (word, count) => count > 4)

scala > filtered.collect().foreach(println)

Experiment - 10

1.4.3 Question:

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.

→ Hadoop MapReduce to sort content in alphabetic order listing only top 10 maximum words

• Mapper.py

```
import sys  
#! /usr/bin/env python3  
for lines in sys.stdin:  
    line = lines.strip().split()  
    for word in line:  
        print(f'{word} {1}')
```

• Reducer.py

```
#!/usr/bin/env python3  
import sys  
from collections import defaultdict  
w = defaultdict(int)  
for lines in sys.stdin:  
    word, count = lines.split()  
    w[word] += int(count)  
w = sorted(w.items(), key=lambda x: (x[1], x[0]), reverse=True)  
j = 0  
for i in range(10):  
    print(i)  
    j += 1  
    if j == 10:  
        break
```

hadoop jar "path/to/hadoop-streaming.jar" 1

- mappers mapper.py
- reducer reducer.py
- input /bda/input.txt
- output /bda/output.txt

1.4.4 Code with Output:

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 = "[\u00a3#\u00a3>\u00a3=\u00a3[\u00a3]\u00a3*\u00a3/\u00a3\\,\u00a3.\u00a3:\u00a3()?\u00a3!""]";
        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));
}
}
}
}

```

```

C:\hadoop-3.3.0\sbin>jps
11072 DataNode
20528 Jps
5620 ResourceManager
15532 NodeManager
6140 NameNode

C:\hadoop-3.3.0\sbin>hdfs dfs -mkdir /input_dir

C:\hadoop-3.3.0\sbin>hdfs dfs -ls /
Found 1 items
drwxr-xr-x - Anusree supergroup          0 2021-05-08 19:46 /input_dir

C:\hadoop-3.3.0\sbin>hdfs dfs -copyFromLocal C:\input.txt /input_dir

C:\hadoop-3.3.0\sbin>hdfs dfs -ls /input_dir
Found 1 items
-rw-r--r--  1 Anusree supergroup         36 2021-05-08 19:48 /input_dir/input.txt

C:\hadoop-3.3.0\sbin>hdfs dfs -cat /input_dir/input.txt
hello
world
hello
hadoop
bye

```

```
C:\hadoop-3.3.0\bin>hadoop jar C:\sort.jar samples.topN /input_dir/input.txt /output_dir
2021-05-08 19:54:54,582 INFO client.DefaultHddsWANFailoverProxyProvider: Connecting to ResourceManager at /0.0.0.0:8032
2021-05-08 19:54:55,291 INFO mapreduce.JobResourceUploader: Disabling Erasure Coding for path: /tmp/hadoop-yarn/staging/Anusree/.staging/job_1620483374279_0001
2021-05-08 19:54:55,821 INFO input.FileInputFormat: Total input files to process : 1
2021-05-08 19:54:56,261 INFO mapreduce.JobSubmitter: number of splits:1
2021-05-08 19:54:56,478 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1620483374279_0001
2021-05-08 19:54:56,552 INFO mapreduce.JobSubmitter: Executing with tokens: []
2021-05-08 19:54:56,843 INFO conf.Configuration: resource-types.xml not found
2021-05-08 19:54:56,843 INFO resource.ResourcesUtils: Unable to find 'resource-types.xml'.
2021-05-08 19:54:57,387 INFO impl.YarnClientImpl: Submitted application application_1620483374279_0001
2021-05-08 19:54:57,587 INFO mapreduce.Job: The url to track the job: http://LAPTOP-JG329FSD:8088/proxy/application_1620483374279_0001/
2021-05-08 19:54:57,588 INFO mapreduce.Job: Running job: job_1620483374279_0001
2021-05-08 19:55:13,792 INFO mapreduce.Job: Job job_1620483374279_0001 running in uber mode : false
2021-05-08 19:55:13,794 INFO mapreduce.Job: map 0% reduce 0%
2021-05-08 19:55:20,828 INFO mapreduce.Job: map 100% reduce 0%
2021-05-08 19:55:27,116 INFO mapreduce.Job: map 100% reduce 100%
2021-05-08 19:55:33,199 INFO mapreduce.Job: Job job_1620483374279_0001 completed successfully
2021-05-08 19:55:33,334 INFO mapreduce.Job: Counters: 54
    File System Counters:
        FILE: Number of bytes read=65
        FILE: Number of bytes written=530397
        FILE: Number of read operations=0
        FILE: Number of large read operations=0
        FILE: Number of write operations=0
        HDFS: Number of bytes read=142
        HDFS: Number of bytes written=31
        HDFS: Number of read operations=8
        HDFS: Number of large read operations=0
        HDFS: Number of write operations=2
        HDFS: Number of bytes read erasure-coded=0
```