

# VISVESVARAYA TECHNOLOGICAL UNIVERSITY

"JnanaSangama", Belgaum -590014, Karnataka.



## LAB REPORT

on

### Big Data Analytics (23CS6PCBDA)

*Submitted by:*

Rachana N (1BM23CS416)

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 **Rachana N (1BM23CS416)**, 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

<b>Sl.no</b>	<b>Program details</b>	<b>Pg no</b>
<b>1</b>	<b>MongoDB- CRUD Operations Demonstration (Practice and Self Study)</b>	<b>1-7</b>
<b>2</b>	<b>Perform the following DB operations using Cassandra.</b>	<b>8-12</b>
<b>3</b>	<b>Perform the following DB operations using Cassandra</b>	<b>13-20</b>
<b>4</b>	<b>Execution of HDFS Commands for interaction with Hadoop Environment. (Minimum 10 commands to be executed)</b>	<b>21-25</b>
<b>5</b>	<b>Implement Wordcount program on Hadoop framework</b>	<b>26-31</b>
<b>6</b>	<b>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.</b>	<b>32-32</b>
<b>7</b>	<b>Write a Scala program to print numbers from 1 to 100 using for loop.</b>	<b>33-34</b>
<b>8</b>	<b>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.</b>	<b>35-40</b>
<b>9</b>	<b>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).</b>	<b>41-41</b>

## **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.

1/11/25

CQ15h

Cells - 4

PAGE: / /  
DATE: / /

## - Working with Cassandra

> Create Keyspace:

CREATE KEYSPACE Students WITH REPLICATION:

{'class': 'SimpleStrategy', 'replication\_factor': 1};

> Describe the existing Keyspaces:

DESCRIBE KEYSPACES;

> For More details on existing Keyspace:

SELECT \* FROM system\_schema.Keyspaces;

> Use the Keyspace "Students":

USE Students;

> To Create table (column family) by name Student\_info:

CREATE TABLE Students\_info (Roll-No int PRIMARY KEY,  
StudentName text, Date Of Joining timestamp, Last\_exam\_Percentage double);

> Lookup the names of all the tables in the current Keyspace:

DESCRIBE TABLES;

> Describe the table information

DESCRIBE TABLE Students\_info;

CRUD Insert:

BEGIN BATCH

INSERT INTO Student\_info (Roll\_No, StudName, DateOfJoining,  
Last\_Exam\_Percent) VALUES (1, 'Asha', '2018-03-12', 79.9)

INSERT INTO ...

APPLY BATCH;

> View data from the table "Student\_info"

SELECT \* FROM Student\_info;

> View data from the "Student\_info" where Roll\_no Column either  
have 1 or 2 or 3

SELECT \* FROM Student\_info WHERE RollNo IN(1,2,3);

> To execute a non primary key - will throw an error  
Select \* from Student\_info where StudName = 'Asha';

→ So Create an INDEX on the Column as below:

~~To create an INDEX on StudName Column of the~~

~~Student\_info column family~~

CREATE INDEX ON Student\_Info (StudName);

"Now execute the above command";

> To Specify the Number of rows returned in the output

~~select \* from ht~~

select Roll-No, StudName from Student\_info LIMIT 2;

> Alias for column:

select Roll-No as USN from Students\_info

> update

UPDATE Students\_info SET StudName='David Sheen'  
where Roll\_no=2;

"In Primary key you cannot perform update  
operation"

> Delete

DELETE lastExamPercent FROM Students\_info  
WHERE Roll-No=2;

> Delete a Row

DELETE FROM Students\_info WHERE RollNo=2;

Set collection:

ALTER TABLE Students\_info ADD hobbies set(text)

List collection:

ALTER TABLE Students\_info ADD Languages list(text);

UPDATE student\_info

SET hobbies = hobbies + { 'chess', 'Table Tennis' }  
where RollNo = 1;

Select \* from student\_info where RollNo = 1;

UPDATE student\_info

SET language = language + [ 'Hindi', 'English' ]  
where Roll No = 1;

"Can remove an element from a set using the  
Subtraction (-) operator.

> using a counter:

CREATE TABLE library\_book ( count int, value count,  
book\_name varchar, stud\_name = 'Big data Analytics' AND  
stud\_name = 'jat' );

> Time no line:

CREATE TABLE userlogin (uid int PRIMARY KEY,  
password text);

INSERT INTO userlogin (uid, password) VALUES (1,  
'injy') USING TTL 30;

SELECT TTL(password) FROM userlogin WHERE uid=1;

ANSWER

### 1.1.2 Code with Output:

```
bmseccse@bmseccse-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: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 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> 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.

PAGE : \_\_\_\_\_  
DATE : \_\_\_\_\_

Lab 5:

1) Create a Keyspace by name Library

→ create Keyspace library with replication =  
{ 'class': 'SimpleStrategy', 'replication\_factor': 1 };

2) Create a column family library\_info

→ We cannot mix counter and non counter columns

create table library.library\_info (  
 stud\_id int Primary Key,  
 stud\_name text,  
 book\_name text,  
 book\_id int,  
 date\_of\_issue date  
);

create table library.Book\_counters (  
 stud\_id int,  
 book\_name text,  
 counter\_value counter  
 PRIMARY KEY (stud\_id, book\_name)

3) Insert the values into the table in batch

→ Begin Batch

Insert into library.library\_info (stud\_id, stud\_name, book\_name, book\_id, date\_of\_issue)  
values (112, 'John', 'BDA', 101, '2025-04-06');  
Insert into library.library\_info (stud\_id, stud\_name, book\_name, book\_id, date\_of\_issue)  
values (113, 'Alice', 'DBMS', 102, '2025-05-07');

Apply Batch

Begin Batch

update library.Book\_Count

set counter\_value = counter\_value + 1

where stud\_id = 112 and Book\_Name = 'BDA'

update library.Book\_Count

set counter\_value = counter\_value + 1

where stud\_id = 112 and Book\_Name = 'BDA'

Apply Batch.

- 4) Display details and increase Counter

select \* from library.library\_info

select \* from library.book\_counter

update library.Book\_Count

set counter\_value = counter\_value + 1

where stud\_id = 112 and Book\_Name = 'BDA'

- 5) Write a query to show that a student <sup>with</sup> id 112 has taken a book "BDA" 2 times.

→ select counter\_value from library.Book\_Count  
where stud\_id = 112 and Book\_Name = 'BDA'

- 6) Export the created column to a csv file

copy library.library\_info To 'library\_info.csv'  
with header = TRUE;

7 Import csv into column family

→ copy library.library\_info from 'library.info.csv'  
with header = TRUE;

(a)

Import csv file to cassandra

copy library.library\_info (stud\_id, stud\_name,  
book\_name, book\_id, Date\_of\_issue)  
from 'library.info.csv' with header = TRUE

08/15/25  
go

## 1.2.2 Code with Output:

```

bnsccse@bnsccse-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 = 864000
    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, 'Rutu', '2023-10-10', 97) insert into Students_info(Roll_no, Studname, DateOfJoining, last_exam_Percent) Values(3, 'Rachana', '2023-10-10', 97.5) 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 | Rutu
  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 | Rutu
  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 | Ritu |
| 3 | 2023-10-09 18:30:00.000000+0000 | 97.5 | Rachana |
+-----+-----+-----+
(3 rows)
cqlsh:students> select * from students_info where Studname='Charu';
InvalideRequest: 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;
+-----+-----+
| roll_no | studname |
+-----+-----+
| 1 | Sadhana |
| 2 | Ritu |
+-----+-----+
(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.

PAGE: \_\_\_\_\_  
DATE: \_\_\_\_\_

IX to set a particular field value to null.

```
db.Students.update({id: 3}, {$set: {Location: null}})
```

\* count the number of documents in Students collection.

```
db.Students.count()
```

\* count the number of documents in Student Collection with grade : VII

```
db.Students.count({grade: "VII"})
```

# retrieve first 3 documents

```
db.Students.find({grade: "VII"}).limit(3).pretty();
```

# sort the document in Ascending order

```
db.Students.find().sort({StudName: 1}).pretty();
```

# for descending order

```
db.Students.find().sort({StudName: -1}).pretty();
```

# to skip the 1<sup>st</sup> two documents from the Students Collection

```
db.Students.find().skip(2).pretty()
```

# Create a collections by name "food" and add to each document add a "fruits" array constitute of "grape", "mango" and "apple".

```
db.food.find({fruits: ['grapes', 'mango', 'apple']}).pretty()
```

# Student Names begins with M

```
db.Student.find({$ StudName: /^M/}).pretty();
```

# Student Names contains 'hi' in any position

```
db.Student.find({$ StudName: /hi/}).pretty();
```

```
db.Student.count();
```

# descending order:

```
db.Student.find().sort({$ StudName: -1}).pretty();
```

# new field in an existing Document  
<sup>Add</sup>

```
db.Student.update({$ id: 2}, {$set: {Location: "Network"}})
```

# Remove field

```
db.Student.update({$ id: 2}, {$unset: {Location: "Network"}})
```

<sup>VII</sup>  
finding document based on search criteria  
Supressing few fields

```
db.Student.find({$ id: 2, $ StudName: 1, Grade: 1, _id: 0})
```

# To find those documents where the Grade is not set to 'VII'

```
db.Student.find({Grade: {$ne: 'VII'}}).pretty();
```

# To find documents from the Students collection where the StudName ends with s.

```
db.Student.find({$ StudName: /s$/}).pretty();
```

## Working with mongoDB

## 1. Create Database in mongoDB

```
use myDB;
```

```
db;
```

```
show dbs;
```

## 2. CRUD (Create, Read, update, delete) operations

```
i) db.createCollection("Student");
```

```
db.Student.drop();
```

```
→ db.Student.insert({_id: 1, StudName: "Rach", Grade: "6th",  
Hobbies: "Nothing"});
```

```
→ db.Student.update({_id: 1, StudName: "Rach", Grade: "6th"},  
{$set: {Hobbies: "Playing", $current: true}});
```

```
db.Student.find({StudName: "Rach"});
```

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

```
db.Student.find({Grade: {$eq: '6th'}}).pretty();
```

```
→ db.Student.insert({_id: 2, StudName: "Monisha", Grade: "6th",  
Hobbies: "Nothing"});
```

```
db.Student.find({Hobbies: {$in: ['Nothing', 'playing']}},  
{pretty: 1});
```

# To find in "fruits" array having "mango" in the first index position

db.food.find({ "fruits.1": "grapes" })

# To find those documents from the "food" collection where the size of the array is two

db.food.find({ "fruits": { "\$size": 2 } })

# To find the documents with a particular id and display the first two elements from the array "fruits"

db.food.find({ "\_id": 1, "fruits": { "\$slice": 2 } })

# To find all the documents from the food collection which have elements mango and grapes in the array "fruits"

db.food.find({ "fruits": { "\$all": ["mango", "grapes"] } })

update on Array:

using particular id replace the element in the 1<sup>st</sup> index position of the fruits array with apple

db.food.update({ "\_id": 3 }, { "\$set": { "fruits.1": "apple" } })

# insert new key value pairs in the fruits array

db.food.update({ "\_id": 2 }, { "\$push": { "price": 80, "mango": 200, "cherry": 100 } })

### 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)

Lab - 6

PAGE : \_\_\_\_\_  
DATE : \_\_\_\_\_

#### HDFS Command in Hadoop file system

1) Start hadoop (must be in hdfsuser)  
\$ start-all.sh

2) creating a directory inside hadoop - mkdir  
\$ hdfs dfs -mkdir /bda\_hadoop

3) listing all content inside hadoop - ls  
\$ hadoop fs -ls /

Output: Found 1 items  
drwxr-xr-x - hadoop supergroup

4) copying files from desktop using put command - put  
\$ hdfs dfs -put /home/hadoop/Desktop/bda-local.txt/  
bda\_hadoop/file.txt

5) cat command (listing the content of file in hadoop)

\$ hdfs dfs -cat /bda\_hadoop/file.txt

6) ~~Copying files from local reference using copyFromLocal  
-l cmd~~

\$ hdfs dfs -copyFromLocal /home/hadoop/Desktop/bda-to-  
local.txt /bda\_hadoop/file\_local.txt cp -local.txt

7) `hdfs dfs -cat /bda_hadoop/file_cp_local.txt`

8) get command

~~`hdfs dfs -get /bda_hadoop/file.txt /home/hadoop/Downloads/downloaded_file.txt`~~

9) Cannot use getmerge for the files  
if both the files are in same (local system)

`hdfs dfs -getmerge /bda_hadoop/ /home/hadoop/Downloads/merged_output.txt`

10) Display Access control list (ACL) permissions  
for a file or directory in HDFS

~~\$ hadoop fs -getfacl /bda\_hadoop/~~

11) Copy To Local

~~+ hadoop~~

~~\$ hdfs dfs -copyToLocal /bda\_hadoop/file.txt /home/hadoop/Desktop~~

12) mv

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

13) ls

~~\$ hadoop fs -ls /abc~~



PAGE :

DATE : / /

13 - copy

\$ hadoop fs -cp /hello/ /hadoop\_ls

15/A

### 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@bmscsecse-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@bmscsecse-HP-Elite-Tower-800-G9-Desktop-PC:~/Desktop$ hdfs dfs -cp /test_Lab05/ /Lab05
hadoop@bmscsecse-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@bmscsecse-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.

24/12/25  
PAGE: / /  
DATE: / /

Lab 7:  
Mapper code

```
// importing Libraries
```

```
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> {
    // Map function
    public void map(LongWritable key, Text value,
                    OutputCollector<Text, IntWritable> output, Reporter reporter) throws IOException {
        String line = value.toString();
        // Splitting the line on spaces
        for (String word : line.split(" ")) {
            if (word.length() > 0)
                output.collect(new Text(word), new IntWritable(1));
        }
    }
}
```

```
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 WordReducer extends MapReduceBase implements Reducer<Text, IntWritable, Text, IntWritable> {
```

// Reduce function

```
public void reduce(Text key, Iterator<IntWritable> value, OutputCollector<Text, IntWritable> output, Reporter reporter) throws IOException {
```

```
    int count = 0;
    while (value.hasNext())
```

```
        IntWritable i = value.next();
        count += i.get();
```

```
    output.collect(key, new IntWritable(count));
```

g

```

import java.io.IOException;
import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapred.FileInputFormat;
import org.apache.hadoop.mapred.FileOutputFormat;
import org.apache.hadoop.JobClient;
import org.apache.hadoop.mapred.JobConf;
import org.apache.hadoop.util.Tool;
import org.apache.hadoop.util.ToolRunner;

```

```

public class WCDriver extends Configuration
    implements Tool {
    public int run(String args[]) throws IOException
    {

```

~~if (args.length < 2)~~

~~System.out.println ("Please give valid inputs");~~  
~~return -1;~~

~~3~~

```

JobConf conf = new JobConf (WCDriver.class);
FileInputFormat.setInputPaths(conf, new Path(args[0]));
FileOutputFormat.setOutputPath(conf, new Path(args[1]));
conf.setReduces(1);
conf.setMapperClass(WCMapper.class);
conf.setReducerClass(WCReducer.class);
conf.setMapOutputKeyClass(Text.class);
conf.setMapOutputValueClass(IntWritable.class);
JobClient.runJob(conf);
return 0;

```

~~• 3~~

public static void main(String args[]) throws Exception

int exitCode = ToolRunner.run(new WCDomains,  
args);

System.out.println(exitCode);

y  
y

80/5/15

```

Activities Terminal May 20 14:47
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]
localhost: namenode is running as process 8499. Stop it first and ensure /tmp/hadoop-hadoop-namenode.pid file is empty before retry.
Starting datanodes
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.
hadoop@bmscsece-HP-Elite-Tower-600-G9-Desktop-PC: $ hadoop dfs -ls /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
15386 ResourceManager
15741 NodeManager
6270 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 2025-05-28 13:48 /Folder1
drwxr-xr-x  - hadoop supergroup          0 2025-05-28 13:48 /Folder2
drwxr-xr-x  - hadoop supergroup          0 2025-05-28 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.MetricsSystemImpl: Scheduled Metric snapshot period at 10 second(s).
2025-05-20 14:45:00,315 INFO impl.MetricsSystemImpl: JobTracker metrics system started
2025-05-20 14:45:00,321 WARN impl.MetricsSystemImpl: JobTracker metrics system already initialized!
2025-05-20 14:45:00,384 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to remedy this.
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: ~

```

```

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=86
  Combine input records=0
  Combiner output records=0
  Reduce input groups=12
  Reduce shuffle 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_REDUCE=0
File Input Format Counters
  Bytes Read=66
  File Output Format Counters
  Bytes Written=86
0
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 2025-05-28 14:45 /rgs/output/_SUCCESS
-rw-r--r-- 1 hadoop supergroup          86 2025-05-28 14:45 /rgs/output/part-00000
hadoop@bmscsece-HP-Elite-Tower-600-G9-Desktop-PC: $ hadoop fs -cat /rgs/output/part-00000
an
are
be
becz
executed
feeling
good
hitit
how
i
program
the
you

```

Activities Terminal May 20 14:48

```

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));
}
}

```

```

C:\hadoop-3.3.0\sbin>hadoop jar C:\avgtemp.jar temp.AverageDriver /input_dir/temp.txt /avgtemp_outputdir
2021-05-15 14:52:50,635 INFO client.DefaultNoHARMFailoverProxyProvider: Connecting to ResourceManager at /0.0.0.0:8032
2021-05-15 14:52:51,085 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-JG329E5D: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 14:53:19,868 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 14:53:26,096 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 MeanMaxDriver /input_dir/temp.txt ./meanmax_output
2021-05-21 20:28:05,258 INFO client.DefaultNoHARFailoverProxyProvider: Connecting to ResourceManager at /0.0.0.0:8032
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_1621608943095_0001
2021-05-21 20:28:08,426 INFO input.FileInputFormat: Total input files to process : 1
2021-05-21 20:28:09,197 INFO mapreduce.JobSubmitter: number of splits:1
2021-05-21 20:28:09,741 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1621608943095_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.ResourceUtil: Unable to find 'resource-types.xml'.
2021-05-21 20:28:10,676 INFO impl.YarnClientImpl: Submitted application application_1621608943095_0001
2021-05-21 20:28:11,005 INFO mapreduce.Job: The url to track the job: http://LAPTOP-7G329ESD:8088/proxy/application_1621608943095_0001/
2021-05-21 20:28:11,006 INFO mapreduce.Job: Running job: job_1621608943095_0001
2021-05-21 20:28:29,385 INFO mapreduce.Job: Job job_1621608943095_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_1621608943095_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=894868
        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)=8077
        Total time spent by all reduces in occupied slots (ms)=7511
        Total time spent by all map tasks (ms)=8077
        Total time spent by all reduce tasks (ms)=7511
        Total vcore-milliseconds taken by all map tasks=8077
        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.

Lab 9:  
Qspark:

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

```
> spark-shell  
scala> for (i <- 1 to 100){  
    println(i)  
}  
1  
2  
3  
:  
99  
100
```

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

```
> nano input.txt  
Hi Hi Hi Hi Hi  
Hello Hello Hello Hello Hello  
in is is my world
```

> spark-shell

```
> val fileRDD = sc.textFile("input.txt")  
val wordsRDD = fileRDD.flatMap(line => line.split(" ")).filter(_.nonEmpty)
```

## 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.

1/25 PAGE: DATE: / /

Lab 9:  
Qspark:

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

```
> spark-shell  
scala> for (i <- 1 to 100){  
    println(i)  
}
```

1  
2  
3  
:  
99  
100

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

```
> nano input.txt  
Hi Hi Hi Hi Hi  
Hello Hello Hello Hello Hello  
in is is my world
```

> spark-shell

```
val fileRDD = sc.textFile("input.txt")  
val wordsRDD = fileRDD.flatMap(line => line.split(" \n\t+").filter(_.nonEmpty))
```

- > val wordCountRDD = wordsRDD.map(word => (word.toLowerCase(), 1)).reduceByKey(\_ + \_)
- > wordCountsRDD.collect().foreach(println)
  - (Hi, 6)
  - (Hello, 5)
  - (in, 3)
  - (my, 1)
  - (world, 1)
- > val filteredWordsRDD = wordCountsRDD.filter(\_.count > 4)
- > val result = filteredWordsRDD.collect()
- > result.foreach(println)
  - (Hi, 6)
  - (Hello, 5)

Q For a given text file, create a map reduce program to sort the content in an alphabetical order listing only top 10 maximum occurrences of words

Mapper:

```
public class WordCountMapper extends Mapper<LongWritable,
Text, Text, IntWritable> {
    private final static IntWritable one = new IntWritable(1);
    public void map(LongWritable key, Text value, Context context)
        throws IOException, InterruptedException {
        String[] words = value.toString().toLowerCase().split(" ").
        for (String word : words) {
            context.write(new Text(word), one);
        }
    }
}
```

## 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.

PAGE : / / DATE : / /

```
> val wordCountsRDD = wordsRDD.map(word =>
  (word.toLowerCase(), 1)).reduceByKey(_ + _)

> wordCountsRDD.collect().foreach(println)
  (Hi, 6)
  (Hello, 5)
  (in, 3)
  (*my, 1)
  (world, 1)

> val filteredWordsRDD = wordCountsRDD
  .filter(_.case (word, count) => count > 4)

> val result = filteredWordsRDD.collect()

> result.foreach(println)
  (Hi, 6)
  (Hello, 5)

? For a given text file, create a map reduce program to sort the content in an alphabetical order listing only top 10 maximum occurrences of words

Mapper:
public class WordCountMapper extends Mapper<LongWritable,
Text, Text, IntWritable> {
  private final static IntWritable one = new IntWritable(1);
  public void map(LongWritable key, Text value, Context context)
    throws IOException, InterruptedException {
    String[] words = value.toString().toLowerCase(),
      .split("\\W+");
    for (String word : words) {
      if (!word.equals("")) {
        context.write(new Text(word), one);
      }
    }
  }
}
```

DATE: / /

for (String word : words) {  
 if (!word.isEmpty()) {  
 context.write(new Text(word), one);  
 }  
}

Reducer:

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

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

```
public void reduce(Text key, Iterable<IntWritable> values,  
Context context) throws  
IOException, InterruptedException {
```

```
int sum = 0;  
for (IntWritable val : values) {  
    sum += val.get();  
}
```

```
countMap.put(new Text(key), sum);
```

3

```
protected void cleanup(Context context) throws IOException,  
InterruptedException {
```

```
List<Map.Entry<Text, Integer>> sorted = new ArrayList<>(count  
Map.entrySet());
```

```
sorted.sort((e1, e2) → e2.getValue().compareTo(e1.get  
Value()));
```

```
int counter = 0;
```

```
for (Map.Entry<Text, Integer> entry : sorted) {
```

```
if (counter++ == 10) break;
```

```
context.write(entry.getKey(), new IntWritable(counter));
```

3

#### 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 = "[_|#<>|^=\\[\\]\\*\\/\\\\;,;.\\:-;?!\""]";
        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\sbin>hadoop jar C:\sort.jar samples.topN /input_dir/input.txt /output_dir
2021-05-08 19:54:54,582 INFO client.DefaultDNHAWNFFailoverProxyProvider: 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 19:54:56,552 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.ResourceUtils: 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-JG329ESD: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 19:55:20,820 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
```

```
C:\hadoop-3.3.0\sbin>hdfs dfs -cat /output_dir/*
hello    2
hadoop   1
world    1
bye      1

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

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).

(Open ended question)

PAGE: \_\_\_\_\_  
DATE: \_\_\_\_\_

lab 10

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 (or

- Preinstall dependencies & download NLTK data:  
text\_cleaning - treat.py

```
from pyspark import SparkContext  
from pyspark.streaming import StreamingContext  
from nltk.corpus import stopwords  
from nltk.stem import WordNetLemmatizer  
import re
```

```
sc = SparkContext("local[2]", "Text Cleaning Stream")
```

```
ssc = StreamingContext(sc, 5)
```

```
stop_words = set(stopwords.words('english'))
```

```
lemmatizer = WordNetLemmatizer()
```

```
def clean_text(line):
```

```
    line = line.lower()
```

```
    line = re.sub(r"\b[a-z]+\b", "", line)
```

~~tokens = [word for word in line.split() if word not in stop\_words]~~

~~lemmatized = [lemmatizer.lemmatize(word) for word in tokens]~~

~~return " ".join(lemmatized)~~

```
lines = ssc.socketTextStream("localhost", 9999)
```

```
cleaned = lines.map(clean_text)
```

```
cleaned.print()
```

```
ssc.start()
```

```
ssc.awaitTermination()
```

```
{entry.getValue()}
```

O/P = The quick brown fox jumps over the lazy dog

→ Quick brown fox jump lazy dog.