

VISVESVARAYA TECHNOLOGICAL UNIVERSITY

“JnanaSangama”, Belgaum -590014, Karnataka.



LAB REPORT

on

Big Data Analytics (23CS6PCBDA)

Submitted by:

Rohit Ramchandra Gandhi (1BM23CS417)

**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 **Rohit Ramchandra Gandhi(1BM23CS417)**, 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. Dr.

Associate Professor
Department of CSE
BMSCE, Bengaluru

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)	1-8
2	Perform the DB operations using Cassandra.	9-13
3	Perform the DB operations using Cassandra	14-16
4	Execution of HDFS Commands for interaction with Hadoop Environment. (Minimum 10 commands to be executed)	17-19
5	Implement Wordcount program on Hadoop framework	20-23
6	a)Create a MapReduce program to find average temperature for each year from the NCDC data set. b) find the mean max temperature for every month.	24-30
7	For a given Text file, Create a Map Reduce program to sort the content in an alphabetic order listing only top 10 maximum occurrences of words.	31-34
8	Write a Scala program to print numbers from 1 to 100 using a for loop.	35
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.	36-37
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. (Open Ended Question).	38-39

Github Link: <https://github.com/RohitGnadi2599/BDALab>

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.

Program 1

MongoDB- CRUD Operations Demonstration (Practice and Self Study)

- Created a database named **myDB** and verified its existence.
- Created and dropped collections like **Student** and **Students**.
- Inserted student data into collections.
- Performed **upsert** to insert or update a student record.
- Used to find queries with various filters: by name, grade, hobbies, regex, etc.
- Retrieved specific fields while suppressing `_id`.
- Counted total documents and documents with specific criteria.
- Sorted records in ascending and descending order.
- Imported data from a CSV file and exported data to a CSV file.
- Used `save()` to insert or replace documents.
- Added, removed, and set fields to `null` in documents.
- Retrieved limited records and skipped initial entries.
- Created a **food** collection with arrays and queried arrays by value, index, size, etc.
- Updated specific elements in an array.
- Practiced query optimizations using `$in`, `$all`, `$ne`, `$regex`, `$slice`, and more.

Observation:

14/03/25

LAB-1

+ creating Database in MongoDB.

1) use mydb;

+ switched to db mydb

2) show dbs;

+ admin 232.00 KiB
+ local 11.59 GiB

3) create Database

1) db.createCollection("student");
+ 80K:13

2) db.students.insertOne({ _id: 1, studName: "Rohit", grade: "VII", hobbies: ["Internet", "Surfing"]});
+ 1 inserted: true, inserted: 1 3

+ find:

1) db.students.find({\$studName: "Rohit"});
2) db.students.find({\$3, \$studName: "Rohit"}, {grade: 1, _id: 0})

+ [{ \$studName: "Rishabh", grade: 'VII' },
{ \$studName: "Rohit", grade: 'VII' }]

3) db.students.find({
 grade : {
 \$eq : "VII"
 }}).pretty();
+ [

{

- id : 1,
 studName : "Rohit";
 grade : 'VII'

4) db.students.find({
 Hobbies : {
 \$in : ["chess", "skating"]
 }}).pretty();
+ [{ id : 3,

 studName : "Rohit",
 Hobbies : "chess" }]

5) db.students.find({
 student : 1713}).pretty();
[

{

- id : 7
 studName : 'Michelle Iactintinha', 3
]

6) db.students.find({
 studentName : 1213}).pretty();
[

{

- id : 1,
 studName : "Michelle" }]

+ count

7) db.students.countDocuments();
+ 2

2) db.students.find().sort({student:-1}).pretty();

id : 2,

studName : 'Richelle Tacinatha',

Grade : 'VII'

Hobbies : 'Internet Surfing'

+ Importing CSV (txt file) to db:

Command: mongorestore --host localhost --port 27017 -d dbalb1 -c users --drop
mongod --port 27017 --db dbalb1 --collection users

+ 2025-03-04:04:28:167 Connected to: mongo+srv://
[++ Redacted ++]@dbalb1:27017/mongo
db.net/test

Code with Output:

```
hadoop@bmscse-HP-Elite-Tower-600-G9-Desktop-PC:~$ mongosh
Current Mongosh Log ID: 6833f9c9126af1945c47586f
Connecting to:      mongodb://127.0.0.1:27017/?directConnection=true&serverSelectionTimeoutMS=2000&appName=mongosh+2.0.1
Using MongoDB:     7.0.2
Using Mongosh:     2.0.1
mongosh 2.5.1 is available for download: https://www.mongodb.com/try/download/shell
For mongosh info see: https://docs.mongodb.com/mongodb-shell/

-----
The server generated these startup warnings when booting
2025-05-26T10:46:48.806+05:30: Using the XFS filesystem is strongly recommended with the WiredTiger storage engine. See http://dochub.mongodb.org/core/prodnotes-filesystem
2025-05-26T10:46:50.937+05:30: Access control is not enabled for the database. Read and write access to data and configuration is unrestricted
-----

test> use MyDB
switched to db MyDB
MyDB> db
MyDB
MyDB> show dbs
admin          40.00 KiB
config         72.00 KiB
local          80.00 KiB
myNewDatabase  72.00 KiB
MyDB> db.createCollection("Student");
{ ok: 1 }
MyDB> db.Student.insert({_id:1,Name:"Preeti",Grade:"V",Hobbies:"Dancing"},{_id:2,Name:"Prajwal",Grade:"V",Hobbies:"Drawing"});
DeprecationWarning: Collection.insert() is deprecated. Use insertOne, insertMany, or bulkWrite.
{ acknowledged: true, insertedIds: { '0': 1 } }
MyDB> db.find();
TypeError: db.find is not a function
MyDB> db.Student.find();
[ { _id: 1, Name: 'Preeti', Grade: 'V', Hobbies: 'Dancing' } ]
MyDB> db.Student.insertMany([ { _id:2,Name:"Rachana",Grade:"V",Hobbies:"Painting"},{_id:3,Name:"Prajwal",Grade:"V",Hobbies:"Drawing"} ]);
MongoInvaldArgumentError: Argument "docs" must be an array of documents
MyDB> db.Student.insertMany([ { _id:2,Name:"Rachana",Grade:"V",Hobbies:"Painting"},{_id:3,Name:"Prajwal",Grade:"V",Hobbies:"Drawing"} ]);
{ acknowledged: true, insertedIds: { '0': 2, '1': 3 } }
MyDB> db.Student.update({ _id:2,Name:"Rachana",Grade:"V"},{$set:[{Hobbies:"Singing"}]}, {upsert:true});
DeprecationWarning: Collection.update() is deprecated. Use updateOne, updateMany, or bulkWrite.
{
  acknowledged: true,
  insertedId: null,
  matchedCount: 1,
  modifiedCount: 1,
  upsertedCount: 0
}
MyDB> db.Student.find();
[
  { _id: 1, Name: 'Preeti', Grade: 'V', Hobbies: 'Dancing' },
  { _id: 2, Name: 'Rachana', Grade: 'V', Hobbies: 'Singing' },
  { _id: 3, Name: 'Prajwal', Grade: 'V', Hobbies: 'Drawing' }
]
```

```
{
  _id: 1, fruits: [ 'grapes', 'banana', 'apple' ] ,
  {
    _id: 2,
    fruits: [ 'grapes', 'mango', 'cherry' ],
    price: [ { grapes: 80, mango: 100, cherry: 200 } ]
  },
  { _id: 3, fruits: [ 'banana', 'mango' ] }
]
MyDB> db.food.find().pretty();
[
  { _id: 1, fruits: [ 'grapes', 'banana', 'apple' ] ,
  {
    _id: 2,
    fruits: [ 'grapes', 'mango', 'cherry' ],
    price: [ { grapes: 80, mango: 100, cherry: 200 } ]
  },
  { _id: 3, fruits: [ 'banana', 'mango' ] }
]
MyDB> db.createCollection("customer");
{ ok: 1 }
MyDB> var val=[{custid:1,accbal:100,acctype:"A"},{custid:1,accbal:100,acctype:"B"},{custid:2,accbal:200,acctype:"B"}];
Uncaught:
SyntaxError: Unexpected token, expected "," (1:144)
> 1 | var val=[{custid:1,accbal:100,acctype:"A"},{custid:1,accbal:100,acctype:"B"},{custid:2,accbal:200,acctype:"B"}];
|   ^
2 |
MyDB> var val=[{custid:1,accbal:100,acctype:"A"},{custid:1,accbal:100,acctype:"B"},{custid:2,accbal:200,acctype:"B"}];
MyDB> db.customer.insert(val);
{
  acknowledged: true,
  insertedIds: [
    '0': ObjectId("683405cb126af1945c475870"),
    '1': ObjectId("683405cb126af1945c475871"),
    '2': ObjectId("683405cb126af1945c475872"),
    '3': ObjectId("683405cb126af1945c475873")
  ]
}
MyDB> db.customer.aggregate([{$group:{_id:'$custid',totalbal:{$sum:'$accbal'}}}]);
[ { _id: 1, totalbal: 200 }, { _id: 2, totalbal: 400 } ]
MyDB> db.Customers.aggregate( { $match:{AcctType:"S"}},{$group : { _id : "$custID",TotAccBal :
Uncaught:
SyntaxError: Unexpected character "'". (1:43)
> 1 | db.Customers.aggregate( { $match:{AcctType:"S"}},{$group : { _id : "$custID",TotAccBal :
```

```

[ { _id: 1, fruits: [ 'grapes', 'banana', 'apple' ] },
  {
    _id: 2,
    fruits: [ 'grapes', 'mango', 'cherry' ],
    price: [ { grapes: 80, mango: 100, cherry: 200 } ]
  },
  { _id: 3, fruits: [ 'banana', 'mango' ] }
]
MyDB> db.food.find().pretty();
[
  { _id: 1, fruits: [ 'grapes', 'banana', 'apple' ] },
  {
    _id: 2,
    fruits: [ 'grapes', 'mango', 'cherry' ],
    price: [ { grapes: 80, mango: 100, cherry: 200 } ]
  },
  { _id: 3, fruits: [ 'banana', 'mango' ] }
]
MyDB> db.createCollection("customer");
{ ok: 1 }
MyDB> var val=[{custid:1,accbal:100,acctype:"A"},{custid:1,accbal:100,acctype:"B"},{custid:2,accbal:200,acctype:"B"}];
Uncaught:
SyntaxError: Unexpected token, expected "," (1:144)

> 1 | var val=[{custid:1,accbal:100,acctype:"A"},{custid:1,accbal:100,acctype:"B"},{custid:2,accbal:200,acctype:"B"}];
|   ^
2 |

MyDB> var val=[{custid:1,accbal:100,acctype:"A"},{custid:1,accbal:100,acctype:"B"},{custid:2,accbal:200,acctype:"B"}];

MyDB> db.customer.insert(val);
{
  acknowledged: true,
  insertedIds: {
    '0': ObjectId("683405cb126af1945c475870"),
    '1': ObjectId("683405cb126af1945c475871"),
    '2': ObjectId("683405cb126af1945c475872"),
    '3': ObjectId("683405cb126af1945c475873")
  }
}
MyDB> db.customer.aggregate({$group:{_id:'$custid',totalbal:{$sum:'$accbal'}}});
[ { _id: 1, totalbal: 200 }, { _id: 2, totalbal: 400 } ]
MyDB> db.Customers.aggregate ( {$match:{AcctType:"S"}},{$group : { _id : "$custID",TotAccBal :
Uncaught:
SyntaxError: Unexpected character '''. (1:43)

> 1 | db.Customers.aggregate ( {$match:{AcctType:"S"}},{$group : { _id : "$custID",TotAccBal :

MyDB> db.customer.aggregate({$match:{acctype:"A"}},{$group:{_id:"$custid",totalbal:{$sum:'$accbal'}}});
[ { _id: 1, totalbal: 100 }, { _id: 2, totalbal: 200 } ]
MyDB> db.customer.aggregate({$match:{acctype:"A"}},{$group:{_id:"$custid",totalbal:{$sum:'$accbal'}}},{$mat$rn

MyDB> db.customer.aggregate({$match:{acctype:"A"}},{$group:{_id:"$custid",totalbal:{$sum:'$accbal'}}},{$match

MyDB> db.customer.aggregate({$match:{acctype:"A"}},{$group:{_id:"$custid",totalbal:{$sum:'$accbal'}}},{$match
[ { _id: 2, totalbal: 200 } ]
MyDB> S

```

Program 2

Perform the following DB operations using Cassandra.

a) Create a keyspace by name Employee

b) Create a column family by name

Employee-Info with attributes Emp_Id Primary Key, Emp_Name, Designation, Date_of_Joining, Salary, Dept_Name

c) Insert the values into the table in batch

d) Update Employee name and Department of Emp-Id 121

e) Sort the details of Employee records based on salary

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

g) Update the altered table to add project names.

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

Observation:

The image shows a page from a handwritten lab exercise notebook. At the top right, there is a logo for 'CLASSEmate' with the text 'Edit', 'Page', and 'Share'. Below the logo, the date '11/03/24' is written. In the center, the text 'LAB-02' is written. To the left of 'LAB-02', there is some very faint, illegible handwriting. Below 'LAB-02', the text 'MongoDB LAB Exercises' is written. Underneath this, there is a numbered list of tasks. Task 1 is partially visible, showing '1) Perform Create a collection by name customers with the following attributes.' and '+ use customers;'. Task 2 is fully visible, showing '2) Insert at least 5 values into table:' and '+ code'. Below task 2, there is a block of code for inserting five documents into a 'customers' collection. The code uses a loop with index variables i and j, and includes fields like '_id', 'acc_bal', and 'acc_type'. Below this code, the word 'Output:' is written, followed by a series of numbers enclosed in quotes: '0', '1', '2', '3', and '4' followed by a closing parenthesis. The rest of the page contains faint, illegible handwriting.

11/03/24 LAB-02

MongoDB LAB Exercises

1) Perform Create a collection by name customers with the following attributes.
+ use customers;

2) Insert at least 5 values into table:
+ code

```
db.customers.insertOne({  
    '_id': 1, output:  
    'acc_bal': 1500, ack: true,  
    'acc_type': '2'3}); insertedId: ObjectId('07c...')  
db.customers.insertMany([  
    {  
        '_id': 1, 'acc_bal': 1500, 'acc_type': '2'3},  
        {  
            '_id': 2, 'acc_bal': 800, 'acc_type': 'x'3},  
            {  
                '_id': 3, 'acc_bal': 2500, 'acc_type': '2'3},  
                {  
                    '_id': 4, 'acc_bal': 1500, 'acc_type': '2'3},  
                    {  
                        '_id': 5, 'acc_bal': 1800, 'acc_type': '2'3});  
Output:  
E acknowledged: true,  
insertedIds: [
```

3) Write a query to display those records whose total account balance is greater than 1200 of account type '2'.

+ db.customers.find({
 ba-Bal: {
 \$gt: 1200},
 Acct-type: '2'})

{

- id: 2,
ba-Bal: 1500,
Acct-type: '2')

}

- id: objId('67cffcfb'),
last-Id: 2,
ba-Bal: 1500,
Acct-type: '2'

},

4) Determine min & max account balance for each customer.

db.customers.aggregate([

 {
 \$group: {
 _id: '\$cust-id',
 min-balance: {
 \$min: "\$ba-Bal"},
 max-bal: {
 \$max: "\$ba-Bal"}
 }
 },
 {
 \$group: {
 _id: '\$cust-id',
 min-balance: {
 \$min: "\$ba-Bal"},
 max-bal: {
 \$max: "\$ba-Bal"}
 }
 }

]);

Output:

{ id : 2; min_balance : 800; max_bal : 800 } ;

{ id : 5; min_balance : 1800; max_bal : 1800 } ;

E-commerce platform

1) Retrieve all products :

db.products.find()

Output:

{
 id : 'p1' ,
 product_id : 'p1' ,
 category : "Electronics" ,
 quantity : 100 ,
 min_balance : 800 ,
 max_balance : 800 } ,

{
 id : 'p2' ,
 product_id : 'p2' ,
 category : "Electronics" ,
 quantity : 100 ,
 min_balance : 1800 ,
 max_balance : 1800 } ,

{
 id : 'p3' ,
 product_id : 'p3' ,
 category : "Electronics" ,
 quantity : 100 ,
 min_balance : 1800 ,
 max_balance : 1800 } ,

2) db.products.find({ category : "Electronics" }) ;

Output:

{
 id : object id : ('62d005') ,
 product_id : 'p1' ,
 category : "Electronics" ,
 quantity : 100 ,
 min_balance : 800 ,
 max_balance : 800 } ,

{
 id : object id : ('67d006') ,
 product_id : 'p2' ,
 category : "Electronics" ,
 quantity : 100 ,
 min_balance : 1800 ,
 max_balance : 1800 } ,

3) Retrieve Products with quantity greater than 0

db.products.find({ quantity : { \$gt : 0 } }) ;

Output:
{
 id : 'p1' ,
 product_id : 'p1' ,
 category : "Electronics" ,
 quantity : 100 ,
 min_balance : 800 ,
 max_balance : 800 } ,

{
 id : 'p2' ,
 product_id : 'p2' ,
 category : "Electronics" ,
 quantity : 100 ,
 min_balance : 1800 ,
 max_balance : 1800 } ,

{
 id : 'p3' ,
 product_id : 'p3' ,
 category : "Electronics" ,
 quantity : 100 ,
 min_balance : 1800 ,
 max_balance : 1800 } ,

4) Retrieved products sorted by price in Ascending order

```
db.products.find().sort({$price: 1});
```

Output:

```
[{"product_id": "P7", "product_name": "Laptop", "product_desc": "Dell XPS 15", "product_price": 1200}, {"product_id": "P4", "product_name": "Smartphone", "product_desc": "Samsung Galaxy S20", "product_price": 800}, {"product_id": "P1", "product_name": "Headphones", "product_desc": "Sony WH-1000XM4", "product_price": 600}, {"product_id": "P6", "product_name": "Monitor", "product_desc": "Dell U2720Q", "product_price": 1000}, {"product_id": "P2", "product_name": "Keyboard", "product_desc": "Logitech G910", "product_price": 500}, {"product_id": "P3", "product_name": "Mouse", "product_desc": "SteelSeries Rival 650", "product_price": 400}, {"product_id": "P5", "product_name": "Charger", "product_desc": "Anker PowerPort Atom III", "product_price": 200}, {"product_id": "P8", "product_name": "Power Bank", "product_desc": "Anker PowerCore 20000", "product_price": 300}, {"product_id": "P9", "product_name": "USB-C Cable", "product_desc": "AmazonBasics 6ft USB-C to USB-A", "product_price": 100}, {"product_id": "P10", "product_name": "MicroSD Card", "product_desc": "SanDisk 128GB MicroSDXC", "product_price": 150}];
```

5) Retrieve products with price less than \$100

```
db.products.find({$price: {$lte: 100}});
```

Output:

```
[{"product_id": "P7", "product_name": "Laptop", "product_desc": "Dell XPS 15", "product_price": 1200}, {"product_id": "P4", "product_name": "Smartphone", "product_desc": "Samsung Galaxy S20", "product_price": 800}, {"product_id": "P1", "product_name": "Headphones", "product_desc": "Sony WH-1000XM4", "product_price": 600}, {"product_id": "P6", "product_name": "Monitor", "product_desc": "Dell U2720Q", "product_price": 1000}, {"product_id": "P2", "product_name": "Keyboard", "product_desc": "Logitech G910", "product_price": 500}, {"product_id": "P3", "product_name": "Mouse", "product_desc": "SteelSeries Rival 650", "product_price": 400}, {"product_id": "P5", "product_name": "Charger", "product_desc": "Anker PowerPort Atom III", "product_price": 200}, {"product_id": "P8", "product_name": "Power Bank", "product_desc": "Anker PowerCore 20000", "product_price": 300}, {"product_id": "P9", "product_name": "USB-C Cable", "product_desc": "AmazonBasics 6ft USB-C to USB-A", "product_price": 100}, {"product_id": "P10", "product_name": "MicroSD Card", "product_desc": "SanDisk 128GB MicroSDXC", "product_price": 150}];
```

6) Retrieved products added to a user's cart

Code:

```
db.user_carts.findone({user_id: "789ghi"}, {cart_items: 1});
```

Output:

```
[{"product_id": "P1", "quantity": 2}];
```

7) Retrieved orders placed by a User

Code:

```
db.orders.find({user_id: "123abc"});
```

Output:

```
[{"id": "order_id_123"}];
```

8) Retriving total count of orders placed by a User

```
db.orders.aggregate([{$match: {user_id: "123abc"}}, {"$group": {"user_id": "123abc", "count": {"$sum": 1}}}], {allowDiskUse: true});
```

classmate

Note

Page

3) Calculate total price of orders closed by each user:

db.user.aggregate([

 {

 \$group: {

 "_id": "user_id",

 "total_price": { \$sum: "\$total_price" } }

 }

]);

2) Find user with highest total price:

{ sort: { total_price: -1 } }

 { limit: 1 }

});

8) Find total avg price:

db.user.aggregate([

 { group: {

 "_id": "user_id",

 "total_price": { \$sum: "\$total_price" } }

 }

]);

Code with Output:

```
...
cqlsh> CREATE KEYSPACE Student WITH REPLICATION= {'class':'SimpleStrategy','replication_factor':1};
cqlsh> describe keyspaces;
'keyspaces' not found in keyspaces
cqlsh> describe keyspaces;

student      system      system_distributed  system_traces  system_virtual_schema
students    system_auth   system_schema       system_views

cqlsh> use students;
cqlsh:students> create table st_info(rollno int primary key, name text, doj timestamp, percent double);
cqlsh:students> describe tables;

library_book  st_info  students_info  userlogin

cqlsh:students> describe table<st_info>;
Improper describe command.
cqlsh:students> describe table st_info;

CREATE TABLE students.st_info (
  rollno int PRIMARY KEY,
  doj timestamp,
  name text,
  percent double
) 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 st_info(rollno, name, doj, percent)
cqlsh:students> select * from st_info;

rollno | doj                  | name    | percent
-----+-----+-----+-----+
  1 | 2010-02-28 18:30:00.000000+0000 | preeti |     90
  2 | 2010-03-19 18:30:00.000000+0000 | prajwal |     89
  4 | 2010-04-22 18:30:00.000000+0000 | rachana |     90

(3 rows)
cqlsh:students> select * from st_info where rollno in(1,2);

rollno | doj                  | name    | percent
-----+-----+-----+-----+
  1 | 2010-02-28 18:30:00.000000+0000 | preeti |     90
  2 | 2010-03-19 18:30:00.000000+0000 | prajwal |     89

(2 rows)
cqlsh:students> select * from st_info where name="preeti";
SyntaxException: line 1:42 no viable alternative at input ';' (...* from st_info where name=[preeti])
cqlsh:students> create index on st_info(name);
cqlsh:students> select * from st_info where name="preeti";
SyntaxException: line 1:42 no viable alternative at input ';' (...* from st_info where name=[preeti])
cqlsh:students> select * from st_info where name='preeti';

rollno | doj                  | name    | percent
-----+-----+-----+-----+
  1 | 2010-02-28 18:30:00.000000+0000 | preeti |     90

(1 rows)
cqlsh:students> select rollno, name, percent from st_info limit 2;

rollno | name    | percent
-----+-----+-----+
  1 | preeti |     90
  2 | prajwal |     89

(2 rows)
cqlsh:students> slect rollno as usn from st_info;
SyntaxException: line 1:0 no viable alternative at input 'slect' ([slect]...)
cqlsh:students> select rollno as usn from st_info;

usn
-----
  1
```

```
usn
-----
1
2
4

(3 rows)
cqlsh:students> create table library(c_val counter,book_name varchar,stud_name varchar,primary key(book_name,stud_name))
cqlsh:students> update library set c_val=c_val+1 where book_name='BDA' and stud_name='preeti';
cqlsh:students> create table userlogin(id int primary key,pass text);
AlreadyExists: Table 'students.userlogin' already exists
cqlsh:students> create table login(id int primary key,pass text);
cqlsh:students> insert into login(id,pass) values(1,'infy')using ttl 30;
cqlsh:students> select ttl(pass) from login where id=1;

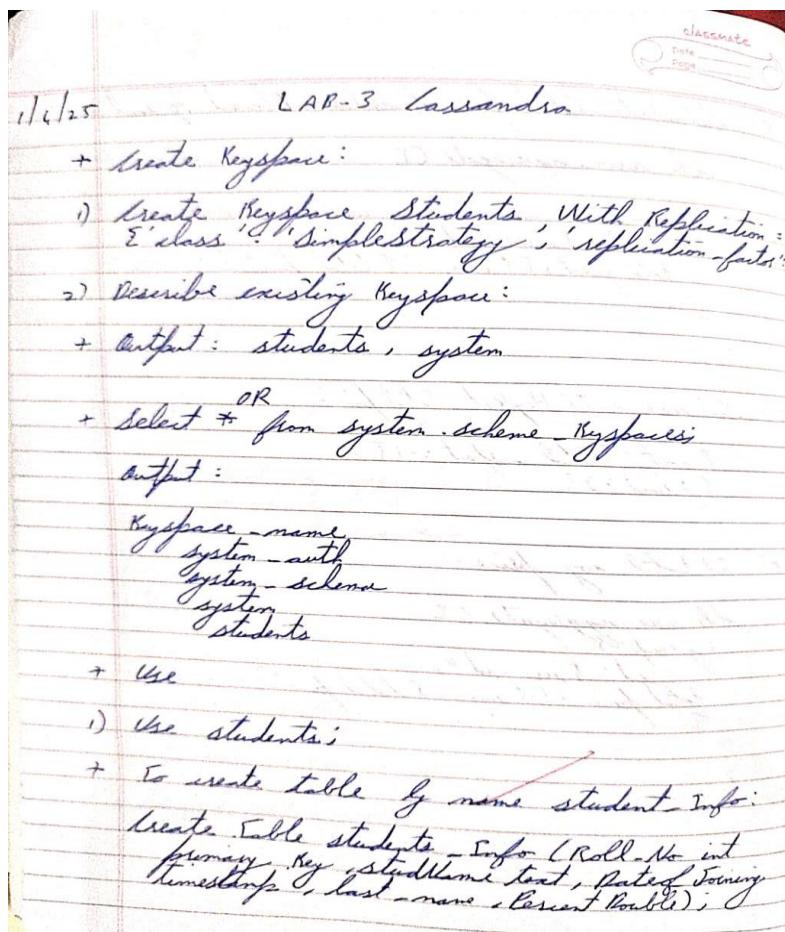
ttl(pass)
-----
3
```

Program 3

Perform the following DB operations using Cassandra.

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

Observation:



+ Describle Table:

i) Describle Table students - Info;

→ Create Table students . students . info (

- rollno int PK,
- dateofjoining,
- last_exam_percent
- studname text

+ Insertion:

i) Begin Batch

Insert into students_info (rollNo, studname,
last_exam_percent)

Values (1, 'Abha', '2012-03-12', 79.9);

Values (2, 'Kiran', '2012-03-12', 78.9);

Values (3, 'Smitha', '2012-03-12', 67.9);

Apply Batch;

+ Display:

Select * from students - Info:

→

roll no	datejoining	last_exam	studname
5	2023	67.9	smitha
2	2024	68.5	Abha

+ Queries:

i) Selection

+ Select * from students - Info where Roll_no (1,2,3);

+ Load data:

Update library book set counter_value:=counter_value+1 where name='Big data analytics'
And stud_name='Test'.

Exercise live:

→ Create Table userlogin (userid int PK, tent);
Insert Into userlogin(userid, pass) Userj, TEL70;
select TEL (pass) userlogin where userid=2;

Q
11/105

Code with Output:

```
hadoop@bmseccse-HP-Elite-Tower-600-G9-Desktop-PC:~$ cqlsh
Connected to Test Cluster at 127.0.0.1:9042
[cqlsh 6.1.0 | Cassandra 4.1.8 | CQL spec 3.4.6 | Native protocol v5]
Use HELP for help.
cqlsh> create keyspace library with replication={'class':'Simplestrategy','replication_factor':1};
ConfigurationException: Unable to find replication strategy class 'org.apache.cassandra.locator.Simplestrategy'
cqlsh> create keyspace library with replication={'class':'Simplestrategy','replication_factor':1};exit
ConfigurationException: Unable to find replication strategy class 'org.apache.cassandra.locator.Simplestrategy'
hadoop@bmseccse-HP-Elite-Tower-600-G9-Desktop-PC:~$ cqlsh
Connected to Test Cluster at 127.0.0.1:9042
[cqlsh 6.1.0 | Cassandra 4.1.8 | CQL spec 3.4.6 | Native protocol v5]
Use HELP for help.
cqlsh> create keyspace library with replication={'class':'SimpleStrategy','replication_factor':1};exit
hadoop@bmseccse-HP-Elite-Tower-600-G9-Desktop-PC:~$ cqlsh
Connected to Test Cluster at 127.0.0.1:9042
[cqlsh 6.1.0 | Cassandra 4.1.8 | CQL spec 3.4.6 | Native protocol v5]
Use HELP for help.
cqlsh> create keyspace library with replication={'class':'SimpleStrategy','replication_factor':1};
AlreadyExists: Keyspace 'library' already exists
cqlsh> exit
hadoop@bmseccse-HP-Elite-Tower-600-G9-Desktop-PC:~$ cqlsh
Connected to Test Cluster at 127.0.0.1:9042
[cqlsh 6.1.0 | Cassandra 4.1.8 | CQL spec 3.4.6 | Native protocol v5]
Use HELP for help.
cqlsh> create keyspace libraries with replication={'class':'SimpleStrategy','replication_factor':1};
cqlsh> keyspaces
...
cqlsh> describe keyspaces;

libraries    students      system_distributed   system_views
library      system       system_schema        system_virtual_schema
student     system_auth  system_traces

cqlsh> use libraries;
cqlsh:libraries> create table l_info(sid int primary key, c_val counter, sname varchar,bname varchar,bid int,doi timestamp);
InvalidRequest: Error from server: code=2200 [Invalid query] message="Cannot mix counter and non counter columns in the same table"
cqlsh:libraries> create table l_info(sid int primary key, sname varchar,bname varchar,bid int,doi timestamp);
cqlsh:libraries> create table count(sid int primary key,c_val counter);
cqlsh:libraries> begin batch
... insert into l_info(sid,sname,bname,bid,doi)
... values(112,'alice','bda',1,'2020-03-03')
... insert into l_info(sid,sname,bname,bid,doi)
... values(113,'preeti','cn',2,'2020-03-04')
... apply batch;
cqlsh:libraries> update l_info
```

```
```
cqlsh:libraries> select * from l_info;

 sid | bid | bname | doi
-----+-----+-----+-----+-----+-----+-----+
 113 | 2 | cn | 2020-03-03 18:30:00.000000+0000 | preeti
 112 | 1 | bda | 2020-03-02 18:30:00.000000+0000 | alice

(2 rows)
cqlsh:libraries> select * from count;

 sid | c_val
-----+-----
 112 | 1

(1 rows)
cqlsh:libraries> □
```

## Program 4

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

Observation:

18/4/25  
LAP-5

classmate  
Date \_\_\_\_\_  
Page \_\_\_\_\_

Hadoop Exercise :

- 1) start-all.sh (start all hadoop files)
- 2) Create a directory inside hadoop-1s  
hadoop fs -ls /
- 3) drwxr-xr-x -hadoop supergroup 0 2025-04-15  
14:23 /ldahadoop
- 4) Copy files from desktop having put command  
hdfs dfs -put /home/hadoop/Desktop/file.txt /hadoop/file.txt
- 5) Copying files using cp from local command  
hdfs dfs -cp /local /home/hadoop/ldahadoop/file.txt
- 6) cat command (-cat)  
hdfs dfs -cat /hadoop/file.txt
- 7) get command  
hdfs dfs -get /hadoop/file.txt /home/hadoop/downloads.txt
- 8) get merge command

lfs. lfs - getmotify /bda/file1.txt /file2.txt

9) ladoop fs - getfacl /bda/local/

# file : /bda/laptops/  
# file owner : hduces  
# group : supergroup  
user : suwe  
group : suwx  
dhes : r-x

10) copy to local

lfs lfs - copy blocal.bda.hadoop/file.txt

11) move command

ladoop fs - mv lab lab

12) copy command

ladoop fs - cp /hello/ ladoop/lab:

## Code with Output:

```
hadoop@bmscecse-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]
Starting datanodes
Starting secondary namenodes [bmscecse-HP-Elite-Tower-600-G9-Desktop-PC]
Starting resourcemanager
Starting nodemanagers
hadoop@bmscecse-HP-Elite-Tower-600-G9-Desktop-PC:~$ hdfs dfs -mkdir /bda_hadoop
hadoop@bmscecse-HP-Elite-Tower-600-G9-Desktop-PC:~$ hdfs dfs -ls /
Found 4 items
drwxr-xr-x - hadoop supergroup 0 2025-04-15 15:07 /abc
drwxr-xr-x - hadoop supergroup 0 2025-05-26 14:13 /bda_hadoop
drwxr-xr-x - hadoop supergroup 0 2025-05-22 16:32 /pqr
drwxr-xr-x - hadoop supergroup 0 2025-05-20 16:36 /rgs
hadoop@bmscecse-HP-Elite-Tower-600-G9-Desktop-PC:~$ hdfs dfs -put /home/hadoop/sample.txt /bda_hadoop/file.txt
hadoop@bmscecse-HP-Elite-Tower-600-G9-Desktop-PC:~$ hdfs dfs -copyFromLocal /home/hadoop/sample.txt /bda_hadoop/local.txt
hadoop@bmscecse-HP-Elite-Tower-600-G9-Desktop-PC:~$ hdfs dfs -cat /bda_hadoop/file.txt
hi how are you
how is your job
how is your family
how is your brother
how is your sister
eof
hadoop@bmscecse-HP-Elite-Tower-600-G9-Desktop-PC:~$ hdfs dfs -get /bda_hadoop/file.txt /home/hadoop/sample.txt
get: '/home/hadoop/sample.txt': File exists
hadoop@bmscecse-HP-Elite-Tower-600-G9-Desktop-PC:~$ hdfs dfs -get /bda_hadoop/file.txt /home/hadoop/get.txt
hadoop@bmscecse-HP-Elite-Tower-600-G9-Desktop-PC:~$ hadoop fs -getfacl /bda_hadoop/
file: /bda_hadoop
owner: hadoop
group: supergroup
user::rwx
group::r-x
other::r-x

hadoop@bmscecse-HP-Elite-Tower-600-G9-Desktop-PC:~$ hdfs dfs -copyToLocal /bda_hadoop/file.txt /home/hadoop/tolocal.txt
hadoop@bmscecse-HP-Elite-Tower-600-G9-Desktop-PC:~$ hadoop fs -cp /bda_hadoop /abc
hadoop@bmscecse-HP-Elite-Tower-600-G9-Desktop-PC:~$ hadoop fs -ls /abc
Found 3 items
drwxr-xr-x - hadoop supergroup 0 2025-05-26 14:28 /abc/bda_hadoop
-rw-r--r-- 1 hadoop supergroup 55 2025-04-15 15:05 /abc/file.txt
-rw-r--r-- 1 hadoop supergroup 55 2025-04-15 15:07 /abc/file_cp_.txt
hadoop@bmscecse-HP-Elite-Tower-600-G9-Desktop-PC:~$
```

## Program 5

### Implement Wordcount program on Hadoop framework

Observation:

1/5/25

discrete

Data

Page

LAB-6

Word Count Program Map reduce:

1) Mappercode.java

```
public class Mappercode Mapper
< long variable Int
 public void map (log
Variable key, Text Values, Context)
throws IOException
String [] words = Values.to string()
split()
for (String Word : words)
Context.write (Text, (1))
```

3

3

3

2) Reducer code.java

```
public class Reducer code extends Reducer
(Text, IntWritable, Text, IntWritable)
Values, Context
throws IOException
int sum = 0
for (Int Writable Val : Values)
sum + Val.get
3
Context.write (Key, new IntWritable (sum)),
3
```

3

→ Raw code :

```
public class RawCode {
 public static void main (String [] args) {
 throw exception
 Configuration conf = new configuration ();
 Job set = Recycles (this class)
 Job.set Reduces class (Reduce class)
 Job.setoutput by class (Text, class)
```

```
Enter Input format.addInput (Job, newpath)
(Lazy []) file output format.setoutput (Job. new
path, [seq []) System.out.println (Job for completion (true))
```

Hadoop code :

```
> hdfs dfs - mddir /text
> hdfs dfs - put, home /hadoop, sample
text input
> hadoop jar word convert Jar, Raw
Code/ input sample.txt /output
> hdfs dfs - cat /output /node 5-000
```

## Code with Output:

```
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 12082. Stop it first and ensure /tmp/hadoop-hadoop-namenode.pid file is empty before retry.
Starting datanodes
localhost: datanode is running as process 12255. 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]
secondarynamenode is running as process 12557. Stop it first and ensure /tmp/hadoop-hadoop-secondarnamenode.pid file is empty before retry.
Starting resourcemanager
resourcemanager is running as process 12845. Stop it first and ensure /tmp/hadoop-hadoop-resourcemanager.pid file is empty before retry.
Starting nodemanagers
localhost: nodemanager is running as process 13014. Stop it first and ensure /tmp/hadoop-hadoop-nodemanager.pid file is empty before retry.
hadoop@bmscsece-HP-Elite-Tower-600-G9-Desktop-PC:~$ jps
12082 NameNode
13014 NodeManager
15100 org.eclipse.equinox.launcher_1.6.1000.v20250227-1734.jar
17036 Jps
12845 ResourceManager
12557 SecondaryNameNode
12255 DataNode
hadoop@bmscsece-HP-Elite-Tower-600-G9-Desktop-PC:~$ hdfs dfs -copyFromLocal /home/hadoop/sample.txt /bda_hadoop/input.txt
hadoop@bmscsece-HP-Elite-Tower-600-G9-Desktop-PC:~$ hadoop jar /home/hadoop/Desktop/hdpwordcount.jar WCDriver /bda_hadoop/input.txt /bda_hadoop/output
Exception in thread "main" java.lang.ClassNotFoundException: WCDriver
 at java.base/java.net.URLClassLoader.findClass(URLClassLoader.java:476)
 at java.base/java.lang.ClassLoader.loadClass(ClassLoader.java:594)
 at java.base/java.lang.ClassLoader.loadClass(ClassLoader.java:527)
 at java.base/java.lang.Class.forName(Native Method)
 at java.base/java.lang.Class.forName(Class.java:398)
 at org.apache.hadoop.util.RunJar.run(RunJar.java:321)
 at org.apache.hadoop.util.RunJar.main(RunJar.java:241)
hadoop@bmscsece-HP-Elite-Tower-600-G9-Desktop-PC:~$ hadoop jar /home/hadoop/Desktop/hdpwordcount.jar hdpwordcount.WCDriver /bda_hadoop/input.txt /bda_hadoop/output
2025-05-26 14:40:01,404 INFO impl.MetricsConfig: Loaded properties from hadoop-metrics2.properties
2025-05-26 14:40:01,440 INFO impl.MetricsSystemImpl: Scheduled Metric snapshot period at 10 second(s).
2025-05-26 14:40:01,440 INFO impl.MetricsSystemImpl: JobTracker metrics system started
2025-05-26 14:40:01,446 WARN impl.MetricsSystemImpl: JobTracker metrics system already initialized!
2025-05-26 14:40:01,501 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to remedy this.
2025-05-26 14:40:01,545 INFO mapred.FileInputFormat: Total input files to process : 1
2025-05-26 14:40:01,567 INFO mapreduce.JobSubmitter: number of splits:1
2025-05-26 14:40:01,624 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_local276129153_0001
2025-05-26 14:40:01,624 INFO mapreduce.JobSubmitter: Executing with tokens: []
2025-05-26 14:40:01,677 INFO mapreduce.Job: The url to track the job: http://localhost:8080/
2025-05-26 14:40:01,679 INFO mapred.LocalJobRunner: OutputCommitter set in config null
2025-05-26 14:40:01,679 INFO mapreduce.Job: Running job: job_local276129153_0001
2025-05-26 14:40:01,680 INFO mapred.LocalJobRunner: OutputCommitter is org.apache.hadoop.mapred.FileOutputCommitter
2025-05-26 14:40:01,682 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 2
2025-05-26 14:40:01,682 INFO output.FileOutputCommitter: FileOutputCommitter skip cleanup _temporary folders under output directory:false, ign
```

```
hadoop@bmscsece-HP-Elite-Tower-600-G9-Desktop-PC:~$ hdfs dfs -ls /bda_hadoop/output
Found 2 items
-rw-r--r-- 1 hadoop supergroup 0 2025-05-26 14:40 /bda_hadoop/output/_SUCCESS
-rw-r--r-- 1 hadoop supergroup 75 2025-05-26 14:40 /bda_hadoop/output/part-00000
hadoop@bmscsece-HP-Elite-Tower-600-G9-Desktop-PC:~$ hdfs dfs -cat /bda_hadoop/output/part-00000
are 1
brother 1
eof 1
family 1
hi 1
how 5
is 4
job 1
sister 1
you 1
your 4
hadoop@bmscsece-HP-Elite-Tower-600-G9-Desktop-PC:~$
```

## Program 6

From the following link extract the weather data

<https://github.com/tomwhite/hadoop-book/tree/master/input/ncdc/all>

- Create a MapReduce program to find average temperature for each year from the NCDC data set.
- find the mean max temperature for every month

Observation:

LAB 06 Find Average Temperature

```
AVDriver.java
package WAverage;
import org.apache.hadoop.fs.Path;
public class AVDriver
{
 public static void main (String[] args)
 {
 if (args.length != 2)
 System.out.println ("Enter");
 System.exit (-1);
 job.setJarByClass (AVDriver.class);
 job.setJobName ('Max temperature');
 FileInputFormat.addInputPath (
 job, new Path (args[0]));
 FileOutputFormat.setOutputPath (
 job, new Path (args[1]));
 job.setOutputKeyClass (Text.class);
 job.setOutputValueClass (IntWritable.class);
 }
}
```

AVMapper.java

package MAvrage;

import java.io.IOException;

public class AVMapper extends Mapper<  
LongWritable, Text, IntWritable>

public static final int Missing = 999;

public void map(LongWritable key,  
Text value)

int temp;

String line = value.toString();

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

if (line.charAt(82) == '+')

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

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

if (temp <= 9999 && quality.matches  
("([01459])"))

context.write(new Text(year), new  
IntWritable(temp));

}

## AVReducer.java

```
package WAverage;

import java.io.IOException;

public class AVReducer extends Reducer
<Text, IntWritable, Text>
{
 public void reduce(Text key, Iterable
 <IntWritable> values,
 Text, IntWritable>.Collector content)
 throws IOException
 {
 int max = 0;
 int count = 0;
 for (IntWritable value : values)
 {
 max = max + value.get();
 count++;
 }
 content.write(key, new IntWritable
 (max / count));
 }
}
```

Output:

1901 46

MNMapper.java.

package mean;

import java.io.IOException;

public class MNMapper extends Mapper  
<LongWritable, Text, Text, IntWritable>  
{

    public static final int miss = 9999;

    public void map(LongWritable key,  
                Text value, Mapper<LongWritable,  
                Text, Text>)

    {

        String line = value.toString();

        String month = line.substring(19, 21);

        if (line.charAt(0) == '+')

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

        String q = line.substring(92, 93);

        if (temp != miss)

            write(new Text(month));

}

}

MNReducer.java

package mean;

import java.io.IOException;

public class MNReducer extends Reducer

public void reduce(Text key, Iterable<IntWritable>

int max = Int. min.

int totalTemp = 0;

int count = 0;

int days = 0;

for (IntW val : vals)

int t = val.get();

if (t > max)

max = temp;

c++

if (c == 3)

st. s = max;

max = Int. min

c = 0

d++;

}

## Code with Output:

### a) Average temperature

```
hadoop@bmscecse-HP-Elite-Tower-600-G9-Desktop-PC:~$ jps
12082 NameNode
17908 Jps
13014 NodeManager
15100 org.eclipse.equinox.launcher_1.6.1000.v20250227-1734.jar
12845 ResourceManager
12557 SecondaryNameNode
12255 DataNode
hadoop@bmscecse-HP-Elite-Tower-600-G9-Desktop-PC:~$ hdfs dfs -copyFromLocal /home/hadoop/Downloads/1901 /bda_hadoop/avinput.txt
hadoop@bmscecse-HP-Elite-Tower-600-G9-Desktop-PC:~$ hadoop jar /home/hadoop/Desktop/WeatherAverage.jar WeatherAverage.AVDriver /bda_h
adoop/avinput.txt /bda_hadoop/avoutput
2025-05-26 14:49:09,290 INFO impl.MetricsConfig: Loaded properties from hadoop-metrics2.properties
2025-05-26 14:49:09,327 INFO impl.MetricsSystemImpl: Scheduled Metric snapshot period at 10 second(s).
2025-05-26 14:49:09,327 INFO impl.MetricsSystemImpl: JobTracker metrics system started
2025-05-26 14:49:09,380 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to remedy this.
2025-05-26 14:49:09,427 INFO input.FileInputFormat: Total input files to process : 1
2025-05-26 14:49:09,452 INFO mapreduce.JobSubmitter: number of splits:1
2025-05-26 14:49:09,510 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_local1313646497_0001
2025-05-26 14:49:09,510 INFO mapreduce.JobSubmitter: Executing with tokens: []
2025-05-26 14:49:09,566 INFO mapreduce.Job: The url to track the job: http://localhost:8080/
2025-05-26 14:49:09,566 INFO mapreduce.Job: Running job: job_local1313646497_0001
2025-05-26 14:49:09,567 INFO mapred.LocalJobRunner: OutputCommitter set in config null
2025-05-26 14:49:09,570 INFO output.PathOutputCommitterFactory: No output committer factory defined, defaulting to FileOutputCommitte
rFactory
2025-05-26 14:49:09,571 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 2
2025-05-26 14:49:09,571 INFO output.FileOutputCommitter: FileOutputCommitter skip cleanup _temporary folders under output directory: f
alse, ignore cleanup failures: false
2025-05-26 14:49:09,571 INFO mapred.LocalJobRunner: OutputCommitter is org.apache.hadoop.mapreduce.lib.output.FileOutputCommitter
2025-05-26 14:49:09,620 INFO mapred.LocalJobRunner: Waiting for map tasks
2025-05-26 14:49:09,620 INFO mapred.LocalJobRunner: Starting task: attempt_local1313646497_0001_m_000000_0
2025-05-26 14:49:09,629 INFO output.PathOutputCommitterFactory: No output committer factory defined, defaulting to FileOutputCommitte
rFactory
2025-05-26 14:49:09,629 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 2
2025-05-26 14:49:09,629 INFO output.FileOutputCommitter: FileOutputCommitter skip cleanup _temporary folders under output directory: f
alse, ignore cleanup failures: false
2025-05-26 14:49:09,635 INFO mapred.Task: Using ResourceCalculatorProcessTree : []
2025-05-26 14:49:09,637 INFO mapred.MapTask: Processing split: hdfs://localhost:9000/bda_hadoop/avinput.txt:0+888190
2025-05-26 14:49:09,666 INFO mapred.MapTask: (EQUATOR) 0 kvi 26214396(104857584)
2025-05-26 14:49:09,666 INFO mapred.MapTask: mapreduce.task.io.sort.mb: 100
2025-05-26 14:49:09,666 INFO mapred.MapTask: soft limit at 83886080
2025-05-26 14:49:09,666 INFO mapred.MapTask: bufstart = 0; bufvoid = 104857600
2025-05-26 14:49:09,666 INFO mapred.MapTask: kvstart = 26214396; length = 6553600
2025-05-26 14:49:09,668 INFO mapred.MapTask: Map output collector class = org.apache.hadoop.mapred.MapTask$MapOutputBuffer
2025-05-26 14:49:09,730 INFO mapred.LocalJobRunner:
2025-05-26 14:49:09,731 INFO mapred.MapTask: Starting flush of map output
2025-05-26 14:49:09,731 INFO mapred.MapTask: Spilling map output
2025-05-26 14:49:09,731 INFO mapred.MapTask: bufstart = 0; bufend = 59076; bufvoid = 104857600
2025-05-26 14:49:09,731 INFO mapred.MapTask: kvstart = 26214396(104857584); kvend = 26188144(104752576); length = 26253/6553600
2025-05-26 14:49:09,739 INFO mapred.MapTask: Finished spill 0
2025-05-26 14:49:09,743 INFO mapred.Task: Task:attempt_local1313646497_0001_m_000000_0 is done. And is in the process of committing
2025-05-26 14:49:09,745 INFO mapred.LocalJobRunner: map
```

```
Merged Map Outputs=1
GC time elapsed (ms)=0
Total committed heap usage (bytes)=1052770304
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=888190
File Output Format Counters
Bytes Written=
hadoop@bmscecse-HP-Elite-Tower-600-G9-Desktop-PC: $ hdfs dfs -ls /bda_hadoop/avoutput
Found 2 items
-rw-r--r-- 1 hadoop supergroup 0 2025-05-26 14:49 /bda_hadoop/avoutput/_SUCCESS
-rw-r--r-- 1 hadoop supergroup 8 2025-05-26 14:49 /bda_hadoop/avoutput/part-r-00000
hadoop@bmscecse-HP-Elite-Tower-600-G9-Desktop-PC: $ hdfs dfs -cat /bda_hadoop/avoutput/part-r-00000
1901 46
hadoop@bmscecse-HP-Elite-Tower-600-G9-Desktop-PC: $
```

## b) Maximum temperature

```
hadoop@bmscecsse-HP-Elite-Tower-600-G9-Desktop-PC:~$ jps
18721 Jps
12082 NameNode
13014 NodeManager
15100 org.eclipse.equinox.launcher_1.6.1000.v20250227-1734.jar
12845 ResourceManager
12557 SecondaryNameNode
12255 DataNode
hadoop@bmscecsse-HP-Elite-Tower-600-G9-Desktop-PC:~$ hdfs dfs -copyFromLocal /home/hadoop/Downloads/1901 /bda_hadoop/minput.txt
hadoop@bmscecsse-HP-Elite-Tower-600-G9-Desktop-PC:~$ hadoop jar /home/hadoop/Desktop/meanTemp.jar Mean.MNDriver /bda_hadoop/minput.txt /bda_hadoop/moutput
2025-05-26 14:54:41,993 INFO impl.MetricsConfig: Loaded properties from hadoop-metrics2.properties
2025-05-26 14:54:42,029 INFO impl.MetricSystemImpl: Scheduled Metric snapshot period at 10 second(s).
2025-05-26 14:54:42,029 INFO impl.MetricSystemImpl: JobTracker metrics system started
2025-05-26 14:54:42,083 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to remedy this.
2025-05-26 14:54:42,131 INFO input.FileInputFormat: Total input files to process : 1
2025-05-26 14:54:42,158 INFO mapreduce.JobSubmitter: number of splits:1
2025-05-26 14:54:42,216 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_local862196817_0001
2025-05-26 14:54:42,216 INFO mapreduce.JobSubmitter: Executing with tokens: []
2025-05-26 14:54:42,272 INFO mapreduce.Job: The url to track the job: http://localhost:8080/
2025-05-26 14:54:42,273 INFO mapreduce.Job: Running job: job_local862196817_0001
2025-05-26 14:54:42,273 INFO mapred.LocalJobRunner: OutputCommitter set in config null
2025-05-26 14:54:42,276 INFO output.PathOutputCommitterFactory: No output committer factory defined, defaulting to FileOutputCommitterFactory
2025-05-26 14:54:42,277 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 2
2025-05-26 14:54:42,277 INFO output.FileOutputCommitter: FileOutputCommitter skip cleanup _temporary folders under output directory:false, ignore cleanup failures: false
2025-05-26 14:54:42,277 INFO mapred.LocalJobRunner: OutputCommitter is org.apache.hadoop.mapreduce.lib.output.FileOutputCommitter
2025-05-26 14:54:42,319 INFO mapred.LocalJobRunner: Waiting for map tasks
2025-05-26 14:54:42,319 INFO mapred.LocalJobRunner: Starting task: attempt_local862196817_0001_m_000000_0
2025-05-26 14:54:42,328 INFO output.PathOutputCommitterFactory: No output committer factory defined, defaulting to FileOutputCommitterFactory
2025-05-26 14:54:42,329 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 2
2025-05-26 14:54:42,329 INFO output.FileOutputCommitter: FileOutputCommitter skip cleanup _temporary folders under output directory:false, ignore cleanup failures: false
2025-05-26 14:54:42,335 INFO mapred.Task: Using ResourceCalculatorProcessTree : []
2025-05-26 14:54:42,336 INFO mapred.MapTask: Processing split: hdfs://localhost:9000/bda_hadoop/minput.txt:0+888190
2025-05-26 14:54:42,366 INFO mapred.MapTask: (EQUATOR) 0 kvi 26214396(104857584)
2025-05-26 14:54:42,366 INFO mapred.MapTask: mapreduce.task.io.sort.mb: 100
2025-05-26 14:54:42,366 INFO mapred.MapTask: soft limit at 83886080
2025-05-26 14:54:42,366 INFO mapred.MapTask: bufstart = 0; bufvoid = 104857600
2025-05-26 14:54:42,366 INFO mapred.MapTask: kvstart = 26214396; length = 6553600
2025-05-26 14:54:42,368 INFO mapred.MapTask: Map output collector class = org.apache.hadoop.mapred.MapTask$MapOutputBuffer
2025-05-26 14:54:42,428 INFO mapred.LocalJobRunner:
2025-05-26 14:54:42,428 INFO mapred.MapTask: Starting flush of map output
2025-05-26 14:54:42,429 INFO mapred.MapTask: Spilling map output
2025-05-26 14:54:42,429 INFO mapred.MapTask: bufstart = 0; bufend = 45948; bufvoid = 104857600
```

```
BAD_ID=0
CONNECTION=0
IO_ERROR=0
WRONG_LENGTH=0
WRONG_MAP=0
WRONG_REDUCE=0
File Input Format Counters
 Bytes Read=888190
File Output Format Counters
 Bytes Written=81
hadoop@bmscecsse-HP-Elite-Tower-600-G9-Desktop-PC:~$ hdfs dfs -ls /bda_hadoop/moutput
Found 2 items
-rw-r--r-- 1 hadoop supergroup 0 2025-05-26 14:54 /bda_hadoop/moutput/_SUCCESS
-rw-r--r-- 1 hadoop supergroup 81 2025-05-26 14:54 /bda_hadoop/moutput/part-r-00000
hadoop@bmscecsse-HP-Elite-Tower-600-G9-Desktop-PC:~$ hdfs dfs -cat /bda_hadoop/moutput/part-r-00000
01 -13
02 -66
03 -15
04 43
05 100
06 168
07 219
08 198
09 141
10 100
11 1
12 -61
hadoop@bmscecsse-HP-Elite-Tower-600-G9-Desktop-PC:~$
```

## Program 7

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

Observation:

1/5/25 Lab - 7

Minimax . Program map reduce  
Minimax Mapper

Public class Minimax Mapper extends Mapper<Text, Text, Doc> {

private DoubleWritable Value = new DoubleWritable();  
private double local min = Double.NEGATIVE\_INFINITY;  
private double max = Double.NEGATIVE\_INFINITY;

throws IOException, InterruptedException {  
double val = Double.parseDouble(value.get());  
local min = Math.min(localmin, val);  
local max = Math.max(localmax, val);  
}

private void Reducer() {  
protected void Reduce(Text key, Iterable<DoubleWritable> values, Context context) throws IOException, InterruptedException {  
double mi = Double.MAX\_VALUE;  
for (DoubleWritable val : values) {  
mi = Math.min(mi, val.get());  
}

content.write (new Text ("Global Min"))  
new Double .writeTable (min)

3

### Min Max Driver

```
public class MinMaxDriver {
 public static void main (String [] args) {
 try {
 if (args.length == 1) {
 System.out ("Usage")
 System.out (-1)
 }
 } catch (Exception e) {
 e.printStackTrace ();
 }
 }
}
```

Configuration conf = new Configuration ()  
Job job = conf.getJob ("job")

MinMaxDriver !

job.setMapperClass (MinMaxMapper.class)

job.setReducerClass (MinMaxReducer.class)

Table.class

File inputFormat.addInput (job, mapots)

System.out.println (job.waitForCompletion  
(true))

0/0 37

Sample Input file:

10

3

55

27

8

91

17

o/p  
Raffer  
Run 3  
Run 55  
Raffer 2  
Run 8  
Run 91

Reducers

Run 3  
Run 8  
Run 5.5  
Run 9

$$\text{Global Min} = \min(3, 8) = 3$$

$$\text{Global Max} = \max(55, 91) = 91$$

## Code with Output:

```
hadoop@bmscecse-HP-Elite-Tower-600-G9-Desktop-PC:~$ jps
12082 NameNode
19238 Jps
13014 NodeManager
15100 org.eclipse.equinox.launcher_1.6.1000.v20250227-1734.jar
12845 ResourceManager
12557 SecondaryNameNode
12255 DataNode
hadoop@bmscecse-HP-Elite-Tower-600-G9-Desktop-PC:~$ hdfs dfs -copyFromLocal /home/hadoop/Desktop/TopN.txt /bda_hadoop/tinput.txt
copyFromLocal: '/bda_hadoop/tinput.txt': No such file or directory: 'hdfs://localhost:9000/bda_hadoop/tinput.txt'
hadoop@bmscecse-HP-Elite-Tower-600-G9-Desktop-PC:~$ hdfs dfs -copyFromLocal /home/hadoop/Desktop/TopN.txt /bda_hadoop/tinput.txt
hadoop@bmscecse-HP-Elite-Tower-600-G9-Desktop-PC:~$ hadoop jar /home/hadoop/Desktop/TopN.jar TopN.TNDriver /bda_hadoop/tinput.txt /bda_hadoop/toutput
2025-05-26 14:59:03,334 INFO impl.MetricsConfig: Loaded properties from hadoop-metrics2.properties
2025-05-26 14:59:03,372 INFO impl.MetricsSystemImpl: Scheduled Metric snapshot period at 10 second(s).
2025-05-26 14:59:03,372 INFO impl.MetricsSystemImpl: JobTracker metrics system started
2025-05-26 14:59:03,426 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. I implement the Tool interface and execute your application with ToolRunner to remedy this.
2025-05-26 14:59:03,472 INFO input.FileInputFormat: Total input files to process : 1
2025-05-26 14:59:03,497 INFO mapreduce.JobSubmitter: number of splits:1
2025-05-26 14:59:03,554 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_local1824101299_0001
2025-05-26 14:59:03,554 INFO mapreduce.JobSubmitter: Executing with tokens: []
2025-05-26 14:59:03,609 INFO mapreduce.Job: The url to track the job: http://localhost:8080/
2025-05-26 14:59:03,610 INFO mapreduce.Job: Running job: job_local1824101299_0001
2025-05-26 14:59:03,610 INFO mapred.LocalJobRunner: OutputCommitter set in config null
2025-05-26 14:59:03,614 INFO output.PathOutputCommitterFactory: No output committer factory defined, defaulting to FileOutputCommitterFactory
2025-05-26 14:59:03,614 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 2
2025-05-26 14:59:03,614 INFO output.FileOutputCommitter: FileOutputCommitter skip cleanup _temporary folders under output directory:false, ignore cleanup failures: false
2025-05-26 14:59:03,614 INFO mapred.LocalJobRunner: OutputCommitter is org.apache.hadoop.mapreduce.lib.output.FileOutputCommitter
2025-05-26 14:59:03,654 INFO mapred.LocalJobRunner: Waiting for map tasks
2025-05-26 14:59:03,655 INFO mapred.LocalJobRunner: Starting task: attempt_local1824101299_0001_m_000000_0
2025-05-26 14:59:03,664 INFO output.PathOutputCommitterFactory: No output committer factory defined, defaulting to FileOutputCommitterFactory
2025-05-26 14:59:03,664 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 2
2025-05-26 14:59:03,664 INFO output.FileOutputCommitter: FileOutputCommitter skip cleanup _temporary folders under output directory:false, ignore cleanup failures: false
2025-05-26 14:59:03,670 INFO mapred.Task: Using ResourceCalculatorProcessTree : []
2025-05-26 14:59:03,672 INFO mapred.MapTask: Processing split: hdfs://localhost:9000/bda_hadoop/tinput.txt:0+95
2025-05-26 14:59:03,701 INFO mapred.MapTask: (EQUATOR) 0 kvi 26214396(104857584)
2025-05-26 14:59:03,701 INFO mapred.MapTask: mapreduce.task.io.sort.mb: 100
2025-05-26 14:59:03,701 INFO mapred.MapTask: soft limit at 83886080
2025-05-26 14:59:03,701 INFO mapred.MapTask: bufstart = 0; bufvoid = 104857600
2025-05-26 14:59:03,701 INFO mapred.MapTask: kvstart = 26214396; length = 6553600
2025-05-26 14:59:03,702 INFO mapred.MapTask: Map output collector class = org.apache.hadoop.mapred.MapTask$MapOutputBuffer
2025-05-26 14:59:03,738 INFO mapred.LocalJobRunner:
2025-05-26 14:59:03,739 INFO mapred.MapTask: Starting flush of map output
```

```

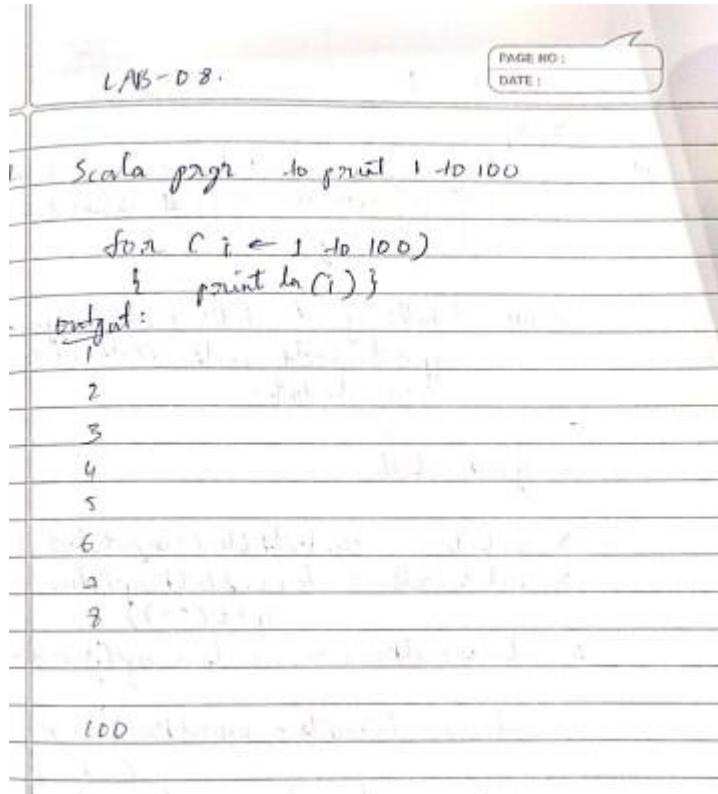
File System Counters
 FILE: Number of bytes read=10682
 FILE: Number of bytes written=1291808
 FILE: Number of read operations=0
 FILE: Number of large read operations=0
 FILE: Number of write operations=0
 HDFS: Number of bytes read=190
 HDFS: Number of bytes written=40
 HDFS: Number of read operations=15
 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=3
 Map output records=15
 Map output bytes=154
 Map output materialized bytes=190
 Input split bytes=108
 Combine input records=0
 Combine output records=0
 Reduce input groups=5
 Reduce shuffle bytes=190
 Reduce input records=15
 Reduce output records=5
 Spilled Records=30
 Shuffled Maps =1
 Failed Shuffles=0
 Merged Map outputs=1
 GC time elapsed (ms)=0
 Total committed heap usage (bytes)=1052770304
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=95
File Output Format Counters
 Bytes Written=40
hadoop@bmscecse-HP-Elite-Tower-600-G9-Desktop-PC:~$ hdfs dfs -ls /bda_hadoop/toutput
Found 2 items
-rw-r--r-- 1 hadoop supergroup 0 2025-05-26 14:59 /bda_hadoop/toutput/_SUCCESS
-rw-r--r-- 1 hadoop supergroup 40 2025-05-26 14:59 /bda_hadoop/toutput/part-r-00000
hadoop@bmscecse-HP-Elite-Tower-600-G9-Desktop-PC:~$ hdfs dfs -cat /bda_hadoop/toutput/part-r-00000
banana 5
apple 4
fruit 3
mango 2
kiwi 1
hadoop@bmscecse-HP-Elite-Tower-600-G9-Desktop-PC:~$

```

# Program 8

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

Observation:



Code with Output:

```
bmscse@bmscsece-HP-Elite-Tower-600-G9-Desktop-PC: $ spark-shell
25/05/26 15:58:23 WARN Utils: Your hostname, bmscsece-HP-Elite-Tower-600-G9-Desktop-PC resolves to a loopback address: 127.0.1.1; using 10.124.5.27 instead (on interface eno1)
25/05/26 15:58:23 WARN Utils: Set SPARK_LOCAL_IP if you need to bind to another address
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
25/05/26 15:58:26 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
Spark context Web UI available at http://10.124.5.27:4040
Spark context available as 'sc' (master = local[*], app id = local-1748255306894).
Spark session available as 'spark'.
Welcome to

 / \
 / \
 / \
/ \
 \ /
 \ /
 _____/ version 3.5.4

Using Scala version 2.12.18 (OpenJDK 64-Bit Server VM, Java 11.0.26)
Type in expressions to have them evaluated.
Type :help for more information.

scala> for(i <- 1 to 100){println(i)};
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
```

## Program 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

Observation:

Lab - 9

Using RDD & flatMap, how many times each word appears in file & write out a list of words whose count is strictly greater than 6 using spark.

Code:

```
val input = "hello world hello spark
hello sofa hello spark
hello hello"
```

```
val inputRDD = sc.parallelize(Seq(input))
```

```
val wordCounts = inputRDD.
 .flatMap(_.split(" ")).
 .map(_.toLowerCase).
 .filter(_.nonEmpty).
 .map(word => (word, 1)).
 .reduceByKey(_ + _).
 .filter((_, count) => count > 4)
```

```
wordCounts.collect().foreach { case word
 count => println(s"$word : $count") }
```

Output:

hello : 6

## Code with Output:

```
bmscecse@bmscecse-HP-Elite-Tower-600-G9-Desktop-PC:~$ echo "code code code code code spark spark spark spark spark hell
o hello hi hi joe ken">input.txt
bmscecse@bmscecse-HP-Elite-Tower-600-G9-Desktop-PC:~$ spark-shell
25/05/26 16:01:15 WARN Utils: Your hostname, bmscecse-HP-Elite-Tower-600-G9-Desktop-PC resolves to a loopback address:
127.0.1.1; using 10.124.5.27 instead (on interface eno1)
25/05/26 16:01:15 WARN Utils: Set SPARK_LOCAL_IP if you need to bind to another address
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
25/05/26 16:01:17 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java c
lasses where applicable
Spark context Web UI available at http://10.124.5.27:4040
Spark context available as 'sc' (master = local[*], app id = local-1748255477930).
Spark session available as 'spark'.
Welcome to

version 3.5.4
Using Scala version 2.12.18 (OpenJDK 64-Bit Server VM, Java 11.0.26)
Type in expressions to have them evaluated.
Type :help for more information.

scala> val lines=sc.textFile("input.txt")
lines: org.apache.spark.rdd.RDD[String] = input.txt MapPartitionsRDD[1] at textFile at <console>:23

scala> val words=lines.flatMap(line => line.split(" "))
<console>:23: error: value flatmap is not a member of org.apache.spark.rdd.RDD[String]
 val words=lines.flatMap(line => line.split(" "))
 ^

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

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

scala> val wordcount = wordParts.reduceByKey(_+_)
wordcount: org.apache.spark.rdd.RDD[(String, Int)] = ShuffledRDD[4] at reduceByKey at <console>:23

scala> val freq = wordcount.filter {case (word,count) => count > 4}
freq: org.apache.spark.rdd.RDD[(String, Int)] = MapPartitionsRDD[5] at filter at <console>:23

scala> freq.collect().foreach(println)
(spark,5)
(code,5)

scala> ■
```

## Program 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. (Open Ended Question).

### Observation:

20/5/25 Lab - 10

classmate  
Date \_\_\_\_\_  
Page \_\_\_\_\_

+ write a simple streaming program in spark text data streams on a particular port, perform basic text cleaning & print the cleaned text on the screen.

Code: scala - spark

```
import org.apache.spark.sql.functions
import org.apache.spark.sql.types
def lemmatize(word: String): String = word match
 case w if w.endsWith("ing") => w.substring(0, w.length - 3)
 case _ => word
val cleanText = udf((line: String) => {
 line.toLowerCase()
 .replaceAll("[^\\a-zA-Z]", " ")
 .split(" ")
 .map(lemmatize)
 .mkString(" ")
})
val line: spark.readStream
val query = "socket"
stream("host", "localhost")
load()
val closed: DataFrame = line.select("value")
```

very - antonymation)

New terminal:

nc - lk. 999

Text = this is a good day.

Output:

Batch 1:

Value                    cleaned  
~~this is a good day~~    good day

OK  
OK















