

# VISVESVARAYA TECHNOLOGICAL UNIVERSITY

“JnanaSangama”, Belgaum -590014, Karnataka.



## LAB REPORT

on

### **Big Data Analytics(23CS6PCBDA)**

*Submitted by*

**MEHUL VALLAMKONDA (1BM22CS153)**

*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**

**Feb-2025 to June-2025**

**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 **Mehul Vallamkonda(1BM22CS153)**, 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 a **Big Data Analytics – (23CS6PCBDA)** work prescribed for the said degree.

Amruta  
Assistant Professor  
Department of CSE  
BMSCE, Bengaluru

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

## Index Sheet

<b>Sl. No.</b>	<b>Experiment Title</b>	<b>Page No.</b>
1	MongoDB	4
2	MongoDB	8
3	Neo4j	17
4	Cassandra	20
5	Cassandra	22
6	Hadoop	27
7	Hadoop	36
8	Hadoop	42
9	Scala	51
8	Spark	53

## Course Outcome

<b>CO1</b>	Apply the concept of NoSQL, Hadoop or Spark for a given task
<b>CO2</b>	Analyse big data analytics mechanisms that can be applied to obtain solution for a given problem.
<b>CO3</b>	Design and implement solutions using data analytics mechanisms for a given problem.

## Lab 1: MongoDB- CRUD Demonstration

Question: Perform basic CRUD (Create, Read, Update, Delete) operations in MongoDB.

## Code with Output:



```
Atlas atlas-ws5rct-shard-0 [primary] mydb> db.createCollection('Customers')
{ ok: 1 }

Atlas atlas-ws5rct-shard-0 [primary] mydb> db.createCollection('Student')
{ ok: 1 }

Atlas atlas-ws5rct-shard-0 [primary] mydb> db.createCollection('Student')
{ ok: 1 }

Atlas atlas-ws5rct-shard-0 [primary] test> use mydb
switched to db mydb

Atlas atlas-ws5rct-shard-0 [primary] mydb> db.Student.deleteMany({Grade:'VII'})
{ acknowledged: true, deletedCount: 3 }

Atlas atlas-ws5rct-shard-0 [primary] mydb> db.Student.deleteOne({StudName:'JacobAdam'})
{ acknowledged: true, deletedCount: 0 }

Atlas atlas-ws5rct-shard-0 [primary] mydb> db.Student.drop()
true

Atlas atlas-ws5rct-shard-0 [primary] mydb> db.dropDatabase()
{ ok: 1, dropped: 'mydb' }

Atlas atlas-ws5rct-shard-0 [primary] mydb> db.Student.remove({StudName:'JacobAdam'})
)
DeprecationWarning: Collection.remove() is deprecated. Use deleteOne, deleteMany, findOneAndDelete, or bulkWrite.
{ acknowledged: true, deletedCount: 0 }
```

```
[Atlas atlas-ws5rct-shard-0 [primary] mydb> db.Customers.find()
[
  {
    _id: ObjectId('67c6c71f812483cc27dd4a64'),
    cust_id: 1,
    balance: 200,
    type: 'S'
  },
  {
    _id: ObjectId('67c6c739812483cc27dd4a65'),
    cust_id: 1,
    balance: 1000,
    type: 'Z'
  },
  {
    _id: ObjectId('67c6c74d812483cc27dd4a66'),
    cust_id: 2,
    balance: 100,
    type: 'Z'
  },
  {
    _id: ObjectId('67c6c75e812483cc27dd4a67'),
    cust_id: 2,
    balance: 1000,
    type: 'C'
  },
  {
    _id: ObjectId('67c6c76e812483cc27dd4a68'),
    cust_id: 2,
    balance: 500,
    type: 'C'
  },
  {
    _id: ObjectId('67c6c781812483cc27dd4a69'),
    cust_id: 2,
    balance: 50,
    type: 'S'
  },
  {
    _id: ObjectId('67c6c795812483cc27dd4a6a'),
    cust_id: 3,
    balance: 500,
    type: 'Z'
  }
]
```

Google Classroom

```
Atlas atlas-ws5rct-shard-0 [primary] mydb> db.Student.find()
[
  {
    _id: ObjectId('67c6c3c3812483cc27dd4a5d'),
    RollNo: 1,
    Age: 21,
    Cont: 9876,
    email: 'antara.de9@gmail.com'
  },
  {
    _id: ObjectId('67c6c3d8812483cc27dd4a5e'),
    RollNo: 1,
    Age: 21,
    Cont: 9876,
    email: 'antara.de9@gmail.com'
  },
  {
    _id: ObjectId('67c6c458812483cc27dd4a5f'),
    RollNo: 2,
    Age: 22,
    Cont: 9976,
    email: 'anushka.de@gmail.com'
  },
  {
    _id: ObjectId('67c6c47f812483cc27dd4a60'),
    RollNo: 3,
    Age: 21,
    Cont: 5576,
    email: 'anubhav.de@gmail.com'
  },
  {
    _id: ObjectId('67c6c4a2812483cc27dd4a61'),
    RollNo: 4,
    Age: 20,
    Cont: 4476,
    email: 'pani.de9@gmail.com'
  },
  {
    _id: ObjectId('67c6c4db812483cc27dd4a62'),
    RollNo: 10,
    Age: 23,
    Cont: 2276,
    email: 'rekha.de9@gmail.com'
  }
]
```

```
Atlas atlas-ws5rct-shard-0 [primary] mydb> db.Student.find()
[[
  {
    _id: ObjectId('67c6c3c3812483cc27dd4a5d'),
    RollNo: 1,
    Age: 21,
    Cont: 9876,
    email: 'antara.de9@gmail.com'
  },
  {
    _id: ObjectId('67c6c3d8812483cc27dd4a5e'),
    RollNo: 1,
    Age: 21,
    Cont: 9876,
    email: 'antara.de9@gmail.com'
  },
  {
    _id: ObjectId('67c6c458812483cc27dd4a5f'),
    RollNo: 2,
    Age: 22,
    Cont: 9976,
    email: 'anushka.de@gmail.com'
  },
  {
    _id: ObjectId('67c6c47f812483cc27dd4a60'),
    RollNo: 3,
    Age: 21,
    Cont: 5576,
    email: 'anubhav.de@gmail.com'
  },
  {
    _id: ObjectId('67c6c4a2812483cc27dd4a61'),
    RollNo: 4,
    Age: 20,
    Cont: 4476,
    email: 'pani.de9@gmail.com'
  },
  {
    _id: ObjectId('67c6c4db812483cc27dd4a62'),
    RollNo: 10,
    Age: 23,
    Cont: 2276,
    email: 'abhinav@gmail.com'
  },
]]
```

```
{  
  _id: ObjectId('67c6c616812483cc27dd4a63'),  
  RollNo: 11,  
  Age: 22,  
  Name: 'FEM',  
  cont: 2276,  
  email: 'rea.de9@gmail.com'  
},  
{  
  _id: 1,  
  StudName: 'Michelle Jacintha',  
  Grade: 'VII',  
  Hobbies: 'InternetSurfing'  
},  
{ _id: 2, StudName: 'Jannie', Grade: 'VIII', Hobbies: 'Music' },  
{ _id: 3, StudName: 'Jacob Adam', Grade: 'VII', Hobbies: 'Swimming' },  
{  
  _id: 4,  
  StudName: 'Amy Jacks',  
  Grade: 'X',  
  Hobbies: 'Dancing',  
  Location: 'Network'  
},  
{ _id: 6, StudName: 'Aryan David', Grade: 'VII', Hobbies: 'Skating' }  
]
```

```

Atlas atlas-ws5rct-shard-0 [primary] mydb> db.Customers.insert({cust_id:1,balance:200,type:'S'})
{
  acknowledged: true,
  insertedIds: { '0': ObjectId('67c6c71f812483cc27dd4a64') }
}
Atlas atlas-ws5rct-shard-0 [primary] mydb> db.Customers.insert({cust_id:1,balance:1000,type:'Z'})
{
  acknowledged: true,
  insertedIds: { '0': ObjectId('67c6c739812483cc27dd4a65') }
}
[Atlas atlas-ws5rct-shard-0 [primary] mydb> db.Customers.insert({cust_id:2,balance:100,type:'Z'})
{
  acknowledged: true,
  insertedIds: { '0': ObjectId('67c6c74d812483cc27dd4a66') }
}
[Atlas atlas-ws5rct-shard-0 [primary] mydb> db.Customers.insert({cust_id:2,balance:1000,type:'C'})
{
  acknowledged: true,
  insertedIds: { '0': ObjectId('67c6c75e812483cc27dd4a67') }
}
[Atlas atlas-ws5rct-shard-0 [primary] mydb> db.Customers.insert({cust_id:2,balance:500,type:'C'})
{
  acknowledged: true,
  insertedIds: { '0': ObjectId('67c6c76e812483cc27dd4a68') }
}
Atlas atlas-ws5rct-shard-0 [primary] mydb> db.Customers.insert({cust_id:2,balance:50,type:'S'})
{
  acknowledged: true,
  insertedIds: { '0': ObjectId('67c6c781812483cc27dd4a69') }
}
[Atlas atlas-ws5rct-shard-0 [primary] mydb> db.Customers.insert({cust_id:3,balance:500,type:'Z'})
{
  acknowledged: true,
  insertedIds: { '0': ObjectId('67c6c795812483cc27dd4a6a') }
}

```

```

Atlas atlas-ws5rct-shard-0 [primary] mydb> db.Student.insert({RollNo:1,Age:21,Cont:t:9876,email:'antara.de9@gmail.com'});
DeprecationWarning: Collection.insert() is deprecated. Use insertOne, insertMany, or bulkWrite.
{
  acknowledged: true,
  insertedIds: { '0': ObjectId('67c6c3c3812483cc27dd4a5d') }
}
Atlas atlas-ws5rct-shard-0 [primary] mydb> db.Student.insertOne({RollNo:1,Age:21,Cont:t:9876,email:'antara.de9@gmail.com'});
{
  acknowledged: true,
  insertedId: ObjectId('67c6c3d8812483cc27dd4a5e')
}
Atlas atlas-ws5rct-shard-0 [primary] mydb> show mydb
MongoshInvalidInputError: [COMMON-10001] 'mydb' is not a valid argument for "show".
Atlas atlas-ws5rct-shard-0 [primary] mydb> db.Student.insert({RollNo:2,Age:22,Cont:t:9976,email:'anushka.de@gmail.com'});
{
  acknowledged: true,
  insertedIds: { '0': ObjectId('67c6c458812483cc27dd4a5f') }
}
Atlas atlas-ws5rct-shard-0 [primary] mydb> db.Student.insert({RollNo:3,Age:21,Cont:t:5576,email:'anubhav.de@gmail.com'});
{
  acknowledged: true,
  insertedIds: { '0': ObjectId('67c6c47f812483cc27dd4a60') }
}
Atlas atlas-ws5rct-shard-0 [primary] mydb> db.Student.insert({RollNo:4,Age:20,Cont:t:4476,email:'pani.de9@gmail.com'});
{
  acknowledged: true,
  insertedIds: { '0': ObjectId('67c6c4a2812483cc27dd4a61') }
}
Atlas atlas-ws5rct-shard-0 [primary] mydb> db.Student.insert({RollNo:10,Age:23,Cont:t:2276,email:'rekha.de9@gmail.com'});
{
  acknowledged: true,
  insertedIds: { '0': ObjectId('67c6c4db812483cc27dd4a62') }
}
Atlas atlas-ws5rct-shard-0 [primary] mydb> db.Student.insertMany([{_id:3,StudName:'Jacob Adam',Grade:'VII',Hobbies:'Swimming'}, {_id:4,StudName:'Amy Jacks',Grade:'X',Hobbies:'Dancing'}])
{ acknowledged: true, insertedIds: { '0': 3, '1': 4 } }
Atlas atlas-ws5rct-shard-0 [primary] mydb> db.Student.insert({_id:1,StudName:'Michelle Jacintha',Grade:'VII',Hobbies:'InternetSurfing'})
{ acknowledged: true, insertedIds: { '0': 1 } }
Atlas atlas-ws5rct-shard-0 [primary] mydb> db.Student.insertOne({_id:2,StudName:'Janine',Grade:'VIII',Hobbies:'Music'})
{ acknowledged: true, insertedId: 2 }

```

```

Atlas atlas-ws5rct-shard-0 [primary] mydb> db.Student.find().pretty()
[
  {
    _id: ObjectId('67c6c3c3812483cc27dd4a5d'),
    RollNo: 1,
    Age: 21,
    Cont: 9876,
    email: 'antara.de9@gmail.com'
  },
  {
    _id: ObjectId('67c6c3d8812483cc27dd4a5e'),
    RollNo: 1,
    Age: 21,
    Cont: 9876,
    email: 'antara.de9@gmail.com'
  },
  {
    _id: ObjectId('67c6c458812483cc27dd4a5f'),
    RollNo: 2,
    Age: 22,
    Cont: 9976,
    email: 'anushka.de@gmail.com'
  },
  {
    _id: ObjectId('67c6c47f812483cc27dd4a60'),
    RollNo: 3,
    Age: 21,
    Cont: 5576,
    email: 'anubhav.de@gmail.com'
  },
  {
    _id: ObjectId('67c6c4a2812483cc27dd4a61'),
    RollNo: 4,
    Age: 20,
    Cont: 4476,
    email: 'pani.de9@gmail.com'
  },
  {
    _id: ObjectId('67c6c4db812483cc27dd4a62'),
    RollNo: 10,
    Age: 23,
    Cont: 2276,
    email: 'rekha.de9@gmail.com'
  }
]
Atlas atlas-ws5rct-shard-0 [primary] mydb> db.Student.save({StudName:'Vamsi',Grade:'VI'})

```

```
Atlas atlas-ws5rct-shard-0 [primary] mydb> db.Student.updateOne({_id:6,StudName:'Aryan David',Grade:'VII'},{$set:{Hobbies:'Skating'}},{upsert:true})
{
  acknowledged: true,
  insertedId: 6,
  matchedCount: 0,
  modifiedCount: 0,
  upsertedCount: 1
}
```

```
Atlas atlas-ws5rct-shard-0 [primary] mydb> db.Student.insert({RollNo:11,Age:22,Name :"ABC",cont:2276,email:"rea.de9@gmail.com"})
{
  acknowledged: true,
  insertedIds: { '0': ObjectId('67c6c616812483cc27dd4a63') }
}
Atlas atlas-ws5rct-shard-0 [primary] mydb> db.Student.update({RollNo:11,Name:"ABC"},{$set:{Name:"FEM"}})
{
  acknowledged: true,
  insertedId: null,
  matchedCount: 1,
  modifiedCount: 1,
  upsertedCount: 0
}
Atlas atlas-ws5rct-shard-0 [primary] mydb> db.Student.updateMany({Grade:'VII'},{$set:{status:'Active'}})
{
  acknowledged: true,
  insertedId: null,
  matchedCount: 3,
  modifiedCount: 2,
  upsertedCount: 0
}
Atlas atlas-ws5rct-shard-0 [primary] mydb> db.Student.updateOne({Grade:'VII'},{$set:{status:'Active'}})
{
  acknowledged: true,
  insertedId: null,
  matchedCount: 1,
  modifiedCount: 1,
  upsertedCount: 0
}
```

```
Atlas atlas-ws5rct-shard-0 [primary] mydb> db.Student.update({_id:4}, {$set:{Location:'Network'}})
{
  acknowledged: true,
  insertedId: null,
  matchedCount: 1,
  modifiedCount: 1,
  upsertedCount: 0
}
Atlas atlas-ws5rct-shard-0 [primary] mydb> db.Student.update({_id:4}, {$unset:{Location:'Network'}})
{
  acknowledged: true,
  insertedId: null,
  matchedCount: 1,
  modifiedCount: 1,
  upsertedCount: 0
}
Atlas atlas-ws5rct-shard-0 [primary] mydb> db.Student.update({RollNo:10}, {$set:{email:'abhinav@gmail.com'}})
DeprecationWarning: Collection.update() is deprecated. Use updateOne, updateMany, or bulkWrite.
{
  acknowledged: true,
  insertedId: null,
  matchedCount: 1,
  modifiedCount: 1,
  upsertedCount: 0
}
```

## LAB-1

Import and Export Command :

```
mongoexport mongodbd +srv : // Mehul @ cluster0 Mehul @  
cluster0 mongodbs.net /  
-- collection = Student --out  
C:\Users\Student\Downloads\DBMS\output.JSON
```

```
mongoimport mongodbd +srv : // Mehul @ Mehul @  
cluster0 Mehul @ cluster0 mongodbs.net /  
-- collection = New_Student --type json  
-- file C:\Users\Student\Downloads\DBMS\  
output.json
```

DAB-2

\* Import in local host :

```
mongoimport --db myDB --collection Student-new --type csv --headerline --file Home/Desktop/Student.csv
```

Export:

```
mongoexport --host localhost --db myDB --collection Student --csv --out Home/Desktop/Student.csv --fields "StudentName", "Grade", "Notes"
```

> use myDB;

Find:

```
db.Student.find({StudentName: "Ayhan David"});
```

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

```
db.Student.find({StudentName: /e/}).pretty();
```

Aggregate:

```
db.Customer.aggregate([{$group: {_id: "$CustomerID", TotalAmount: {$sum: "$Amount"}, Count: {$sum: 1}}}, {"$sort": {"Count": -1}}]);
```

db.Customer.aggregate([{\$match: {Category: "J" || "S", \$group: {\_id: "Count", TotalAmount: {\$sum: "\$Amount"}, \$sum: "1"}, \$sort: {"Count": -1}}});

Update:

db.Student.update({\_id: 3, StudentName: "Ayhan David", Grade: "VII"}, {\$set: {Notes: "Skating", Amount: 33, upsert: true}});

## **Lab 2:** Cassandra

Question: Perform the following DB operations using Cassandra.

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

Code with Output:

```

cqlsh> CREATE KEYSPACE Employee WITH REPLICATION = { 'class' : 'SimpleStrategy', 'replication_factor' : 1 };
cqlsh> CREATE TABLE Employee.Employee_Info (
...     Emp_Id int,
...     Salary DECIMAL,
...     EMP_Name TEXT,
...     Designation TEXT,
...     Date_of_Joining DATE,
...     Dept_Name TEXT,
...     PRIMARY KEY (Emp_Id, Salary)
... ) WITH CLUSTERING ORDER BY (Salary ASC);
cqlsh> BEGIN BATCH
...     INSERT INTO Employee.Employee_Info (Emp_Id, Salary, EMP_Name, Designation, Date_of_Joining, Dept_Name) VALUES (121, 60000, 'John Doe', 'Developer', '2023-01-15', 'IT');
...     INSERT INTO Employee.Employee_Info (Emp_Id, Salary, EMP_Name, Designation, Date_of_Joining, Dept_Name) VALUES (122, 80000, 'Jane Smith', 'Manager', '2022-05-20', 'HR');
...     INSERT INTO Employee.Employee_Info (Emp_Id, Salary, EMP_Name, Designation, Date_of_Joining, Dept_Name) VALUES (123, 55000, 'Alice Johnson', 'Analyst', '2021-11-10', 'Finance');
...     APPLY BATCH;
cqlsh> UPDATE Employee.Employee_Info SET Emp_Name = 'Johnathan Doe', Dept_Name = 'Engineering' WHERE Emp_Id = 121 AND Salary = 60000;
cqlsh> SELECT * FROM Employee.Employee_Info WHERE Emp_Id = 121 ORDER BY Salary;
emp_id | salary | date_of_joining | dept_name | designation | emp_name
-----+-----+-----+-----+-----+
121 | 60000 | 2023-01-15 | Engineering | Developer | Johnathan Doe
(1 rows)
cqlsh> ALTER TABLE Employee.Employee_Info ADD Projects SET<TEXT>;
cqlsh> UPDATE Employee.Employee_Info SET Projects = {'Project A', 'Project B'} WHERE Emp_Id = 121 AND Salary = 60000;
cqlsh> INSERT INTO Employee.Employee_Info (Emp_Id, Salary, EMP_Name, Designation, Date_of_Joining, Dept_Name) VALUES (124, 30000, 'Temp Employee', 'Intern', '2023-10-01', 'Temp Dept') USING TTL 15;
cqlsh> SELECT * FROM Employee.Employee_Info;
emp_id | salary | date_of_joining | dept_name | designation | emp_name | projects
-----+-----+-----+-----+-----+-----+-----+
123 | 55000 | 2021-11-10 | Finance | Analyst | Alice Johnson | null
122 | 80000 | 2022-05-20 | HR | Manager | Jane Smith | null
121 | 60000 | 2023-01-15 | Engineering | Developer | Johnathan Doe | ['Project A', 'Project B']
(3 rows)
cqlsh> []

```

Lab-4

Cambridge (AP) → Availability → Partition Tolerance

→ Not care controller

> create keyspace Student with replication = { 'class': 'SimpleStrategy', 'replication\_factor': 2 },

> describe keyspace;

student1 system system\_info

> use student1;

> create table student\_info (rollno int primary key, timestamp timestamp, lastmarks double, studentname text);

> describe table;

→ student info

> describe table student\_info;

> begin batch

insert into student1(rollno, studentname, timestamp, lastmarks)

values (1, 'A', '2023-03-25', 80)

insert into student1(rollno, studentname, timestamp, lastmarks)

values (2, 'B', '2023-04-15', 90)

values (3, 'C', '2023-03-15', 79)

values (4, 'D', '2023-03-20', 82)

apply batch;

rollno	timestamp	lastmarks	studentname
1	2023-03-25 08:30:00	80	A
2	2023-04-15 08:30:00	90	B
4	2023-03-19 08:30:00	82	D
3	2023-03-14 08:30:00	79	C

Mark \* from student\_info where rollno in (1,2,5).

> select \* from student\_info where studentname = 'A';

Some old concern as studentname is not a primary key so makes no order on value part.

> update student\_info set hobbies = hobbies + 'hobbies' where rollno = 1;

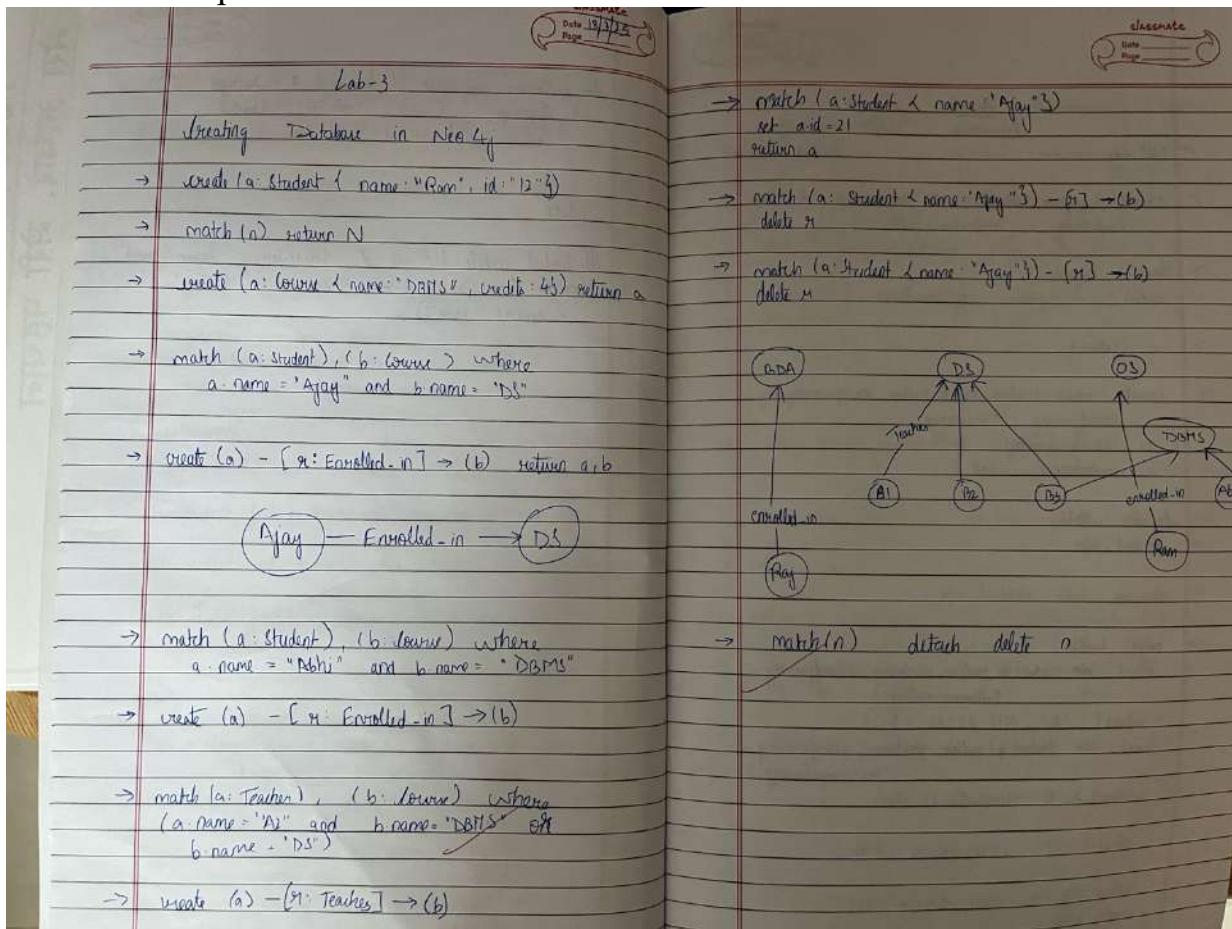
> select \* from student\_info;

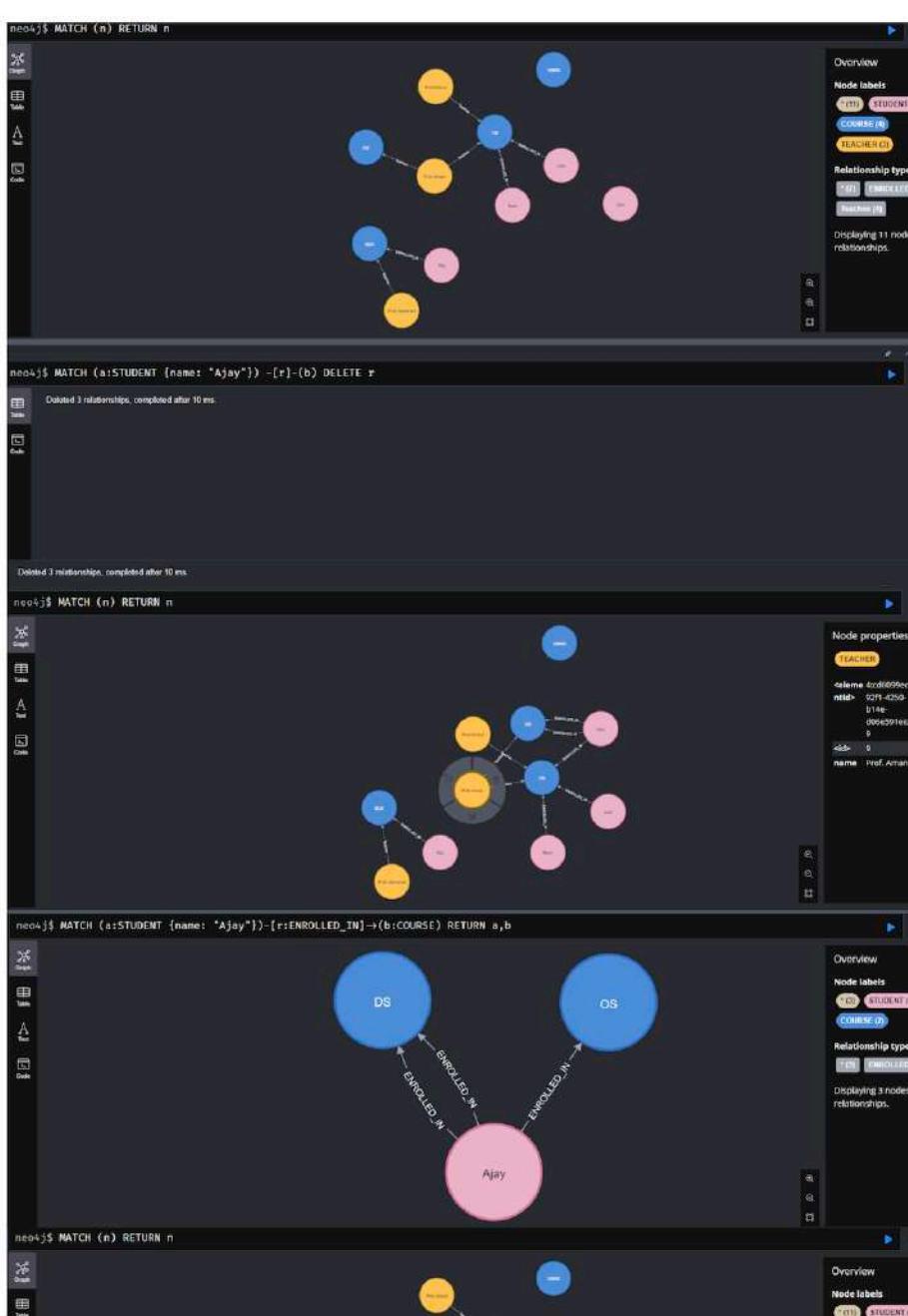
rollno	timestamp	hobbies	lastmarks	studentname
1	2023-03-25	hobbies	80	A
4	2023-03-19	hobbies	82	D
3	2023-03-14	hobbies	79	C

## Lab 3: Neo4J

Question: Create a graph database in Neo4J

Code with output:





## Lab 4: Cassandra

Question: Perform the following DB operations using Cassandra.

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

Code with Output:

```
cqlsh> CREATE KEYSPACE Library WITH REPLICATION = { 'class' : 'SimpleStrategy', 'replication_factor' : 1 };
cqlsh> CREATE TABLE Library.Library_Info (
    ...     Stud_Id int,
    ...     Book_Name TEXT,
    ...     Book_Id int,
    ...     Date_of_issue DATE,
    ...     PRIMARY KEY (Stud_Id, Book_Name, Date_of_issue)
    ... );
cqlsh> BEGIN BATCH
    ...     INSERT INTO Library.Library_Info (Stud_Id, Book_Name, Book_Id, Date_of_issue) VALUES (112, 'BDA', 1, '2023-09-01');
    ...     INSERT INTO Library.Library_Info (Stud_Id, Book_Name, Book_Id, Date_of_issue) VALUES (112, 'BDA', 1, '2023-09-05');
    ...     INSERT INTO Library.Library_Info (Stud_Id, Book_Name, Book_Id, Date_of_issue) VALUES (113, 'ML', 2, '2023-09-02');
    ...     INSERT INTO Library.Library_Info (Stud_Id, Book_Name, Book_Id, Date_of_issue) VALUES (114, 'AI', 3, '2023-09-03');
    ...     INSERT INTO Library.Library_Info (Stud_Id, Book_Name, Book_Id, Date_of_issue) VALUES (115, 'DBMS', 4, '2023-09-04');
    ...     APPLY BATCH;
cqlsh> SELECT * FROM Library.Library_Info;

stud_id | book_name | date_of_issue | book_id
-----+-----+-----+-----+
  114 |      AI | 2023-09-03 |      3
  113 |      ML | 2023-09-02 |      2
  112 |      BDA | 2023-09-01 |      1
  112 |      BDA | 2023-09-05 |      1
  115 |    DBMS | 2023-09-04 |      4

(5 rows)
cqlsh> SELECT COUNT(*) FROM Library.Library_Info WHERE Stud_Id = 112 AND Book_Name = 'BDA';

count
-----
  2

(1 rows)
```

```

cqlsh> COPY Library.Library_Info TO 'library_info.csv' WITH HEADER = TRUE;
Using 16 child processes

Starting copy of library.library_info with columns [stud_id, book_name, date_of_issue, book_id].
Processed: 5 rows; Rate: 96 rows/s; Avg. rate: 96 rows/s
5 rows exported to 1 files in 0.089 seconds.
cqlsh> COPY Library.Library_Info FROM 'library_info.csv' WITH HEADER = TRUE;
Using 16 child processes

Starting copy of library.library_info with columns [stud_id, book_name, date_of_issue, book_id].
Processed: 5 rows; Rate: 9 rows/s; Avg. rate: 13 rows/s
5 rows imported from 1 files in 0.375 seconds (0 skipped).
cqlsh> 

```

Date: 8/1/23  
Page: 1

classmate

Nb-5 Cassandra

Q

8. Keyspace employee and perform operations like inserting, updating, delete, TR.

> Create keyspace employee with replication =  
<del> l'strategy class = 'SimpleStrategy', replication\_factor = 1 </del>

> use employee

> Create table employee\_info (emp\_id int primary key, emp\_name text, designation text, dep\_date, salary int, dept text);

> begin batch

insert into employee\_info (emp\_id, emp\_name, designation, dep, salary, dept)  
values (1, 'Rahul', 'Manager', '2020-03-23', 10000,  
'CRO')

insert into employee\_info ( )  
values (2, 'Sugan', 'SDE-2', '2021-01-25', 15000, 'Tx')

insert into employee\_info ( )  
values (3, 'Sabil', 'TA', '2021-03-19', 30000, 'Tx')

insert into employee\_info ( )  
values (4, 'Zeshan', 'Marketing Head', '2021-10-15',  
20000, 'Sales')

apply batch;

> select \* from employee\_info;

emp_id	dept	designation	dep	emp_name	salary
1	CRO	Manager	2020-03-23	Rahul	10000
2	Tx	SDE-2	2021-01-25	Sugan	15000
4	Sales	Marketing Head	2021-10-15	Zeshan	20000
3	Tx	CFO	2021-03-19	Sabil	30000

Date: 8/1/23  
Page: 2

classmate

> alter table employee\_info add projects set <key>;

> update employee\_info set projects = projects +  
<del> 2 </del> to do list '3'  
where emp\_id = 1;

> insert into employee\_info (emp\_id, emp\_name, designation, dep, salary, dept)  
values (5, 'Monish', 'Officer', '2020-10-20',  
50000, 'Management') using TTL 15;

> Select TTL (designation) from employee\_info where emp\_id = 5;  
Output: TTL (designation)

Q

9. Keyspace library and use counter, insert, update.

> Create keyspace library with replication =  
<del> l'strategy class = 'SimpleStrategy', replication\_factor = 1 </del>

> Create table library\_info (stud\_id int primary key, counter\_value counter, stud\_name text, book\_name text, book\_id int, dep\_date);

> begin batch

insert into library\_info (stud\_id, counter\_value, stud\_name, book\_name, book\_id, dep)

values (.....)

insert into library\_info (.....)

values (.....)

apply batch;

4. update library\_info counter\_table  
set counter\_value = counter\_value + 1  
where stud\_id = 1;

5. Export:

copy library\_info (stud\_id, counter) to  
'home\Desktop'\name.txt';

6. Import:

copy library\_info (stud\_id, counter) from  
'home\Desktop' ;

7. select count(\*) from library\_info where  
stud\_id = 1 and bookname = 'BDA'  
allow filtering;

24

## Lab 5: HDFS

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

Code with Output:

```
hadoop@bmscsecse-HP-Elite-Tower-800-G9-Desktop-PC:~$ cd ./Desktop/
hadoop@bmscsecse-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 [bmscsecse-HP-Elite-Tower-800-G9-Desktop-PC]
Starting resourcemanager
Starting nodemanagers
hadoop@bmscsecse-HP-Elite-Tower-800-G9-Desktop-PC:~/Desktop$ hdfs dfs -mkdir /Lab05
```

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

```
hadoop@bmscsecse-HP-Elite-Tower-800-G9-Desktop-PC:~/Desktop$ touch test.txt
hadoop@bmscsecse-HP-Elite-Tower-800-G9-Desktop-PC:~/Desktop$ nano text.txt
hadoop@bmscsecse-HP-Elite-Tower-800-G9-Desktop-PC:~/Desktop$ hdfs dfs -put ./text.txt /Lab05/text.txt
hadoop@bmscsecse-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@bmscsecse-HP-Elite-Tower-800-G9-Desktop-PC:~/Desktop$ hdfs dfs -cat /Lab05/text.txt
Hello
How are you?
```

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

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

Lab - 6  
Hadoop

> login into Ubuntu , Hadoop User : password: bmsce

> Navigate to terminal and execute  
→ start-all.sh

namenodes [localhost]  
datanodes

secondary namenodes [bmscse - HP - like . . . ]  
resource manager  
node managers

→ jps

(1) Creating a directory

→ hdfs dfs -mkdir /mehul153

(2) Checking the hdfs directory structure

→ hdfs dfs -ls /mehul153

(3) Copying a local file into HDFS.

i) → hdfs dfs -put /home /hadoop /Documents /sample.txt  
/mehul153 /wc.txt

ii) hdfs dfs -copyFromLocal /home /"

The difference in (i) and (ii) is that (i) can ~~copy~~ copy from anywhere but (ii) can only copy from local machine.

P.T.O

④ To get files from HDFS to local machine  
→ hdfs dfs -get mehwil53/abc.txt  
/home/hduser/Documents/def.txt

⑤ → hdfs dfs -getmerge mehwil53/abc.txt /mehwil53  
/def.txt /home/hadoop/Documents/merged.txt

⑥ → hdfs dfs -cat /mehwil53/abc.txt

⑦ hdfs fs -mv /mehwil53/FFF

⑧ hadoop fs -cp /FFF /LLL

-cp → copies from one directory to another in same.

Q/A G

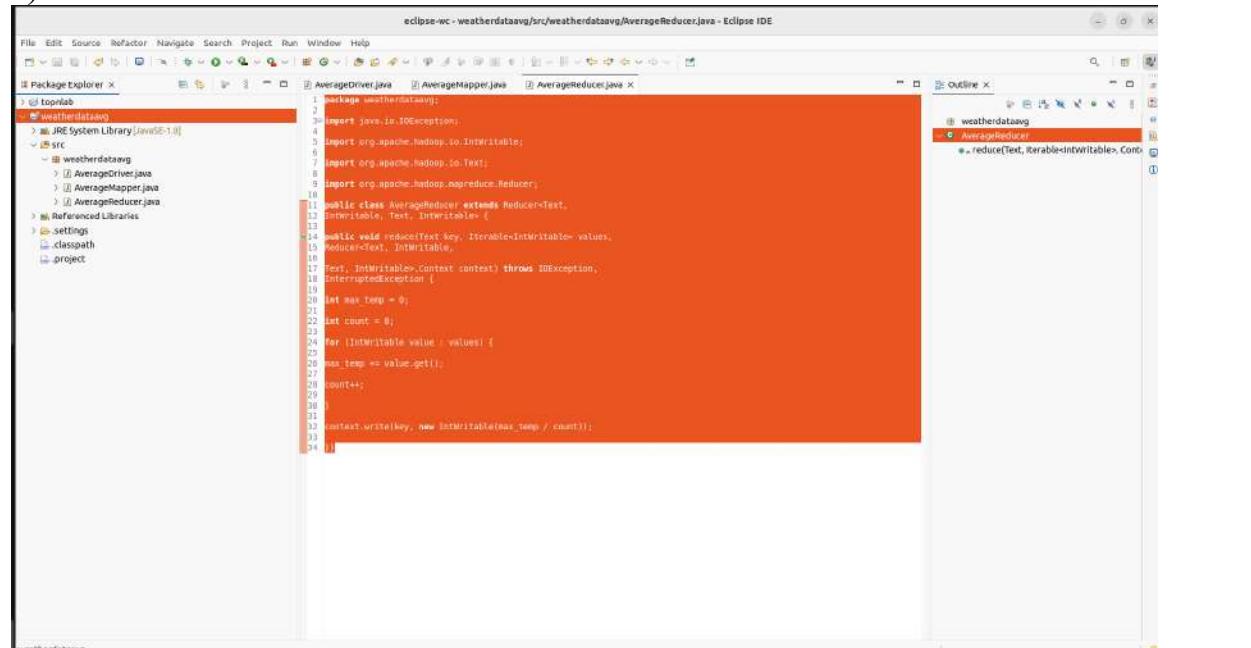
## Lab 6: Hadoop

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

- find average temperature for each year from NCDC data set.
- find the mean max temperature for every month

Code with Output:

a)



```
eclipse-wc - weatherdataavg/src/weatherdataavg/AverageReducer.java - Eclipse IDE
File Edit Source Refactor Navigate Search Project Run Window Help
Package explorer X
src
  AverageDriver.java
  AverageMapper.java
  AverageReducer.java
weatherdataavg
  AverageDriver.java
  AverageMapper.java
  AverageReducer.java
  settings
  classpath
  project
Outline X
weatherdataavg
  AverageReducer.java
    reduce(Text, Iterable<IntWritable>, Context)
File Edit Source Refactor Navigate Search Project Run Window Help
weatherdataavg
hadoop@bnseccecsa-HP-Elite-Tower-800-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.
Aborting.
Starting namenodes on [localhost]
localhost: namenode is running as process 9745. Stop it first and ensure /tmp/hadoop-hadoop-namenode.pid file is empty before retry.
Starting datanodes
localhost: datanode is running as process 9928. Stop it first and ensure /tmp/hadoop-hadoop-datanode.pid file is empty before retry.
Starting secondary namenodes [bnseccecsa-HP-Elite-Tower-800-G9-Desktop-PC]
bnseccecsa-HP-Elite-Tower-800-G9-Desktop-PC:~$ secondarynamenode is running as process 10221. Stop it first and ensure /tmp/hadoop-hadoop-secondarynamenode.pid file is empty before retry.
Starting resourcemanager
Please enter the input and output parameters
hadoop@bnseccecsa-HP-Elite-Tower-800-G9-Desktop-PC:~$ ^C
hadoop@bnseccecsa-HP-Elite-Tower-800-G9-Desktop-PC:~$ hadoop fs -rm -r /renout
rm: /renout: No such file or directory
hadoop@bnseccecsa-HP-Elite-Tower-800-G9-Desktop-PC:~$ hadoop fs -ls /
Found 9 items
drwxr-xr-x - hadoop supergroup 0 2025-04-15 15:11 /CSE
drwxr-xr-x - hadoop supergroup 0 2025-04-15 15:10 /FFF
drwxr-xr-x - hadoop supergroup 0 2025-04-15 15:09 /LLL
drwxr-xr-x - hadoop supergroup 0 2024-05-14 15:39 /abc
drwxr-xr-x - hadoop supergroup 0 2025-05-20 14:30 /klm
drwxr-xr-x - hadoop supergroup 0 2025-05-20 13:52 /mno
drwxr-xr-x - hadoop supergroup 0 2025-05-20 13:58 /res
drwxr-xr-x - hadoop supergroup 0 2025-04-29 15:38 /rgs
drwxr-xr-x - hadoop supergroup 0 2024-05-21 15:37 /wordcount
hadoop@bnseccecsa-HP-Elite-Tower-800-G9-Desktop-PC:~$ hadoop fs -ls /klm/
Found 1 item
-rw-r--r-- 1 hadoop supergroup 888198 2025-05-20 14:30 /klm/1901
hadoop@bnseccecsa-HP-Elite-Tower-800-G9-Desktop-PC:~$ hadoop fs -cat /klm/1901
```

```

hadoop@bmseccce-HP-ELITE-Tower-800-CR-Desktop-MC:~$ hadoop jar /home/hadoop/Desktop/weatherdataavg.jar weatherdataavg.AverageDriver /kln/1901 /rem
2025-05-20 14:34:21,074 INFO impl.MetricsConfig: Loaded properties from hadoop-metrics2.properties
2025-05-20 14:34:21,116 INFO Iml.MetricsSystemImpl: Jobtracker metrics system started
2025-05-20 14:34:21,176 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 14:34:21,238 INFO impl.FileInputFormat: Total input files to process : 1
2025-05-20 14:34:21,250 INFO mapreduce.JobResourceUploader: Number of split(s) = 1
2025-05-20 14:34:21,334 INFO mapreduce.JobSubmission: Submitting tokens for Job: job_local261815800_0001
2025-05-20 14:34:21,408 INFO mapreduce.Job: The url to track the job: http://localhost:8080/
2025-05-20 14:34:21,408 INFO mapreduce.Job: Running job: job_local261815800_0001
2025-05-20 14:34:21,408 INFO mapred.LocalJobRunner: OutputCommitter set in config null
2025-05-20 14:34:21,408 INFO mapred.LocalJobRunner: No output committer factory defined, defaulting to FileOutputCommitterFactory
2025-05-20 14:34:21,408 INFO mapred.FileOutputCommitter: FileOutputCommitter algorithm version is 2
2025-05-20 14:34:21,408 INFO mapred.FileOutputCommitter: FileOutputCommitter skip cleanup _temporary folders under output directory:false, ignore cleanup failures: false
2025-05-20 14:34:21,440 INFO mapred.LocalJobRunner: Waiting for map tasks
2025-05-20 14:34:21,447 INFO mapred.LocalJobRunner: Starting task: attempt_local261815800_0001_m_000000_0
2025-05-20 14:34:21,458 INFO output.PathOutputCommitterFactory: No output committer factory defined, defaulting to FileOutputCommitterFactory
2025-05-20 14:34:21,468 INFO output.FileOutputCommitter: FileOutputCommitter skip cleanup _temporary folders under output directory:false, ignore cleanup failures: false
2025-05-20 14:34:21,468 INFO output.FileOutputCommitter: FileOutputCommitter skip cleanup _temporary folders under output directory:false, ignore cleanup failures: false
2025-05-20 14:34:21,472 INFO mapred.MapTask: Processing split: hdfs://localhost:9000/kln/1901/6+888190
2025-05-20 14:34:21,511 INFO mapred.MapTask: (EQUATOR) 0 kv 26214390(104857504)
2025-05-20 14:34:21,511 INFO mapred.MapTask: mapreduce.task.io.sort.mb: 100
2025-05-20 14:34:21,511 INFO mapred.MapTask: sort limit at B3880600
2025-05-20 14:34:21,511 INFO mapred.MapTask: bufstart = 0; bufvoid = 104857600
2025-05-20 14:34:21,511 INFO mapred.MapTask: Kvstart = 20214390; length = 6553600
2025-05-20 14:34:21,588 INFO mapred.MapTask: Output collector class = org.apache.hadoop.mapred.MapTask$MapOutputBuffer
2025-05-20 14:34:21,588 INFO mapred.LocalJobRunner:
2025-05-20 14:34:21,592 INFO mapred.MapTask: Starting flush of map output
2025-05-20 14:34:21,597 INFO mapred.MapTask: Spilling map output
2025-05-20 14:34:21,597 INFO mapred.MapTask: bufstart = 0; bufend = 59076; bufvoid = 104857600
2025-05-20 14:34:21,597 INFO mapred.MapTask: kvstart = 26214390(104857504); kvend = 26188144(104752576); length = 26253/6553600
2025-05-20 14:34:21,597 INFO mapred.MapTask: Finished spill 0
2025-05-20 14:34:21,600 INFO mapred.Task: attempt_local261815800_0001_m_000000_0 is done. And is in the process of committing
2025-05-20 14:34:21,602 INFO mapred.LocalJobRunner: map
2025-05-20 14:34:21,602 INFO mapred.Task: Task 'attempt_local261815800_0001_m_000000_0' done.
2025-05-20 14:34:21,605 INFO mapred.Task: Final Counters for attempt_local261815800_0001_m_000000_0: Counters: 23
File System Counters
FILE: Number of bytes read=4430
FILE: Number of bytes written=713998
FILE: Number of read operations=0
FILE: Number of large read operations=0
FILE: Number of write operations=0
HDFS: Number of bytes read=888190
HDFS: Number of bytes written=0
HDFS: Number of read operations=5
HDFS: Number of large read operations=1
HDFS: Number of write operations=1
HDFS: Number of bytes read erasure-coded=0
Map-Reduce Framework
  Map input records=6565
  Map output records=564
  HDFS: Number of bytes written=8
  HDFS: Number of read operations=18
  HDFS: Number of large read operations=0
  HDFS: Number of write operations=3
  HDFS: Number of bytes read erasure-coded=0
Map-Reduce Framework
  Combine input records=0
  Combine output records=0
  Reduce input records=0
  Reduce output records=0
  Reduce shuffle bytes=72210
  Reduce input records=6564
  Reduce output records=1
  Spilled Records=6564
  Shuffled Maps=1
  Failed Shuffles=0
  Merged Map outputs=1
  GC time elapsed (ms)=0
  Total committed heap usage (bytes)=633339904
Shuffle Errors
  BAD_ID=0
  CONNECTION=0
  IO_ERROR=0
  WRONG_LENGTH=0
  WRONG_TYPE=0
  WRONG_MAP=0
  WRONG_REDUCE=0
File Output Format Counters
  Bytes Written=8
2025-05-20 14:34:21,803 INFO mapred.LocalJobRunner: Finishing task: attempt_local261815800_0001_r_000000_0
2025-05-20 14:34:21,803 INFO mapred.LocalJobRunner: reduce Task executor complete.
2025-05-20 14:34:22,007 INFO mapreduce.Job: Job: job_local261815800_0001 running in uber mode : false
2025-05-20 14:34:22,007 INFO mapreduce.Job:  map 0% reduce 100%
2025-05-20 14:34:22,007 INFO mapreduce.Job: Job job_local261815800_0001 completed successfully
2025-05-20 14:34:22,414 INFO mapreduce.Job: Counters: 36
File System Counters
FILE: Number of bytes read=15312
FILE: Number of bytes written=1500200
FILE: Number of read operations=0
FILE: Number of large read operations=0
FILE: Number of write operations=8
HDFS: Number of bytes read=1776380
HDFS: Number of bytes written=8
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=6565
  Map output records=564
  Map output bytes=59076
  Map output materialized bytes=72210
  Input split bytes=95
  Combine input records=0
  Combine output records=0
  Reduce input groups=1
  Reduce shuffle bytes=72210
  Reduce input records=6564

```



b)

The screenshot shows the Eclipse IDE interface with the following details:

- Project Explorer:** Shows the project structure for "weatherdatameanmax". It includes a "src" folder containing "MeanMaxDriver.java", "MeanMaxMapper.java", and "MeanMaxReducer.java".
- Code Editor:** Displays the content of "MeanMaxReducer.java". The code implements a Reducer that processes temperature data to find the maximum value.
- Outline View:** Shows the class hierarchy, including "MeanMaxReducer" which extends "Reducer<Text, IntWritable, Text, IntWritable>".
- Terminal Window:** Shows a terminal session on a Linux system (hadoop@bmseccse-HP-Elite-Tower-800-G9-Desktop-PC) running the command "start-all.sh". The output indicates the start of various Hadoop daemons.
- Bottom Help:** A detailed list of Hadoop command-line options and their descriptions, such as "archive", "checknative", "classpath", "conftest", "createtoken", "distcp", "distutil", "envvars", "fs", "gridmix", "jar", and "rmr".

```
hadoop@hmsccse-HP-Elite-Tower-800-C9-Desktop-PC: ~

jpath applications, not this command.
prints the Java.library.path
kdiag Diagnose Kerberos Problems
kerbname show auth_to_local principal conversion
key manage keys via the KeyProvider
runenfolde scale a runen input trace
runentrace convert Logs into a runen trace
s3cmd S3 Commands
trace view and modify Hadoop tracing settings
version print the version

  Daemon Commands:

kms run KMS, the Key Management Server
registrydns run the registry DNS server

SUBCOMMAND may print help when invoked w/ parameters or with -h,
hadoop@hmsccse-HP-Elite-Tower-800-C9-Desktop-PC:~$ hadoop fs -mkdir /local /home/hadoop/Desktop/1901 /onnn/1901
hadoop@hmsccse-HP-Elite-Tower-800-C9-Desktop-PC:~$ hadoop fs -copyFromLocal /home/hadoop/Desktop/1901 /onnn/1901
hadoop@hmsccse-HP-Elite-Tower-800-C9-Desktop-PC:~$ hadoop jar /home/hadoop/Desktop/weatherdatameanmax.jar weatherdatameanmax.MeanMaxDriver /onnn/1901 /ren
2025-05-20 14:53:15.576 INFO Impl.MetricsConfig: Loaded properties from hadoop-metrics2.properties
2025-05-20 14:53:15.515 INFO LogUtil.MetricsSystemImpl: Scheduled Metric snapshot period at 10 second(s).
2025-05-20 14:53:15.515 INFO LogUtil.MetricsSystemImpl: Metrics system has been started
2025-05-20 14:53:15.579 INFO Mapred.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to remedy this.
2025-05-20 14:53:15.737 INFO Input.FileInputFormat: Total input files to process : 1
2025-05-20 14:53:15.764 INFO mapreduce.JobSubmissionWriter: number of splits:1
2025-05-20 14:53:15.838 INFO mapreduce.JobSubmissionWriter: Submitting tokens for job: job_local12143084439_0001
2025-05-20 14:53:15.896 INFO mapreduce.Job: Job: url to track the job: http://localhost:8080/
2025-05-20 14:53:15.901 INFO mapreduce.Job: Job ID: job_local12143084439_0001
2025-05-20 14:53:15.905 INFO mapreduce.JobSubmitter: JobConf needs cleanupCommitter = true, config null
2025-05-20 14:53:15.905 INFO output.PathOutputCommitterFactory: No output Committer factory defined, defaulting to FileOutputCommitterFactory
2025-05-20 14:53:15.905 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 2
2025-05-20 14:53:15.905 INFO output.FileOutputCommitter: FileOutputCommitter skip cleanup _temporary folders under output directory:false, ignore cleanup failures: false
2025-05-20 14:53:15.905 INFO mapreduce.LocalJobRunner: OutputCommitter is org.apache.hadoop.mapreduce.lib.output.FileOutputCommitter
2025-05-20 14:53:15.981 INFO mapreduce.LocalJobRunner: Waiting for no tasks
2025-05-20 14:53:15.982 INFO mapreduce.LocalJobRunner: Starting Task: attempt_local12143084439_0001_m_000000_0
2025-05-20 14:53:16.000 INFO mapreduce.Task: Using ResourceCalculatorProcessTree : []
2025-05-20 14:53:16.000 INFO mapreduce.Task: Using FileOutputCommitter: FileOutputCommitter Algorithm version is 2
2025-05-20 14:53:16.004 INFO mapreduce.Task: Using ResourceCalculatorProcessTree : []
2025-05-20 14:53:16.007 INFO mapred.MapTask: Processing split: hdfs://localhost:9000/onn/1901:0+868190
2025-05-20 14:53:16.044 INFO mapred.MapTask: (EQUATOR) 0 kvl 26214396(184857594)
2025-05-20 14:53:16.044 INFO mapred.MapTask: mapred.task.io.sort.mb: 100
2025-05-20 14:53:16.044 INFO mapred.MapTask: soft limit at 83886080
2025-05-20 14:53:16.044 INFO mapred.MapTask: hard limit at 104857600
2025-05-20 14:53:16.044 INFO mapred.MapTask: kvstart = 26214396; length = 6553600
2025-05-20 14:53:16.046 INFO mapred.MapTask: Map output collector class = org.apache.hadoop.mapred.MapTask$MapOutputBuffer
2025-05-20 14:53:16.118 INFO mapred.LocalJobRunner
2025-05-20 14:53:16.119 INFO mapred.MapTask: Starting flush of map output
2025-05-20 14:53:16.119 INFO mapred.MapTask: Spilling map output
2025-05-20 14:53:16.119 INFO mapred.MapTask: bufstart = 0; buflen = 45948; bufvoid = 104857600
2025-05-20 14:53:16.119 INFO mapred.MapTask: kvstart = 26214396(104857584); kvend = 26188144(164752576); length = 26253/6553600
2025-05-20 14:53:16.126 INFO mapred.MapTask: Finished split B
2025-05-20 14:53:16.133 INFO mapred.Task: Task:attempt_local12143084439_0001_m_000000_0 is done. And is in the process of committing
2025-05-20 14:53:16.133 INFO mapred.Task: LocalJobRunner is interrupting

hadoop@hmsccse-HP-Elite-Tower-800-C9-Desktop-PC: ~

2025-05-20 14:53:16.145 INFO output.FileOutputCommitter: FileOutputCommitter skip cleanup _temporary folders under output directory:false, ignore cleanup failures: false
2025-05-20 14:53:16.145 INFO mapred.Task: Using ResourceCalculatorProcessTree : []
2025-05-20 14:53:16.146 INFO mapred.ReduceTask: Using ShuffleConsumerPlugin: org.apache.hadoop.mapreduce.task.reduce.ShuffleJobConsumer
2025-05-20 14:53:16.147 MARNI.MetricSystemImpl: JobTracker metrics system already initialized!
2025-05-20 14:53:16.155 INFO Reduce.MergeManagerImpl: MergerManager: memoryLimit=5029453312, maxSingleShuffleLimit=1457363328, mergeThreshold=3047439360, ioSortFactor=10, memToMemMergeOutputsThreshold=10
2025-05-20 14:53:16.156 INFO Reduce.EventFetcher: attempt_local12143084439_0001_r_000000_0 Thread started: EventFetcher for fetching Map Completion Events
2025-05-20 14:53:16.178 INFO Reduce.LocalFetcher: localfetcher#1 about to shuffle output of map attempt_local12143084439_0001_n_000000_0 decomp: 59082 len: 59082 to MEMORY
2025-05-20 14:53:16.178 INFO Reduce.LocalFetcher: localfetcher#1 about to shuffle output of map attempt_local12143084439_0001_n_000000_0 decomp: 59082 len: 59082 to MEMORY
2025-05-20 14:53:16.178 INFO Reduce.MergeManagerImpl: closing temporary file > map-output or size: 59087, inMemoryMapOutputs.size() -> 1, commitMemory -> 0, usedMemory -> 59087
2025-05-20 14:53:16.178 INFO reduce.EventFetcher: EventFetcher is interrupted.. Returning
2025-05-20 14:53:16.178 INFO mapred.LocalJobRunner: 1 / 1 copied.
2025-05-20 14:53:16.178 INFO reduce.MergeManagerImpl: finalMerge called with 1 in-memory map-outputs and 0 on-disk map-outputs
2025-05-20 14:53:16.176 INFO mapred.Merger: Merging 1 sorted segments
2025-05-20 14:53:16.176 INFO mapred.Merger: Down to the last merge-pass, with 1 segments left of total size: 59073 bytes
2025-05-20 14:53:16.186 INFO Reduce.MergeManagerImpl: Merger: 1 segments, 59078 bytes to disk to satisfy reduce memory lmt
2025-05-20 14:53:16.186 INFO Reduce.MergeManagerImpl: Reducer: 1 segments, 59082 bytes from disk
2025-05-20 14:53:16.191 INFO reduce.MergeManagerImpl: Merging 0 segments, 0 bytes from memory into reduce
2025-05-20 14:53:16.191 INFO mapred.Merger: Merging 1 sorted segments
2025-05-20 14:53:16.192 INFO mapred.Merger: Down to the last merge-pass, with 1 segments left of total size: 59073 bytes
2025-05-20 14:53:16.192 INFO mapred.LocalJobRunner: 1 / 1 copied.
2025-05-20 14:53:16.212 INFO Configuration.deprecation: mapred.skip.on is deprecated. Instead, use mapreduce.job.skippedrecords
2025-05-20 14:53:16.273 INFO mapred.Task: Task:attempt_local12143084439_0001_r_000000_0 is done. And is in the process of committing
2025-05-20 14:53:16.275 INFO mapred.LocalJobRunner: 1 / 1 copied.
2025-05-20 14:53:16.275 INFO mapred.Task: Task attempt_local12143084439_0001_r_000000_0 is allowed to commit now
2025-05-20 14:53:16.289 INFO output.FileOutputCommitter: Saved output of task 'attempt_local12143084439_0001_r_000000_0' to hdfs://localhost:9000/ren
2025-05-20 14:53:16.296 INFO mapred.LocalJobRunner: Reduce > reduce
2025-05-20 14:53:16.296 INFO mapred.Task: Task 'attempt_local12143084439_0001_r_000000_0' done.
2025-05-20 14:53:16.296 INFO mapred.Task: Final Counters for attempt_local12143084439_0001_r_000000_0: Counters: 30
File System Counter:
FILE: Number of bytes read=122769
FILE: Number of bytes written=59073
FILE: Number of read operations=0
FILE: Number of large read operations=0
FILE: Number of write operations=0
HDFS: Number of bytes read=888190
HDFS: Number of bytes written=81
HDFS: Number of read operations=18
HDFS: Number of large read operations=0
HDFS: Number of write operations=1
HDFS: Number of bytes read erasure-coded=0
Number of bytes read erasure-coded=0

Map-Reduce Framework
Combine input records=0
Combine output records=0
Reduce input groups=12
Reduce shuffle bytes=59082
Reduce input records=6564
Reduce output records=12
Spilled Records=4564
Shuffled Maps =1
Failed Shuffles=0
Merged Map outputs=1
GC time elapsed (ms)=0
Total committed heap usage (bytes)=526385152
Shuffle sort bytes=0
Bad ID=0
CONNECTION=0
```

Screenshot captured  
You can paste the image from the clipboard.

```

2025-05-20 14:53:15,982 INFO mapred.LocalJobRunner: Starting task attempt_local2143084439_0001_m_000000_0
2025-05-20 14:53:15,996 INFO output.PathOutputCommitterFactory: No output committer factory defined, defaulting to FileOutputCommitterFactory
2025-05-20 14:53:15,996 INFO output.FileOutputCommitter: FileOutputCommitter Algorithm version is 2
2025-05-20 14:53:15,996 INFO output.FileOutputCommitter: skip cleanup _temporary folders under output directory:false, ignore cleanup failures: false
2025-05-20 14:53:15,996 INFO output.FileOutputCommitter: max number of tasks using same output committer:1
2025-05-20 14:53:16,007 INFO mapred.MapTask: processing split(s) : []
2025-05-20 14:53:16,044 INFO mapred.MapTask: MapTask@ kvt=26214396(i104857584)
2025-05-20 14:53:16,044 INFO mapred.MapTask: mapreduce.task.io.sort.mb: 100
2025-05-20 14:53:16,044 INFO mapred.MapTask: soft limit at B3886080
2025-05-20 14:53:16,044 INFO mapred.MapTask: bufstart = 0; bufvoid = 104857600
2025-05-20 14:53:16,044 INFO mapred.MapTask: kvstart = 26214396; length = 6553600
2025-05-20 14:53:16,044 INFO mapred.MapTask: map output collector class = org.apache.hadoop.mapred.MapTask$MapOutputBuffer
2025-05-20 14:53:16,119 INFO mapred.LocalJobRunner: Map Task attempt_local2143084439_0001_m_000000_0 started
2025-05-20 14:53:16,119 INFO mapred.MapTask: Starting flush of map output
2025-05-20 14:53:16,119 INFO mapred.MapTask: Spilling map output
2025-05-20 14:53:16,119 INFO mapred.MapTask: bufstart = 0; bufend = 45948; bufvoid = 104857600
2025-05-20 14:53:16,128 INFO mapred.MapTask: kvstart = 26214396(i104857584); kvend = 26188144(104752576); length = 26253/6553600
2025-05-20 14:53:16,128 INFO mapred.MapTask: Finished spilt #0
2025-05-20 14:53:16,139 INFO mapred.Task: attempt_local2143084439_0001_m_000000_0 is done. And is in the process of committing
2025-05-20 14:53:16,139 INFO mapred.Task: Task 'attempt_local2143084439_0001_m_000000_0' done.
2025-05-20 14:53:16,139 INFO mapred.Task: Final Counters For attempt_local2143084439_0001_m_000000_0: Counters: 23
  File System Counter
    FILE: Number of bytes read=4573
    FILE: Number of bytes written=704111
    FILE: Number of read operations=0
    FILE: Number of large read operations=0
    FILE: Number of write operations=0
    HDFS: Number of bytes read=888190
    HDFS: Number of bytes written=0
    HDFS: Number of read operations=5
    HDFS: Number of large read operations=0
    HDFS: Number of write operations=1
    HDFS: Number of bytes read erasure-coded=0
  Map-Reduce Framework
    Map Input records=6565
    Map Output records=6564
    Map output bytes=45948
    Map output materialized bytes=59082
    Input split bytes=95
    Combine input records=0
    Spilled Records=6564
    Failed Shuffles=0
    Merged Map outputs=0
    GC time elapsed (ms)=0
    Total committed heap usage (bytes)=526385152
  File Input Format Counters
    Bytes Read=888190
2025-05-20 14:53:16,139 INFO mapred.LocalJobRunner: Finishing task: attempt_local2143084439_0001_m_000000_0
2025-05-20 14:53:16,139 INFO mapred.LocalJobRunner: map task executor complete.
2025-05-20 14:53:16,140 INFO mapred.LocalJobRunner: Waiting for reduce tasks
2025-05-20 14:53:16,140 INFO mapred.LocalJobRunner: Starting task: attempt_local2143084439_0001_r_000000_0
2025-05-20 14:53:16,145 INFO output.PathOutputCommitterFactory: No output committer factory defined, defaulting to FileOutputCommitterFactory
2025-05-20 14:53:16,145 INFO output.FileOutputCommitter: FileOutputCommitter Algorithm version is 2
2025-05-20 14:53:16,145 INFO output.FileOutputCommitter: skip cleanup _temporary folders under output directory:false, ignore cleanup failures: false
2025-05-20 14:53:16,145 INFO output.FileOutputCommitter: max number of tasks using same output committer:1
  File System Counter
    FILE: Number of bytes written=703193
    FILE: Number of read operations=0
    FILE: Number of large read operations=0
    FILE: Number of write operations=0
    HDFS: Number of bytes read=888190
    HDFS: Number of bytes written=81
    HDFS: Number of read operations=10
    HDFS: Number of large read operations=0
    HDFS: Number of write operations=1
    HDFS: Number of bytes read erasure-coded=0
  Map-Reduce Framework
    Combine input records=0
    Combine output records=0
    Reduce input groups=12
    Reduce shuffle bytes=59082
    Reducer input records=6564
    Reducer output records=12
    Spilled Records=6564
    Shuffled Maps =1
    Failed Shuffles=0
    Merged Map outputs=1
    GC time elapsed (ms)=0
    Total committed heap usage (bytes)=526385152
  Shuffle Errors
    IO_ERROR=0
    CONNECTION=0
    ID_ERROR=0
    WRONG_LENGTH=0
    WRONG_MAP=0
    WRONG_REDUCE=0
  File Output Format Counters
    Bytes Written=81
2025-05-20 14:53:16,290 INFO mapred.LocalJobRunner: Finishing task: attempt_local2143084439_0001_r_000000_0
2025-05-20 14:53:16,291 INFO mapred.LocalJobRunner: reduce task executor complete.
2025-05-20 14:53:16,895 INFO mapreduce.Job: Job job_local2143084439_0001 running in uber mode : false
2025-05-20 14:53:16,897 INFO mapreduce.Job: map 100% reduce 100%
2025-05-20 14:53:16,899 INFO mapreduce.Job: Job job_local2143084439_0001 completed successfully
2025-05-20 14:53:16,899 INFO mapreduce.Job: Counters: 36
  File System Counter
    FILE: Number of bytes read=127342
    FILE: Number of bytes written=1467304
    FILE: Number of read operations=0
    FILE: Number of large read operations=0
    FILE: Number of write operations=0
    HDFS: Number of bytes read=1776388
    HDFS: Number of bytes written=91
    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=6565
    Map Output records=6564
    Map output bytes=45948
    Map output materialized bytes=59082
    Input null+ header=0

```

Screenshot captured  
You can paste the image from the clipboard.

```

hadoop@bmsccse-HP-Elite-Tower-800-G9-Desktop-PC: ~
FILE: Number of bytes written=1467304
FILE: Number of read operations=0
FILE: Number of large read operations=0
FILE: Number of write operations=0
HDFS: Number of bytes read=177888
HDFS: Number of bytes written=81
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=6055
Map output records=5564
Map output bytes=45948
Map output materialized bytes=59002
Input split bytes=95
Combine input records=0
Combine output records=0
Combine output bytes=0
Reduce shuffle groups=3
Reduce shuffle bytes=59802
Reduce input records=5564
Reduce output records=12
Spilled Records=13128
Shuffled Maps =1
Failed Shuffles=0
Merged Map outputs=1
GC time elapsed (ms)=8
Total committed heap usage (bytes)=1052770304
Shuffle Errors
  BAD_ID=0
  CONNECTION=0
  IO_ERROR=0
  WRONG_LENGTH=0
  WRONG_MAGIC=0
  WRONG_NUMBER=0
File Input Format Counters
  Bytes Read=888190
File Output Format Counters
  Bytes Written=81
hadoop@bmsccse-HP-Elite-Tower-800-G9-Desktop-PC: ~$ hadoop fs -cat /omn/part-r-00000
cat: /omn/part-r-00000: No such file or directory
hadoop@bmsccse-HP-Elite-Tower-800-G9-Desktop-PC: ~$ hadoop fs -cat /ren/part-r-00000
01      13
02      -66
03      -15
04      43
05      108
06      168
07      210
08      198
09      141
10      100
11      1
12      -81
hadoop@bmsccse-HP-Elite-Tower-800-G9-Desktop-PC: ~$ ■

```

## Average Temperature

- ① Import average temperature files
- ② Copy shov and map reduce to external jars.
- ③ Change java version to Java - SE - 1.8.
- ④ Change mapped-site.xml.

### Commands

```
> start-all.sh  
> jps  
> hadoop fs -mkdir weather  
> hadoop fs -copyFromLocal /home/hadoop/tool.txt  
/weather/test.txt  
  
> hadoop jar /home/hadoop/AverageTemperature.jar  
AverageDriver // Name of Driver  
/weather/test.txt // Iteration in HDFS  
/weather/output // output location (new)  
  
> hadoop fs -cat /weather/output/part-r-00000  
  
Output :-  
190 | 46 ✓  
  
(OR)  
> hdfs dfs -cat /weather/output/*
```

P.T.O

## Lab 7: Hadoop

Question: Implement Wordcount program on Hadoop framework

Code with Output:

The screenshot shows the Eclipse IDE interface with the following details:

- Project Explorer:** Shows the project structure with packages `word_count` and `wordcount`. Inside `wordcount`, there are three files: `WCDriver.java`, `WCMapper.java`, and `WCReducer.java`.
- Code Editor:** Displays the `WCDriver.java` file content:1 import java.io.IOException;
2 import org.apache.hadoop.conf.Configuration;
3 import org.apache.hadoop.fs.Path;
4 import org.apache.hadoop.io.IntWritable;
5 import org.apache.hadoop.io.Text;
6 import org.apache.hadoop.mapred.FileInputFormat;
7 import org.apache.hadoop.mapred.FileOutputFormat;
8 import org.apache.hadoop.mapred.JobClient;
9 import org.apache.hadoop.mapred.JobConf;
10 import org.apache.hadoop.util.Tool;
11 import org.apache.hadoop.util.ToolRunner;
12
13 public class WCDriver extends Configuration implements Tool {
14
15 public int run(String[] args) throws IOException {
16 if (args.length < 2) {
17 System.out.println("Please provide input and output paths");
18 return -1;
19 }
20
21 JobConf conf = new JobConf(WCDriver.class);
22 conf.setJobName("WordCount");
23 conf.setJarByClass(WCDriver.class); // Ensures job runs from correct JAR
24
25 FileInputFormat.setInputPaths(conf, new Path(args[0]));
26 FileOutputFormat.setOutputPath(conf, new Path(args[1]));
27
28 conf.setMapperClass(WCMapper.class);
29 conf.setReducerClass(WCReducer.class);
30
31 conf.setMapOutputKeyClass(Text.class);
32 }
- Console:** Shows the terminal output of the Hadoop command execution:hadoop@bmseccse-HP-Elite-Tower-800-G9-Desktop-PC:~\$ hadoop fs -mkdir /rsg
hadoop@bmseccse-HP-Elite-Tower-800-G9-Desktop-PC:~\$ hadoop fs -copyFromLocal /home/hadoop/Desktop/sample.txt /rsg/sample.txt
hadoop@bmseccse-HP-Elite-Tower-800-G9-Desktop-PC:~\$ hadoop jar /home/hadoop/Desktop/wordcount.jar WCDriver /rsg/sample.txt /result
2025-05-06 15:05:01,261 INFO impl.MetricsConfig: Loaded properties from hadoop-metrics2.properties
2025-05-06 15:05:01,299 INFO impl.MetricsSystemImpl: Scheduled Metric snapshot period at 10 second(s).
2025-05-06 15:05:01,299 INFO impl.MetricsSystemImpl: JobTracker metrics system started
2025-05-06 15:05:01,305 WARN impl.MetricsSystemImpl: JobTracker metrics system already initialized!
2025-05-06 15:05:01,365 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to remedy this.
2025-05-06 15:05:01,414 INFO mapred.FileInputFormat: Total input files to process : 1
2025-05-06 15:05:01,445 INFO mapreduce.JobSubmitter: number of splits:1
2025-05-06 15:05:01,511 INFO mapreduce.JobSubmitter: Submitting tokens for job: job\_local90897529\_0001
2025-05-06 15:05:01,511 INFO mapreduce.JobSubmitter: Executing with tokens: []
2025-05-06 15:05:01,565 INFO mapreduce.Job: The url to track the job: http://localhost:8080/
2025-05-06 15:05:01,566 INFO mapred.LocalJobRunner: OutputCommitter set in config null
2025-05-06 15:05:01,567 INFO mapred.LocalJobRunner: OutputCommitter is org.apache.hadoop.mapred.FileOutputCommitter
2025-05-06 15:05:01,569 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 2

```

hadoop@bmscse-HP-Elite-Tower-800-G9-Desktop-PC: ~
2025-05-06 15:05:01,299 INFO impl.MetricsSystemImpl: Scheduled Metric snapshot period at 10 second(s).
2025-05-06 15:05:01,299 INFO impl.MetricsSystemImpl: JobTracker metrics system started
2025-05-06 15:05:01,305 WARN impl.MetricsSystemImpl: JobTracker metrics system already initialized!
2025-05-06 15:05:01,365 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to remedy this.
2025-05-06 15:05:01,414 INFO mapred.FileInputFormat: Total input files to process : 1
2025-05-06 15:05:01,445 INFO mapreduce.JobSubmitter: number of splits:1
2025-05-06 15:05:01,511 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_local90897529_0001
2025-05-06 15:05:01,511 INFO mapreduce.JobSubmitter: Executing with tokens: []
2025-05-06 15:05:01,565 INFO mapreduce.Job: The url to track the job: http://localhost:8080/
2025-05-06 15:05:01,566 INFO mapred.LocalJobRunner: Running job: job_local90897529_0001
2025-05-06 15:05:01,566 INFO mapred.LocalJobRunner: OutputCommitter set in config null
2025-05-06 15:05:01,567 INFO mapred.LocalJobRunner: OutputCommitter is org.apache.hadoop.mapred.FileOutputCommitter
2025-05-06 15:05:01,569 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 2
2025-05-06 15:05:01,569 INFO output.FileOutputCommitter: FileOutputCommitter skip cleanup _temporary folders under output directory:false , ignore cleanup failures: false
2025-05-06 15:05:01,606 INFO mapred.LocalJobRunner: Waiting for map tasks
2025-05-06 15:05:01,607 INFO mapred.LocalJobRunner: Starting task: attempt_local90897529_0001_m_000000_0
2025-05-06 15:05:01,618 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 2
2025-05-06 15:05:01,618 INFO output.FileOutputCommitter: FileOutputCommitter skip cleanup _temporary folders under output directory:false , ignore cleanup failures: false
2025-05-06 15:05:01,624 INFO mapred.Task: Using ResourceCalculatorProcessTree : []
2025-05-06 15:05:01,631 INFO mapred.Task: Processing split: hdfs://localhost:9000/rsg/sample.txt:0+89
2025-05-06 15:05:01,640 INFO mapred.MapTask: numReduceTasks: 1
2025-05-06 15:05:01,671 INFO mapred.MapTask: (EQUATOR) 0 kv1=26214396(104857584)
2025-05-06 15:05:01,671 INFO mapred.MapTask: mapreduce.task.io.sort.mb: 100
2025-05-06 15:05:01,671 INFO mapred.MapTask: soft limit at 83886000
2025-05-06 15:05:01,671 INFO mapred.MapTask: bufstart = 0; bufvoid = 104857600
2025-05-06 15:05:01,671 INFO mapred.MapTask: kvstart = 26214396; length = 6553600
2025-05-06 15:05:01,673 INFO mapred.MapTask: Map output collector class = org.apache.hadoop.mapred.MapTask$MapOutputBuffer
2025-05-06 15:05:01,742 INFO mapred.LocalJobRunner:
2025-05-06 15:05:01,742 INFO mapred.MapTask: Starting flush of map output
2025-05-06 15:05:01,742 INFO mapred.MapTask: Spilling map output
2025-05-06 15:05:01,742 INFO mapred.MapTask: bufstart = 0; bufend = 169; bufvoid = 104857600
2025-05-06 15:05:01,742 INFO mapred.MapTask: kvstart = 26214396(104857584); kvend = 26214320(104857280); length = 77/6553600
2025-05-06 15:05:01,745 INFO mapred.MapTask: Finished spill 0
2025-05-06 15:05:01,751 INFO mapred.Task: Task:attempt_local90897529_0001_m_000000_0 is done. And is in the process of committing
2025-05-06 15:05:01,753 INFO mapred.LocalJobRunner: hdfs://localhost:9000/rsg/sample.txt:0+89
2025-05-06 15:05:01,753 INFO mapred.Task: Task 'attempt_local90897529_0001_m_000000_0' done.
2025-05-06 15:05:01,753 INFO mapred.Task: Final Counters for attempt_local90897529_0001_m_000000_0: Counters: 23
File System Counters
FILE: Number of bytes read=4273
FILE: Number of bytes written=639534
FILE: Number of read operations=0
hadoop@bmscse-HP-Elite-Tower-800-G9-Desktop-PC: ~
BAD_ID=0
CONNECTION=0
IO_ERROR=0
WRONG_LENGTH=0
WRONG_MAP=0
WRONG_REDUCE=0
File Output Format Counters
Bytes Written=69
2025-05-06 15:05:01,897 INFO mapred.LocalJobRunner: Finishing task: attempt_local90897529_0001_r_000000_0
2025-05-06 15:05:01,897 INFO mapred.LocalJobRunner: reduce task executor complete.
2025-05-06 15:05:02,569 INFO mapreduce.Job: Job job_local90897529_0001 running in uber mode : false
2025-05-06 15:05:02,572 INFO mapreduce.Job: map 100% reduce 100%
2025-05-06 15:05:02,574 INFO mapreduce.Job: Job job_local90897529_0001 completed successfully
2025-05-06 15:05:02,584 INFO mapreduce.Job: Counters: 36
File System Counters
FILE: Number of bytes read=9008
FILE: Number of bytes written=1279283
FILE: Number of read operations=0
FILE: Number of large read operations=0
FILE: Number of write operations=0
HDFS: Number of bytes read=178
HDFS: Number of bytes written=69
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=5
Map output records=20
Map output bytes=169
Map output materialized bytes=215
Input split bytes=88
Combine input records=0
Combine output records=0
Reduce input groups=10
Reduce shuffle bytes=215
Reduce input records=20
Reduce output records=10
Spilled Records=40
Shuffled Maps =1
Failed Shuffles=0
Merged Map outputs=1
GC time elapsed (ms)=0
Total committed heap usage (bytes)=1052770304

```

```

hadoop@bmscsece-HP-Elite-Tower-800-G9-Desktop-PC: ~
FILE: Number of write operations=0
HDFS: Number of bytes read=89
HDFS: Number of bytes written=0
HDFS: Number of read operations=5
HDFS: Number of large read operations=0
HDFS: Number of write operations=1
HDFS: Number of bytes read erasure-coded=0
Map-Reduce Framework
  Map input records=5
  Map output records=20
  Map output bytes=169
  Map output materialized bytes=215
  Input split bytes=88
  Combine input records=0
  Spilled Records=20
  Failed Shuffles=0
  Merged Map outputs=0
  GC time elapsed (ms)=0
  Total committed heap usage (bytes)=526385152
File Input Format Counters
  Bytes Read=89
2025-05-06 15:05:01,756 INFO mapred.LocalJobRunner: Finishing task: attempt_local90897529_0001_m_000000_0
2025-05-06 15:05:01,757 INFO mapred.LocalJobRunner: map task executor complete.
2025-05-06 15:05:01,758 INFO mapred.LocalJobRunner: Waiting for reduce tasks
2025-05-06 15:05:01,759 INFO mapred.LocalJobRunner: Starting task: attempt_local90897529_0001_r_000000_0
2025-05-06 15:05:01,762 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 2
2025-05-06 15:05:01,762 INFO output.FileOutputCommitter: FileOutputCommitter skip cleanup _temporary folders under output directory:false , ignore cleanup failures: false
2025-05-06 15:05:01,762 INFO mapred.Task: Using ResourceCalculatorProcessTree : []
2025-05-06 15:05:01,763 INFO mapred.ReduceTask: Using ShuffleConsumerPlugin: org.apache.hadoop.mapreduce.task.reduce.shuffle@636a90e9
2025-05-06 15:05:01,764 WARN impl.MetricsSystemImpl: JobTracker metrics system already initialized!
2025-05-06 15:05:01,771 INFO reduce.MergeManagerImpl: MergerManager: memoryLimit=5827985408, maxSingleShuffleLimit=1456996352, mergeThreshold=3846470400, loSortFactor=10, memToMemMergeOutputsThreshold=10
2025-05-06 15:05:01,772 INFO reduce.EventFetcher: attempt_local90897529_0001_r_000000_0 Thread started: EventFetcher for fetching Map Completion Events
2025-05-06 15:05:01,785 INFO reduce.LocalFetcher: localfetcher#1 about to shuffle output of map attempt_local90897529_0001_m_000000_0 decmp: 211 len: 215 to MEMORY
2025-05-06 15:05:01,787 INFO reduce.InMemoryMapOutput: Read 211 bytes from map-output for attempt_local90897529_0001_m_000000_0
2025-05-06 15:05:01,788 INFO reduce.MergeManagerImpl: closeInMemoryFile -> map-output of size: 211, inMemoryMapOutputs.size() -> 1, commitMemory -> 0, usedMemory -> 211
2025-05-06 15:05:01,788 INFO reduce.EventFetcher: EventFetcher is interrupted.. Returning
2025-05-06 15:05:01,789 INFO mapred.LocalJobRunner: 1 / 1 copied.
2025-05-06 15:05:01,789 INFO reduce.MergeManagerImpl: finalMerge called with 1 in-memory map-outputs and 0 on-disk map-outputs
2025-05-06 15:05:01,792 INFO mapred.Merger: Merging 1 sorted segments
hadoop@bmscsece-HP-Elite-Tower-800-G9-Desktop-PC: ~
2025-05-06 15:05:01,792 INFO reduce.MergeManagerImpl: Merged 1 segments, 211 bytes to disk to satisfy reduce memory limit
2025-05-06 15:05:01,793 INFO reduce.MergeManagerImpl: Merging 1 files, 215 bytes from disk
2025-05-06 15:05:01,793 INFO reduce.MergeManagerImpl: Merging 0 segments, 0 bytes from memory into reduce
2025-05-06 15:05:01,793 INFO mapred.Merger: Merging 1 sorted segments
2025-05-06 15:05:01,793 INFO mapred.Merger: Down to the last merge-pass, with 1 segments left of total size: 205 bytes
2025-05-06 15:05:01,793 INFO mapred.LocalJobRunner: 1 / 1 copied.
2025-05-06 15:05:01,867 INFO mapred.Task: Task:attempt_local90897529_0001_r_000000_0 is done. And is in the process of committing
2025-05-06 15:05:01,869 INFO mapred.LocalJobRunner: 1 / 1 copied.
2025-05-06 15:05:01,869 INFO mapred.Task: Task attempt_local90897529_0001_r_000000_0 is allowed to commit now
2025-05-06 15:05:01,894 INFO output.FileOutputCommitter: Saved output of task 'attempt_local90897529_0001_r_000000_0' to hdfs://localhost:9000/result
2025-05-06 15:05:01,896 INFO mapred.LocalJobRunner: reduce > reduce
2025-05-06 15:05:01,896 INFO mapred.Task: Task 'attempt_local90897529_0001_r_000000_0' done.
2025-05-06 15:05:01,897 INFO mapred.Task: Final Counters for attempt_local90897529_0001_r_000000_0: Counters: 30
  File System Counters
    FILE: Number of bytes read=4735
    FILE: Number of bytes written=639749
    FILE: Number of read operations=0
    FILE: Number of large read operations=0
    FILE: Number of write operations=0
    HDFS: Number of bytes read=89
    HDFS: Number of bytes written=69
    HDFS: Number of read operations=10
    HDFS: Number of large read operations=0
    HDFS: Number of write operations=3
    HDFS: Number of bytes read erasure-coded=0
  Map-Reduce Framework
    Combine input records=0
    Combine output records=0
    Reduce input groups=10
    Reduce shuffle bytes=215
    Reduce input records=20
    Reduce output records=10
    Spilled Records=20
    Shuffled Maps =1
    Failed Shuffles=0
    Merged Map outputs=1
    GC time elapsed (ms)=0
    Total committed heap usage (bytes)=526385152
  Shuffle Errors
    BAD_ID=0
    CONNECTION=0
    IO_ERROR=0
    WRONG_LENGTH=0

```

```
hadoop@bmscsecse-HP-Elite-Tower-800-G9-Desktop-PC: ~
HDFS: Number of write operations=4
HDFS: Number of bytes read erasure-coded=0
Map-Reduce Framework
  Map input records=5
  Map output records=20
  Map output bytes=169
  Map output materialized bytes=215
  Input split bytes=88
  Combine input records=0
  Combine output records=0
  Reduce input groups=10
  Reduce shuffle bytes=215
  Reduce input records=20
  Reduce output records=10
  Spilled Records=40
  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=89
File Output Format Counters
  Bytes Written=69
Exit Code: 0
hadoop@bmscsecse-HP-Elite-Tower-800-G9-Desktop-PC: ~$ hadoop fs -cat /result/part-00000
are 1
brother 1
family 1
hi 1
how 5
is 4
job 1
sister 1
you 1
your 4
hadoop@bmscsecse-HP-Elite-Tower-800-G9-Desktop-PC: ~$ []
```

machine

Date 6/5/25  
Page

dab - 7

Jar File creation :- To perform map reduce program for word count using eclipse by executing jar file.

① Open eclipse → new → java project → change java version to 1.8

② Import core jar file in hadoop mapreduce and common jar file in hadoop common to path.

③ Create 3 files for driver, Mapper and Reducer in src and save.

in source.

### Execution

> start-all.sh

> jps

> hadoop fs -ls /  
(Gives all groups in hadoop)

> hadoop fs -mkdir /rigs  
(Create a directory with name rigs)

> hadoop fs -copyFromLocal /home/hadoop/Desktop/file1.txt /rigs/test.txt  
(Copy the input file from local system into hadoop file system)

=> Get the jar file's driver class to test on input and produce output file

> hadoop jar /home/hadoop/wordcount.jar  
Driver /tmp/test.txt /tmp/output.txt

\* Read the output

> hadoop fs -cat /tmp/output.txt/part-00000  
(output.txt will be a directory)

(cat displays all file contents)

Output :-

an = 1

brother = 1

family = 1

hi = 1

how = 5

is = 4

job = 4

sister = 1

you = 1

your = 4

① 9  
② 1  
③ 3  
④ 4

>  
>  
>

>

>

>

o  
h

## Lab 8: Hadoop

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.

## Code with Output:

The screenshot shows the Eclipse IDE interface with the following details:

- File Bar:** File, Edit, Source, Refactor, Navigate, Search, Project, Run, Window, Help.
- Project Explorer:** Shows the project structure under toplab:
  - src
    - toplab
      - TopN.java
      - TopNReducer.java
      - TopNMapper.java
      - TopNCombiner.java
  - JRE System Library [JavaSE-1.8]
  - Referenced Libraries
  - settings
  - .classpath
  - .project- Editor:** The code editor displays the `TopNMapper.java` file. The code implements the `Mapper<Object, Text, Text, IntWritable>` interface. It includes methods for cleaning text, tokenizing, and emitting tokens with counts.
- Outline View:** Shows the class structure and the `_map` method implementation.

```
hadoop@bmscsece-HP-Elite-Tower-800-G9-Desktop-PC:~$ start -all.sh
Command 'start' not found, did you mean:
  command 'stars' from snap stars (2.7jrc3)
  command 'rstart' from deb x11-session-utils (7.7+4build2)
  command 'kstart' from deb kde-cli-tools (4:5.24.4-0ubuntu1)
  command 'startx' from deb xinit (1.4.1-0ubuntu4)
  command 'stat' from deb coreutils (8.32-4.1ubuntu1.2)
  command 'tart' from deb tart (3.10-1build1)
See 'snap info <snapname>' for additional versions.
hadoop@bmscsece-HP-Elite-Tower-800-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 [bmscsece-HP-Elite-Tower-800-G9-Desktop-PC]
Starting resourcemanager
Starting nodemanagers
hadoop@bmscsece-HP-Elite-Tower-800-G9-Desktop-PC:~$ hadoop fs -mkdir /mno
mkdir: Cannot create directory /mno. Name node is in safe mode.
hadoop@bmscsece-HP-Elite-Tower-800-G9-Desktop-PC:~$ ^C
hadoop@bmscsece-HP-Elite-Tower-800-G9-Desktop-PC:~$ ^C
hadoop@bmscsece-HP-Elite-Tower-800-G9-Desktop-PC:~$ hdfs dfsadmin -safemode get
Safe mode is OFF
hadoop@bmscsece-HP-Elite-Tower-800-G9-Desktop-PC:~$ hadoop fs -mkdir /mno
hadoop@bmscsece-HP-Elite-Tower-800-G9-Desktop-PC:~$ hadoop fs -copyFromLocal /home/hadoop/Desktop/sample.txt /mno/sample.txt
hadoop@bmscsece-HP-Elite-Tower-800-G9-Desktop-PC:~$ hadoop jar /home/hadoop/Desktop/topn.jar TopN /mno/sample.txt /res
Exception in thread "main" java.lang.ClassNotFoundException: TopN
        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.forName0(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-800-G9-Desktop-PC:~$ hadoop jar /home/hadoop/Desktop/topn.jar topnlab /mno/sample.txt /res
Exception in thread "main" java.lang.ClassNotFoundException: topnlab
        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.forName0(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-800-G9-Desktop-PC:~$ hadoop jar /home/hadoop/Desktop/topn.jar TopN /mno/sample.txt /res
Exception in thread "main" java.lang.NoClassDefFoundError: topnlab/TopN (wrong name: TopN)
        at java.base/java.lang.ClassLoader.defineClass1(Native Method)
```

```
hadoop@bmsccse-HP-Elite-Tower-800-G9-Desktop-PC: ~
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.forName0(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@bmsccse-HP-Elite-Tower-800-G9-Desktop-PC:~$ hadoop jar /home/hadoop/Desktop/topn.jar topnlab /mn0/sample.txt /res
Exception in thread "main" java.lang.ClassNotFoundException: topnlab
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.Class.forName0(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@bmsccse-HP-Elite-Tower-800-G9-Desktop-PC:~$ hadoop jar /home/hadoop/Desktop/topn.jar TopN /mn0/sample.txt /res
Exception in thread "main" java.lang.NoClassDefFoundError: topnlab/TopN (wrong name: TopN)
at java.base/java.lang.ClassLoader.defineClass1(Native Method)
at java.base/java.lang.ClassLoader.defineClass(ClassLoader.java:1022)
at java.base/java.security.SecureClassLoader.defineClass(SecureClassLoader.java:174)
at java.base/java.net.URLClassLoader.defineClass(URLClassLoader.java:555)
at java.base/java.net.URLClassLoader$1.run(URLClassLoader.java:458)
at java.base/java.net.URLClassLoader$1.run(URLClassLoader.java:452)
at java.base/java.security.AccessController.doPrivileged(Native Method)
at java.base/java.net.URLClassLoader.findClass(URLClassLoader.java:451)
at java.base/java.lang.ClassLoader.loadClass(ClassLoader.java:594)
at java.base/java.lang.ClassLoader.loadClass(ClassLoader.java:527)
at java.base/java.lang.Class.forName0(Native Method)
at java.base/java.lang.Class.forName(Class.java:398)
at org.apache.hadoop.util.RunJar.run(RunJar.java:321)
at org.apache.hadoop.util.RunJar.main(RunJar.java:241)
hadoop@bmsccse-HP-Elite-Tower-800-G9-Desktop-PC:~$ hadoop jar /home/hadoop/Desktop/topn.jar topnlab.TopN /mn0/sample.txt /res
2025-05-20 13:58:09,506 INFO impl.MetricsConfig: Loaded properties from hadoop-metrics2.properties
2025-05-20 13:58:09,545 INFO impl.MetricSystemImpl: Scheduled Metric snapshot period at 10 second(s).
2025-05-20 13:58:09,545 INFO impl.MetricsSystemImpl: JobTracker metrics system started
2025-05-20 13:58:09,658 INFO input.FileInputFormat: Total input files to process : 1
2025-05-20 13:58:09,709 INFO mapreduce.JobSubmitter: number of splits:1
2025-05-20 13:58:09,777 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_local408680812_0001
2025-05-20 13:58:09,778 INFO mapreduce.JobSubmitter: Executing with tokens: []
2025-05-20 13:58:09,836 INFO mapreduce.Job: The url to track the job: http://localhost:8080/
2025-05-20 13:58:09,837 INFO mapreduce.Job: Running job: job_local408680812_0001
2025-05-20 13:58:09,838 INFO mapred.LocalJobRunner: OutputCommitter set in config null
2025-05-20 13:58:09,841 INFO output.PathOutputCommitterFactory: No output committer factory defined, defaulting to FileOutputCommitterFactory
2025-05-20 13:58:09,842 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 2
2025-05-20 13:58:09,842 INFO output.FileOutputCommitter: FileOutputCommitter skip cleanup _temporary folders under output directory:false, ignore cleanup failures: false
```

```
+ hadoop@bmscerse-HP-Elite-Tower-800-G9-Desktop-PC: ~
2025-05-20 13:58:09,842 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 2
2025-05-20 13:58:09,842 INFO output.FileOutputCommitter: FileOutputCommitter skip cleanup _temporary folders under output directory:false, ignore cleanup failures: false
2025-05-20 13:58:09,842 INFO mapred.LocalJobRunner: OutputCommitter is org.apache.hadoop.mapreduce.lib.output.FileOutputCommitter
2025-05-20 13:58:09,884 INFO mapred.LocalJobRunner: Waiting for map tasks
2025-05-20 13:58:09,885 INFO mapred.LocalJobRunner: Starting task: attempt_local408680812_0001_m_000000_0
2025-05-20 13:58:09,895 INFO output.PathOutputCommitterFactory: No output committer factory defined, defaulting to FileOutputCommitterFactory
2025-05-20 13:58:09,895 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 2
2025-05-20 13:58:09,895 INFO output.FileOutputCommitter: FileOutputCommitter skip cleanup _temporary folders under output directory:false, ignore cleanup failures: false
2025-05-20 13:58:09,903 INFO mapred.Task: Using ResourceCalculatorProcessTree : []
2025-05-20 13:58:09,906 INFO mapred.MapTask: Processing split: hdfs://localhost:9000/mno/sample.txt:0+75
2025-05-20 13:58:09,945 INFO mapred.MapTask: (EQUATOR) 0 kvi 26214396(104857584)
2025-05-20 13:58:09,945 INFO mapred.MapTask: mapreduce.task.io.sort.mb: 100
2025-05-20 13:58:09,945 INFO mapred.MapTask: soft limit at 83886080
2025-05-20 13:58:09,945 INFO mapred.MapTask: bufstart = 0; bufvoid = 104857600
2025-05-20 13:58:09,945 INFO mapred.MapTask: kvstart = 26214396; length = 6553600
2025-05-20 13:58:09,947 INFO mapred.MapTask: Map output collector class = org.apache.hadoop.mapred.MapTask$MapOutputBuffer
2025-05-20 13:58:10,006 INFO mapred.LocalJobRunner:
2025-05-20 13:58:10,007 INFO mapred.MapTask: Starting flush of map output
2025-05-20 13:58:10,007 INFO mapred.MapTask: Spilling map output
2025-05-20 13:58:10,007 INFO mapred.MapTask: bufstart = 0; bufend = 135; bufvoid = 104857600
2025-05-20 13:58:10,007 INFO mapred.MapTask: kvstart = 26214396(104857584); kvend = 26214340(104857360); length = 57/6553600
2025-05-20 13:58:10,010 INFO mapred.MapTask: Finished spill 0
2025-05-20 13:58:10,014 INFO mapred.Task: Task:attempt_local408680812_0001_m_000000_0 is done. And is in the process of committing
2025-05-20 13:58:10,016 INFO mapred.LocalJobRunner: map
2025-05-20 13:58:10,017 INFO mapred.Task: Task 'attempt_local408680812_0001_m_000000_0' done.
2025-05-20 13:58:10,020 INFO mapred.Task: Final Counters for attempt_local408680812_0001_m_000000_0: Counters: 23
    File System Counters
        FILE: Number of bytes read=751
        FILE: Number of bytes written=645435
        FILE: Number of read operations=0
        FILE: Number of large read operations=0
        FILE: Number of write operations=0
        HDFS: Number of bytes read=75
        HDFS: Number of bytes written=0
        HDFS: Number of read operations=5
        HDFS: Number of large read operations=0
        HDFS: Number of write operations=1
        HDFS: Number of bytes read erasure-coded=0
    Map-Reduce Framework
        Map input records=2
        Map output records=15
        Map output bytes=135
```

```
hadoop@bmscerse-HP-Elite-Tower-800-G9-Desktop-PC: ~
2025-05-20 13:58:10,168 INFO mapred.Task: Task:attempt_local408680812_0001_r_000000_0 is done. And is in the process of committing.
2025-05-20 13:58:10,169 INFO mapred.LocalJobRunner: 1 / 1 copied.
2025-05-20 13:58:10,169 INFO mapred.Task: Task attempt_local408680812_0001_r_000000_0 is allowed to commit now
2025-05-20 13:58:10,194 INFO output.FileOutputCommitter: Saved output of task 'attempt_local408680812_0001_r_000000_0' to hdfs://localhost:9000/res
2025-05-20 13:58:10,195 INFO mapred.LocalJobRunner: reduce > reduce
2025-05-20 13:58:10,196 INFO mapred.Task: Task 'attempt_local408680812_0001_r_000000_0' done.
2025-05-20 13:58:10,197 INFO mapred.Task: Final Counters for attempt_local408680812_0001_r_000000_0: Counters: 30
    File System Counters
        FILE: Number of bytes read=7887
        FILE: Number of bytes written=645606
        FILE: Number of read operations=0
        FILE: Number of large read operations=0
        FILE: Number of write operations=0
        HDFS: Number of bytes read=75
        HDFS: Number of bytes written=105
        HDFS: Number of read operations=10
        HDFS: Number of large read operations=0
        HDFS: Number of write operations=3
        HDFS: Number of bytes read erasure-coded=0
    Map-Reduce Framework
        Combine input records=0
        Combine output records=0
        Reduce input groups=15
        Reduce shuffle bytes=171
        Reduce input records=15
        Reduce output records=15
        Spilled Records=15
        Shuffled Maps =1
        Failed Shuffles=0
        Merged Map outputs=1
        GC time elapsed (ms)=0
        Total committed heap usage (bytes)=526385152
    Shuffle Errors
        BAD_ID=0
        CONNECTION=0
        IO_ERROR=0
        WRONG_LENGTH=0
        WRONG_MAP=0
        WRONG_REDUCE=0
    File Output Format Counters
        Bytes Written=105
2025-05-20 13:58:10,197 INFO mapred.LocalJobRunner: Finishing task: attempt_local408680812_0001_r_000000_0
2025-05-20 13:58:10,197 INFO mapred.LocalJobRunner: reduce task executor complete.
2025-05-20 13:58:10,840 INFO mapreduce.Job: Job job_local408680812_0001 running in uber mode : false
2025-05-20 13:58:10,842 INFO mapreduce.Job: map 100% reduce 100%
2025-05-20 13:58:10,843 INFO mapreduce.Job: Job job_local408680812_0001 completed successfully
```

```
hadoop@bmscerse-HP-Elite-Tower-800-G9-Desktop-PC: ~
2025-05-20 13:58:10,168 INFO mapred.Task: Task:attempt_local408680812_0001_r_000000_0 is done. And is in the process of committing.
2025-05-20 13:58:10,169 INFO mapred.LocalJobRunner: 1 / 1 copied.
2025-05-20 13:58:10,169 INFO mapred.Task: Task attempt_local408680812_0001_r_000000_0 is allowed to commit now
2025-05-20 13:58:10,194 INFO output.FileOutputCommitter: Saved output of task 'attempt_local408680812_0001_r_000000_0' to hdfs://localhost:9000/res
2025-05-20 13:58:10,195 INFO mapred.LocalJobRunner: reduce > reduce
2025-05-20 13:58:10,196 INFO mapred.Task: Task 'attempt_local408680812_0001_r_000000_0' done.
2025-05-20 13:58:10,197 INFO mapred.Task: Final Counters for attempt_local408680812_0001_r_000000_0: Counters: 30
    File System Counters
        FILE: Number of bytes read=7887
        FILE: Number of bytes written=645606
        FILE: Number of read operations=0
        FILE: Number of large read operations=0
        FILE: Number of write operations=0
        HDFS: Number of bytes read=75
        HDFS: Number of bytes written=105
        HDFS: Number of read operations=10
        HDFS: Number of large read operations=0
        HDFS: Number of write operations=3
        HDFS: Number of bytes read erasure-coded=0
    Map-Reduce Framework
        Combine input records=0
        Combine output records=0
        Reduce input groups=15
        Reduce shuffle bytes=171
        Reduce input records=15
        Reduce output records=15
        Spilled Records=15
        Shuffled Maps =1
        Failed Shuffles=0
        Merged Map outputs=1
        GC time elapsed (ms)=0
        Total committed heap usage (bytes)=526385152
    Shuffle Errors
        BAD_ID=0
        CONNECTION=0
        IO_ERROR=0
        WRONG_LENGTH=0
        WRONG_MAP=0
        WRONG_REDUCE=0
    File Output Format Counters
        Bytes Written=105
2025-05-20 13:58:10,197 INFO mapred.LocalJobRunner: Finishing task: attempt_local408680812_0001_r_000000_0
2025-05-20 13:58:10,197 INFO mapred.LocalJobRunner: reduce task executor complete.
2025-05-20 13:58:10,840 INFO mapreduce.Job: Job job_local408680812_0001 running in uber mode : false
2025-05-20 13:58:10,842 INFO mapreduce.Job: map 100% reduce 100%
2025-05-20 13:58:10,843 INFO mapreduce.Job: Job job_local408680812_0001 completed successfully
```

```
+ hadoop@bmscerse-HP-Elite-Tower-800-G9-Desktop-PC: ~
Map input records=2
Map output records=15
Map output bytes=135
Map output materialized bytes=171
Input split bytes=101
Combine input records=0
Spilled Records=15
Failed Shuffles=0
Merged Map outputs=0
GC time elapsed (ms)=0
Total committed heap usage (bytes)=526385152
File Input Format Counters
  Bytes Read=75
2025-05-20 13:58:10,020 INFO mapred.LocalJobRunner: Finishing task: attempt_local408680812_0001_m_000000_0
2025-05-20 13:58:10,020 INFO mapred.LocalJobRunner: map task executor complete.
2025-05-20 13:58:10,021 INFO mapred.LocalJobRunner: Waiting for reduce tasks
2025-05-20 13:58:10,022 INFO mapred.LocalJobRunner: Starting task: attempt_local408680812_0001_r_000000_0
2025-05-20 13:58:10,027 INFO output.PathOutputCommitterFactory: No output committer factory defined, defaulting to FileOutputCommitterFactory
2025-05-20 13:58:10,027 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 2
2025-05-20 13:58:10,027 INFO output.FileOutputCommitter: FileOutputCommitter skip cleanup _temporary folders under output directory:false, ignore cleanup failures: false
2025-05-20 13:58:10,027 INFO mapred.Task: Using ResourceCalculatorProcessTree : [ ]
2025-05-20 13:58:10,028 INFO mapred.ReduceTask: Using ShuffleConsumerPlugin: org.apache.hadoop.mapreduce.task.reduce.Shuffle@6622090a
2025-05-20 13:58:10,029 WARN impl.MetricsSystemImpl: JobTracker metrics system already initialized!
2025-05-20 13:58:10,037 INFO reduce.MergeManagerImpl: MergerManager: memoryLimit=5829453312, maxSingleShuffleLimit=1457363328, mergeThreshold=3847439360, ioSortFactor=10, memToMemMergeOutputsThreshold=10
2025-05-20 13:58:10,038 INFO reduce.EventFetcher: attempt_local408680812_0001_r_000000_0 Thread started: EventFetcher for fetching Map Completion Events
2025-05-20 13:58:10,053 INFO reduce.LocalFetcher: localfetcher#1 about to shuffle output of map attempt_local408680812_0001_m_000000_0 decomp: 167 len: 171 to MEMORY
2025-05-20 13:58:10,054 INFO reduce.InMemoryMapOutput: Read 167 bytes from map-output for attempt_local408680812_0001_m_000000_0
2025-05-20 13:58:10,055 INFO reduce.MergeManagerImpl: closeInMemoryFile -> map-output of size: 167, inMemoryMapOutputs.size() -> 1, commitMemory -> 0, usedMemory ->167
2025-05-20 13:58:10,056 INFO reduce.EventFetcher: EventFetcher is interrupted.. Returning
2025-05-20 13:58:10,056 INFO mapred.LocalJobRunner: 1 / 1 copied.
2025-05-20 13:58:10,056 INFO reduce.MergeManagerImpl: finalMerge called with 1 in-memory map-outputs and 0 on-disk map-outputs
2025-05-20 13:58:10,059 INFO mapred.Merger: Merging 1 sorted segments
2025-05-20 13:58:10,059 INFO mapred.Merger: Down to the last merge-pass, with 1 segments left of total size: 162 bytes
2025-05-20 13:58:10,060 INFO reduce.MergeManagerImpl: Merged 1 segments, 167 bytes to disk to satisfy reduce memory limit
2025-05-20 13:58:10,060 INFO reduce.MergeManagerImpl: Merging 1 files, 171 bytes from disk
2025-05-20 13:58:10,060 INFO reduce.MergeManagerImpl: Merging 0 segments, 0 bytes from memory into reduce
2025-05-20 13:58:10,060 INFO mapred.Merger: Merging 1 sorted segments
2025-05-20 13:58:10,061 INFO mapred.Merger: Down to the last merge-pass, with 1 segments left of total size: 162 bytes
2025-05-20 13:58:10,061 INFO mapred.LocalJobRunner: 1 / 1 copied.
```

```
hadoop@bmscse-HP-Elite-Tower-800-G9-Desktop-PC: ~
2025-05-20 13:58:10,840 INFO mapreduce.Job: Job job_local408680812_0001 running in uber node : false
2025-05-20 13:58:10,842 INFO mapreduce.Job: map 100% reduce 100%
2025-05-20 13:58:10,843 INFO mapreduce.Job: Job job_local408680812_0001 completed successfully
2025-05-20 13:58:10,854 INFO mapreduce.Job: Counters: 36
    File System Counters
        FILE: Number of bytes read=15400
        FILE: Number of bytes written=1291041
        FILE: Number of read operations=0
        FILE: Number of large read operations=0
        FILE: Number of write operations=0
        HDFS: Number of bytes read=150
        HDFS: Number of bytes written=105
        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=2
        Map output records=15
        Map output bytes=135
        Map output materialized bytes=171
        Input split bytes=101
        Combine input records=0
        Combine output records=0
        Reduce input groups=15
        Reduce shuffle bytes=171
        Reduce input records=15
        Reduce output records=15
        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=75
    File Output Format Counters
        Bytes Written=105
hadoop@bmscsecse-HP-Elite-Tower-800-G9-Desktop-PC:~$ hadoop fs -cat /res/part-00000
cat: '/res/part-00000': No such file or directory
hadoop@bmscsecse-HP-Elite-Tower-800-G9-Desktop-PC:~$ hadoop fs -ls /res
```

```
hadoop@bmscsece-HP-Elite-Tower-800-G9-Desktop-PC: ~
Input split bytes=101
Combine input records=0
Combine output records=0
Reduce input groups=15
Reduce shuffle bytes=171
Reduce input records=15
Reduce output records=15
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=75
File Output Format Counters
Bytes Written=105
hadoop@bmscsece-HP-Elite-Tower-800-G9-Desktop-PC:~$ hadoop fs -cat /res/part-00000
cat: '/res/part-00000': No such file or directory
hadoop@bmscsece-HP-Elite-Tower-800-G9-Desktop-PC:~$ hadoop fs -ls /res
Found 2 items
-rw-r--r-- 1 hadoop supergroup      0 2025-05-20 13:58 /res/_SUCCESS
-rw-r--r-- 1 hadoop supergroup  105 2025-05-20 13:58 /res/part-r-00000
hadoop@bmscsece-HP-Elite-Tower-800-G9-Desktop-PC:~$ ^C
hadoop@bmscsece-HP-Elite-Tower-800-G9-Desktop-PC:~$ hadoop fs -cat /res/part-r-00000
college 1
in      1
bms     1
hi      1
i       1
inna    1
am      1
m       1
bhuvana 1
how    1
are    1
avyukth 1
of     1
you    1
engineering 1
hadoop@bmscsece-HP-Elite-Tower-800-G9-Desktop-PC:~$
```

## Lab 9: Scala

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

## Code with Output:

Scala:-

To run Scala in spark

> \$ spark-shell

> // creating file  
name pointnumbers.scala

object Main

def Main (args : Array [String])

for (w < 0 to 100)

pointn (w)

}

ctrl + O → save

enter

ctrl + X → exit nano.

scala >: load pointnumbers.scala.

scala > pointnumbers.main (Array(1))

100

## Lab 10: Spark

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

## Code with Output:

```
bmscsece@bmscsece-HP-Elite-Tower-800-G9-Desktop-PC:~$ spark-shell
25/05/20 15:32:38 WARN Utils: Your hostname, bmscsece-HP-Elite-Tower-800-G9-Desktop-PC resolves to a loopback address: 127.0.1.1
: using 10.124.2.8 instead (on interface eno1)
25/05/20 15:32:38 WARN Utils: Set SPARK_LOCAL_IP if you need to bind to another address
WARNING: An illegal reflective access operation has occurred
WARNING: Illegal reflective access by org.apache.spark.unsafe.Platform (file:/opt/spark/jars/spark-unsafe_2.12-3.0.3.jar) to con-
structor java.nio.DirectByteBuffer(long,int)
WARNING: Please consider reporting this to the maintainers of org.apache.spark.unsafe.Platform
WARNING: Use --illegal-access=warn to enable warnings of further illegal reflective access operations
WARNING: All illegal access operations will be denied in a future release
25/05/20 15:32:38 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes wh-
ere applicable
Using Spark's default log4j profile: org/apache/spark/log4j-defaults.properties
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
Spark context Web UI available at http://10.124.2.8:4040
Spark context available as 'sc' (master = local[*], app id = local-1747735361481).
Spark session available as 'spark'.
Welcome to

    / \ \ / \ / \ \ / \ / \
   / \ \ / \ / \ / \ / \ / \ / \
  / \ \ / \ / \ / \ / \ / \ / \ / \
 / \ \ / \ / \ / \ / \ / \ / \ / \ / \
version 3.0.3

Using Scala version 2.12.10 (OpenJDK 64-Bit Server VM, Java 11.0.26)
Type in expressions to have them evaluated.
Type :help for more information.

scala> val textFile = sc.textFile("/home/bmscsece/Desktop/sparkdata.txt")
textFile: org.apache.spark.rdd.RDD[String] = /home/bmscsece/Desktop/sparkdata.txt MapPartitionsRDD[1] at textFile at <console>:2
4

scala>

scala> val counts = textFile
counts: org.apache.spark.rdd.RDD[String] = /home/bmscsece/Desktop/sparkdata.txt MapPartitionsRDD[1] at textFile at <console>:24

scala> .flatMap(line => line.split(" "))
res0: org.apache.spark.rdd.RDD[String] = MapPartitionsRDD[2] at flatMap at <console>:26

scala> .map(word => (word, 1))
scala> val data = sc.textFile("sparkdata.txt")
data: org.apache.spark.rdd.RDD[String] = sparkdata.txt MapPartitionsRDD[1] at textFile at <console>:25

scala> val splitdata = data.flatMap(line => line.split(" "))
splitdata: org.apache.spark.rdd.RDD[String] = MapPartitionsRDD[2] at flatMap at <console>:26

scala> val mapdata = splitdata.map(word => (word, 1))
mapdata: org.apache.spark.rdd.RDD[(String, Int)] = MapPartitionsRDD[3] at map at <console>:26

scala> val reducedata = mapdata.reduceByKey(_ + _)
reducedata: org.apache.spark.rdd.RDD[(String, Int)] = ShuffledRDD[4] at reduceByKey at <console>:26

scala> reducedata.collect.foreach(println)
(),1
(hello,2)
(world,1)
(spark,1)
```

```

scala> val textFile = sc.textFile("/home/bmscsece/Desktop/WC.txt")
textFile: org.apache.spark.rdd.RDD[String] = /home/bmscsece/Desktop/WC.txt MapPartitionsRDD[31] at textFile at <console>:31

scala> val words = textfile.flatMap(line => line.split(" "))
words: org.apache.spark.rdd.RDD[String] = MapPartitionsRDD[32] at flatMap at <console>:32

scala>

scala> val pairs = words.map(word => (word, 1))
pairs: org.apache.spark.rdd.RDD[(String, Int)] = MapPartitionsRDD[33] at map at <console>:32

scala>

scala> val counts = pairs.reduceByKey(_ + _)
counts: org.apache.spark.rdd.RDD[(String, Int)] = ShuffledRDD[34] at reduceByKey at <console>:32

scala> val countsArray = counts.collect() // This is Array[(String, Int)]
countsArray: Array[(String, Int)] = Array(("",1), (hello,6), (world,1), (spark,1))

scala> val sorted = ListMap(countsArray.sortWith(_._2 > _._2): _*)
sorted: scala.collection.immutable.ListMap[String,Int] = ListMap(hello -> 6, "" -> 1, world -> 1, spark -> 1)

scala> for ((k, v) <- sorted) {
|   if (v > 4) println(s"$k, $v")
| }
hello, 6

scala>

```

